

Meta-analysis of dichotomous and ordinal tests with an imperfect gold standard

Enzo Cerullo ^{*,1,2}, Hayley E. Jones ³, Olivia Carter⁴, Terry J. Quinn ⁵, Nicola J. Cooper ^{1,2}, and Alex J. Sutton ^{1,2}

¹Biostatistics Research Group, Department of Health Sciences, University of Leicester, Leicester, UK

²Complex Reviews Support Unit, University of Leicester & University of Glasgow, Glasgow, UK

³Population Health Sciences, Bristol Medical School, University of Bristol, Bristol, UK

⁴No affiliation

⁵Institute of Cardiovascular and Medical Sciences, University of Glasgow, Glasgow, UK

Abstract

Standard methods for the meta-analysis of medical tests, without assuming a gold standard, are limited to dichotomous data. Multivariate probit models are used to analyze correlated dichotomous data, and can be extended to model ordinal data. Within the context of an imperfect gold standard, they have previously been used for the analysis of dichotomous and ordinal test data from a single study, and for the meta-analysis of dichotomous tests. However, they have not previously been used for the meta-analysis of ordinal tests.

In this paper, we developed a Bayesian multivariate probit latent class model for the simultaneous meta-analysis of ordinal and dichotomous tests without assuming a gold standard, which also allows one to obtain summary estimates of joint test accuracy. We fitted the models using the software Stan, which uses a state-of-the-art Hamiltonian Monte Carlo algorithm, and we applied the models to a dataset in which studies evaluated the accuracy of tests, and test combinations, for deep vein thrombosis. We demonstrate the issues with dichotomising ordinal test accuracy data in the presence of an imperfect gold standard, before applying and comparing several variations of our proposed model which do not require the data to be dichotomised.

The models proposed will allow researchers to more appropriately meta-analyse ordinal and dichotomous tests without a gold standard, potentially leading to less biased estimates of test accuracy. This may lead to a better understanding of which tests, and test combinations, should be used for any given medical condition.

Keywords

Meta-Analysis, test accuracy, multivariate probit, latent class, imperfect gold, ordinal tests

*Corresponding Author

Email address: enzo.cerullo@bath.edu

1 Introduction

Medical tests are used to screen, monitor and diagnose medical conditions. In order to evaluate their accuracy, we can carry out test accuracy studies - studies which estimate the accuracy of a test by comparing its results to some existing test assumed to be perfect (i.e. 100% sensitive and specific). The tests under evaluation, and those assumed to be perfect, are referred to as *index* tests and *reference* (or *gold standard*) tests, respectively. Index tests often have a lower sensitivity and/or specificity than the gold standard; however, they may be quicker, less invasive, and/or less costly. Unfortunately, the fact that gold standard tests are often imperfect is ignored in routinely used methods^{1,2,3} to meta-analyse studies of test accuracy, which can lead to misleading results⁴.

The results between tests are usually *conditionally dependent* - that is, they are correlated within each disease class (diseased and non-diseased individuals). Models which account for this dependency, in addition to an imperfect gold standard have been proposed^{5,6,7,8,9}. These models - which we will refer to as *traditional latent class models* (TLCMs) - assume that all tests are measuring the same latent disease, and each individual is modelled as belonging in either disease class. Since they can model imperfect gold standards, they also allow one to compare the accuracy between gold standard and index tests.

Proposed models which can model an imperfect gold standard based on TLCMs have some limitations, which motivates the proposal of more flexible latent class models^{10,11,8,12,13,14,9}. For instance, multivariate probit latent class (MVP-LC) models^{10,11,8,12,13,14,15}, which are a type of regression model. Unlike TLCMs, MVP-LC models can be extended to model ordinal test accuracy data^{12,16} without forcing the user to dichotomise it, whilst simultaneously modelling conditional dependence. For example, Xu et al¹⁰ presented an MVP-LC model to analyse primary studies evaluating multiple dichotomous tests without assuming a gold standard, which they later extended¹¹ to model ordinal tests with two cutpoints. The latent trait model proposed by Qu et al¹³ is a variation of the MVP-LC model which is defined by specifying a series of univariate regressions with a common subject-specific latent variable. This model was later expanded upon by Sadatsafavi et al⁸ to the meta-analysis setting - to analyse studies evaluating up to three dichotomous tests using direct comparisons - whilst allowing the test accuracy to vary between studies. However, it cannot appropriately model ordinal tests, nor can it model between-study variation for the conditional dependence parameters.

In clinical practice, tests are rarely used in isolation. The accuracy of two or more tests used in combination is often referred to as the *joint test accuracy*. Few meta-analytical methods have been proposed which can simultaneously calculate summary joint test accuracy and incorporate ordinal tests, all of which assume a perfect gold standard. For instance, Novielli et al¹⁷ proposed a model based on conditional probabilities, in which studies evaluated up to one ordinal test with two cutpoints and two dichotomous tests. This model can estimate summary test accuracy at each cutpoint whilst modelling conditional dependence.

To address the gaps in the literature discussed above, we developed a Bayesian model for the meta-analysis of studies evaluating both ordinal and dichotomous tests without assuming a perfect gold standard. The model also enables the estimation of summary joint test accuracy, whilst allowing the conditional dependence parameters to vary between studies. The proposed model is an extension of previous MVP-LC models which have been developed to analyse multiple tests in a single study^{10,11}. In section 2, we describe the case study dataset which will serve to motivate our proposed model, which we will describe in section 3. Then, we apply several variations of our proposed model to this dataset in section 4. Finally, in section 5 we discuss the benefits and limitations of the model, as well as possible extensions.

2 Motivating example

Deep vein thrombosis (DVT) is the formation of a blood clot in a deep (i.e. not superficial) vein. DVT can occur in the upper (proximal) or lower (distal) part of the leg, with the former more likely

to be life-threatening. A potential complication of DVT occurring in up to a third of patients¹⁸ is pulmonary embolism (PE). PE occurs when a blood vessel in the lungs becomes blocked by a blood clot (formed as a result of DVT) which has migrated from the legs to the lungs. Contrast venography is generally considered to be a gold standard for DVT, as it is almost 100% sensitive and specific^{19,20}. However, it is not commonly used in clinical practice because it is time consuming and invasive^{19,20}. Instead, ultrasound is often used to diagnose DVT, since it is non-invasive and cost-effective^{18,21,22,23}. However, it is less accurate than contrast venography^{24,25} for both distal and proximal DVT, with its sensitivity being lower for distal DVT²⁴. Furthermore, although ultrasound is known to have a very high specificity, it is still nonetheless imperfect^{24,25}. A commonly used¹⁸ screening tool for DVT is a questionnaire called the Wells score²⁶, which groups patients into one of three risk categories - 'low', 'intermediate', or 'high'. Another DVT test is the D-Dimer assay: a blood test measuring the amount of a protein fragment called D-Dimer, higher concentrations of which are indicative of DVT. Despite being considered to be generally more accurate than the Wells²⁷, the D-Dimer assay is intended to be used for screening as opposed to diagnosis²⁷, since a number of other conditions can elevate serum D-dimer concentrations¹⁸.

Investigating the joint test accuracy of the aforementioned tests for DVT is important for a variety of reasons. The Wells and D-Dimer are both relatively cheap, quick and non-invasive to carry out, particularly the Wells test. A combined screening approach utilising the Wells and D-Dimer may be more cost-effective and reduce test burden for patients compared to using either alone. Furthermore, despite the fact that neither the Wells nor the D-Dimer alone are generally considered to be diagnostic tools for DVT, they may have diagnostic potential when combined^{17,24,28}. An example of a potential screening strategy is to use the Wells prior to the D-dimer in the diagnostic pathway as a pre-screening tool to rule out individuals at low risk for DVT. Following this, individuals who scored as intermediate or high risk are subsequently screened using the D-Dimer assay, and only patients who also test positive on the D-Dimer undertake ultrasound. Another potential strategy is to refer patients scoring as high risk on the Wells score directly to ultrasound. Both of the aforementioned joint testing strategies are examples of '*believe the negatives*' (BTN) strategies¹⁷. This is a testing strategy where only those patients who test positive on an initial test go on to receive a second test, then only individuals who also test positive on the second test are considered positive. Conversely, '*Believe the positives*' (BTP) is a testing strategy where only those patients who test negative on the first test go on to receive a second test, with only those patients who also test negative on this test being considered negative. Joint testing strategies are important across clinical areas besides DVT, for example for depression screening and for COVID-19 - see discussion section 5.1 for more details.

Novielli et al¹⁷ proposed a statistical model in order to conduct a meta-analysis of studies investigating the D-dimer, Wells score and ultrasound for DVT. The proposed model allowed them to model the Wells score without dichotomising the data whilst modelling the conditional dependence between tests, enabling them to estimate summary-level joint test accuracy. However, their model assumes that ultrasound is a perfect gold standard, which could have led to biased estimates of the performance of other tests under evaluation. Novielli et al¹⁷ carried out several analyses based on different datasets – for instance, one based on the 11 studies which directly compared the D-dimer, Wells' score via the gold standard (ultrasound), and another which also included studies which only analysed one of Wells or D-dimer tests, and utilised indirect comparisons. In section 4 of this paper, we re-analyse the direct comparisons data (see table 1) from Novielli et al¹⁷ without assuming a perfect gold standard, using a variety of models we propose in section 3; namely, models which dichotomised the Wells score and those which modelled it as an ordinal test, those which assumed conditional independence and dependence between tests, as well as models which assumed ultrasonography was perfect or imperfect. This dataset consisted of 11 studies, with a total of 4096 individuals and 12,288 observations, with all 11 studies evaluating all three tests.

Table 1: Sample of case study dataset

Study	Ultrasound -'ve						Ultrasound +'ve					
	D-Dimer -'ve			D-Dimer +'ve			D-Dimer -'ve			D-Dimer +'ve		
	Wells score ¹						Wells score ¹					
	L	M	H	L	M	H	L	M	H	L	M	H
1	32	20	5	8	18	2	0	0	2	1	6	8
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
11	243	16	3	233	104	29	1	0	0	28	117	109

Note: All test results are modelled at the individual level. We show the aggregate data in this table for ease of presentation.

¹ The Wells score is classified as L = Low, M = Moderate, H = High

3 Methods

Before describing our proposed Bayesian MVP-LC model, we will first define some terminology and notation in section 3.1. For a formal model specification, please refer to the full technical model specification (supplementary material 1).

