

Challenges and opportunities of real-world data: Statistical analysis plan for the Optimise:MS multicentre prospective cohort pharmacovigilance study

1 Dr Ed Waddingham^{1*}, Dr Aleisha Miller¹, Dr Ruth Dobson², Prof Paul M. Matthews¹

2 1. Department of Brain Sciences and UK Dementia Research Institute, Imperial College
3 London, Hammersmith Campus, Du Cane Road, London, W12 0NN

4 2. Preventive Neurology Unit, Wolfson Institute of Preventive Medicine, Queen Mary University
5 of London

6 *Keywords: Real-world data, multiple sclerosis, statistical analysis plan, pharmacovigilance, cohort
7 study, signal detection*

8 Number of figures: 3, Number of tables: 5, Word count: 4,668

9 1 Abstract

10 *Introduction:* Optimise:MS is an observational pharmacovigilance study aimed at characterising the
11 safety profile of disease-modifying therapies (DMTs) for multiple sclerosis (MS) in a real world
12 population. The study will categorise and quantify the occurrence of serious adverse events (SAEs)
13 in a cohort of MS patients recruited from clinical sites around the UK.

14 The study was motivated particularly by a need to establish the safety profile of newer DMTs, but
15 will also gather data on outcomes among treatment-eligible but untreated patients and those receiving
16 established DMTs (interferons and glatiramer acetate). It will also explore the impact of treatment
17 switching.

18 *Methods:* Causal pathway confounding between treatment selection and outcomes, together with the
19 variety and complexity of treatment and disease patterns observed among MS patients in the real
20 world, present statistical challenges to be addressed in the analysis plan. We developed an approach
21 for analysis of the OPTIMISE:MS data that will include disproportionality-based signal detection
22 methods adapted to the longitudinal structure of the data and a longitudinal time-series analysis of a
23 cohort of participants receiving second-generation DMT for the first time. The time-series analyses
24 will use a number of exposure definitions in order to identify temporal patterns, carryover effects and
25 interactions with prior treatments. Time-dependent confounding will be allowed for via inverse-
26 probability-of-treatment weighting (IPTW). Additional analyses will examine rates and outcomes of
27 pregnancies and explore interactions of these with treatment type and duration.

28 *Results:* To date 13 hospitals have joined the study and over 2000 participants have been recruited.
29 A statistical analysis plan has been developed and is described here.

30 *Conclusion:* Optimise:MS is expected to be a rich source of data on the outcomes of DMTs in real-
31 world conditions over several years of follow-up in an inclusive sample of UK MS patients.
32 Analysis is complicated by the influence of confounding factors including complex treatment
33 histories and a highly variable disease course, but the statistical analysis plan includes measures to

34 mitigate the biases such factors can introduce. It will enable us to address key questions that are
35 beyond the reach of randomised controlled trials.

36 **2 Introduction**

37 OPTIMISE:MS is a prospective observational cohort study lasting at least 7 years (with the
38 possibility of extension depending on funding), focused on evaluating the safety profile of MS DMTs
39 in the real-world setting. A sample size of around 4,000 multiple sclerosis (MS) patients is
40 anticipated, to be recruited from several sites (MS treatment centres) around the UK. This sample
41 size is based on the recruitment level that is expected to be achievable in practice, rather than on
42 considerations relating to statistical power. The study is open to all MS patients (as defined by the
43 2017 McDonald criteria(1)), of any MS subtype, attending a participating site and eligible for
44 treatment based on current UK guidelines(2) , regardless of their actual treatment history. The study
45 has been recruiting since May 2019, and as of 2022 remains open to new recruits. The length of the
46 recruitment window, coupled with the introduction of remote consenting, should ensure that the
47 sample is not heavily skewed towards those attending clinics most frequently. Details of the study
48 design and protocol have already been published (3). The study is academically initiated and led, but
49 is guided by a public-private partnership between academic clinical investigators and pharmaceutical
50 companies with marketing authorisations for DMTs.

51 Subjects taking second-generation DMTs will be the main focus of investigation, and controls will
52 include those eligible but not receiving treatment and those receiving first-generation DMTs (see
53 Table 1 for a current list of first- and second-generation DMTs; any new DMTs becoming available
54 for use by patients in the UK during the course of the study will be classed as second-generation).

55 The primary objective of the study is to establish the incidence of serious adverse events (SAEs)
56 among MS patients receiving any second-generation DMT, and compare it with that observed in
57 untreated but treatment-eligible patients and those receiving first-generation DMT.

58 Secondary objectives are:

- 59
- 60 • to measure and compare SAE rates for individual DMTs;
 - 61 • to assess associations between second-generation DMT therapy and incidence of
62 lymphopaenia;
 - 63 • to assess associations between second-generation DMT therapy and moderately and severely
64 abnormal liver function, as indicated by blood tests for alanine transaminase or aspartate
65 transaminase;
 - 66 • to assess the impact of sequential DMT therapy on the incidence of SAEs;
 - 67 • to assess the relative efficacy of DMT classes with regard to suppression of relapses,
68 disability progression and new lesion formation on MRI; and
 - 69 • to measure the frequencies of pregnancies and their outcomes.

70 SAEs are defined as adverse events resulting in death, persistent or significant disability/incapacity,
71 or hospitalisation (or extension of a hospital stay for an inpatient). These are classified according to
72 the following categories: Opportunistic infections, infections requiring hospitalization, MS relapses,
73 deaths, COVID-19 infections, other SAEs deemed to be related to treatment (eg malignancies), and
74 other SAEs.

