



This is a repository copy of *Subgroup analysis and interpretation for phase 3 confirmatory trials: White paper of the EFSPI/PSI working group on subgroup analysis*.

White Rose Research Online URL for this paper:
<https://eprints.whiterose.ac.uk/140633/>

Version: Accepted Version

Article:

Dane, A. orcid.org/0000-0002-2997-4074, Spencer, A. orcid.org/0000-0002-4194-3561, Rosenkranz, G. et al. (3 more authors) (2019) Subgroup analysis and interpretation for phase 3 confirmatory trials: White paper of the EFSPI/PSI working group on subgroup analysis. *Pharmaceutical Statistics*, 18 (2). pp. 126-139. ISSN 1539-1604

<https://doi.org/10.1002/pst.1919>

This is the peer reviewed version of the following article: Dane A, Spencer A, Rosenkranz G, Lipkovich I, Parke T, on behalf of the PSI/EFSPI Working Group on Subgroup Analysis. Subgroup analysis and interpretation for phase 3 confirmatory trials: White paper of the EFSPI/PSI working group on subgroup analysis. *Pharmaceutical Statistics*. 2018, which has been published in final form at <https://doi.org/10.1002/pst.1919>. This article may be used for non-commercial purposes in accordance with Wiley Terms and Conditions for Self-Archiving.

Reuse

Items deposited in White Rose Research Online are protected by copyright, with all rights reserved unless indicated otherwise. They may be downloaded and/or printed for private study, or other acts as permitted by national copyright laws. The publisher or other rights holders may allow further reproduction and re-use of the full text version. This is indicated by the licence information on the White Rose Research Online record for the item.

Takedown

If you consider content in White Rose Research Online to be in breach of UK law, please notify us by emailing eprints@whiterose.ac.uk including the URL of the record and the reason for the withdrawal request.



eprints@whiterose.ac.uk
<https://eprints.whiterose.ac.uk/>



**Subgroup Analysis and Interpretation for Phase 3
Confirmatory Trials: White paper of the EFSP/PSI Working
Group on Subgroup Analysis**

Journal:	<i>Pharmaceutical Statistics</i>
Manuscript ID	PST-16-0105.R4
Wiley - Manuscript type:	Main Paper
Date Submitted by the Author:	24-Sep-2018
Complete List of Authors:	Dane, Aaron; DaneStat Consulting Limited Spencer, Amy; Statistical Services Unit, University of Sheffield Rosenkranz, Gerd; Medizinische Universität Wien Lipkovich, Ilya; Quintiles Innovation, Parke, Tom; Berry Consultants LLC
Key Words:	subgroup analysis, regulatory labelling, bias adjustment, late-phase clinical programs
Abstract:	<p>Subgroup by treatment interaction assessments are routinely performed when analysing clinical trials and are particularly important for Phase 3 trials where the results may affect regulatory labelling. Interpretation of such interactions is particularly difficult, as on one hand the subgroup finding can be due to chance, but equally such analyses are known to have a low chance of detecting differential treatment effects across subgroup levels, so may overlook important differences in therapeutic efficacy. EMA have therefore issued draft guidance on the use of subgroup analyses in this setting. Although this guidance provided clear proposals on the importance of pre-specification of likely subgroup effects and how to use this when interpreting trial results, it is less clear which analysis methods would be reasonable, and how to interpret apparent subgroup effects in terms of whether further evaluation or action is necessary.</p> <p>A PSI/EFSP/PSI Working Group has therefore been investigating a focused set of analysis approaches to assess treatment effect heterogeneity across subgroups in confirmatory clinical trials which take account of the number of subgroups explored, and also investigating the ability of each method to detect such subgroup heterogeneity. This evaluation has shown that the plotting of standardised effects, bias-adjusted bootstrapping method and SIDES method all perform more favourably than traditional approaches such as investigating all subgroup-by-treatment interactions individually or applying a global test of interaction. Therefore, these approaches should be considered to aid interpretation and provide context for observed results from subgroup analyses conducted for Phase 3 clinical trials.</p>

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

SCHOLARONE™
Manuscripts

For Peer Review

1
2
3
4 Subgroup Analysis and Interpretation for Phase 3 Confirmatory Trials: White paper of the
5
6 EFSPi/PSI Working Group on Subgroup Analysis
7
8
9

10 Authors: Aaron Dane¹, Amy Spencer², Gerd Rosenkranz³, Ilya Lipkovich⁴ and Tom Parke⁵ on behalf
11
12 of the PSI/EFSPi Working Group on Subgroup Analysis
13
14
15

16
17 ¹DaneStat Consulting, Macclesfield, UK, ²Statistical Services Unit, University of Sheffield, UK,
18

19 ³Medical University of Vienna, ⁴Quintiles, USA, ⁵Berry Consultants, UK.
20
21
22

23 Corresponding author: Aaron Dane, Email: aarondane@danestat.com, +44 7821 720 631
24
25
26
27
28

29 Running head: PSI WG evaluation of Subgroup Analysis approaches
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55

Abstract

Subgroup by treatment interaction assessments are routinely performed when analysing clinical trials and are particularly important for Phase 3 trials where the results may affect regulatory labelling. Interpretation of such interactions is particularly difficult, as on one hand the subgroup finding can be due to chance, but equally such analyses are known to have a low chance of detecting differential treatment effects across subgroup levels, so may overlook important differences in therapeutic efficacy. EMA have therefore issued draft guidance on the use of subgroup analyses in this setting. Although this guidance provided clear proposals on the importance of pre-specification of likely subgroup effects and how to use this when interpreting trial results, it is less clear which analysis methods would be reasonable, and how to interpret apparent subgroup effects in terms of whether further evaluation or action is necessary.

A PSI/EFSPi Working Group has therefore been investigating a focused set of analysis approaches to assess treatment effect heterogeneity across subgroups in confirmatory clinical trials which take account of the number of subgroups explored, and also investigating the ability of each method to detect such subgroup heterogeneity. This evaluation has shown that the plotting of standardised effects, bias-adjusted bootstrapping method and SIDES method all perform more favourably than traditional approaches such as investigating all subgroup-by-treatment interactions individually or applying a global test of interaction. Therefore, these approaches should be considered to aid interpretation and provide context for observed results from subgroup analyses conducted for Phase 3 clinical trials.

Key words: Subgroup analysis, regulatory labelling, bias adjustment, late-phase clinical programs

2

Introduction

It is common to conduct subgroup analyses in clinical trials during all phases of drug development, and to try to understand heterogeneity in treatment effect across various levels of subgroup factors of interest. This is of relevance in early phase studies in order to better understand which patients to study in future trials, whilst for late stage clinical trials is applied in order to understand the effects of a given treatment across a range of baseline factors. Consistency of treatment effect across trial subgroups indicates that the conclusions made regarding treatment benefit are applicable across various baseline characteristics and associated subpopulations, whilst substantial heterogeneity in treatment effect may suggest clinically relevant differential treatment effects across subpopulations. At the extreme, the presence of substantial heterogeneity may imply that the conclusion of beneficial treatment effect is only relevant for a subset of the population.

Interpretation of subgroup analyses is difficult ^[1] as any apparent heterogeneity can be due to chance, which is particularly likely when a large number of subgroup analyses are undertaken. Conversely, clinical trials are generally not designed for detecting subgroup heterogeneity, so statistical tests may miss important interactions due to low power ^[2]. This is particularly challenging as approaches need to account for both issues simultaneously, as illustrated by Gonnermann ^[3] who also concluded that further work was needed to understand alternative statistical approaches in this area. Despite these issues, potential subgroup differences cannot be ignored, so it is still necessary to understand potential subgroup heterogeneity ^[4], and as such, consideration needs to be given for how such effects are to be analysed and interpreted.