3.1 Terminology & notation

The model is for a meta-analysis dataset with a total of S studies, with each study having a total of N_s individuals, where s is an index for study - so s can be used to denote anything between the first ($s = 1$) and the last ($s = S$) study. Each study is assessing the same number of tests, T . We will use t as an index for test, which can be between 1 and T , and n as an index for individual, which can be between 1 and N_s for study s .

For the n th individual from study s , we will denote the vector of observed test responses as $\mathbf{y}_{s,n} = (y_{s,n,1}, \dots, y_{s,n,T})'$. Each test is either dichotomous or ordinal. For dichotomous tests, each observed test response, $y_{s,n,t}$, is coded as 0 and 1 for negative and positive results, respectively. For ordinal tests, each test t has K_t categories (hence $K_t - 1$ cutpoints). We will use k as an index to refer to any given cutpoint, which can be between 1 and $K_t - 1$. Ordinal test responses are coded according to the category that the individuals' test result falls in: in other words, $y_{s,n,t} = k$ if the test result falls in the k th category for test t . Since we are assuming an imperfect gold standard, the true disease status of each individual, $d_{s,n}$ is not defined by the results of the gold standard. Instead, it is modelled as an unknown (i.e. *latent*) variable, and belongs to one of two classes - 'diseased' ($d_{s,n} = 1$) or 'non-diseased' ($d_{s,n} = 0$).

3.2 Within-study model

We will now define the MVP-LC model within each study. For dichotomous data, MVP-LC models work by transforming the observed, discrete test result data into a *continuous* variable - a statistical technique known as "data augmentation"¹⁵. This augmented continuous data, which we will denote as $\mathbf{Z}_{s,n}$, is an unobserved latent variable, similarly to the disease status, $d_{s,n}$. This data augmentation process allows us to assume that, conditional on $d_{s,n}$, the unobserved test accuracy data, $\mathbf{Z}_{s,n}$, can be modelled by a multivariate normal (MVN) distribution with mean vector $\boldsymbol{\nu}_s^{[d]}$ and variance-covariance matrix $\boldsymbol{\Psi}_s^{[d]}$.

More specifically, $\mathbf{Z}_{s,n} \sim \text{MVN}(\boldsymbol{\nu}_s^{[d]}, \boldsymbol{\Psi}_s^{[d]})$, where:

$$\mathbf{Z}_{s,n} = \begin{pmatrix} Z_{s,n,1} \\ \vdots \\ Z_{s,n,T} \end{pmatrix}, \boldsymbol{\nu}_s^{[d]} = \begin{pmatrix} \nu_{s,1}^{[d]} \\ \vdots \\ \nu_{s,T}^{[d]} \end{pmatrix}, \boldsymbol{\Psi}_s^{[d]} = \begin{pmatrix} (\tau_{s,1}^{[d]})^2 & \cdots & \epsilon_{s,1,T}^{[d]} \cdot \tau_{s,1}^{[d]} \cdot \tau_{s,T}^{[d]} \\ \vdots & \ddots & \vdots \\ \epsilon_{s,T,1}^{[d]} \cdot \tau_{s,T}^{[d]} \cdot \tau_{s,1}^{[d]} & \cdots & (\tau_{s,T}^{[d]})^2 \end{pmatrix} \quad (1)$$

Where $\nu_{s,t}^{[d]}$ and $\tau_{s,t}^{[d]}$, denote the study-specific means and standard deviations, respectively. Each $\epsilon_{s,1,t}^{[d]}$ denotes the study-specific correlations between each test-pair (or 'test-pair' - denoted as 't and t') for the augmented data ($Z_{s,n,t}$ and $Z_{s,n,t'}$) - the *polychoric correlation*^{16,29} - which is not the same as the correlation between the observed data $y_{s,n,t}$ and $y_{s,n,t'}$. Each $\epsilon_{s,1,t}^{[d]}$ models the conditional dependence between each test-pair. However, we can assume *conditional independence* by setting $\epsilon_{s,1,t}^{[d]} = 0$ and $\epsilon_{s,1,t'}^{[d]} = 0$. We must ensure that the number of parameters being estimated from our model is not greater than what is possible for the given dataset; otherwise it may be *non-identifiable* - which means that the model will give misleading results. For example, it may estimate the sensitivity for a test to be equal to both *both* 0.20 and 0.80. To ensure that our model is identifiable^{12,16}, as in Xu et al¹¹, we set each $\tau_{s,t}^{[d]} = 1$ (i.e., set all $\boldsymbol{\Psi}_s^{[d]}$ to be correlation matrices). Please see supplementary material 6.

For dichotomous tests, the augmented data ($Z_{s,n,t}$) will be less than 0 for negative results ($y_{s,n,t} = 0$) or greater than 0 for positive results ($y_{s,n,t} = 1$), and the measures of test accuracy for a given study s are given by,

$$\begin{aligned} Se_{s,t} &= \Phi(\nu_{s,t}^{[1]}) \\ Sp_{s,t} &= 1 - \Phi(\nu_{s,t}^{[0]}) \end{aligned} \quad (2)$$

Where $\Phi(\cdot)$ denotes the cumulative density function (CDF) of the standard normal distribution - that is, a normal distribution with mean 0 and standard deviation 1. For ordinal tests, the augmented data ($Z_{s,n,t}$) will belong to an interval defined by strictly increasing latent cutpoint parameters ($\{C_{1,s,t}^{[d]}, \dots, C_{K_t-1,s,t}^{[d]}\}$, where $C_{k-1,s,t}^{[d]} < C_{k,s,t}^{[d]}$, and k between 2 and $K_t - 1$). This interval will depend on the observed test result as follows - if the test result is below the first cutpoint (i.e. in the first category), then the augmented data will be less than the first cutpoint parameter; if it is above the last cutpoint (i.e., in the last category), then the augmented data will be greater than the last cutpoint parameter; otherwise, if the test result falls between two cutpoints (i.e., the test result belongs to any other category), then the augmented data will fall between the corresponding cutpoint parameters. The measures of test accuracy are given by,

$$\begin{aligned} Se_{s,t,k} &= 1 - \Phi(\nu_{s,t}^{[1]} - C_{k,s,t}^{[1]}) \\ Sp_{s,t,k} &= \Phi(\nu_{s,t}^{[0]} - C_{k,s,t}^{[0]}) \end{aligned} \quad (3)$$

3.3 Between-study model

Now we will explain how we will model the variation in test accuracy between studies - called the *between-study heterogeneity*, as well as the correlation between the sensitivities and specificity between studies - called the *between-study correlation*. It is important to bear in mind the distinction from the within-study correlations (defined in section 3.2), which model the conditional dependence between tests. For each test t , we will assume that the study-specific means ($\nu_{s,t}^{[d]}$ - defined in equation 1 in section 3.2) arise from a bivariate normal (BVM) distribution with means $\mu_t^{[d]}$, between-study standard deviations $\sigma_t^{[d]}$, and between-study correlations ρ_t .

More specifically, $\boldsymbol{\nu}_{s,t} \sim \text{BVN}(\boldsymbol{\mu}_t, \boldsymbol{\Sigma}_t)$, where,

$$\boldsymbol{\mu}_t = \begin{pmatrix} \mu_t^{[1]} \\ \mu_t^{[0]} \end{pmatrix}, \boldsymbol{\Sigma}_t = \begin{pmatrix} (\sigma_t^{[1]})^2 & \rho_t \cdot \sigma_t^{[1]} \cdot \sigma_t^{[0]} \\ \rho_t \cdot \sigma_t^{[1]} \cdot \sigma_t^{[0]} & (\sigma_t^{[0]})^2 \end{pmatrix} \quad (4)$$

The model described in equation 4 is known as a *partial pooling* model (using the terminology from Gelman & Hill³⁰ - otherwise known as '*random-effects*'). These models allow the study-specific accuracy parameters across studies to inform one another, without assuming full homogeneity like a full pooling (i.e., "fixed-effects") would - which would allow no between-study variation in the means $\nu_{s,t}^{[d]}$. The disease prevalence's in each study, p_s , are modelled independently of each other, known as a *no pooling* model. There are several differences between partial pooling and no pooling models³¹. For example, the former uses less parameters than no pooling, which means that there is less likelihood of encountering parameter identifiability issues. An advantage of our partial pooling model is that allows us to summarise the results using the parameters which are shared across studies (see section 3.3), allowing us to more easily summarise test accuracy as well as the heterogeneity in accuracy between studies and correlation between sensitivities and specificities. We can incorporate meta-regression covariates into the model by extending the partial pooling model defined in equation 4 above - see supplementary material 1, meta-regression section (section 1.2.1) for details. We can assume that a given test is a perfect gold standard by setting $\mu_t^{[0]} = -5$ and $\mu_t^{[1]} = 5$, which correspond to approximately 100% sensitivity and specificity, respectively, and, by assuming a complete pooling model (i.e. setting $\sigma_t^{[d]} = 0$).

3.3.1 Within-study correlations

We will model the within-study correlation matrices ($\Psi_s^{[d]}$) defined in equation (1) in section 3.2 using a partial pooling model. As suggested by Goodrich³², this can be achieved by specifying each $\Psi_s^{[d]}$ as a weighted linear combination of a global 'average' correlation matrix across studies ($\Psi_G^{[d]}$), and a matrix of study-level deviations from this global matrix ($\Psi_s^{[d]\Delta}$), with weight $\beta^{[d]}$ which is between 0 and 1. More specifically, $\Psi_s^{[d]} = (1 - \beta^{[d]}) \cdot \Psi_G^{[d]} + \beta^{[d]} \cdot \Psi_s^{[d]\Delta}$. Note that we can also model the conditional dependence between only certain pairs of tests by setting the relevant terms for the other tests in $\Psi_G^{[d]}$ and $\Psi_s^{[d]\Delta}$ to zero.

3.3.2 Cutpoints

The cutpoint parameters can be modelled using an induced Dirichlet model, an approach which has been proposed by Betancourt³³. By taking advantage of the properties of a type of statistical distribution called a Dirichlet distribution, this model is able to map the latent cutpoint parameters in each study ($\{C_{1,s,t}^{[d]}, \dots, C_{K_t-1,s,t}^{[d]}\}$) defined in section 3.2 to a simplex (i.e., a vector whose elements sum to 1) of ordinal probabilities ($P_{1,s,t}^{[d]}, \dots, P_{K_t,s,t}^{[d]}$). Each probability $P_{k,s,t}^{[d]}$ corresponds to the probability that an individual's test result for test t falls in category k for study s . In this paper, we used a partial pooling model for the cutpoints across studies, enabling us to model the between-study heterogeneity in the cutpoints. We can also obtain 'average' cutpoints, ($C_t^{[d]}$) by using the posterior distribution of the induced Dirichlet partial pooling model, enabling us to obtain summary accuracy measures for ordinal tests. For the full details of this model, see supplementary material 1 and supplementary material 4.