75 **3 Methods and Analysis**

76 **3.1 Study sites, data entry and storage**

77 Participants are recruited at participating MS clinics at hospitals around the UK. Currently there are
78 13 participating hospital sites and over 2000 individuals have been enrolled in the study.

79 At each site, study data is entered onto a local secure database held on a dedicated PC. These
80 machines connect securely to the Optimise:MS server (hosted by the Data Science Institute at
81 Imperial College London) and automatically upload (“push”) the data to the central database at
82 regular intervals. Regular quality checks on the data central data are performed centrally through
83 monitoring data completeness, internal consistency, concordance with expected ranges, and
84 harmonization of units; queries are fed back to the site staff for resolution.

85 Participants’ data is managed in line with the requirements of the General Data Protection
86 Regulation, Imperial College London’s policies and the study’s own Standard Operating Procedures.
87 Personally identifiable data is kept to a minimum; names and contact details are accessible only by
88 local site staff and are not stored on the central study database.

89 **3.2 Longitudinal cohort structure and outcome assessment**

90 MS patients may join the study if they are eligible for treatment with DMT, regardless of whether or
91 not they actually receive DMT. Upon enrolment the patient’s basic demographic and clinical data
92 (including their MS diagnosis and any comorbidities) are entered onto the study database by site
93 staff. Retrospective data is also collected at enrolment, including disability assessments and relapses,
94 lab test results, a full history of DMT use, and any past serious infections or malignancies.

95 Whenever a participant attends a clinic visit while under observation in the study, the database is
96 updated with the reason for the visit, date of the visit, and details of any other changes in the
97 participant’s data (such as disease progression, new comorbidities, any treatment changes, SAEs, test
98 results, or MRI scan results) since the previous visit. Exact dates for all such events are recorded
99 whenever possible. No additional clinic visits or procedures are required as part of the study.
100 Participants are under observation from their enrolment visit until they withdraw consent, leave a
101 participating clinic, die, or until the end of the study, whichever is the earliest.

102 SAEs (including MS relapses), pregnancies and their outcomes, and any new/enlarging lesions
103 revealed by clinically indicated interval MRI are recorded on the Optimise database by local site staff
104 accessing medical records. Disability is assessed by local clinical staff using the Expanded
105 Disability Status Scale (EDSS) (4) and the total score is recorded on the database; a disability
106 progression outcome is defined as an EDSS measurement scoring at least 1 point higher than the
107 most recent measurement at or after baseline. Laboratory test results (eg blood cell and liver enzyme
108 counts) also are recorded on the database. Abnormal liver function is assessed using blood alanine
109 aminotransferase (ALT) or aspartate aminotransferase (AST) levels. For each, moderate and severe
110 elevation are respectively defined as exceeding 2.5x and 5x , respectively, of the upper limit of the
111 normal ranges established by Imperial North West London Pathology. Lymphopaenia is defined
112 based on absolute lymphocyte count (ALC) according to the following grades:

- 113 ○ Grade 1: Lower limit of normal range \geq ALC \geq 800/mm³
- 114 ○ Grade 2: 800/mm³ \geq ALC \geq 500/mm³
- 115 ○ Grade 3: 500/mm³ $>$ ALC \geq 200/mm³
- 116 ○ Grade 4: 200/mm³ $>$ ALC

117

118 **3.3 Statistical principles**

119 Due to selection effects, MS patients receiving different treatments are likely to have different
120 underlying characteristics and to experience outcomes at different rates even before allowing for the
121 effects of treatment. Thus, confounding is expected between treatment selection and outcomes,
122 leading to biased treatment effect estimates. Confounding variables may include demographics,
123 disease and treatment history, and time variables representing period and cohort effects (5).

124 Controlling for the effects of confounders can be particularly difficult in longitudinal studies where
125 past treatment exposures and covariates may influence future exposures and/or covariates as well as
126 future outcomes. This is known as time-varying causal pathway confounding, and the bias it
127 introduces may not be adequately controlled by the standard multivariate covariate adjustment
128 approach (6, 7) (8). This type of confounding is expected to occur in the Optimise:MS cohort, given
129 the nature of MS as a chronic progressive disease and the factors that are suspected to influence
130 treatment decisions. Methods for controlling confounders have been chosen to mitigate this problem
131 (further details below).

132 The statistical analyses fall into three classes: cohort analyses, signal detection analyses, and
133 pregnancy analyses. These are described under the headings below.

134 **3.4 Cohort analyses**

135 A “new user” cohort of those study subjects who have never received second-generation DMT prior
136 to study enrolment will be the subject of longitudinal analyses. These will examine the effects of
137 DMTs on relapse, disability progression, abnormal liver function, lymphopenia, new lesion
138 formation and SAE rates. The temporal relationship between exposures and outcomes will also be
139 explored.

140 The primary cohort analysis aims to investigate the effectiveness and safety of DMTs using a
141 relatively simple model. Participants will be separated into two strata according to whether or not
142 they have ever received first-generation DMT prior to second-generation DMT initiation (or prior to
143 the end of follow-up, if second-generation DMT is never initiated). Within each stratum, outcomes
144 occurring while exposed to second-generation DMT will be compared to outcomes occurring while
145 unexposed. Follow-up is censored upon cessation of second-generation DMT. Subjects who
146 commence second-generation DMT while under observation will contribute an initial unexposed
147 episode and a subsequent exposed episode of follow-up time to the analysis, as illustrated for two
148 hypothetical patients in Figure 1.