Given these issues, EMA proposed draft guidance in February 2014 ^[5] and held a Workshop in November 2014 to address these considerations. This guidance was generally well received, but a key area of further work related to providing information on methodological approaches and criteria to conclude consistency of effect. Two European Pharmaceutical industry statistical organizations, the European Federation for Statisticians in the Pharmaceutical Industry (EFSPI) and Statisticians in the Pharmaceutical Industry (PSI) therefore formed a working group to further explore possible analysis approaches in order to

1
2
3 inform subgroup evaluation in a regulatory environment. As this work is prompted by the
4 EMA draft guidance and possible regulatory labelling, approaches have focused on assessing
5 heterogeneity of treatment effect across subgroup levels in a Phase 3 trial. As such, this
6 paper assumes the situation where the analysis of the overall population showed a positive
7 effect before any subgroup effects were explored.
8
9

10
11
12 This paper summarises the key activities of the working group, and their possible
13 implications for future regulatory evaluation. As the focus is to understand whether
14 heterogeneity of treatment effect across subgroup levels is real, any reference to subgroup
15 analysis or subgroup effects will be in relation to this heterogeneity. This paper discusses
16 predictive effects only (factors driving a differential effect across treatments) and does not
17 consider prognostic factors (factors predictive of outcome, regardless of treatment). However,
18 factors that are known to be prognostic for the disease of interest would be a more plausible
19 set of candidate variables to explore for treatment effect heterogeneity.
20
21
22
23
24
25

26 27 28 **Pre-specification of subgroup effects**

29
30 When conducting subgroup analyses it is often necessary to assess the heterogeneity of
31 treatment effect across various levels of a number of subgroup factors. When conducting late
32 stage clinical trials there are a variety of reasons for this, ranging from a standard regulatory
33 authority requirement, their importance in a particular clinical setting, or to address key
34 reimbursement questions. Indeed, FDA requires that subgroups defined by gender, age and
35 race are to be analysed and dosage modifications to be identified for specific subgroups^[6],
36 and in 2014 the agency issued an action plan “to enhance the collection and availability of
37 demographic subgroup data”^[7]. Given the multiplicity challenges associated with
38 interpreting a large number of subgroup analyses, the first objective for regulators and
39 sponsors alike is to attempt to minimize the number of subgroups analysed wherever
40 possible, and regulators at both EMA and FDA have recommended that the focus should be
41 on the subgroups of primary interest^[5, 8].
42
43
44
45
46
47
48
49

50
51 All subgroups to be analysed in this way must be pre-specified, and it will often still be
52 necessary to choose a fairly large number of subgroups within a confirmatory clinical trial.
53
54

55 Therefore, it is also important to provide a clear justification for the plausibility of a
56
57

58
59
60

1
2
3 particular subgroup by treatment interaction. This will ease interpretation and will also make
4 for a stronger argument in relation to any subgroup effect found. Any prior evaluation of
5 plausibility should be based on historical data, external trial data obtained from literature or
6 the labels of approved products, or data from earlier trials in the clinical program (Phase I or
7 II) for the project under study. Further, it is important to specify the nature of the interaction
8 (direction of impact on overall outcome), and to state which changes are clinically relevant
9 and would impact labelling. From such pre-specification, it is possible to consider the
10 subgroups in 3 categories. The following categories below are similar to those outlined in the
11 EMA guidance:
12
13
14
15
16
17
18

19 Confirmatory (strong reason to expect a significant, clinically relevant heterogeneous
20 response):
21

22 In this case it is necessary to address the differential levels of effect for a specific factor
23 explicitly and is addressed in the design and interpretation of the trial (for example, through
24 type I error adjustment). This category will not be covered further in this paper.
25
26
27

28 Biologically plausible (existing biological rationale or external evidence of heterogeneous
29 response):
30

31 In this case, a heterogeneous response would be a differential treatment effect across
32 subgroups levels considered clinically meaningful, and as such would impact product
33 labelling. There are likely very few subgroups in this category, and EMA guidance states
34 that it is relevant to discuss and plan for an assessment of consistency of effects. This could
35 well be a situation where subgroups may have some prognostic effect, but there has been no
36 evidence of a treatment by subgroup interaction to date.
37
38
39
40
41
42

43 Hypothesis generating only (no prior evidence to expect a heterogeneous response):
44

45 This would include all other subgroups necessary for study based upon standard regulatory
46 requirements, or clinical practice. As these subgroups are not anticipated to show any
47 differential treatment effect, approaches to understand the effects in the context of the large
48 number of subgroups evaluated can be considered.
49
50
51
52

53 This categorization of subgroup analyses is helpful both when planning subgroup analyses
54 and when interpreting the results. For a subgroup to have some plausibility there must be a
55
56
57
58
59
60

1
2
3 prior hypothesis regarding a potential treatment-by-subgroup interaction. In particular, the
4 “biologically plausible” category represents hypothesized subgroup factors or baseline
5 covariates which may govern differential treatment response. These factors may be simple
6 (e.g., based on a single baseline characteristic) or complex (i.e., based on multiple baseline
7 characteristics simultaneously.) The number of proposed factors in this category is expected
8 to be small and identification of these factors, along with specific hypotheses regarding
9 direction and magnitude of anticipated differences are key to evaluating the robustness of any
10 subgroup findings once the results are available.
11
12
13
14
15
16
17

18 **Current approaches**

19
20 At present, it is common to conduct a statistical test for treatment-by-subgroup interaction by
21 adding these terms to the primary analysis model. This can involve assessing subgroup by
22 treatment interactions for a large number of subgroups, which inflates the chances of a false
23 positive finding ^[9].
24
25
26
27

28 Further, interaction tests have low power, a point raised within the EMA guidance which
29 states that lack of statistical significance of an interaction is not sufficient to conclude
30 subgroup by treatment homogeneity. An approach to improve the power when testing for an
31 interaction is to use arbitrary criteria such as $p < 0.10$. Whilst this makes it easier to detect an
32 interaction, this also increases the chances of false positive findings. Indeed, if 10
33 independent subgroup factors are explored in a trial, none of which have an actual effect, the
34 probability of observing at least one significant interaction at $p < 0.10$ is $1 - (1 - 0.10)^{10}$, or
35 approximately 65%. It is acknowledged that it could equally be argued that a more stringent
36 type I error rate should be applied as a result of the multiple testing of many subgroups ^[10]. It
37 is for this reason that the simulation exercise described has explored the type I error rate,
38 defined as the number of times a subgroup by treatment interaction is incorrectly identified
39 under the “null” that no such interaction is present.
40
41
42
43
44
45
46
47
48
49

50 An alternative approach is to perform variable selection in regression models, where
51 parameters associated with various interaction effects are added or removed from the model
52 depending upon variable importance. However, interpretation of this approach can be
53 challenging as very different conclusions can be reached dependent upon the selection
54
55
56
57
58
59
60

6

1
2
3 method applied (e.g., backward or forward selection). Furthermore, variable selection is
4 unstable where small perturbations in the data may lead to selection of different variables.
5
6 Inference has often been done without consideration of model selection, i.e., as if the selected
7
8 model had been pre-specified. Only recently exact methods have been developed for type I
9
10 error control or confidence intervals for parameter estimates after selection ^[11]. These
11
12 methods are still under debate. Therefore, this approach has not been explored further.
13

14
15 One way of providing some type I error control when investigating multiple subgroups is to
16
17 perform a global test of interaction (comparing model fit of models with no interaction terms,
18
19 and all interaction terms of interest), and this is justified when there is no biological rationale
20
21 to expect an interaction for the subgroups studied.