3.3.3 Test accuracy summaries

For dichotomous tests, the summary sensitivity and specificity estimates for test t are given by evaluating equation (2) at the means of the between-study model (see equation (1)). More specifically, $Se_{G,t} = \Phi(\mu_t^{[1]})$, and $Sp_{G,t} = 1 - \Phi(\mu_t^{[0]})$. Similarly, for ordinal tests, the summary measures for test t at cutpoint k are given by evaluating equation (3) evaluated at the means of the partial pooling model (see equation (1)), and at the global (summary) cutpoints ($C_{k,t}^{[1]}$). That is, $Se_{G,t,k} = 1 - \Phi(C_{k,t}^{[1]} - \mu_t^{[1]})$, and $Sp_{G,t,k} = \Phi(C_{k,t}^{[0]} - \mu_t^{[0]})$. We can generate predictions for a 'new' ($S+1$)-th study by simulating a draw (at each iteration of the parameter sampler) from the posterior predictive distributions of the between-study normal hierarchical model, (see (4)), $\nu_{S+1,t}$, and a

new vector of cutpoints from the between-study cutpoint model (for more details, see section 1.2.4 in supplementary material 1). $\mathbf{C}_{S+1,t}^{[d]}$. Then, the predicted sensitivities and specificities for an $(S+1)$ -th study are given by $Se_{S+1,t} = \Phi(\nu_{S+1,t}^{[1]})$, $Sp_{S+1,t} = 1 - \Phi(\nu_{S+1,t}^{[0]})$ for dichotomous tests, and $Se_{S+1,t,k} = 1 - \Phi(C_{S+1,k,t}^{[1]} - \nu_{S+1,t}^{[1]})$, $Sp_{S+1,t,k} = \Phi(C_{S+1,k,t}^{[0]} - \nu_{S+1,t}^{[0]})$ for ordinal tests.

3.3.4 Joint test accuracy summaries

The summary estimates for the joint test accuracy of tests t and t' at cutpoints k and k' are given by:

$$\begin{aligned}
Se_{G,tt',kk'}^{BTN} &= Se_{G,t,k} * Se_{G,t',k'} + cov_{G,tt',kk'}^{[1]} \\
Sp_{G,tt',kk'}^{BTN} &= 1 - ((1 - Sp_{G,t,k}) * (1 - Sp_{G,t',k'}) + cov_{G,tt',kk'}^{[0]}) \\
Se_{G,tt',kk'}^{BTP} &= 1 - ((1 - Se_{G,t,k}) * (1 - Se_{G,t',k'}) + cov_{G,tt',kk'}^{[1]}) \\
Sp_{G,tt',kk'}^{BTP} &= Sp_{G,t,k} * Sp_{G,t',k'} + cov_{G,tt',kk'}^{[0]}
\end{aligned} \tag{5}$$

With BTN and BTP as defined in section 2. Note that for ordinal tests, the order of the tests can affect joint test accuracy estimates²⁸. However, for dichotomous tests it does not²⁸, although this is often still important for clinical practice. For example, the first test may be cheaper to carry out. The parameter $cov_{G,tt',kk'}^{[d]}$ is the global conditional covariance between all possible test-pairs. Obtaining these covariances requires us to calculate the global conditional correlations between each test-pair, which we will denote as $\rho_{G,tt',kk'}^{[0]}$; this assumes that the test results are the same form as the observed data. However, our model is parameterised in terms of the polychoric correlations ($\epsilon_{G,tt',kk'}^{[d]}$ - see equation (equation 1) in section 3.2). Therefore, in order to be able to estimate the joint test accuracy estimates in equation 5, we will need to convert from $\epsilon_{G,tt',kk'}^{[d]}$ to $\rho_{G,tt',kk'}^{[d]}$. For details on how this is achieved, please refer to section 1.2.4 of supplementary material 1.

3.4 Assessing model fit & model comparison

For our MVP-LC model, We can check how well our model predicts the data by using a technique called *posterior predictive checking* - where we generate data from our model, and compare it to the observed data. For example, we can plot the model-predicted test results against the observed test results for each test-pair^{5,4,6}. We also assessed model fit by plotting the model-predicted within-study correlations against the observed within-study correlations, using the correlation residual plot proposed by Qu et al¹³. For model comparison, we used leave-one-out (LOO) cross-validation³⁴ - an iterative procedure which removes part of the data and re-fits the model, and sees how well the model predicts the missing data. For more details on model comparison and posterior predictive checking, including relevant formulae, please refer to section 1.3 in supplementary material 1.

3.5 Model implementation

We implemented the models in R³⁵ using the probabilistic programming language Stan^{36,37} via the R package CmDStanR³⁸ using a PC with 32GB of RAM and an AMD Ryzen 3900X 12-core CPU with Linux Mint OS. To code the model in Stan, we extended the code for a standard binary multivariate probit model³². This is described in detail in Goodrich 2017³⁹ and is summarised in supplementary material 5. We implemented the between-study partial pooling model for the within-study correlations described in section 3.3.1 in Stan by using the function provided by Stephen Martin and Ben Goodrich⁴⁰. For the cutpoint between-study model, we used Betancourt's induced Dirichlet model³³ described in section 3.3.2; this is described in more detail in supplementary material 4, and this was implemented using code by Betancourt³³.

We ran all models using 4 chains until the split R-hat statistic was less than 1.05 for all parameters and the number of effective samples was satisfactory for all parameters⁴¹. We only reported results

when we obtained no warnings for divergent transitions or energy fraction of missing information (E-FMI), important diagnostics for geometric ergodicity³⁷. We used the CmDStanR diagnostic utility to check all of the aforementioned model diagnostics³⁸. We also inspected trace plots and plotted the posterior distributions to check they were not bimodal. Rather than using $\Phi(\cdot)$, which is prone to numerical instability, we can use the closely resembling logistic function, $\Phi'(x) = \frac{1}{1+e^{-1.702 \cdot x}}$, which has an absolute maximum deviation from $\Phi(\cdot)$ of 0.0095. This is the same probit approximation used for the meta-analysis of dichotomous tuberculosis tests using latent trait models in Sadatsafavi et al⁸. The data, Stan model code, and R code to reproduce the results and figures for the case study application in section 4 is provided at .

4 Application to case study

Since our model is Bayesian, we must formulate a *prior model* - that is, specify prior distributions for the model parameters defined in section 3. We describe this prior model in 4.1. When faced with the task of analysing a dataset with an imperfect gold standard which contains test accuracy data from an ordinal test, in order to be able to apply proposed methods for meta-analysis without assuming a gold standard^{5,7}, one must first dichotomise the data at each cutpoint and conduct a series of stratified analyses. We applied a priori dichotomisation technique using our proposed MVP-model in section 4.2. Finally, in section 4.3, we applied the models proposed in section 3, but without dichotomising the Well's score.

In section 4.1, we will index the gold standard (ultrasound), the D-Dimer, and the Wells' score by $t = 1, t = 2$, and $t = 3$, respectively. In sections 4.2 and 4.3 we will denote summary estimates as " X [Y, Z]", where X is the posterior median and [Y, Z] is the 95% posterior interval.

4.1 Prior distributions

For the summary-level accuracy parameter for the gold standard test (ultrasound - i.e. $\mu_{t=1}^{[d]}$), we constructed informative priors using subject-matter knowledge for ultrasound, based on meta-analyses from the literature^{24,25} (see supplementary material 2 for more details). These priors correspond to 95% prior intervals of (0.49, 0.94) and (0.82, 0.99) for the sensitivity and specificity, respectively. On the other hand, for the summary-level accuracy parameters for the D-Dimer and the Wells score (i.e., $\mu_{t=2}^{[d]}$ and $\mu_{t=3}^{[d]}$), we specified priors conveying very little information - equivalent to assuming a 95% prior interval of (0.04, 0.96) for the sensitivities and specificities.

For the between study deviation parameters for all three tests (i.e., for $\sigma_t^{[d]}$ for $t = \{1, 2, 3\}$ - see equation (4) in section 3.3), we used weakly informative priors corresponding to a 95% prior interval of (0.02, 1.09). The priors are weakly informative since they weakly pull the study-specific sensitivities and specificities towards each other, whilst allowing for large between-study heterogeneity if the data demands. For example, if 0.8 is the value found for the summary sensitivity, and the data suggests a standard deviation equal to 2 (corresponding to a high degree of between-study heterogeneity), then these priors would allow the study-specific sensitivities and specificities to be in the interval (0.44, 0.97) with 95% probability. We also used weak priors for the between-study correlation parameters for all tests (i.e., ρ_t for $t = \{1, 2, 3\}$ - see equation (4) in section 3.3), corresponding 95% prior probability interval of (-0.82, 0.82). Finally, for conditional dependence models, for the within-study correlation parameters (see section 3.2 and equation 1), we used priors which correspond to 95% prior intervals of (-0.65, 0.65) for both the global 'average' correlation matrices ($\Omega_G^{[d]}$) and the study-specific deviation matrices ($\Omega_s^{[d]\Delta}$), respectively. These are weakly informative and allow a moderately large between-study deviation in the strength of the conditional dependence between tests. For more detail on these prior distributions, please see supplementary material 2.