149 To control for confounding in the primary analysis, propensity score weighting will be used; each
150 exposure episode will be weighted in inverse proportion to the estimated propensity (probability) of
151 the observed treatment exposure. The propensity score is based on time-varying covariates measured
152 at the start of the exposure episode (7). The effect of the weighting is to construct a pseudo-
153 population which is effectively “randomised” in the sense that the covariates at the start of exposure
154 episodes are balanced across exposure categories. The propensities are estimated using a pooled
155 logistic regression model.

156 The secondary cohort analyses are aimed at exploring the temporal relationship between DMT use
157 and outcomes, including whether the effects of DMTs persist after treatment cessation/switch.
158 Follow-up is not censored upon cessation of second-generation DMT; instead, participants can
159 contribute multiple periods of exposure to the analysis as they move between treatment classes. This
160 is illustrated in Figure 2 for the two hypothetical patients described in Figure 1. The secondary
161 analyses thus make use of all observed data for the new user cohort and, owing to the more complex

162 longitudinal exposure patterns involved, observations will be weighted using time-varying inverse
163 probability weights (IPTW) to estimate a marginal structural model (MSM) (9). This is similar to the
164 propensity score method described above, but the weights are updated at regular (6-month) intervals
165 based on the latest covariate values and reflect the probability of observing the participant's full
166 treatment history up until that timepoint (6). For details of how these probabilities are modelled, and
167 the formulae for the weights, see the Supplementary Material. This method aims to create a
168 dynamically weighted pseudo-population that is longitudinally balanced, i.e. with covariates equally
169 balanced across all possible treatment histories at every 6-month timepoint. This construction relies
170 on an assumption that the probabilities lie strictly between 0 and 1 for each possible level of the
171 covariates (the positivity assumption). Provided that this condition is met and all confounders are
172 measured at sufficiently frequent intervals, this method can fully control for time-varying causal
173 pathway confounding and generate unbiased estimates of the marginal treatment effects. A three-
174 category treatment variable will be used (no treatment, first-generation DMT or second-generation
175 DMT) instead of the stratified approach of the primary analysis. Parallel analyses will use different
176 exposure models to examine the temporal patterns of treatment effects:

- 177 (a) Outcomes associated with current treatment class (categorical exposure variable)
- 178 (b) Outcomes associated with current treatment class plus carryover effect of any other treatment
179 class in the past 6 months (categorical exposure variables)
- 180 (c) Outcomes associated with cumulative exposures (continuous exposure variable for each
181 treatment category)
- 182 (d) Outcomes associated with time-weighted cumulative exposure, i.e. historic exposures
183 downweighted relative to recent exposures (continuous exposure variable for each treatment
184 category)

185 The tertiary cohort analysis extends exposure model (b) to examine whether there is an interaction
186 effect associated with treatment switching, i.e. whether the carryover effect of previous treatment is
187 dependent on current treatment exposure.

188 Further cohort analyses will examine the effects of second-generation DMTs individually rather than
189 as a collective treatment class. The principle analysis method for all cohort analyses will be time-
190 varying Cox proportional hazards regression (10).

191 **3.5 Signal detection analyses**

192 The signal detection analyses will examine whether the rate of SAEs (excluding MS relapses)
193 occurring for any individual DMT is disproportionate to the overall rate of SAEs in the study sample.
194 SAEs will be analysed according to their classification in the Optimise database as:

- 195 • Infections
- 196 • Opportunistic Infections
- 197 • Malignancies and other SAEs likely related to treatment
- 198 • Deaths (all causes)
- 199 • Covid-19
- 200 • Other SAEs

202 Infections, opportunistic infections and Covid-19 will be further analysed according to the subtypes
203 recorded on the database, currently including the following categories:

- 204 • *Infections*: urinary tract infections, bronchitis, sinusitis, gastroenteritis, thinea, sepsis,
205 bacterial, viral, abscess, other
- 206 • *Opportunistic Infections*: progressive multifocal leukencephalopathy, herpes zoster,
207 herpes simplex, varicella, viral hepatitis, listeria, mycosis, abscess, other
- 208 • *Covid-19*: suspected, confirmed by test, hospitalised, ventilated

209 Classifications based on MedDRA codings or free-text descriptions may also be used.

210 Patient-months will be assigned to treatments according to three different definitions of exposure:

- 211 • Exposure within the month of interest or the previous month
- 212 • Exposure within the preceding 6 months
- 213 • Exposure at any prior time in the patient’s treatment history

214 Only incident events (i.e. the first recorded occurrence in a given study participant) will be analysed;
215 follow-up is censored upon occurrence of the event of interest.

216 A minimum report criterion is also imposed in order to avoid statistical noise in the
217 disproportionality statistics when event counts are too low. For a signal to be triggered, an event
218 must be reported in at least 3 study participants for second-generation DMTs and 5 participants for
219 first-generation DMTs. The higher threshold in the latter case results in fewer false positives and
220 more precise risk estimates, but with reduced sensitivity (11), reflecting the fact that the safety profile
221 of first-generation DMTs is relatively well understood and early detection of signals is less of a
222 priority than for the newer treatments.

223 3.5.1 Signal detection methodologies/measures

224 The key disproportionality methods used in this study, the Reporting Odds Ratio and Bayesian
225 Confidence Propagation Neural Network, were originally developed in the context of spontaneous
226 report databases. In this original context the methods would be used to evaluate whether an event is
227 cited more frequently in AE reports for the treatment of interest than in reports for other treatments.