22
23 It is also common practice to produce univariate forest plots for each subgroup to help
24
25 understand the effect of an individual subgroup on overall response, and forest plots
26
27 including a reference line indicating the point estimate for the overall effect ^[12] will provide
28
29 information on any subgroups whose 95% CI excludes this point estimate for the overall
30
31 treatment effect. When interpreting such forest plots, those subgroup levels with a 95% CI
32
33 looking very different from the overall result would often require further discussion and
34
35 investigation of whether the overall result applies equally to all levels of this subgroup. EMA
36
37 guidance specifically addresses the use of forest plots but suggests caution as a formal rule
38
39 for interpretation that is “both sensitive enough to detect heterogeneity ...and specific
40
41 [enough to detect true subgroups] is not available”. Further, as such plots do not include
42
43 information on confounding of key covariates, they should generally be interpreted alongside
44
45 other methods to understand the likelihood of such findings in the presence of other important
46
47 factors.

48 49 50 **Alternative methods considered**

51
52 There is a wealth of literature describing possible approaches to subgroup analysis.
53
54 Lipkovich et al ^[13] have recently provided a review of the literature, much of which is
55
56 focused on the exploratory setting. Similarly, given the importance of subgroup analysis, the
57
58 whole of issue 24 of the Journal of Biopharmaceutical Statistics in 2014 was dedicated to
59
60 subgroup analysis in clinical trials ^[14] and prominent authors have provided a tutorial on the

7

1
2
3 statistical considerations of such analyses ^[15]. There have also been other papers which have
4 described ideas related to the use of permutation distributions in other settings, such as that to
5 assess regression to the mean ^[16,17], which although not directly related do show the approach
6 has a wide applicability for generating “null” distributions in the context of data driven
7 subgroup evaluation.
8
9

10
11
12 Given the purpose of this paper has been to identify approaches that could be used in the
13 setting of a confirmatory clinical trial and for regulatory labelling, approaches were explored
14 that were consistent with current practice, but also addressed the concerns regarding the
15 power to detect true subgroup heterogeneity, whilst controlling the rate of incorrect
16 identification of a subgroup due to multiplicity. As such, we consider a new method which
17 produces a plot of standardised effects and will provide context for the magnitude of
18 subgroup heterogeneity seen in relation to what would be expected by chance. This approach
19 is motivated by the desire to find an approach akin to a forest plot, but which takes account of
20 the number of subgroups analysed and the correlation between these subgroups. In addition,
21 we present approaches that either control type I error or provide estimates of treatment effect
22 after adjusting for subgroup selection, both of which are important when providing
23 information to regulatory agencies, treating physicians or for reimbursement.
24
25
26
27
28
29
30
31
32

33 This section will outline these approaches, whilst the following sections will describe a
34 simulation study used to evaluate the operating characteristics of these methods in different
35 scenarios.
36
37
38
39
40

41 (a) Plot of Standardised effects

42 A resampling based graphical method to present Standardised Effects Adjusted for Multiple
43 Overlapping Subgroups (SEAMOS) is proposed by Dane ^[18]. SEAMOS is intended to
44 provide a graphical presentation of the results for all pre-specified subgroups, and to illustrate
45 how extreme the results are expected to be by chance, given the number of subgroups
46 analysed and the correlation between the subgroups. This approach can be used when it is
47 possible to clearly divide a subgroup into categories, as would be expected in an evaluation
48 used to assess pre-specified subgroups likely to impact regulatory labelling. Briefly,
49 observed standardised effects for the difference between the subgroup level and overall effect
50
51
52
53
54

1
2
3 are ordered from largest to smallest to highlight the statistically most extreme standardised
4 effects in either direction. These are then compared to a “null” distribution of standardised
5 effects, or the distribution of effects expected by chance when there are no subgroup-by-
6 treatment interactions. This is calculated by using a resampling algorithm and is achieved by
7 randomly permuting the rows of covariates many times against the vector made up of the
8 response variable and the treatment identifier. This has the effect of removing any subgroup
9 by treatment effects but preserves the overall treatment effect and the correlation between
10 subgroups. The smallest and largest standardised effects from each permutation are taken
11 and used to produce a probability interval of extreme effects. If there were no subgroup
12 heterogeneity, we would expect the largest standardised effect to lie within this probability
13 interval 95% of the time. Similarly, we would expect the smallest standardised effect to lie
14 within the interval 95% of the time. These intervals can be used to assess the highest and
15 lowest observed values, with points lying outside these probability intervals requiring further
16 evaluation, as shown in Figure 1, controlling the chances of incorrectly identifying a
17 subgroup when there is no subgroup heterogeneity at $\leq 10\%$. Similarly, probability intervals
18 can be constructed for all ordered subgroups as described in Dane ^[18].

19
20
21
22
23
24
25
26
27
28
29
30 Given there are often challenges in using a formal method which makes firm conclusions
31 regarding statistical assessment criteria, the intention is to provide a method which places the
32 results into context in addition to providing a basis for further exploration of subgroups.
33 There are some approaches available that take a similar approach, but only work when the
34 subgroup levels are all fully independent (for example, when exploring the effects by country
35 in the clinical trial) ^[19,20,21,22].

41 42 (b) Adjusted subgroup estimates via bootstrapping

43 This bootstrapping method proposed by Rosenkranz ^[23,24] is primarily based on estimation.
44 The method fits a reference model containing treatment and all factors as main effects. For
45 each factor defining a subgroup, the method fits a model containing in addition the treatment
46 by factor interaction effect and picks the factor/model leading to the best fit to the data. As it
47 is acknowledged that the estimates provided for subgroup are exaggerated, particularly when
48 a large number of subgroups are investigated, this selection process is repeated on a series of
49 bootstrap samples of the original data to account and adjust for such bias, with the degree of
50 adjustment being dependent upon the amount of times the factor in question is selected from
51
52
53
54
55
56
57
58
59
60

1
2
3 the bootstrap samples (for example, if a factor is picked in every bootstrap sample very little
4 adjustment would occur, whilst much more adjustment would occur if the factor were
5 selected less frequently). The bootstrap samples are then used to calculate an interaction
6 effect estimator (and its standard deviation) adjusted for selection bias and model uncertainty.
7
8
9

10
11 The method allows for all data types and in principle also for continuous variables. Instead of
12 selecting only the factor with the best fit in most bootstrap samples, all factors selected more
13 often than expected by chance could be selected. Furthermore, the full population (i.e., the
14 model containing no treatment by factor interaction term) is part of the selection process such
15 that all subgroups compete with the full population in terms of goodness-of-fit, meaning it is
16 possible to select the model with no treatment effect heterogeneity.
17
18
19
20
21

22 (c) SIDES

23
24 The SIDES method (Subgroup Identification based on Differential Effect Search)^[25,26]
25 applies the following simple search strategy combined with a resampling approach in order to
26 adjust any significance level (or p-value) for the model selection approach. This method was
27 originally developed for controlling type I error in a more exploratory setting, (with
28 subgroups possibly defined as “signatures” of up to 3 continuous variables with data-driven
29 cut-offs) but could equally be applied (as a special case) to the setting where the search is
30 restricted to subgroups based on a single binary factor selected from a candidate set. This
31 method can also be used for all response types and is applied as follows:
32
33
34
35
36

- 37 • All candidate subgroups are evaluated and the covariate with the larger differential effect
38 across subgroup levels is selected. The subgroup level requiring further evaluation is
39 that which has the larger treatment effect.
- 40
41 • Then the null reference distribution for the treatment effect p-value of the selected
42 subgroup is constructed by randomly permuting the treatment labels and applying the
43 above procedure to the resulting data.
- 44
45 • This is repeated 1000 times and the adjusted p-value is determined as the proportion of
46 reference sets where the selected subgroup’s p-value is less than or equal to that found in
47 the observed data.
- 48
49 • For the purposes of this evaluation, a subgroup is selected if the adjusted p-value is below
50 the pre-specified cut-off of 0.1. Thus, the procedure ensures that the type I error of the
51 entire selection strategy is within 10%.
- 52
53
54
55

56 10
57
58
59
60

- To make the null distributions of p-values insensitive to the overall treatment effect the analysis data set is standardised before applying the procedure. This is achieved by subtracting the respective treatment arm's mean and dividing by the arm-specific standard deviation.

