EUnetHTA 21

**EUnetHTA 21 – Individual Practical Guideline Document**

**D4.5 – APPLICABILITY OF EVIDENCE – PRACTICAL GUIDELINE ON MULTIPLICITY, SUBGROUP, SENSITIVITY AND POST HOC ANALYSES**

**Version 0.3, 04/07/2022**
Template version 1.0, 03/03/2022

## 29    DOCUMENT HISTORY AND CONTRIBUTORS

| Version | Date | Description |
|---------|------|-------------|
| V0.1 | 23/03/2022 | First draft for CSCQ and NC-HTAb review |
| V0.2 | 25/05/2022 | Second draft for CSCQ and NC-HTAb review |
| V0.3 | 04/07/2022 | Third draft for public consultation |

### 30    Disclaimer

31   This Practical Guideline was produced under the Third EU Health Programme through a service contract
32   with the European Health and Digital Executive Agency (HaDEA) acting under mandate from the
33   European Commission. The information and views set out in this Practical Guideline are those of the
34   author(s) and do not necessarily reflect the official opinion of the Commission/Executive Agency. The
35   Commission/Executive Agency do not guarantee the accuracy of the data included in this study. Neither
36   the Commission /Executive Agency nor any person acting on the Commission's/Executive Agency's
37   behalf may be held responsible for the use which may be made of the information contained herein.

### 38    Participants

| Hands-on Group | Gemeinsamer Bundesausschuss [G-BA], Germany<br>Haute Autorité de Santé [HAS], France<br>Institut für Qualität und Wirtschaftlichkeit im Gesundheitswesen [IQWIG], Germany<br>Norwegian Medicines Agency [NOMA], Norway |
|---|---|
| Project Management | Zorginstituut Nederland [ZIN], the Netherlands |
| CSCQ<br>CEB | Agencia Española de Medicamentos y Productos Sanitarios [AEMPS], Spain<br>Austrian Institute for Health Technology Assessment [AIHTA], Austria<br>Belgian Health Care Knowledge Centre [KCE], Belgium<br>Gemeinsamer Bundesausschuss [G-BA], Germany<br>Haute Autorité de Santé [HAS], France<br>Institut für Qualität und Wirtschaftlichkeit im Gesundheitswesen, [IQWIG], Germany<br>Italian Medicines Agency [AIFA], Italy<br>National Authority of Medicines and Health Products [INFARMED], Portugal<br>National Centre for Pharmacoeconomics [NCPE], Ireland<br>National Institute of Pharmacy and Nutrition [NIPN], Hungary<br>Norwegian Medicines Agency [NOMA], Norway<br>The Dental and Pharmaceutical Benefits Agency [TLV], Sweden<br>Zorginstituut Nederland [ZIN], The Netherlands |

39   The work in EUnetHTA 21 is a collaborative effort. While the agencies in the Hands-on Group will be actively writing the
40   deliverable, the entire EUnetHTA 21 consortium is involved in its production throughout various stages. This means that the
41   Committee for Scientific Consistency and Quality (CSCQ) will review and discuss several drafts of the deliverable before
42   validation. The Consortium Executive Board (CEB) will then endorse the final deliverable before publication.

### 43    Copyright

45

# TABLE OF CONTENTS

78    **LIST OF ACRONYMS - INITIALISMS**

| | |
|---|---|
| CEB | Consortium executive board |
| CER | Comparisonwise error rate |
| CSCQ | Committee for scientific consistency and quality |
| EUnetHTA | European network of health technology assessment |
| FWER | Familywise error rate |
| HaDEA | European Health and Digital Executive Agency |
| HOG | Hands-on group |
| HTA | Health technology assessment |
| HTAb | Health technology assessment body |
| HTD | Health technology developer |
| ICE | Intercurrent event |
| ICH | International Council for Harmonisation of Technical Requirements for Pharmaceuticals for Human Use |
| JCA | Joint clinical assessment |
| NMA | Network meta-analysis |
| PICO | Population, intervention, comparator, outcome |
| RCT | Randomised controlled trial |
| SAP | Statistical analysis plan |

79

80

81

# 1   INTRODUCTION

## 1.1   *Problem statement, scope and objective of the guideline*

The design and conduct of a clinical study such as a randomised controlled trial (RCT) and analysis of its results are aimed at answering specific research question(s) defined by a health technology developer (HTD). Guidelines on statistical principles that should be used in this context have been developed by the International Council for Harmonisation of Technical Requirements for Pharmaceuticals for Human Use (ICH) [1].

During health technology assessment (HTA) under Regulation (EU) 2021/2282 (the EU HTA regulation), joint clinical assessment (JCA) starts with the assessment scope (Article 8(6) of the regulation; see EUnetHTA 21 Practical Guideline D.4.2.1 *Scoping process*). The assessment scope, expressed as one or more PICO (population, intervention, comparator, and outcome) questions, defines the needs of member states in terms of evidence that should be submitted by the HTD. Thus, these PICO questions are research questions defined by the member states for HTA purposes.

When drafting a JCA report, the assessor and co-assessor should not include any value judgement or conclusion on the clinical added value of the health technology assessed (Article 9(1)). The report should be limited to a factual assessment of the effectiveness and the certainty of results, considering the strengths and limitations of the available evidence submitted by the HTD. The outcome of the JCA should not affect the discretion of member states to draw conclusions regarding the clinical added value of the health technology assessed (Article 9).