228 Longitudinal cohort data also covers periods when no adverse events occur, which provides
229 additional information regarding the relative frequencies of exposures and outcomes. When applying
230 the disproportionality approach in the longitudinal setting it is appropriate to make use of this
231 additional data by altering the methods so that they do not simply count AE reports occurring on
232 treatments, but also take into account periods with no exposure and/or no events(12). This is
233 achieved by treating each patient-month of follow-up as a unit of observation and evaluating whether
234 events occur more frequently during patient-months exposed to the treatment of interest than during
235 all other patient-months. The methods are described under the headings below in accordance with
236 this longitudinal formulation.

237 3.5.1.1 Simple disproportionality measures

238 The reporting odds ratio (ROR) (13) compares the odds of an adverse event occurring during
239 exposed patient-months to the odds of occurrence during unexposed patient-months. For a given
240 drug-event combination the ROR is calculated as follows:

241
$$ROR = \frac{n_{11}n_{00}}{n_{01}n_{10}}$$

242 where n_{00} = number of patient-months without exposure to drug or occurrence of event
 243 n_{01} = number of patient-months without exposure to drug but with occurrence of event
 244 n_{10} = number of patient-months with exposure to drug but without occurrence of event
 245 n_{11} = number of patient-months with exposure to drug and occurrence of event
 246

247 Another simple disproportionality measure is the proportional reporting ratio (PRR), which is
 248 calculated not as an odds ratio, but rather a relative risk in exposed vs unexposed months:

249
$$PRR = \frac{n_{11}n_{0\cdot}}{n_{01}n_{1\cdot}}$$

250 where the dot symbol \cdot indicates summation over the index values 0 and 1 (14). A third measure is
 251 the relative reporting ratio (RRR), a relative risk in exposed vs all months:

252
$$RRR = \frac{n_{11}n_{\cdot\cdot}}{n_{\cdot 1}n_{1\cdot}}$$

253 In practice the PRR, RRR and ROR give near-identical results when used for signal detection (15)
 254 (12).

255 The incidence rate ratio (IRR) is a standard relative measure of incidence in epidemiology and
 256 medical statistics, often estimated by Poisson regression. It is calculated as the incidence of an event
 257 among treated participants divided by its incidence among untreated participants, where the incidence
 258 is the number of events divided by the total amount of follow-up time. It can easily be seen that the
 259 IRR is equivalent to the longitudinal formulation of the PRR described above. This observation
 260 allows us to calculate a confounder-controlled estimate of the PRR via weighted Poisson regression,
 261 using the marginal structural approach described under “Cohort Analyses” above. Indeed, the same
 262 weighted PRR estimate can be obtained by directly substituting weighted equivalents of n_{01} , n_{11} , n_{00}
 263 and n_{10} in the formula above (for details see the Supplementary Material). The latter approach can
 264 be extended to calculate a weighted version of the RRR, which will be used in the “weighted analysis
 265 pathway” (see “Signal Generation Procedure” section below).

266 **3.5.1.2 Shrinkage (Bayesian Confidence Propagation Neural Network)**

267 Owing to the discrete nature of count data, simple disproportionality measures are very unstable
 268 when event rates are low. Chance occurrences of a rare event can easily generate spurious false
 269 positive signals.

270 The Bayesian Confidence Propagation Neural Network (BCPNN) method (16) is designed to reduce
 271 the rate of false positives by using a Bayesian model to express the joint distribution of the
 272 probabilities of drug exposure and event occurrence, with conjugate beta priors that favour an
 273 independent relationship (i.e. no association between drug and event). This achieves a “shrinkage”
 274 effect that pulls the disproportionality estimates back towards the null when event counts are low.

275 The model’s key measure of disproportionality is the Information Component, which is the base-2
 276 logarithm of the RRR. A posterior estimate of the False Discovery Rate (FDR) for each signal, i.e.
 277 the probability of no association between drug and event, can also be obtained (17).

278 **3.5.1.3 Controlling for protopathic bias (LEOPARD)**

279 Signal detection methods are often prone to generating false positives due to protopathic bias, which
 280 occurs if an event is mistakenly ascribed to initiation of a new treatment when both shared a common
 281 cause such as an underlying disease exacerbation(18). LEOPARD is a signal filtering method aimed
 282 at eliminating this bias. The method works by examining the rate of treatment initiations before and
 283 after adverse event incidence; protopathic bias is inferred if treatment follows the event more often
 284 than it precedes it (15). To address this, we will employ a one-sided binomial test of the distribution
 285 of treatment initiation events, with the null hypothesis that treatment initiation is equally likely before
 286 an AE as after it, and the alternative hypothesis that the probability is higher after the AE. This test
 287 will be carried out at the 50% significance level (19); signals where the null hypothesis is rejected
 288 will be discarded.

289 **3.5.2 Signal generation procedure**

290 For each treatment of interest and exposure definition, the analysis will follow the process set out
 291 Figure 3. As the first step in the analysis, a list of events fulfilling the minimum report criterion is
 292 generated (the Level 1 list). Thereafter, three parallel analysis pathways are used: a crude
 293 (unadjusted) disproportionality analysis, and two analyses aimed at controlling for potential
 294 confounding covariates: a subgrouped analysis and a weighted analysis (IPTW).

295 Within each pathway, a Level 2 list is produced containing all signals identified by the Reporting
 296 Odds Ratio or, equivalently, the incidence rate ratio. Signals are triggered when the lower 95%
 297 confidence bound for the disproportionality measure exceeds 1 (for the subgrouped analysis, this
 298 must be observed in at least one subgroup; this approach has been reported to provide better
 299 performance than using a pooled odds ratio (11)).