4.2 The pitfalls of a priori dichotomisation in the presence of an imperfect gold standard

We consider two dichotomisations of the Wells score. For the first, we dichotomised the Wells' score by grouping together those patients who obtain a score of 'low' or 'moderate' as a negative result and those who scored 'high' as positive. On the other hand, for the second dichotomisation, we grouped together patients who scored 'moderate' or 'high' and considered this as a positive result, and those who scored 'low' as a negative result. We will refer to the former dichotomisation as "low + moderate vs high" and the latter as "low vs moderate + high". We applied this technique to this dataset, to allow comparison with our "full" model, using the models proposed in section 3, fitting both conditional independence (CI) and dependence (CD) models, the results of which are shown in in section

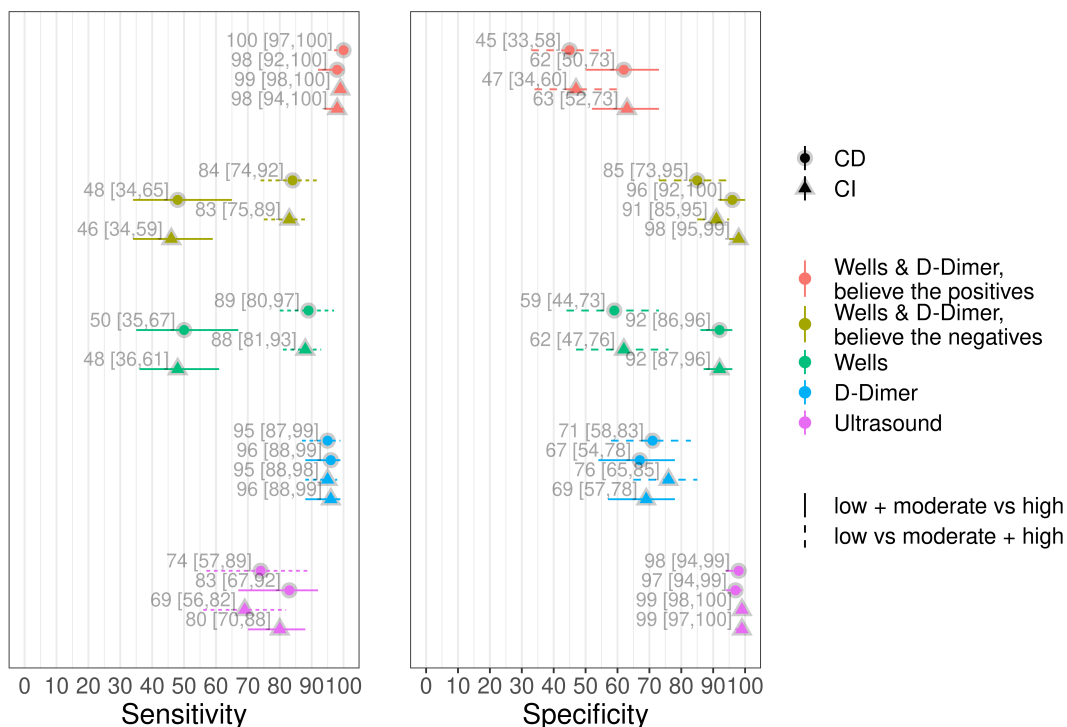


Figure 1: Posterior medians and 95% posterior intervals for models dichotomising the Well's score. Note: CD = Conditional Dependence; CI = Conditional Independence

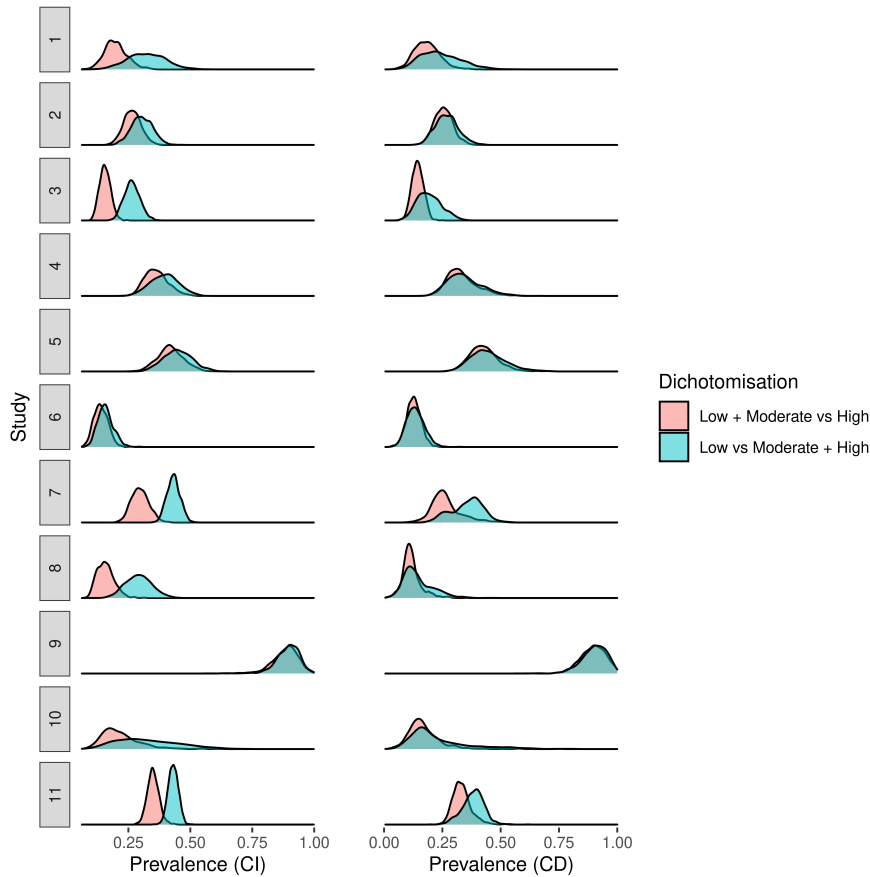


Figure 2: Posterior density plots for disease prevalence parameters. Note: CD = Conditional Dependence; CI = Conditional Independence

When assuming conditional independence between all three tests, we see that (figure 1) some of the estimates of the accuracy of the other two tests change substantially depending on whether we dichotomise the Wells score as low+moderate vs high, or as low vs moderate+high. For the former dichotomisation, the sensitivity of ultrasound was estimated as 0.80 [0.70, 0.88] whereas for the latter it was 0.69 [0.56, 0.82]. The specificity of ultrasound and the sensitivity of the D-Dimer were similar between both dichotomisations. However, there was a notable difference in the specificities of the D-Dimer test, where we obtained specificities of 0.69 [0.57, 0.78] and 0.76 [0.65, 0.85] for the low+moderate vs high and low vs moderate+high dichotomisations, respectively.

The differences in the results were similar when modelling conditional dependence between the three tests (see figure 1). In the low+moderate vs high dichotomisation, for the ultrasound sensitivity we obtained 0.83 [0.67, 0.92] and for the low vs moderate + high dichotomisation 0.74 [0.57, 0.89]. For the D-Dimer specificities, we obtained 0.67 [0.54, 0.78] and 0.71 [0.58, 0.83] for the low+moderate vs high and low vs moderate+high dichotomisations, respectively. As with conditional independence, the specificity of the ultrasound and the sensitivity of the D-Dimer were similar between the two dichotomisations. We can also see the estimates of disease prevalence increase for most studies for the low vs moderate + high dichotomisation relative to the low + moderate vs high dichotomisation, for both conditional independence (left panel of figure 2) and dependence models (right panel of figure 2).

Overall, regardless of whether we assume conditional independence or dependence, some of the accuracy estimates change notably depending on how we dichotomise the Wells score. This is not surprising, since imperfect gold standard models based on latent class analysis utilise the full distribution of test responses from all tests to estimate accuracy and disease prevalence⁴. This simple example demon-

strates the importance of modelling all the available data for ordinal non-dichotomous tests, such as the Wells score, in the presence of an imperfect gold standard, as opposed to simply conducting multiple stratified analyses at each cutpoint of the ordinal test using simpler methods. This observation serves to motivate the implementation of ordinal regression into the models to appropriately model the ordinal nature of the Wells score.

4.3 Modelling the Wells score as an ordinal test

Now we fit the models without dichotomising the Wells score, by simultaneously modelling all three categories. For these models, we used weakly informative priors of $\mu_3 \sim N(0, 1)$ for the mean parameters for the Wells test. We used the partial pooling model on the Wells score cutpoint parameters (see section 3.2, equation 3). For the Dirichlet population parameters, we used a weakly informative prior $\kappa^{[d]} \sim N_{\geq 0}(0, 50)$. This allows considerable asymmetry in the Dirichlet population vector α_k , as can be seen from the prior predictive check (see figure 1 in section 1.3 in supplementary material 1). The rest of the priors were the same as those discussed in section 4.1.

We fit the following models: one assuming that ultrasound is a perfect gold standard and conditional independence between all three tests (**M1**); the same model but modelling the conditional dependence between the Well’s score and D-Dimer (**M2**); a model assuming ultrasound to be an imperfect gold standard and conditional independence between all three tests (**M3**); and a variation of M3 which modelled the conditional dependence between all three tests (**M4**).

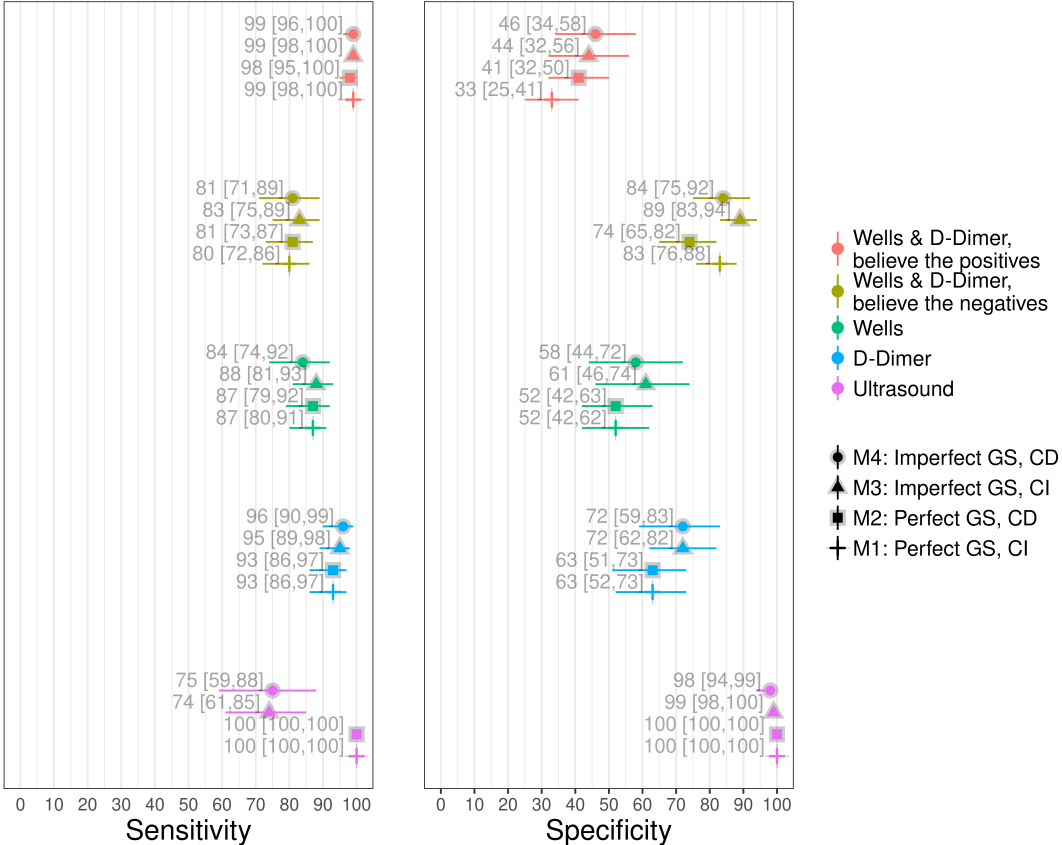


Figure 3: Posterior medians and 95% posterior intervals for summary sensitivities and specificities, for models 1 - 4. Note: The Wells score summary estimates are dichotomised as low vs moderate + high. CD = Conditional Dependence; CI = Conditional Dependence; GS= Gold Standard.