Simulation exercise to evaluate performance of the methods

Simulation scenarios used to evaluate operating characteristics

When investigating the performance of various methods, two broad simulation scenarios have been used, the first (“simple set”) where all subgroup factors were simulated independently and the second (“complex set”) where some correlation higher than that expected by chance was generated. For both approaches the operating characteristics will focus on the situation with a total of 10 binary subgroup factors. Note that we have not considered scenarios where the number of subgroup factors is much less than this (2 or 3 factors, for example), as this is not felt to be a realistic scenario for Phase 3 confirmatory trials where it is necessary to address a number of regulatory and clinical questions regarding subgroups. Similarly, we have not considered scenarios of an interaction between two subgroups and with treatment, as we aimed to simulate a situation where a Phase 3 trial is to be conducted and no subgroup effects are anticipated, so such a scenario seems unlikely.

Given all of these approaches have the potential to identify more than one subgroup as part of the procedure this may involve identifying both a correct and an incorrect subgroup. For the purposes of operating characteristic calculation, a second step is performed which explores which subgroup has the most compelling result - this “primary subgroup” is used to define whether a correct (or incorrect) subgroup identification has occurred. The only exception is for SEAMOS. As this approach presents information on both levels of each subgroup factor there is a small chance that two different subgroup factors will be identified as the most extreme positive and most extreme negative effects. In this case they have both exhibited results beyond those expected by chance and would both need to be explored further. As such, it is possible to identify both a correct and incorrect subgroup in this case.

1
2
3 The methods used for the simulation scenarios are similar to those described with full details
4 in Dane ^[18]. Briefly, this involved performing 1,000 simulations for each scenario and
5 assessing how often a true subgroup was identified (in scenarios with true treatment effect
6 heterogeneity for one subgroup factor), or how often a subgroup was incorrectly identified
7 when there was no treatment effect heterogeneity across subgroups. This number of
8 simulations was chosen as a balance between time taken for the simulation exercise and
9 achieving sufficient precision, and is considered reasonable as the standard deviation around
10 a 10% error rate and 1,000 simulations is estimated to be $0.95\% = \sqrt{\frac{(0.1*0.9)}{1000}}$. This was felt
11 sufficient for ensuring the error rates were reasonably well controlled.
12
13
14
15
16
17
18
19

20 The simulations used assumptions regarding an overall treatment effect and an enhanced
21 effect for one level of the subgroup (*S+* group), defined as k , where k took possible values of
22 1-, 2- or 4-times the overall treatment effect). Individual patient outcomes were simulated
23 from a normal distribution using the mean based on the allocation to treatment arm and to
24 subgroup *S+* or *S-*. In the “simple set” of simulations, for all other subgroups, defined by the
25 other 9 factors, the overall effect, θ , is assumed. Using these assumptions, 4 different sized
26 trials were considered by varying the magnitude of effect and size of study, and within each,
27 separate simulations were performed according to the proportion of patients in the subgroups.
28 Subjects had either a 0.1, 0.25 or 0.5 probability of being in the *S+* subgroup based on the
29 first factor, and this same proportional split was used for all other factors. This magnitude of
30 effect was assessed under “high power” and “normal power” scenarios, where “normal
31 power” is defined as 90% power to detect $p < 0.05$ (2-sided) for the overall effect, and “high
32 power” as 90% power to detect an effect with $p < 0.01$. Further, each method was assessed
33 when a smaller or larger trial was required for these scenarios to understand whether the
34 ability to detect heterogeneity was affected by having a smaller number of patients within the
35 subgroup. This resulted in sample sizes of 200 and 270 patients per group for smaller sample
36 size and the normal and high-power scenarios respectively, and 760 and 1080 patients per
37 group for the larger sample sizes.
38
39
40
41
42
43
44
45
46
47
48
49
50

51 The previous text describes the simple simulations, in which the 10 subgroup factors were all
52 generated independently. The “complex set” of simulations were produced which introduced
53 some inter-relationship between the true factor and one correlated “noise” factor. In this case
54
55

56 12
57
58
59
60

1
2
3 the correlation of these two binary variables was calibrated in order to achieve 50% overlap
4 (See details in the Appendix).
5
6

7
8 Results for the “simple” and “complex” set of simulation scenarios were very similar.
9
10 Therefore, the following section presents plots of the operating characteristics for only the
11 complex scenario as we feel this is likely to be closer to real data observed in a clinical trial.
12
13 Similar plots using a “subject level” evaluation which estimates how well the identified
14 subgroup from the data captures the patients belonging to the “ideal” subgroup (see Appendix
15 for further details) gave similar conclusions to those described in this paper.
16
17

18
19 When considering these scenarios, it is also important to consider the size of treatment effect
20 in the remaining population (θ_{S-}), and hence whether the interaction is qualitative or
21 quantitative. This is important as it can have implications for regulatory labelling and will
22 depend upon the proportion of subjects in the subgroup positive group (P_{S+}) and the
23 magnitude of differential effect in this positive subgroup level (k). Some examples are
24 presented when the overall treatment effect is 5 units (see Table 1). As can be seen in Table
25
26 1, when one subgroup level has a very large positive effect (e.g., $k=4$), the effect in the
27 remainder of the population may be either zero (for $P_{S+}=0.25$) or negative (for $P_{S+}=0.5$).
28
29 This must be considered alongside a review of the performance of the analysis methods, as
30 there are key implications regarding benefit-risk and whether a new treatment should be
31 approved in all patients.
32
33
34
35
36
37
38

39 **Results of simulation exercise**

40
41
42 In order to understand the performance of the methods outlined in this paper it is necessary to
43 assess how often a given approach correctly identifies subgroup heterogeneity and how often
44 a subgroup is incorrectly identified when there is no heterogeneity. Using the simulation
45 scenarios described, operating characteristics were calculated for traditional interaction
46 testing, the global interaction test and the methods outlined in the “Alternative methods”
47 section. Because traditional interaction testing is the only one of these methods which does
48 not attempt to correct for multiple testing, the results for this method are presented separately
49 in the “Performance of traditional interaction testing” section, whilst the evaluation of the
50
51
52
53
54

1
2
3 remaining methods are included in the section summarising “Performance of methods which
4 account for the evaluation of multiple subgroups”.