300 The Level 2 list is expected to contain some false positives due to (i) volatility of disproportionality
 301 measures associated with low event counts, and (ii) protopathic bias. The Level 3 list tackles these
 302 problems by (i) applying Bayesian shrinkage to pull disproportionality estimates back towards the
 303 null (the Bayesian Confidence Propagation Neural Network Method) and (ii) verifying that
 304 prescriptions tend to precede rather than follow events (the LEOPARD filter). Signals with an FDR
 305 estimate below 5% which are not rejected by the LEOPARD filter will be included on the Level 3
 306 list. In the subgrouped- analysis, these conditions must be achieved in at least one sub-group; in the
 307 weighted analysis, the BCPNN calculations are based on the weighted event counts described in the
 308 Supplementary Material.

309 Pooled lists at levels 2 and 3 will be produced in which signals will be ranked according to the
 310 number of pathways in which the signal was observed and the associated disproportionality statistics
 311 (level 2) or estimated false discovery rates (level 3) (17).

312 Sensitivity analyses may explore the use of alternative decision rules, such as varying the minimum
 313 report or FDR thresholds, and alternative methodologies, such as replacing BCPNN with the
 314 Gamma-Poisson Shrinker (15, 20) or Information Component Temporal Pattern Discovery (21)).

315 After drug-event signals have been identified, the data will be further examined for evidence of drug-
 316 drug-event signals, i.e. adverse events associated with treatment interactions. These analyses will
 317 also proceed using the procedure set out in Figure 3, with different exposure definitions and
 318 background rates depending on the context (these are set out in the full Statistical Analysis Plan).

319 An additional paediatric signal detection analysis be carried out in participants under 18 years old.
320 For this purpose the threshold for the minimum report criterion will be reduced to 2 cases, and only
321 the crude analysis pathway will be used.

322 **3.6 Pregnancy analyses**

323 The average rate of pregnancy per person-year of follow-up will be estimated, both among all
324 females aged 18 to 50 in the study population and according to DMT class and specific DMT being
325 received at the date of conception.

326 Multinomial or binomial logistic regression will be used to estimate the effect of the treatment
327 received at conception on the eventual outcome of pregnancy.

328 **3.7 Planned interim analyses**

329 The study is in a position to reveal previously unobserved adverse drug reactions, particularly in
330 connection with the more novel second-generation DMTs. To facilitate timely detection of such
331 signals, a simplified set of analyses will be performed on an annual basis while data is being accrued.
332 These will consist of the signal detection analyses (crude analysis pathway and single-drug-event
333 associations only), and simple (constant-hazard) unadjusted Poisson regressions of the occurrence of
334 any SAE according to current treatment received.

335 **4 Discussion**

336 Optimise:MS is being carried out in a routine sub-specialty referral care setting, and will thus provide
337 “real-world” data on outcomes occurring under the sort of treatment and clinical monitoring regimes
338 that patients typically experience, rather than the idealised conditions of a randomised controlled trial
339 (RCT) (22). The study participants should be more representative of the general population of MS
340 patients in the UK than would be the case in a typical RCT, since the inclusion criteria are less
341 restrictive and the study does not burden the participants with additional procedures or impose any
342 new treatment regimes. This also facilitates recruitment, and over a long period of follow-up, despite
343 the lack of additional investigations or procedures, enables a comprehensive set of clinical data to be
344 gathered. The use of electronic consent forms and remote/virtual clinic visits has also helped in this
345 regard, particularly during the COVID-19 pandemic.

346 The sample size and length of follow-up thus exceed most RCTs and, together with the detailed data
347 gathered on participants’ DMT and disease histories, will enable the estimation of washout,
348 switching and subgroup effects that often lie beyond the scope and capabilities of trials.

349 Of course, observational studies have well-known drawbacks compared to RCTs - chiefly the
350 absence of randomisation, which leaves treatment selection potentially subject to the influence of
351 prognostic factors and therefore vulnerable to confounding with outcomes. The likely existence of
352 time-varying causal pathway confounding in the MS context makes this problem particularly
353 challenging to address analytically, but the marginal structural modelling approach (IPTW) has shown
354 that it has the capability to produce unbiased estimates - at least under ideal conditions when
355 positivity is satisfied, probability models are specified correctly and there are few extreme weights(8,
356 23, 24). The comprehensive longitudinal data collection in Optimise should facilitate MSM
357 estimation, which will be particularly important for the secondary cohort analyses investigating the
358 effect of longitudinal treatment trajectories. The estimation of probability weights in itself may

359 provide useful insight into the prevalence of DMT use in particular subgroups, and other factors
360 influencing treatment decisions.

361 We have also specified a simpler cross-sectional propensity-score weighting approach, as this
362 improves the chances of positivity and reduces the potential for extreme weights. Although this
363 model may not fully control for the influence of prior treatment history on outcomes, this is less
364 likely to be a major concern in the primary cohort analysis since exposure histories are relatively
365 simple (Figure 1) compared to the more complex exposure histories in the secondary analysis (Figure
366 2).