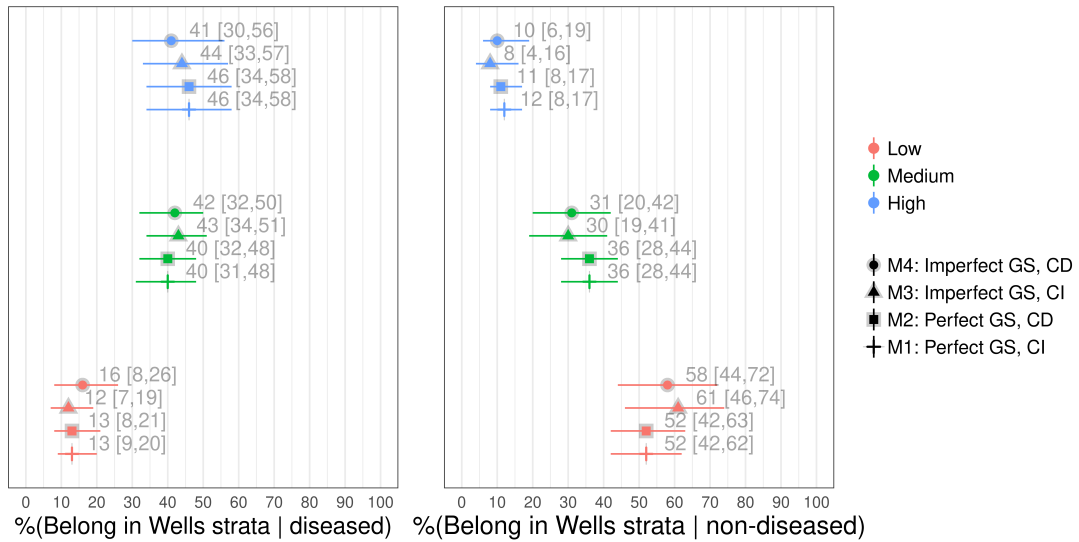


Figure 4: Posterior medians and 95% posterior intervals for the Well’s score stratum, for models 1 - 4. Note: CD = Conditional Dependence; CI = Conditional Dependence; GS= Gold Standard.

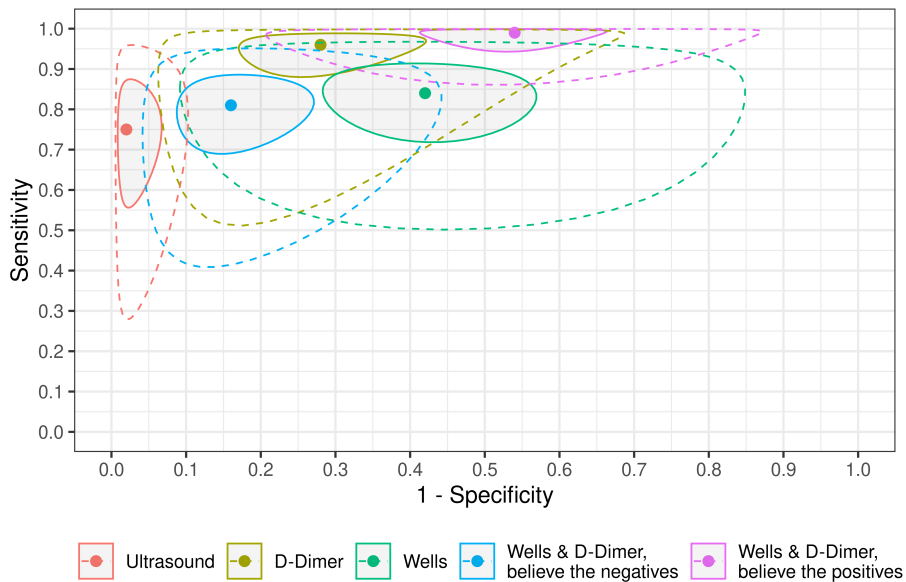


Figure 5: Summary Receiver Operating Characteristic (sROC) plot for M4. Shaded regions represent 95% posterior regions and regions surrounded by dashed lines represent 95% prediction regions. Note: The Wells score summary estimates are dichotomised as low vs moderate + high

The results for the summary sensitivity and specificity estimates for the four models are shown in figure 3, and the results for each of the Wells score strata are shown in figure 4. The estimates for the two models assuming a perfect gold standard (M1 and M2) are within 2% of those obtained from Novielli et al¹⁷. The similarity of the results is not surprising, since despite using different models and different link functions (logit vs approximate probit), both models assume that ultrasound is perfect.

For both the conditional independence (M1) and dependence (M2) models which assumed that ultrasound is a perfect reference test, the results we obtained for the accuracy of the BTP and BTN testing strategies for the Wells & D-Dimer tests were similar to those obtained by Novielli et al¹⁷. More specifically, for the BTP testing strategy, we found that the summary specificity estimates for M1 (33 [25, 41]) were around 8% lower than M2 (41 [32, 50]). For the BTN strategy, we found that

the estimates for M1 (74 [65, 82]) were around 9% higher (83 [76, 88]) than M2.

When we modelled ultrasound as an imperfect test, the summary estimate for the sensitivity of the Wells test for the model assuming conditional independence (M3 - 88 [81, 93]) was around 4% higher than the model which modelled conditional dependence (M4 - 84 [74, 92]), and around 5% higher for the sensitivity of the Wells & D-Dimer BTN testing strategy (89 [83, 94] and 84 [75, 92] for M3 and M4, respectively). The other differences between M3 and M4 were 3% or less (see figure 3).

Whilst assuming conditional dependence, the model assuming ultrasound is perfect estimated the specificity of the Wells score to be around 6% lower than the conditional dependence model (52 [42, 63] and 58 [44, 72] for M2 and M4, respectively). We also found that the specificity of the D-Dimer was around 9% lower (63 [51, 73] and 72 [59, 83] for M2 and M4, respectively), that the specificity was around 5% lower for the Wells & D-Dimer BTP testing strategy (41 [32, 50] and 46 [34, 58] for M2 and M4, respectively), and that the specificity of the BTN testing strategy was around 10% lower (74 [65, 82] and 84 [75, 92] for M2 and M4, respectively).

The summary receiver operating characteristic plot for M4 is shown in figure 5. The prediction regions suggest that there is substantial between-study heterogeneity for the sensitivity and specificity for most estimates. However, we found relatively narrow prediction regions for the specificity of ultrasound and for the Wells and D-Dimer BTP testing strategy, so we can be more confident in generalising our inferences for these estimates.

The LOO-CV results for all of the models are shown in table 2. The results suggested that M1 has the poorest fit. Modelling the dependency between the D-Dimer and Wells tests (M2) improved the fit (LOO-IC = 16038.6 and 15819.0 for M1 and M2, respectively). Out of the two models not assuming a perfect gold standard, conditional independence model gave a worse fit than the conditional dependence model (difference in ELPD between M3 and M4 = -31.4, se = 6.4). The two conditional dependence models were the two best fitting models, with the conditional dependence model giving the best fit (difference in ELPD between M2 and M4 = -20.8, se = 6.2). The posterior predictive checks for this model are shown in figure 6 (correlation residual plot) and figure 2 in supplementary material 3 (2x2 table count residual plot). Both plots show that the model fits the data well.

Table 2: Leave-One-Out Cross Validation (LOO-CV) for comparison of model fit for case study 1 dataset

Model ¹	LOO-IC ²	$ELPD_{M4} - ELPD_{M_i}$	$^{3,4}se(ELPD_{M4} - ELPD_{M_i})$ ⁴
4 (Imperfect ultrasound + CD)	15,777.4	0	0
2 (Perfect ultrasound + CD)	15,819.0	-20.8	6.2
3 (Imperfect ultrasound + CI)	15,840.1	-31.4	6.4
1 (Perfect ultrasound + CI)	16,038.6	-130.6	15.4

¹ Models are ordered from best to worst fitting

² LOO-IC = Leave-One-Out Information Criterion; note that LOO-IC is on the deviance scale

³ ELPD = Estimated Log pointwise Predictive density for a new Dataset

⁴ M_i denotes the i th model

CI = Conditional Independence; CD = Conditional Dependence

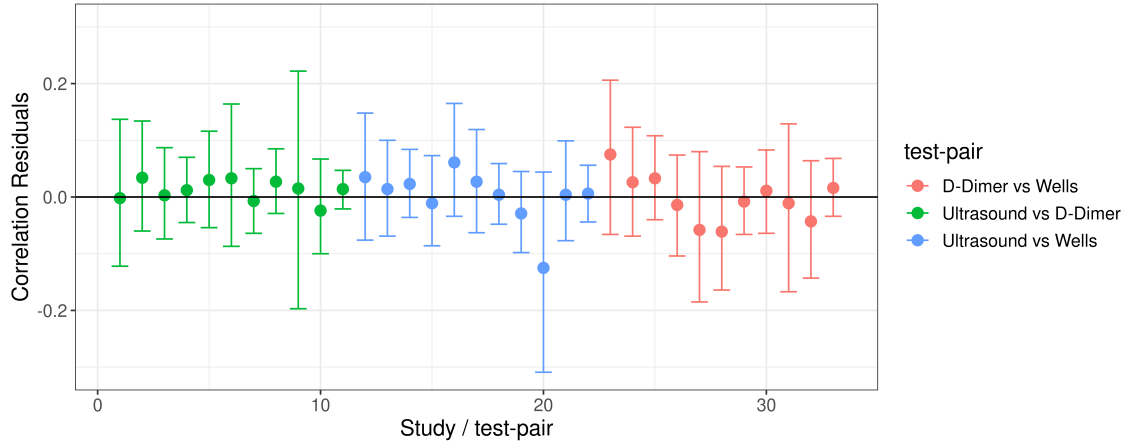


Figure 6: Posterior predictive check for model 4; correlation residual plot

5 Discussion

5.1 Summary

Our proposed MVP-LC model addresses the novel problem of carrying out meta-analysis of two or more conditionally dependent tests when there is no perfect gold standard, for the case where there are both ordinal and dichotomous test(s) under evaluation, and estimation of joint test accuracy is of interest.