5 6 7 Performance of traditional interaction testing

8
9
10 The performance of the more traditional approach of fitting a multivariate model including all
11 subgroups involved defining a subgroup of interest if any of the subgroup-by-treatment
12 interaction terms from the statistical model were statistically significant using the criteria
13 $p < 0.1$. Table 2, below demonstrates that, as expected, when 10 subgroups are explored and
14 criteria of $p < 0.1$ is applied the type I error is very high. Table 2 also demonstrates that even
15 when there is true heterogeneity within one subgroup, other subgroups without this
16 heterogeneity are incorrectly identified on many occasions. This characteristic is concerning
17 when no subgroups interactions are anticipated, particularly when decisions are required
18 regarding labelling for such subgroups. The performance characteristics for the alternative
19 methods are presented in the next section and show that these methods perform more
20 favourably than the approach of assessing all interactions individually in this way.
21
22
23
24
25
26
27
28

29 Performance of methods accounting for the evaluation of multiple subgroups

30
31 The operating characteristics for the scenarios with 10 subgroup factors, where one of the
32 “unimportant” factors has an overlap of 50% with the one for which the “true” effect was
33 generated are given in Figure 2. The panels a-d show the results for four sample size and
34 overall effect size combinations.
35
36
37
38

39 The evaluation of the various approaches for the scenarios with 10 subgroup factors, both
40 with independent subgroups (not shown) and with an enhanced correlation between the true
41 subgroup and one other subgroup showed that the global interaction (GI) test controlled the
42 type I error at approximately 10%. SEAMOS also controlled the type I error at 10% for
43 subgroup sizes 0.1 and 0.25, and at a slightly lower rate for a subgroup size of 0.5. The
44 bootstrapping method tended to result in higher type I error with the smaller sample sizes of
45 200 or 270 patients, but this higher type I error was not seen with larger sample sizes.
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

14

1
2
3 For each method, the power to detect true subgroup heterogeneity increased with true
4 subgroup size (P_{S+}), subgroup effect (k) and sample size (for the same true effect size), as
5 would be expected. Comparing the power of the different methods, for the scenarios where
6 the true treatment effect in one subgroup level was twice as large as the treatment effect for
7 the population as a whole ($k=2$), the bootstrapping method and then SEAMOS had the
8 highest power, whilst SIDES had lower power along with the lower type I error stated
9 previously. However, SIDES performs more preferably than the GI test at the lower
10 subgroups sizes of 0.1 and 0.25. When the true treatment effect within the better performing
11 level of the subgroup was 4-times the treatment effect in the population as a whole ($k=4$), all
12 three of the newer methods had similar power to detect the subgroup by treatment interaction
13 and performed better than the GI test.
14
15
16
17
18
19
20
21

22 When many subgroups are analysed simultaneously and there is a subgroup with a truly
23 heterogeneous effect it is also possible to incorrectly identify another subgroup as exhibiting
24 treatment effect heterogeneity. As a result, we also explored the number of times an incorrect
25 subgroup was identified within the simulation scenarios when there was true subgroup
26 heterogeneity for one subgroup factor. The results from this exploration were similar to those
27 presented previously, in that the SIDES method tended to give the lowest rates of incorrect
28 identification and SEAMOS also tended to give lower values than the GI test. Regarding the
29 bootstrapping method, the degree of incorrect identification again tended to depend on the
30 parameters of the simulation scenario. In the scenario with subgroup effects twice that of the
31 overall effect, where the overall effect size was 0.33 and subgroup size was 0.1 or 0.25, this
32 method gives comparatively high incorrect identification rates, whilst for other simulations
33 the results are similar to SEAMOS.
34
35
36
37
38
39
40
41
42
43

44 **Conclusions and discussion**

45 The EMA draft guidance is a great advancement in terms of the pre-specification and
46 interpretation of subgroup analyses. The work adopted by this Working Group is looking to
47 add guidance on how to best analyse the data when a reasonably large number of subgroups
48 are considered, and when the results have the potential to impact regulatory labelling and
49 patient access to medicines.
50
51
52
53
54
55

1
2
3 With any subgroup analysis, pre-specification of the more plausible subgroup-by-treatment
4 interactions is key to effective interpretation. This is a point EMA have stated very clearly in
5 their guidance ^[5], and a paper by FDA authors ^[8] also talks of reducing the number of
6 subgroups. Our paper considers the situation where effects are, a priori, assumed to be
7 consistent (or sufficiently consistent to apply the overall result to that subgroup), and this
8 premise would ideally be agreed with regulatory agencies at the design stage of a trial. By
9 this we mean that the subgroups to be explored can be pre-specified in the “Biologically
10 plausible” or “Hypothesis generating only” categories. This investigation would be
11 performed separately within these two categories. Situations where there is a stronger belief
12 regarding subgroup heterogeneity have not been considered, as this should lead to a different
13 study design, and not be addressed at the analysis stage.

14
15
16
17
18
19
20
21
22 When exploration of the “Biologically plausible” or “Hypothesis generating only” categories
23 is required, methods such as SEAMOS, SIDES and the adjusted effect size estimates via
24 bootstrapping provide useful context when interpreting subgroup analyses and compare
25 favourably with traditional approaches which test all subgroups independently and have very
26 high error rates.

27
28
29
30
31 The graphical SEAMOS method and bias-adjustment methods presented should be
32 considered key tools to help quantify whether subgroup findings are real. They have benefits
33 over traditional methods of analysis such as assessing all interactions separately, stepwise
34 regression or a global test of interaction, as they appropriately account for the number of
35 subgroups analysed and the correlation between subgroups. It should be noted that we do not
36 advocate the use of any of these methods in isolation, but rather suggest presenting their
37 results along with the observed effects to provide context when assessing the likelihood that
38 any subgroup heterogeneity is real.

39
40
41
42
43
44
45
46 The situation when each of these three methods will be most useful will depend upon the key
47 aims for the specific analysis in question. For example, the number of subgroups being
48 explored, or the degree of focus on type I error control will affect when each method should
49 be adopted. Similarly, whether the subgroups are categorical or continuous (or whether there
50 is an established method of categorizing continuous variables) will be critical. This is
51 particularly true of SEAMOS which, although providing very useful context for the observed
52

1
2
3 results, can only be undertaken when a clear method of dividing a subgroup into categories is
4 possible. However, such categorization is most likely in the case of the late phase
5 confirmatory clinical trials considered here. Finally, the most appropriate method will also
6 depend upon whether the aim of additional subgroup work is to provide context for the
7 observed results (SEAMOS), provide adjusted estimates of treatment effect (bias-adjusted
8 bootstrapped estimates) or provide adjusted p-values of subgroup effects (SIDES).
9
10
11
12
13

14 The simulation scenarios investigated have included a set of independent subgroups and those
15 with a degree of overlap above that expected by chance. It is acknowledged that there are
16 many other possible scenarios with respect to the overlap (or correlation) between subgroups,
17 and it has not been possible to assess all of these. However, we have no reason to believe we
18 would see differences with regard to the conclusions regarding type I error and power.
19 Additional simulation scenarios could be the subject of further work, but it would also be
20 important for anybody applying these approaches to demonstrate the type I error control
21 based upon observed subgroups and the relationship between them.
22
23
24
25
26
27

28 The approach to error control in this setting is a balance between controlling the rate of
29 incorrect subgroup identification (or incorrect labelling restrictions), with incorrectly
30 allowing a broad label. The approaches explored have looked to control the incorrect
31 subgroup identification error at ~10%. Other values could be used in discussion with
32 regulatory agencies, and would depend upon the magnitude of overall effect, the therapeutic
33 index and number of other therapeutic options in that particular situation.
34
35
36
37
38