367 The use of weighted event counts in the disproportionality-based signal detection methods is, to our
368 knowledge, novel, but is well-founded (see the Supplementary Material). This is the only method we
369 are aware of that can control for time-varying causal pathway confounding when using
370 disproportionality methods such as the ROR, BCPNN or GPS. However, it can only be used when
371 these methods are applied to longitudinal cohort data, rather than to the spontaneous report data for
372 which such methods were originally developed. Linking cases to their treatment histories, and hence
373 examining drug-drug-event signals involving washout effects of prior treatments, is also more
374 straightforward in the longitudinal setting. These considerations favour the Optimise cohort-based
375 design for future signal detection databases. Another reason, of course, is the additional data gained
376 from periods with no treatment exposure or adverse events, which may improve the performance of
377 disproportionality methods (15). Without this additional data, disproportionality analyses of
378 spontaneous reports can unfairly penalise drugs with low overall AE rates if any one AE occurs more
379 often than others (an example is shown in the Supplementary Material). Alongside the novel
380 weighted analysis, a parallel subgrouped- analysis provides another means of controlling for
381 confounders and is better established in signal detection (11, 20) - although this method may still be
382 vulnerable to time-varying causal pathway confounding. since the subgroups are based on cross-
383 sectional covariate values rather than full exposure and covariate histories.

384 A disadvantage of using the Optimise study for signal detection purposes, as opposed to a
385 spontaneous report registry, is the relatively small sample size. This exacerbates the known problem
386 of volatility in disproportionality statistics when event counts are low - hence the importance of using
387 a shrinkage methodology such as BCPNN. Protopathic bias presents another significant problem for
388 pharmacovigilance in MS patients, as false signals may easily be generated by both the
389 relapsing/remitting and progressive aspects of the disease, and the wide range of symptoms it can
390 produce. Direct comparisons between safety profiles of different DMTs - in particular between first-
391 and second-generation DMTs - may also be biased due to the fact that exposure and follow-up time
392 are more limited for newer drugs, and so treatment effects that manifest over the longer term cannot
393 be observed. Finally, the potential for differences in the intensity of follow up on different treatments
394 to bias event detection is not specifically accounted for in the analysis. The impact of this varies
395 greatly by outcome; for example, it would be expected to be greater for imaging measures of disease
396 activity such as new or enlarging lesions than for SAEs. Although imaging results may also be
397 affected by the use of different scanners, acquisition protocols and schedules, this is not expected to
398 be strongly related to treatment.

399 In summary, OPTIMISE is observational, inclusive, and does not impose any fixed timelines on
400 those taking part. Participants can be enrolled at any stage of their MS or treatment history; there is
401 no unifying milestone marking for the start of follow-up, and no set course of treatment to be
402 followed thereafter. This inclusivity makes recruitment easier, enhances data collection and may
403 increase the population representativeness and generalisability of results, but it presents major

404 challenges from a statistical perspective. We have tried to address these and realise opportunities
405 arising from the design. Our approach to signal detection analyses will ensure a healthy mix of data
406 from as wide a population as possible, although care has been needed to plan the analysis in a way
407 that controls for treatment selection and protopathic bias. For longitudinal cohort analyses, the lack
408 of fixed timelines for participants is a complicating factor, but also creates the potential for a wealth
409 of useful data if handled appropriately. Our cohort analyses simplify the structure of the data by
410 focusing on a sub-population of participants initiating second-generation DMT for the first time, as it
411 is the safety profile of these drugs that is the primary outcome of interest. Further analytical choices
412 have been made to either mitigate the confounding influence of variability in patient
413 characteristics/histories (eg marginal structural modelling) or exploit this variability to gain
414 additional insights (eg the analyses of washout/cumulative/switch effects).

415 **5 Author Contributions**

416 EW developed the Statistical Analysis Plan with support and input from RD, AM and PM, and
417 drafted the manuscript. RD and PM contributed to conception and design of the study and revisions
418 of the manuscript. RD drafted the study protocol with contributions from PM. AM co-ordinated the
419 activation and ongoing management of the study. All authors contributed to reviewing the
420 manuscript and approved the submitted version.

421 **6 Funding**

422 PM acknowledges generous personal and research support from the Edmond J Safra Foundation and
423 Lily Safra, an NIHR Senior Investigator Award, the UK Dementia Research Institute and the NIHR
424 Biomedical Research Centre at Imperial College London. The study is funded by a partnership
425 comprising:

- 426 • Biogen IDEC Ltd (grant reference P76049, CrossRef Funder ID: 10.13039/100006314)
- 427 • Merck Serono Ltd., Feltham, UK, an affiliate of Merck KGaA, Darmstadt, Germany (grant
428 reference WMCN_P74840, CrossRef Funder ID: 10.13039/100009945)
- 429 • Celgene Ltd (Bristol-Myers Squibb) (grant ref P76049, CrossRef Funder ID:
430 10.13039/100006436)

431 The funding companies are represented on the study's Steering Committee but are not involved in
432 day-to-day administration of the study, data collection or analysis.

433 Manuscript publication fees are paid by Imperial College London.

434 **7 Conflicts of Interest**

435 The study received funding from Biogen IDEC Limited, Merck Serono Ltd and Celgene Ltd. The
436 statistical analysis plan and this manuscript were developed with input from the funders in a
437 reviewing capacity.

438 PM acknowledges consultancy fees from Novartis, Bristol Myers Squibb, Celgene and Biogen. He
439 has received honoraria or speakers' honoraria from Novartis, Biogen and Roche and has received
440 research or educational funds from Biogen, Novartis, GlaxoSmithKline and Nodthera.