Using the case study as a demonstrative aid for the model (see section 4), we showed why treating ordinal tests as dichotomous in the context of an imperfect gold standard is suboptimal (see section 4.2). When we modelled the Wells test as ordinal and treated ultrasound as a perfect gold standard (see section 4.3), the summary estimates from Novielli et al¹⁷ are replicated in our findings. However, we found that the most complex model - which treated ultrasound as an imperfect gold standard in addition to modelling the conditional dependence between tests - had the best fit to the data. For this model, our estimates of test accuracy differed considerably compared to other models we fit (which gave worse fit to the data) and compared to the results obtained in the analysis conducted by Novielli et al¹⁷. In particular, we obtained considerably different estimates of specificity for both the D-Dimer and the Wells score tests when used alone, and for the joint specificity of the Wells and D-Dimer BTN testing strategy. However, the large between-study heterogeneity limited the generalisability of our results.

5.2 Potential applications

The methods we have developed in this paper have a wide scope of applicability in clinical practice, further than just DVT. For instance, Hamza et al⁴² re-analysed a meta-analysis⁴³, which assessed the accuracy of the CAGE questionnaire⁴⁴ - a 4-category ordinal test used as a screening tool to detect individuals who may be suffering from alcoholism. However, their model assumed a perfect gold standard¹. Our proposed MVP-LC model could be used to more appropriately estimate the accuracy of the CAGE questionnaire, since we would not need to assume that the reference test in each study is perfect.

The methods could also be used to more appropriately assess joint testing strategies. For instance, current UK Health Security Agency guidance⁴⁵ states that individuals who have symptoms suggestive of the severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2) and test positive using Lateral flow tests (LFTs) should be considered as positive and require no subsequent testing. On the other hand, it states that individuals with negative LFTs should be assessed with a polymerase chain reaction (PCR), with only those who also test negative on PCR being considered negative. Our methods could

be used to investigate this joint test accuracy strategy without modelling PCR as an imperfect gold standard, particularly with respect to its sensitivity. For depression screening, one potential BTN testing strategy is one in which individuals undertake a very brief 2-item version of the Patient Health Questionnaire (PHQ-2⁴⁶) followed by the 9-item version (PHQ-9⁴⁷). This was investigated recently by Levis et al⁴⁸; however, they assumed perfect gold standards, and they only used around half of the available studies, since they discarded those studies which used inferior gold standards. Our MVP-LC model could be used to analyse these data without assuming a perfect gold standard whilst accounting for differences between reference tests with meta-regression, and using all of the available data. Furthermore, we would be able to model the differences in gold standards between studies using meta-regression (see section 3.3).

5.3 Advantages

Our proposed Bayesian MVP-LC model addresses some important limitations which are present in models based on TLCMs^{4,49,5,6,7}. For example, although TLCMs have fast run times due to being computationally inexpensive, and they can model the conditional dependence between tests⁵⁰, an important limitation is that, unlike our proposed MVP-LC model, they cannot appropriately model ordinal tests. For example, if one wishes to simultaneously model ordinal tests whilst modelling conditional dependence, they would first need to dichotomise the data a priori. As we showed in section 4.2, this is suboptimal in the context of an imperfect gold standard, since the test accuracy and disease prevalence estimates were varied depending on which cutpoint we dichotomise the data at. A limitation of TLCMs which have been proposed for meta-analysis^{5,6,7} is that, unless one assumes a complete pooling model between studies, it is not possible estimate summary correlation parameters - parameters which are required to estimate summary-level joint test accuracy. This is due to the fact that, in contrast to our MVP-LC model (which uses the within-study correlations), TLCMs model the conditional dependence using the within-study covariances, making it difficult to construct a partial pooling model for the within-study conditional dependence parameters. These covariances have bounds based on the sensitivity and specificity parameters in each study^{49,51}. Therefore, any summary-level covariance parameters obtained would be questionable. Our MVP-LC model also has advantages over more advanced models for meta-analysis of test accuracy, such as the model proposed by Sadatsafavi et al⁸, which is also based on multivariate probit regression and is an extension of the latent trait model¹³. Two important limitations of this model - not present in our MVP-LC model - is that it can only model dichotomous data, and it assumes that the within-study correlations are fixed across studies. Furthermore, since our proposed MVP-LC model is an extension of the model for single studies proposed by¹¹, another benefit over the model by Sadatsafavi et al⁸ is that it can also be used to specify more general correlation structures (by setting certain correlations to zero - see section 3.2). The fact that our model is Bayesian means that one can incorporate subject-matter knowledge into the model, as we did for our case study. Furthermore, the Induced Dirichlet partial pooling model³⁷ (see section 3.3.2) and supplementary material 4) for the ordinal tests makes it possible to specify priors for ordinal tests and obtain summary estimates.

5.4 Limitations

When applying the model to our case study dataset (see section 2), we used available subject-matter knowledge²⁴ to construct informative prior distributions for the gold standard test (ultrasound), and weakly informative priors for other parameters (see section 4.1 and supplementary material 1). Attempts to conduct sensitivity analysis using more diffuse priors led to diagnostic errors. This is likely due to the fact that Stan is quite sensitive at detecting non-identifiabilities in the posterior distributions³⁶, and non-identifiability is more likely to occur with less informative priors, particularly for latent class models due to the large number of parameters relative to the data. Another limitation of our case study analysis is that, although our model can easily incorporate meta-regression coefficients (see supplementary material 1), the case study dataset did not contain any study-level covariates, since primary studies did not report sufficient data. In an ideal world where such data were available, a more principled analysis could be carried out by using a meta-regression covariate for the proportion

of patients who have proximal versus distal DVT, which would have enabled us to model the variation of ultrasound sensitivity that exists between the two DVT groups in Novielli et al's¹⁷ data.

A limitation of our model, which is present across all imperfect gold standard methods based on latent class models (including TLCMs), is that full cross-classification tables (i.e. the full distribution of test results) are required for each study. This is a potential barrier to the uptake of our proposed MVP-LC model, as this data is frequently not reported for studies evaluating 3 or more tests and/or studies assessing ordinal tests. One way in which we could have assessed the general performance of our MVP-LC model is by running a simulation study⁵². A simulation study comparing our proposed MVP-LC model to other models would also be very useful. However, it is important to note that, at the time of writing, no other models have been proposed to simultaneously meta-analyse both dichotomous and ordinal tests without assuming a perfect gold standard. That being said, a simulation study would still be useful, since we could compare the performance of our model to other proposed models which do assume a perfect gold standard (e.g., Novielli et al¹⁷) under a variety of different scenarios.

Although our proposed MVP-LC model offers considerable benefits in comparison to the more commonly used TLCM models^{5,6,7,50} (see section 5.3), we found that our proposed model was considerably less time efficient than TLCM models. Although this was not prohibitive for the case study used in this paper², our MVP-LC model may be intractable for larger sample sizes. Speeding up models based on augmented continuous data, such as our MVP-LC model, is an active area of research^{53,54,55,56,57,58,59,60}. An important area for future research would be to apply the models developed in this paper using these more efficient algorithms, which would make our proposed MVP-LC model more suitable for general use with larger meta-analyses, and it would also make it easier to conduct more meaningful simulation studies.

5.5 Future work

Models for the meta-analysis of test accuracy which can incorporate patient-level covariates - otherwise known as individual patient data (IPD) - have been proposed⁶¹, but only for dichotomous data and they assume a perfect gold standard⁶¹. Modelling IPD can lead to results which are more applicable to clinical practice as they can more easily be applied to patients when there is between-study heterogeneity, rather than only providing summary estimates which relate to some "average" patient. Extending our model to incorporate IPD would be relatively straightforward, since our model uses the patient-level data (as reconstructed from the reported contingency tables) as opposed to aggregated data for each study. It is straightforward to extend our model to the case where not all studies are assessing the same number of tests, using direct comparisons only. This could be further extended to allow indirect comparisons (network meta-analysis [NMA]), by assuming tests are missing at random (MAR)⁶², and extending the between-study model described in section 3.3 to an arm-based network-meta analysis model^{63,64}. Another straightforward modelling extension would be to incorporate data from ordinal tests which have missing data for some categories.

Our model could also be extended to synthesize data from ordinal tests for the case where some (or all) studies do not report data for every cutpoint - which is common in research. One could formulate such a 'missing cutpoint' version of our MVP-LC model by extending the partial pooling between-study cutpoint model (see section 3.3.2 and supplementary material 4), and viewing the cutpoints as MAR. Another possible 'missing cutpoint' model could be constructed by modelling the cutpoint parameters as the same in the diseased and non-diseased classes, and assume that they are fixed between studies by using a no pooling model. Then, as opposed to our MVP-LC model, in which the within-study variances are set to 1 to ensure parameter identifiability (see section 3.2), the no pooling cutpoint model would allow us to introduce within-study variance parameters and model them using a partial pooling model without encountering significant identifiability issues. These within-study parameters could be set to vary between the two latent classes, which would result in a smooth, non-symmetric receiver operating characteristic (ROC) curve. Another possible 'missing cutpoint' approach would be one based on the model proposed by Dukic et al⁶⁵, which assumes a perfect gold standard. This model also results in a smooth, non-symmetric ROC curve, since it assumes that the cutpoints vary

between studies and are the same in the diseased and non-diseased class. However, it would be more parsimonious since it assumes that the sensitivity is some location-scale change of the false positive rate.

For the case where studies report thresholds at explicit numerical cutpoints (as is sometimes reported for continuous tests, such as biomarkers), some 'missing threshold' methods which assume a perfect gold standard have been proposed^{66,67}. Rather than modelling the cutpoints as parameters, these methods assume that the cutpoints are constants, equal to the value of the numerical cutpoint, and they estimate separate location and scale parameters in each study and disease class. Our MVP-LC model could be extended to achieve this without assuming a gold standard. An important area for future research would be to construct other models which can be used for the same purposes as our proposed MVP-LC model. For instance, a multivariate logistic regression model could be constructed by using the Bayesian multivariate logistic distribution proposed by O' Brien et al⁶⁸. Such a model would use logistic link functions as opposed to probit (or approximate probit) links like our MVP-LC model, which are more numerically stable than probit and may give better fit to some datasets. Another multivariate regression approach would be to use copulas^{69,70,71}. Besides multivariate regression based on augmented data, another approach to modelling conditionally dependent ordinal diagnostic tests without assuming a perfect gold standard is log-linear models¹¹. These models can account for higher-order correlations¹¹. However, this requires estimation of additional parameters, so it is likely to introduce identifiability issues. Similarly to the multivariate probit models utilised in this paper, it may be possible to extend these models to meta-analyse multiple, imperfect diagnostic tests with multiple cutpoints.