39 Bayesian methods have been considered extensively when evaluating subgroups
40 [27,28,29,30,31,32] and are appealing as they incorporate prior beliefs regarding subgroup
41 heterogeneity into the evaluation. In our preliminary work we considered the “Simple
42 Regression” and “Dixon and Simon” methods reviewed in Jones [28] as these methods
43 appeared promising in the setting we have outlined in this paper. We encountered issues
44 related to model fitting and did not to pursue this further at that time, but additional work is
45 ongoing exploring a range of Bayesian approaches to provide recommendations on how such
46 methods may be applied to the regulatory setting with pre-defined subgroups.
47
48
49
50
51
52

53 This paper has focused on the situation when only one confirmatory Phase 3 trial has been
54 conducted, as any unexpected subgroup heterogeneity is most challenging to interpret in that
55
56
57
58
59
60

1
2
3 setting. When a number of studies have been conducted, this will provide additional, critical
4 information in terms of whether the heterogeneity of treatment effect is replicated, or whether
5 the magnitude of differential treatment effect is consistent across trials. When heterogeneity
6 is not replicated this casts more doubt on the likelihood of a differential effect and should be
7 used alongside an assessment of biological plausibility when interpreting whether the
8 heterogeneity is more likely to be real, or a chance finding.
9
10
11

12
13 These methods could also be used for post-hoc evaluation of efficacy, or when the overall
14 treatment effect is not positive. This special case has not been addressed in this paper as this
15 work is in the context of the regulatory setting where an overall treatment effect would be
16 necessary before investigating consistency in subgroups. Further, post-hoc subgroup analysis
17 and/or analysis in the setting of a negative treatment effect can only be considered hypothesis
18 generating as the implications for type I error are less clear, and the lack of pre-specification
19 of likely subgroup effects (and hence plausibility) makes such results very difficult to
20 interpret.
21
22
23
24
25
26

27
28 In summary, the approaches investigated in this paper regarding SEAMOS, SIDES and the
29 adjusted effect size estimates via bootstrapping all provide useful context when investigating
30 and interpreting subgroup heterogeneity and, when presented alongside the observed data and
31 used in conjunction with a clear pre-specification of more plausible subgroups, are useful
32 tools for interpreting apparent subgroup effects. Ongoing work is defining how these
33 approaches can be used within the regulatory review setting, whether they can contribute in
34 the review of reimbursement dossiers, and how such statistical approaches can be used to
35 inform an assessment of the benefit risk profile of a new treatment.
36
37
38
39
40
41

42 **Acknowledgments.**

43
44 The conclusions described within this document represent the work of the European
45 Federation for Statisticians on the Pharmaceutical Industry (EFSPI) and Statisticians in the
46 Pharmaceutical Industry (PSI) Working Group on Subgroup analysis. The work of the
47 Project Team is gratefully acknowledged, especially because it depends upon the generous
48 contributions of personal time. The EFSP/PSI Working Group members are as follows:
49 Aaron Dane (WG lead, DaneStat Consulting), Chrissie Fletcher (Amgen), Heiko Goette
50 (Merck), Necdet Gunsoy (GSK), Ilya Lipkovich (Quintiles), Henrik Loft (Lundbeck), Brian
51
52
53
54
55
56
57
58
59
60

1
2
3 Millen (Lily), Tom Parke (Berry Consultants), Arne Ring (Medac), Gerd Rosenkranz
4 (Medical University of Vienna), Amy Spencer (University of Sheffield), David Svensson
5 (AstraZeneca).
6
7

8
9 The initial ideas for this work were developed when Aaron Dane and Amy Spencer were
10 employed by AstraZeneca, and Gerd Rosenkranz was employed by Novartis and under
11 funding provided by the UK Medical Research Council, Project No. MR/M005755/1.
12
13

14 Therefore, we would like to thank AstraZeneca and Novartis for their support in this research.
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55

References

- [1] Wittes, J. On looking at subgroups. *Circulation*. 2009; 119: 912-915.
- [2] Wang, R., Lagakos, S. W., Ware, J. H., Hunter, D. J., & Drazen, J. M. Statistics in medicine—reporting of subgroup analyses in clinical trials. *New England Journal of Medicine*, 2007; 357(21), 2189-2194.
- [3] Gonnermann A, Kottas M, Koch A. Biometrische Entscheidungsunterstützung in Zulassung und Nutzenbewertung am Beispiel der Implikationen von heterogenen Ergebnissen in Untergruppen der Studienpopulation. *Bundesgesundheitsbl - Gesundheitsforschung - Gesundheitsschutz*. 2015; 58: Issue 3, pp 274-282
- [4] Koch, A., Framke, T. Reliably basing conclusions on subgroups of randomized clinical trials. *Journal of Biopharmaceutical Statistics*. 2014; 24:1, 42-57.
- [5] EMA. Guideline on the investigation of subgroups in confirmatory clinical trials. European Medicines Agency/Committee for Medicinal Products for Human Use. 2014; EMA/CHMP/539146/2013.
- [6] FDA, Content and format of a new drug application. 1985; 21 CFR 314.50[1]
- [7] FDA, FDA action plan to enhance the collection and availability of demographic subgroup data. 2014.
- [8] Alosch M, Fritsch K, Huque M, Mahjoob K, Pennello G, Rothmann M, Russek-Cohen E, Smith F, Wilson S, Yue L. Statistical Considerations on Subgroup Analysis in Clinical Trials. *Statistics in Biopharmaceutical Research*; 2015; 7(4): 286-303.
- [9] Assmann SF, Pocock SJ, Enos LE, Kasten LE. Subgroup analysis and other (mis)uses of baseline data in clinical trials. *Lancet* 2000; 355(9209): 1064-69
- [10] Lagakos SW. The challenge of subgroup analyses - reporting without distorting. *New England Journal of Medicine* 2006; 354;16: 1667-69.
- [11] Tibshirani RJ, Taylor T, Lockhart R, Tibshirani R. Exact PostSelection Inference for Sequential Regression Procedures, *JASA* 2016, 111, 600-620.
- [12] Cuzick, J. Forest plots and the interpretation of subgroups. *The Lancet*, 2005; 365(9467), 1308.
- [13] Lipkovich I, Dmitrienko A, D'Agostino Snr R. Tutorial in biostatistics: data-driven subgroup identification and analysis in clinical trials. *Statistics in Medicine*; 2017; 36(1); 136-196.
- [14] Wang, Sue-Jane, and Alex Dmitrienko. "Guest Editors' Note: Special Issue on Subgroup Analysis in Clinical Trials." *Journal of biopharmaceutical statistics*. 2014; 24(1): 1-3.
- [15] Alosch M, Huque MF, Bretz F, D'Agostino RB Sr. tutorial on statistical considerations on subgroup analysis in confirmatory clinical trials. *Statistics in medicine* 2017 36(8): 1334-1360.
- [16] Farlow MR, Small GW, Quarg P, Krause A. Efficacy of Rivastigmine in Alzheimer's Disease Patients with Rapid Disease Progression: Results of a Meta-Analysis (Subgroup analysis with corrected p-values for a regression to the mean effect). *Dementia and Geriatric Cognitive Disorders*, 20(2-3): 192-197, 2005.
- [17] Krause A, Pinheiro J. Modeling and Simulation to adjust p-values in Presence of a Regression to the Mean Effect. *The American Statistician*, 61 (4): 302-307, 2007.
- [18] Dane A, Spencer A, Stone A, Svensson D. The use of Standardised Effect Plots when Interpreting Analyses of Overlapping Subgroups. Submitted to *Statistics in Biopharmaceutical Research*, 2017.