441 RD works within the PNU, which is funded by Barts Charity. She receives grant support from the
442 UK MS Society, BMA foundation, NIHR, MRC, NMSS, Horne Family Charitable Trust, Biogen and

443 Merck. She has received honoraria for Advisory boards and/or educational activities from Biogen,
444 Teva, Sanofi, Merck, Janssen, Novartis, and Roche.

445 The authors declare no other commercial or financial relationships that could be seen as potential
446 conflicts of interest.

447 **8 Acknowledgements**

448 We acknowledge the contributions of those who have been involved with the study and provided
449 feedback on the analysis plan and manuscript, particularly Dr Matt Craner (Frimley Park and John
450 Radcliffe Hospitals) and the statistical teams at the funding companies.

451 **9 Contribution to the Field**

452 Observational data gathered during routine clinical care has the potential to improve our
453 understanding of the effects of treatments in “real-world” populations rather than the idealised
454 conditions of randomised controlled trials. The Optimise:MS pharmacovigilance study seeks to
455 make use of routine care data on multiple sclerosis patients, recruited from clinical sites around the
456 UK, to examine the safety and effectiveness of disease-modifying therapies over a period of 7 years.
457 The study may provide important information to support patients’ and clinicians’ treatment decisions.

458 The use of real-world data will enable the study to explore factors that clinical trials frequently
459 cannot, such as the impact of prior treatment history. This type of data also presents statistical
460 challenges, however, not least due to the extensive confounding that is expected. A robust analysis
461 plan is therefore critical to interpretation of the study results. This manuscript describes the statistical
462 analysis plan that has been developed for Optimise:MS. The plan aims to ensure that the study meets
463 its objectives using methods that minimise problems relating to the observational nature of the data,
464 while exploiting the insights such data may provide.

465 **10 References**

- 466 1. Thompson AJ, Banwell BL, Barkhof F, Carroll WM, Coetzee T, Comi G, et al. Diagnosis of multiple
467 sclerosis: 2017 revisions of the McDonald criteria. *The Lancet Neurology* (2018) 17(2):162-73. doi:
468 10.1016/S1474-4422(17)30470-2.
- 469 2. Excellence NifHaC. NICE Pathways: Multiple Sclerosis (2021) [27/09/2021]. Available from:
470 <https://pathways.nice.org.uk/pathways/multiple-sclerosis>.
- 471 3. Dobson R, Craner M, Waddingham E, Miller A, Cavey A, Webb S, et al. OPTIMISE:MS - a
472 pragmatic, prospective observational study to address the need for, and challenges with, real world
473 pharmacovigilance in multiple sclerosis. *BMJ Open* (2021) 11:e050176. doi:10.1136/bmjopen-2021-050176
- 474 4. Kurtzke JF. Rating neurologic impairment in multiple sclerosis. *Neurology* (1983) 33(11):1444. doi:
475 10.1212/WNL.33.11.1444.
- 476 5. Diggle PJ, Heagerty P, Liang K-Y, Heagerty PJ, Zeger S. *Analysis of longitudinal data*. Oxford:
477 Oxford University Press (2002).
- 478 6. Cole SR, Hernán MA. Constructing Inverse Probability Weights for Marginal Structural Models.
479 *American Journal of Epidemiology* (2008) 168(6):656-64. doi: 10.1093/aje/kwn164.
- 480 7. Ray WA, Liu Q, Shepherd BE. Performance of time-dependent propensity scores: a
481 pharmacoepidemiology case study. *Pharmacoepidemiology and drug safety* (2015) 24(1):98-106. Epub
482 2014/11/18. doi: 10.1002/pds.3727. PubMed PMID: 25408360.