Highlights

What is already known?: Standard, well-established methods exist for the synthesising estimates (i.e., conducting a meta-analysis) of test accuracy. These methods estimate test accuracy by comparing test results to some test which is assumed to be perfect - which is referred to as a 'gold standard' test. However, in clinical practice, these tests are often imperfect, which can cause estimates of the tests being evaluated to be biased and potentially lead to the wrong test being used in clinical practice. Meta-analytic methods, which do not assume a gold standard, have previously been proposed, but only for dichotomous tests.

What is new?: We developed a model which allows one to simultaneously meta-analyse ordinal and dichotomous tests without assuming a gold standard. The model also allows one to obtain summary estimates of the accuracy of two tests used in combination (i.e. joint test accuracy).

Potential impact for Research Synthesis Methods readers outside the authors field: The methods are widely applicable. For instance, psychometric measures and radiologic tests are typically ordinal, and the studies assessing these tests often do not use a gold standard; hence, applying standard models to these datasets may lead to misleading conclusions. The methods we proposed may lead to less biased accuracy estimates, and hence potentially a better understanding of which tests, and test combinations, should be used for these conditions.

Acknowledgements

The authors would like to thank Elpida Vounzoulaki for proofreading the manuscript. The authors would also like to thank various members of the Stan community forums (see <https://discourse.mc-stan.org/>) including Ben Goodrich, Michael Betancourt, Stephen Martin, Staffan Betnér, Martin Modrák, Niko Huurre, Bob Carpenter and Aki Vehtari and for providing functions which were utilised in the models and for useful discussions.

Funding: The work was carried out whilst EC was funded by a National Institute for Health Research (NIHR) Complex Reviews Support Unit (project number 14/178/29) and by a National Institute for Health Research Systematic Review Fellowship (project number RM-SR-2017-09-023). The views and opinions expressed herein are those of the authors and do not necessarily reflect those of the NIHR, NHS or the Department of Health. The NIHR had no role in the design of the study and collection, analysis, and interpretation of data and in writing the manuscript. This project is funded by the NIHR Applied Research Collaboration East Midlands (ARC EM). The views expressed are those of the authors and not necessarily those of the NIHR or the Department of Health and Social Care.

Data availability statement

The Data, R and Stan code to reproduce the results and figures from section 4 is available on Github at:

<https://github.com/CerulloE1996/dta-ma-mvp-1>.

References

- [1] Johannes B. Reitsma et al. "Bivariate analysis of sensitivity and specificity produces informative summary measures in diagnostic reviews". In: *Journal of Clinical Epidemiology* (2005). ISSN: 08954356. DOI: 10.1016/j.jclinepi.2005.02.022.
- [2] Carolyn M. Rutter and Constantine A. Gatsonis. "A hierarchical regression approach to meta-analysis of diagnostic test accuracy evaluations". In: *Statistics in Medicine* (2001). ISSN: 02776715. DOI: 10.1002/sim.942.
- [3] Roger M. Harbord et al. "A unification of models for meta-analysis of diagnostic accuracy studies". In: *Biostatistics* (2007). ISSN: 14654644. DOI: 10.1093/biostatistics/kxl004.

- [4] S. L. Hui and S. D. Walter. “Estimating the Error Rates of Diagnostic Tests”. In: *Biometrics* (1980). ISSN: 0006341X. DOI: 10.2307/2530508.
- [5] Haitao Chu, Sining Chen, and Thomas A. Louis. “Random effects models in a meta-analysis of the accuracy of two diagnostic tests without a gold standard”. In: *Journal of the American Statistical Association* (2009). ISSN: 01621459. DOI: 10.1198/jasa.2009.0017.
- [6] J. Menten, M. Boelaert, and E. Lesaffre. “Bayesian meta-analysis of diagnostic tests allowing for imperfect reference standards”. In: *Statistics in Medicine* (2013). ISSN: 10970258. DOI: 10.1002/sim.5959.
- [7] Nandini Dendukuri et al. “Bayesian Meta-Analysis of the Accuracy of a Test for Tuberculous Pleuritis in the Absence of a Gold Standard Reference”. In: *Biometrics* (2012). ISSN: 0006341X. DOI: 10.1111/j.1541-0420.2012.01773.x.
- [8] Mohsen Sadatsafavi et al. “A statistical method was used for the meta-analysis of tests for latent TB in the absence of a gold standard, combining random-effect and latent-class methods to estimate test accuracy”. In: *Journal of Clinical Epidemiology* (2010). ISSN: 08954356. DOI: 10.1016/j.jclinepi.2009.04.008.
- [9] Jian Kang, Rollin Brant, and William A. Ghali. “Statistical methods for the meta-analysis of diagnostic tests must take into account the use of surrogate standards”. In: *Journal of Clinical Epidemiology* (2013). ISSN: 18785921. DOI: 10.1016/j.jclinepi.2012.12.008.
- [10] Huiping Xu and Bruce A. Craig. “A probit latent class model with general correlation structures for evaluating accuracy of diagnostic tests”. In: *Biometrics* (2009). ISSN: 0006341X. DOI: 10.1111/j.1541-0420.2008.01194.x.
- [11] Huiping Xu, Michael A. Black, and Bruce A. Craig. “Evaluating accuracy of diagnostic tests with intermediate results in the absence of a gold standard”. In: *Statistics in Medicine* (2013). ISSN: 02776715. DOI: 10.1002/sim.5695.
- [12] John S. Uebersax. “Probit Latent Class Analysis with Dichotomous or Ordered Category Measures: Conditional Independence/Dependence Models”. In: *Applied Psychological Measurement* 23.4 (1999), pp. 283–297. DOI: 10.1177/01466219922031400.
- [13] Yinsheng Qu, Ming Tan, and Michael H. Kutner. “Random Effects Models in Latent Class Analysis for Evaluating Accuracy of Diagnostic Tests”. In: *Biometrics* (1996). ISSN: 0006341X. DOI: 10.2307/2533043.
- [14] Yinsheng Qu and Alula Hadgu. “A model for evaluating sensitivity and specificity for correlated diagnostic tests in efficacy studies with an imperfect reference test”. In: *Journal of the American Statistical Association* (1998). ISSN: 1537274X. DOI: 10.1080/01621459.1998.10473748.
- [15] James H. Albert and Siddhartha Chib. “Bayesian Analysis of Binary and Polychotomous Response Data”. In: *Journal of the American Statistical Association* (1993). ISSN: 01621459. DOI: 10.2307/2290350.
- [16] WH William H . Greene. *Econometric analysis 7th Ed.* 2012. ISBN: 978-0-273-75356-8.
- [17] Nicola Novielli, Alexander J. Sutton, and Nicola J. Cooper. “Meta-analysis of the accuracy of two diagnostic tests used in combination: Application to the ddimer test and the wells score for the diagnosis of deep vein thrombosis”. In: *Value in Health* (2013). ISSN: 10983015. DOI: 10.1016/j.jval.2013.02.007.
- [18] Jonathan Stone et al. “Deep vein thrombosis: pathogenesis, diagnosis, and medical management”. In: *Cardiovascular Diagnosis and Therapy* 7 (Suppl 3 Dec. 2017), S276–S284. ISSN: 2223-3652. DOI: 10.21037/cdt.2017.09.01. URL: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5778510/> (visited on 08/01/2021).
- [19] Clive Tovey and Suzanne Wyatt. “Diagnosis, investigation, and management of deep vein thrombosis”. In: *BMJ* 326.7400 (2003), pp. 1180–1184. ISSN: 0959-8138. DOI: 10.1136/bmj.326.7400.1180.
- [20] Paul A. Kyrle and Sabine Eichinger. “Deep vein thrombosis”. In: *The Lancet* 365.9465 (2005), pp. 1163–1174. ISSN: 0140-6736. DOI: 10.1016/S0140-6736(05)71880-8.
- [21] C. Kearon et al. “Noninvasive diagnosis of deep venous thrombosis. McMaster Diagnostic Imaging Practice Guidelines Initiative”. In: *Annals of Internal Medicine* 128.8 (Apr. 15, 1998), pp. 663–677. ISSN: 0003-4819. DOI: 10.7326/0003-4819-128-8-199804150-00011.