20

- 1
2
3 [19] Anzures-Cabrera J. and Higgins, JPT. Graphical displays for meta-analysis: An overview with suggestions for
4 practice. *Res. Synth. Methods*, 2010; 1(1): 66–80.
- 5
6 [20] Galbraith, Rex. Graphical display of estimates having differing standard errors. *Technometrics*, 1988; 30 (3):
7 271–281.
- 8
9 [21] Chen J, Zheng H, Quan H, Li G, Gallo P, Soo PO, Binkowitz B, Ting N, Tanaka Y, Luo X, Ibia E, and for the Society for
10 Clinical Trials (SCT). Multi-regional Clinical Trial (MRCT) Consistency Working Group. Graphical assessment of
11 consistency in treatment effect among countries in multi-regional clinical trials. *Clin Trials*, 2013; 10(6): 842-
12 851.
- 13
14 [22] Schou I. M., and C. Marschner I. Methods for exploring treatment effect heterogeneity in subgroup analysis: an
15 application to global clinical trials, *Pharmaceut. Statist*, 2015; 14(1), 44–55.
- 16
17 [23] Rosenkranz, GK. Bootstrap corrections of treatment effect estimates following selection. *Computational*
18 *Statistics and Data Analysis*, 2014; 69, 220–227.
- 19
20 [24] Rosenkranz, GK. Exploratory subgroup analysis in clinical trials by model selection. *Biometrical Journal*, 2016;
21 58(5), 1217-1228.
- 22
23 [25] Lipkovich, I., Dmitrienko, A., Denne, J., Enas, G.. Subgroup Identification based on Differential Effect Search
24 (SIDES): A recursive partitioning method for establishing response to treatment in patient subpopulations.
25 *Statistics in Medicine*, 2011; 30(21): 2601–2621.
- 26
27 [26] Lipkovich, I., Dmitrienko, A. Strategies for Identifying Predictive Biomarkers and Subgroups with Enhanced
28 Treatment Effect in Clinical Trials Using SIDES, *Journal of Biopharmaceutical Statistics*, 2014; 24:1, 130-153.
- 29
30 [27] Millen BA, Dmitrienko A, Song G. Bayesian Assessment of the Influence and Interaction Conditions in
31 Multipopulation Tailoring Clinical Trials. *Journal of Biopharmaceutical Statistics*, 2014; 24:1, 94-109.
- 32
33 [28] Jones HE, Ohlssen DI, Neuenschwander B, Racine A, Branson M. Bayesian models for subgroup analysis in
34 clinical trials. *Clinical Trials*, 2011; 8(2): 129–143.
- 35
36 [29] Berger JO, Wang X, Shen L. A Bayesian Approach to Subgroup Identification. *Journal of Biopharmaceutical*
37 *Statistics*, 2014; 24: 110–129.
- 38
39 [30] Henderson NC, Louis TA, Rosner GL, Varadhan V. Individualized Treatment Effects with Censored Data via Fully
40 Nonparametric Bayesian Accelerated Failure Time Models, 2017. arXiv:1706.06611v1 [stat.ME], available at
41 <https://arxiv.org/abs/1706.06611>
- 42
43 [31] Gu X, Chen N, Wei C, Liu S, Vassilike A, Herbst RS, Lee JJ. *Statistics in Biosciences*, 2016; 8:1; 99-128.
- 44
45 [32] Hsu Y-Y, Zalkikar J, Tiwari RC. Hierarchical Bayes approach for subgroup analysis. *Statistical Methods in Medical*
46 *Research*, July 2017. <https://doi.org/10.1177%2F0962280217721782>
- 47
48
49
50
51
52
53
54
55

Table 1 Effect sizes in subgroup positive group and in complement of that subgroup for the effect sizes and subgroup sizes outlined in the simulation scenarios.

Overall treatment effect (θ)	Differential Effect in positive subgroup level (k)	Treatment Effect for subgroup positive subgroup level (θ_{S+})	Size of positive subgroup level (P_{S+})	Treatment effect in remaining population (θ_{S-})
5 units	4 x overall	20 units	0.1	3.33 units
			0.25	0 units
			0.5	-2.5 units
	2 x overall	10 units	0.1	4.44 units
			0.25	3.33 units
			0.5	0 units

Table 2. Operating characteristics for traditional analysis assessing treatment-by-factor interaction terms at $p < 0.10$

Overall Treatment Effect/standard deviation	N per arm	Subgroup size (% of overall population)	Type I error, %	Probability of detecting true subgroup (power) [Probability of detecting incorrect subgroup], %	
				2x overall effect	4x overall effect
0.167	760	10%	64.3	20.7 [60.0]	81.1 [60.3]
		25%	67.5	45.5 [62.0]	99.9 [61.2]
		50%	65.7	90.9 [62.2]	100 [61.0]
	1080	10%	64.8	26.7 [62.9]	91.3 [60.4]
		25%	64.4	60.2 [63.3]	100 [61.4]
		50%	63.5	97.8 [60.7]	100 [59.4]
0.33	200	10%	66.4	21.9 [58.8]	76.4 [62.1]
		25%	63.2	46.3 [57.5]	99.9 [60.8]
		50%	67.5	91.7 [58.6]	100 [60.6]
	270	10%	63.0	25.8 [60.7]	90.2 [61.3]
		25%	66.4	56.7 [59.9]	100 [60.6]
		50%	63.3	97.6 [59.2]	100 [64.1]

Notes:

1. Estimates in this table are approximate, and are based upon 1000 simulations; Simulations present number of significant treatment x covariate interaction effects ($p < 0.10$) in traditional regression analysis based on 10 covariates (1 factor with "real" treatment interaction and 9 factors with no interaction).
2. The probability of at least 1 of 9 factors with no interaction showing a significant interaction effect is $1 - ((1 - 0.10)^9) \approx 61\%$ (similar to squared brackets in columns "2x overall effect" and "4x overall effect").
3. The probability of at least 1 of 10 factors with no interaction showing a significant interaction effect is $1 - ((1 - 0.10)^{10}) \approx 65\%$ (similar to column "Type I error, %").

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

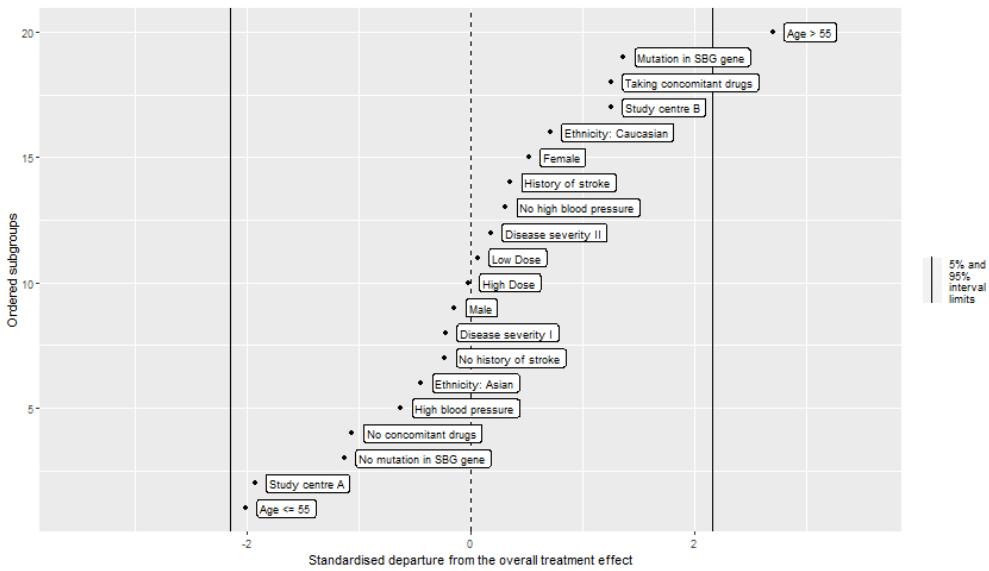


Figure 1: Example SEAMOS plot with probability intervals for the most extreme effects.