- 483 8. Robins JM, Hernán MÁ, Brumback B. Marginal Structural Models and Causal Inference in
484 Epidemiology. *Epidemiology* (2000) 11(5).
- 485 9. Hernán MA, Brumback B, Robins JM. Marginal Structural Models to Estimate the Joint Causal Effect
486 of Nonrandomized Treatments. *Journal of the American Statistical Association* (2001) 96(454):440-8. doi:
487 10.1198/016214501753168154.
- 488 10. Cox DR. Regression Models and Life-Tables. *Journal of the Royal Statistical Society: Series B*
489 (*Methodological*) (1972) 34(2):187-202. doi: <https://doi.org/10.1111/j.2517-6161.1972.tb00899.x>.
- 490 11. Wisniewski AFZ, Bate A, Bousquet C, Brueckner A, Candore G, Juhlin K, et al. Good Signal
491 Detection Practices: Evidence from IMI PROTECT. *Drug Saf* (2016) 39(6):469-90. doi: 10.1007/s40264-016-
492 0405-1.
- 493 12. Zorych I, Madigan D, Ryan P, Bate A. Disproportionality methods for pharmacovigilance in
494 longitudinal observational databases. *Statistical Methods in Medical Research* (2011) 22(1):39-56. doi:
495 10.1177/0962280211403602.
- 496 13. Rothman KJ, Lanes S, Sacks ST. The reporting odds ratio and its advantages over the proportional
497 reporting ratio. *Pharmacoepidemiology and Drug Safety* (2004) 13(8):519-23. doi:
498 <https://doi.org/10.1002/pds.1001>.
- 499 14. Curtis JR, Cheng H, Delzell E, Fram D, Kilgore M, Saag K, et al. Adaptation of Bayesian data mining
500 algorithms to longitudinal claims data: coxib safety as an example. *Med Care* (2008) 46(9):969-75. doi:
501 10.1097/MLR.0b013e318179253b. PubMed PMID: 18725852.
- 502 15. Schuemie MJ. Methods for drug safety signal detection in longitudinal observational databases: LGPS
503 and LEOPARD. *Pharmacoepidemiology and Drug Safety* (2011) 20(3):292-9. doi: 10.1002/pds.2051.
- 504 16. Bate A, Lindquist M, Edwards IR, Olsson S, Orre R, Lansner A, et al. A Bayesian neural network
505 method for adverse drug reaction signal generation. *European Journal of Clinical Pharmacology* (1998)
506 54(4):315-21. doi: 10.1007/s002280050466.
- 507 17. Ahmed I, Haramburu F, Fourrier-Réglat A, Thiessard F, Kreft-Jais C, Miremont-Salamé G, et al.
508 Bayesian pharmacovigilance signal detection methods revisited in a multiple comparison setting. (2009)
509 28(13):1774-92. doi: <https://doi.org/10.1002/sim.3586>.
- 510 18. Horwitz RI, Feinstein AR. The problem of “protopathic bias” in case-control studies. *The American*
511 *Journal of Medicine* (1980) 68(2):255-8. doi: [https://doi.org/10.1016/0002-9343\(80\)90363-0](https://doi.org/10.1016/0002-9343(80)90363-0).
- 512 19. Schuemie MJ, Madigan D, Ryan PB. Empirical Performance of LGPS and LEOPARD: Lessons for
513 Developing a Risk Identification and Analysis System. *Drug Saf* (2013) 36(1):133-42. doi: 10.1007/s40264-
514 013-0107-x.
- 515 20. Dumouchel W. Bayesian Data Mining in Large Frequency Tables, with an Application to the FDA
516 Spontaneous Reporting System. *The American Statistician* (1999) 53(3):177-90. doi:
517 10.1080/00031305.1999.10474456.
- 518 21. Norén GN, Hopstadius J, Bate A, Star K, Edwards IR. Temporal pattern discovery in longitudinal
519 electronic patient records. *Data Mining and Knowledge Discovery* (2010) 20(3):361-87. doi: 10.1007/s10618-
520 009-0152-3.
- 521 22. Cohen JA, Trojano M, Mowry EM, Uitdehaag BM, Reingold SC, Marrie RA. Leveraging real-world
522 data to investigate multiple sclerosis disease behavior, prognosis, and treatment. *Mult Scler* (2020) 26(1):23-
523 37. Epub 2019/11/28. doi: 10.1177/1352458519892555. PubMed PMID: 31778094.
- 524 23. Austin PC, Stuart EA. Moving towards best practice when using inverse probability of treatment
525 weighting (IPTW) using the propensity score to estimate causal treatment effects in observational studies.
526 (2015) 34(28):3661-79. doi: <https://doi.org/10.1002/sim.6607>.

527 24. Xiao Y, Abrahamowicz M, Moodie EE. Accuracy of conventional and marginal structural Cox model
 528 estimators: a simulation study. *Int J Biostat* (2010) 6(2):Article 13. Epub 2010/01/01. doi: 10.2202/1557-
 529 4679.1208. PubMed PMID: 21969997.

530

531

532

533 **Tables**

534 Table 1 – Classification of DMTs in the Optimise:MS study

| FIRST-GENERATION DMTs | | |
|------------------------------|--------------------|--|
| Drug | Product name(s) | Mode and frequency of delivery |
| Glatiramer acetate | Brabio, Copaxone | Subcutaneous, 3-7x weekly |
| Interferon beta-1a | Avonex | Intramuscular, weekly |
| Interferon beta-1a | Rebif | Subcutaneous, 3x weekly |
| Pegylated interferon beta-1a | Plegridy | Subcutaneous or intramuscular, every 2 weeks |
| Interferon beta-1b | Betaferon, Extavia | Subcutaneous, every 2 days |
| SECOND-GENERATION DMTs | | |
| Drug | Product name(s) | Mode and frequency of delivery |
| Alemtuzumab | Lemtrada | Intravenous infusion, 5 consecutive days followed by 3 consecutive days 1 year later |
| Cladribine | Mavenclad | Oral, up to 5 consecutive days per month for 2 months, repeated 1 year later |
| Daclizumab | Zinbryta | Subcutaneous, monthly |
| Dimethyl fumarate | Tecfidera | Oral, 2x daily |
| Fingolimod | Gilenya | Oral, daily |
| Natalizumab | Tysabri | Intravenous infusion, monthly |
| Ocrelizumab | Ocrevus | Intravenous infusion, 2x yearly |
| Ofatumumab | Kesimpta | Subcutaneous, monthly |
| Rituximab | Mabthera, Truxima | Intravenous infusion, up to 2x yearly |
| Siponimod | Mayzent | Oral, daily |
| Teriflunomide | Aubagio | Oral, daily |

535

536

537 **Figure Captions**

538 Figure 1 – Illustrations of the determination of exposure and control periods in the primary cohort
 539 analysis for two hypothetical patients, one in each stratum. The filled blocks represent the treatment
 540 received by the patient; the labels below indicate the periods of follow-up that contribute to the
 541 analysis.

542 Figure 2. Illustration of the determination of exposure and control periods in the secondary and
 543 tertiary cohort analysis for the two hypothetical patients shown in Figure 1. The filled blocks
 544 represent the treatment being received by the patient; the labels below indicate the periods of follow-
 545 up that contribute to the analysis.

546 Figure 3 – signal generation procedure. ROR = Reporting Odds Ratio; IRR = Incidence Rate Ratio;
 547 BCPNN = Bayesian Confidence Propagation Neural Network