- [22] Vincent B. Ho et al. “ACR Appropriateness Criteria(®) on suspected lower extremity deep vein thrombosis”. In: *Journal of the American College of Radiology: JACR* 8.6 (June 2011), pp. 383–387. ISSN: 1558-349X. DOI: 10.1016/j.jacr.2011.02.016.
- [23] Seung-Kee Min et al. “Diagnosis and Treatment of Lower Extremity Deep Vein Thrombosis: Korean Practice Guidelines”. In: *Vascular Specialist International* 32.3 (Sept. 2016), pp. 77–104. ISSN: 2288-7970. DOI: 10.5758/vsi.2016.32.3.77.
- [24] Steve Goodacre et al. “Systematic review and meta-analysis of the diagnostic accuracy of ultrasonography for deep vein thrombosis”. In: *BMC Medical Imaging* (2005). ISSN: 14712342. DOI: 10.1186/1471-2342-5-6.
- [25] M. Di Nisio et al. “Accuracy of diagnostic tests for clinically suspected upper extremity deep vein thrombosis: A systematic review”. In: *Journal of Thrombosis and Haemostasis* (2010). ISSN: 15387933. DOI: 10.1111/j.1538-7836.2010.03771.x.
- [26] P. S. Wells et al. “Value of assessment of pretest probability of deep-vein thrombosis in clinical management”. In: *Lancet (London, England)* 350.9094 (Dec. 20, 1997), pp. 1795–1798. ISSN: 0140-6736. DOI: 10.1016/S0140-6736(97)08140-3.
- [27] Faizan Khan et al. “Venous thromboembolism”. In: *The Lancet* 398.10294 (July 3, 2021), pp. 64–77. ISSN: 0140-6736. DOI: 10.1016/S0140-6736(20)32658-1. URL: <https://www.sciencedirect.com/science/article/pii/S0140673620326581> (visited on 08/01/2021).
- [28] N. Novielli, N. J. Cooper, and A. J. Sutton. “Evaluating the cost-effectiveness of diagnostic tests in combination: is it important to allow for performance dependency?” In: 16.4 (2013), pp. 536–41.
- [29] Joakim Ekström. “A Generalized Definition of the Polychoric Correlation Coefficient”. In: *Department of Statistics, UCLA* (2011).
- [30] Andrew Gelman and Jennifer Hill. *Data Analysis Using Regression and Multilevel/Hierarchical Models*. 2006. DOI: 10.1017/cbo9780511790942.
- [31] Michael Betancourt. *Hierarchical Modeling*. https://github.com/betanalphabet/knitr_case_studies/tree/master/hierarchical_modeling. commit 27c1d260e9ceca710465dc3b02f59f59b729ca43. 2020.
- [32] Ben Goodrich. *A better parameterization of the multivariate probit model*. <https://github.com/stan-dev/example-models/commit/d6f0282d64382b627dfddca6b7f9a551bda3f537>. 2016.
- [33] Michael Betancourt. *Ordinal Regression*. https://github.com/betanalphabet/knitr_case_studies/tree/master/ordinal_regression. commit 23eb263be4cfb44278d0dfb8ddb593a4b142506. 2019.
- [34] Aki Vehtari, Andrew Gelman, and Jonah Gabry. “Practical Bayesian model evaluation using leave-one-out cross-validation and WAIC”. In: *Statistics and Computing* (2017). ISSN: 15731375. DOI: 10.1007/s11222-016-9696-4. arXiv: 1507.04544.
- [35] R Core Team. “R: A Language and Environment for Statistical Computing. R Foundation for Statistical Computing”. In: (2021). URL: <https://www.R-project.org>.
- [36] Bob Carpenter et al. “Stan: A probabilistic programming language”. In: *Journal of Statistical Software* (2017). ISSN: 15487660. DOI: 10.18637/jss.v076.i01.
- [37] Michael Betancourt. *A Conceptual Introduction to Hamiltonian Monte Carlo*. <https://arxiv.org/abs/1701.02434>. 2018. arXiv: 1701.02434 [stat.ME].
- [38] Jonah Gabry. Author Rok Češnovar et al. *CmDStanR: A lightweight interface to Stan for R users*. R package version 0.3.0. 2021. URL: <https://mc-stan.org/cmdstanr/>.
- [39] Ben Goodrich. *Truncated Multivariate Normal Variates in Stan*. <https://groups.google.com/g/stan-users/c/GuWUJogum1o/m/Lvxj1UBnBwAJ?pli=1>. 2017.
- [40] Stephen Martin. *Hierarchical prior for partial pooling on correlation matrices*. <https://discourse.mc-stan.org/t/hierarchical-prior-for-partial-pooling-on-correlation-matrices/4852/27>. 2018.
- [41] *Stan Modeling Language Users Guide and Reference Manual*. https://mc-stan.org/docs/2_25/reference-manual/. 2020.

- [42] Taye H. Hamza et al. “Multivariate random effects meta-analysis of diagnostic tests with multiple thresholds”. In: *BMC Medical Research Methodology* (2009). ISSN: 14712288. DOI: 10.1186/1471-2288-9-73.
- [43] B Aertgeerts, F Buntinx, and A Kester. “The value of the CAGE in screening for alcohol abuse and alcohol dependence in general clinical populations: a diagnostic meta-analysis”. In: *Journal of Clinical Epidemiology* 57.1 (Jan. 1, 2004), pp. 30–39. ISSN: 0895-4356. DOI: 10.1016/S0895-4356(03)00254-3. URL: <https://www.sciencedirect.com/science/article/pii/S0895435603002543> (visited on 08/01/2021).
- [44] J. A. Ewing. “Detecting alcoholism. The CAGE questionnaire”. In: *JAMA* 252.14 (Oct. 12, 1984), pp. 1905–1907. ISSN: 0098-7484. DOI: 10.1001/jama.252.14.1905.
- [45] Em Wilkinson-Brice et al. *Updated UK Health Security Agency guidance – confirmatory PCR tests to be temporarily suspended for positive lateral flow test results*. 2021.
- [46] Kurt Kroenke, Robert L Spitzer, and Janet B W Williams. “The Patient Health Questionnaire-2: validity of a two-item depression screener.” In: *Medical care* 41.11 (2003), pp. 1284–92.
- [47] K Kroenke, R L Spitzer, and J B Williams. “The PHQ-9: validity of a brief depression severity measure.” In: *Journal of general internal medicine* 16.9 (2001), pp. 606–13.
- [48] Brooke Levis et al. “Accuracy of the PHQ-2 Alone and in Combination With the PHQ-9 for Screening to Detect Major Depression: Systematic Review and Meta-analysis”. In: *JAMA* 323.22 (June 2020). Number: 22, pp. 2290–2300. ISSN: 0098-7484. DOI: 10.1001/jama.2020.6504. URL: <https://doi.org/10.1001/jama.2020.6504> (visited on 05/29/2021).
- [49] Pamela M. Vacek. “The Effect of Conditional Dependence on the Evaluation of Diagnostic Tests”. In: *Biometrics* (1985). ISSN: 0006341X. DOI: 10.2307/2530967.
- [50] Zhuoyu Wang et al. “Modeling conditional dependence among multiple diagnostic tests”. In: *Statistics in Medicine* (2017). ISSN: 10970258. DOI: 10.1002/sim.7449.
- [51] Nandini Dendukuri and Lawrence Joseph. “Bayesian Approaches to Modeling the Conditional Dependence Between Multiple Diagnostic Tests”. In: *Biometrics* 57.1 (2001), pp. 158–167. DOI: <https://doi.org/10.1111/j.0006-341X.2001.00158.x>.
- [52] Tim P. Morris, Ian R. White, and Michael J. Crowther. “Using simulation studies to evaluate statistical methods”. In: *Statistics in Medicine* 38.11 (May 20, 2019), pp. 2074–2102. ISSN: 1097-0258. DOI: 10.1002/sim.8086.
- [53] David M. Blei, Alp Kucukelbir, and Jon D. McAuliffe. “Variational Inference: A Review for Statisticians”. In: *Journal of the American Statistical Association* 112.518 (2017), pp. 859–877. DOI: 10.1080/01621459.2017.1285773.
- [54] Leo L. Duan, James E. Johndrow, and David B. Dunson. “Scaling up Data Augmentation MCMC via Calibration”. In: *J. Mach. Learn. Res.* 19.1 (Jan. 2018), 2575–2608. ISSN: 1532-4435.
- [55] Leo L. Duan. *Transport Monte Carlo*. 2020. arXiv: 1907.10448 [stat.CO].
- [56] Charles C. Margossian et al. *Hamiltonian Monte Carlo using an adjoint-differentiated Laplace approximation: Bayesian inference for latent Gaussian models and beyond*. 2020. arXiv: 2004.12550 [stat.CO].
- [57] Laurent Dinh, Jascha Sohl-Dickstein, and Samy Bengio. *Density estimation using Real NVP*. 2017. arXiv: 1605.08803 [cs.LG].
- [58] George Papamakarios, Theo Pavlakou, and Iain Murray. *Masked Autoregressive Flow for Density Estimation*. 2018. arXiv: 1705.07057 [stat.ML].
- [59] Danilo Jimenez Rezende and Shakir Mohamed. *Variational Inference with Normalizing Flows*. 2016. arXiv: 1505.05770 [stat.ML].
- [60] Akash Kumar Dhaka et al. “Robust, Accurate Stochastic Optimization for Variational Inference”. In: 1 (2020), pp. 1–15. arXiv: 2009.00666. URL: <http://arxiv.org/abs/2009.00666>.
- [61] Richard D. Riley et al. “Meta-analysis of diagnostic test studies using individual patient data and aggregate data”. In: *Statistics in Medicine* 27.29 (2008), pp. 6111–6136. DOI: <https://doi.org/10.1002/sim.3441>.
- [62] Donald B. Rubin. “Inference and Missing Data”. In: *Biometrika* (1976). ISSN: 00063444. DOI: 10.2307/2335739.

- [63] Xiaoye Ma et al. “A Bayesian hierarchical model for network meta-analysis of multiple diagnostic tests”. In: *Biostatistics* (2018). ISSN: 14684357. DOI: 10.1093/biostatistics/kxx025.
- [64] Victoria N. Nyaga, Marc Aerts, and Marc Arbyn. “ANOVA model for network meta-analysis of diagnostic test accuracy data”. In: *Statistical Methods in Medical Research* (2018). ISSN: 14770334. DOI: 10.1177/0962280216669182. arXiv: 1604.02018.
- [65] V. Dukic and C. Gatsonis. “Meta-analysis of Diagnostic Test Accuracy Assessment Studies with Varying Number of Thresholds”. In: *Biometrics* 59.4 (2003), pp. 936–946. DOI: <https://doi.org/10.1111/j.0006-341X.2003.00108.x>.
- [66] Hayley E. Jones et al. “Quantifying how diagnostic test accuracy depends on threshold in a meta-analysis”. In: *Statistics in Medicine* (2019). ISSN: 10970258. DOI: 10.1002/sim.8301.
- [67] Susanne Steinhauser, Martin Schumacher, and Gerta Rücker. “Modelling multiple thresholds in meta-analysis of diagnostic test accuracy studies”. In: *BMC Medical Research Methodology* 16.1 (2016), p. 97. ISSN: 1471-2288. DOI: 10.1186/s12874-016-0196-1.
- [68] Sean M. O’Brien and David B. Dunson. “Bayesian Multivariate Logistic Regression”. In: *Biometrics* 60.3 (2004), pp. 739–746. DOI: <https://doi.org/10.1111/j.0006-341X.2004.00224.x>.
- [69] Rainer Winkelmann. “COPULA BIVARIATE PROBIT MODELS: WITH AN APPLICATION TO MEDICAL EXPENDITURES”. In: *Health Economics* 21.12 (2012), pp. 1444–1455. DOI: <https://doi.org/10.1002/hec.1801>.
- [70] Michael Eichler, Hans Manner, and Dennis Turk. *Dynamic copula based multivariate discrete choice models with applications*. https://wisostat.uni-koeln.de/sites/statistik/user_upload/DCMDC.pdf. 2017.
- [71] Christian Meyer. “The Bivariate Normal Copula”. In: *Communications in Statistics - Theory and Methods* 42.13 (2013), 2402–2422. ISSN: 1532-415X. DOI: 10.1080/03610926.2011.611316.