296x169mm (72 x 72 DPI)

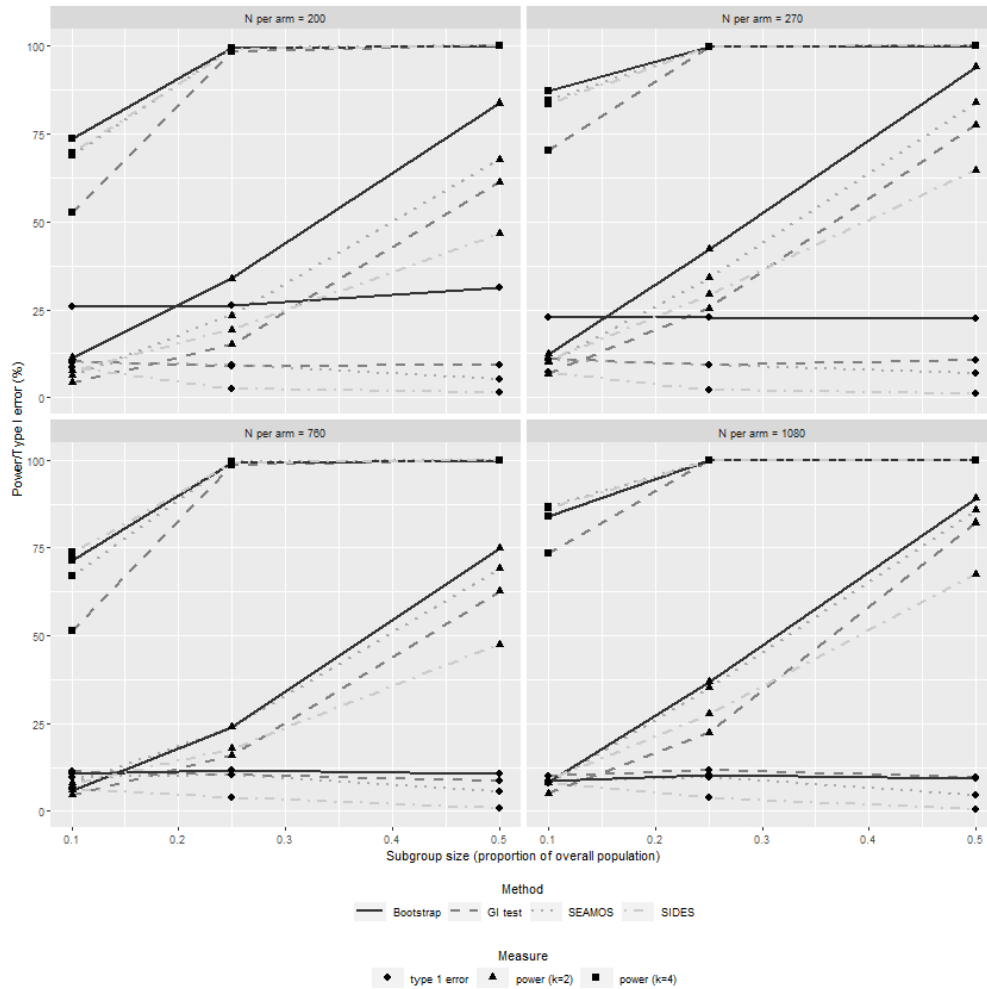


Figure 2: Summary of operating characteristics comparing the subgroup analysis approaches under investigation.

Panels present four sets of simulation results for sample sizes of 200 per group and 270 per group, both with an overall large treatment size; and 720 per group and 1080 per group, both with an overall small treatment size.

296x296mm (72 x 72 DPI)

APPENDIX

Framework for simulation scenarios when simulating binary subgroups with desired overlap

Summary

Following the simple simulation scenarios with 10 subgroup factors simulated independently, further work regarding more complex scenarios was explored. Here we describe the methods used for simulating these slightly more complex scenarios. The basic set-up is as follows:

- A number of subgroup factors is again 10, two of which explain a degree of the subgroup by treatment interaction.
- Simulate data assuming one of these subgroups has an important interaction as before.
- Simulate the second subgroup on the basis of the correlation/overlap between subgroups 1 and 2. This assumes an overlap of 0.5, which is higher than expected by chance for all the simulation scenarios.
- All other subgroups are simulated independently of these first two, as previously.
- Reapply for all previous scenarios in terms of sample size, subgroup size, overall treatment effect and subgroup treatment effect.
- Define performance criteria on identifying either the “correct” subgroup (1) or the “correct” patients (those in subgroup 1).

Simulation of correlated variables

Let X_1, \dots, X_{10} be binary covariates defining “elementary” subgroups. We define the “overlap” measure between two binary covariates X_1 and X_2 using Jaccard coefficient as the ratio of the probability that a subject belongs to both subgroups to the probability that s/he belongs to either of them.

$$J_{12} = \frac{\Pr(\{X_1 = 1\} \cap \{X_2 = 1\})}{\Pr(\{X_1 = 1\} \cup \{X_2 = 1\})}$$

The amount of overlap is a measure of correlation between binary variables, and can be calculated using the underlying multivariate normal distribution.

The correlation can be easily evaluated numerically for the case of bivariate normal.

Criteria of subgroup identification performance evaluation

It is possible to consider two criteria that measure performance of the different analysis methods at:

- Subgroup level. How well was the predictive subgroup identified?
- Subject level. How closely are the subjects who are be in the “true” subgroup approximated with the identified subgroup?

The methods presented in this paper have focussed on the subgroup level identification rate. The subject level identification rate shows higher success rates for all methods, but does not change the conclusions regarding the relative performance of the methods.

Subgroup level

Here we evaluate how well a subgroup identification method can identify true predictors. This is essentially the same analysis used in the simple scenario, but considers which (if any) of the subgroups are picked out by the analysis and uses the following measures:

- Type I error: In the case where $k=1$ (no subgroup effect), if a subgroup is identified in the analysis, this is a type I error.
- Power: Where there is a subgroup effect ($k>1$), if the correct subgroup (S_1) is identified, this contributes to power.
- Incorrect ID: Where there is a subgroup effect, if any subgroup other than the correct one (not S_1) is identified, this is an incorrect ID.

Subject level

Here we evaluate how close the identified subgroup is to some ideal one in the sense that the identified subgroup captures most of the patients defined by the ideal subgroup (even if it may fail to capture the correct subgroup factor). We use standard measures that are used in evaluating the quality of predicting a target binary outcome (here the true subgroup, S) by the identified subgroup, \hat{S} . In this case, we define the true subgroup S as subgroup 1, so this method will give excellent results when subgroup 1 is identified, good results when subgroup 2 (the highly correlated subgroup) is identified and poor results when another subgroup (based on one of the markers uncorrelated with the true subgroup) is identified. The performance measures we use are:

- sensitivity = $|\hat{S} \cap S|/|S|$
- specificity = $|\hat{S}^c \cap S^c|/|S^c|$, where S^c is complement of S .