When assessing the clinical added value of a health technology at a national level, member states are required to give due consideration to the JCA reports published (Article 13(1)). However, different member states can consider certain methodological aspects differently, especially because they can approach consistency or mismatches between the research question(s) as investigated by the HTD and the PICO question(s) differently. These differences in approach can be related to methodological issues pertaining to multiple hypothesis testing, subgroup, sensitivity, or post hoc analyses. Therefore, adequate assessment and reporting of the necessary elements of methods and results within the JCA report are needed to allow member states to draw their own conclusions regarding the clinical added value of a health technology. Moreover, studies submitted by HTDs as evidence can be individual clinical studies and/or evidence synthesis studies. The way in which different member states consider the aforementioned methodological aspects can be impacted by this distinction. Therefore, multiplicity, subgroup, sensitivity, and post hoc analyses are discussed in this guideline from the perspective of assessing individual clinical studies and assessing evidence synthesis studies in separate sections.

The intent for this practical guideline is not to endorse a particular approach regarding appraisal of the aforementioned methodological issues when member states draw their own conclusions at the national level. Rather, the objective is to guide assessors and co-assessors in assessing and reporting all the necessary elements that member states need to carry out national assessment of the clinical added value of the health technology regarding multiple hypothesis testing and subgroup, sensitivity, and post hoc analyses. Thus, all the requirements for assessment and reporting mentioned in this guideline assume that HTDs present the necessary elements in the submission dossier.

Multiplicity, subgroup, sensitivity, and post hoc analyses are not the only methodological aspects to consider when assessing the clinical added value of a health technology. Complementary elements in the reporting and assessment of the certainty of results (internal validity, statistical precision (e.g., confidence intervals), applicability) for individual clinical studies are described in EUnetHTA 21 Practical Guideline D4.6.1 *Validity of clinical studies*, while the validity of evidence synthesis studies is covered in EUnetHTA 21 Practical Guideline D4.3.1 *Direct and indirect comparisons* (along with EUnetHTA 21 Methodological Guideline D4.3.2 *Direct and indirect comparisons*). Additional considerations regarding the definition of clinically relevant outcomes and assessment of their validity, reliability and interpretability are covered in EUnetHTA 21 Practical Guideline D4.4.1 *Endpoints*.

130 This guideline predominantly deals with methodological issues related to inferential statistical analyses.
131 RCTs are the gold standard for answering clinical research questions in accordance with a hypothetico-
132 deductive approach. Therefore, while recommendations in this guideline may be better suited for RCTs,
133 they can apply to various study designs. For simplicity, effectiveness is the common term used here to
134 describe efficacy or effectiveness throughout the rest of the document. Effectiveness also includes
135 safety within the context of this document. Furthermore, treatment is used as a common term for any
136 health technology that can be assessed.

137 ## 1.2 Relevant articles in Regulation (EU) 2021/2282

138 Articles from Regulation (EU) 2021/2282 directly relevant to the content of this practical guideline are:

139    o   Article 8: initiation of joint clinical assessments,

140    o   Article 9: joint clinical assessment reports and the dossier of the health technology developer,

141    o   Article 13: member states' rights and obligations.

150 # 2 DEFINITIONS

143 In the context of this document, the terms "planned" and "prespecified" refer to a given statistical analysis
144 as planned according to a study protocol and/or statistical analysis plan (SAP) of a study submitted as
145 evidence by a HTD.

146 Mirroring the previous definition, the term "post hoc analysis" can be understood, unless stated
147 otherwise, as a synonym for any statistical analysis that was not planned according to a study protocol
148 and/or SAP of a study submitted as evidence by a HTD. More details are provided in Sections 9 and 10
149 of this document.

150 # 3 MULTIPLE STATISTICAL HYPOTHESIS TESTING IN INDIVIDUAL CLINICAL
151 STUDIES

152 ## 3.1 Purposes, definitions and general methodological considerations

153 In general, statistical analyses are performed for sample(s) of patients from a population of interest, as
154 collection of data for all patients in the population is usually not feasible. Thus, the statistics that are
155 produced are estimates and not true values for the population. Therefore, even if a clinical study is free
156 from bias, a difference observed for an **endpoint** of interest between groups (e.g., a difference in
157 mortality observed in an RCT) does not necessarily equate to a true difference in the population of
158 interest because of the sampling hazard (i.e., a form of random error). Thus, the risk of wrongly claiming
159 the existence of treatment effectiveness needs to be controlled at an acceptable level, which is achieved
160 via **statistical hypothesis testing**.

161 Under the **frequentist approach**, statistical hypothesis testing involves testing of two competing
162 hypotheses: the **null hypothesis** (usually denoted $H_0$) and the **alternative hypothesis** (usually denoted
163 $H_1$). In a superiority setting, the null hypothesis usually involves postulating the true absence of
164 difference in the population of interest, while the alternative hypothesis involves postulating a true
165 difference (two-sided tests) or true superiority or inferiority (one-sided tests). Statistical hypothesis
166 testing relies on estimating the **p value**, which is the probability of the occurrence of a difference at least
167 as large as the one observed if the null hypothesis is true. In RCTs, statistical test results are usually
168 interpreted under the **Neyman-Pearson approach**: the p value of a test is compared to a prespecified
169 risk level – the **α level** – and if the p value is less than the α level, the alternative hypothesis is accepted.
170 In biomedical research, the consensus is usually to set the α level of a (two-sided) single test at 0.05
171 (5%). Any statistical test can lead to two errors: rejecting the null hypothesis when it is actually true (i.e.,
172 the **type-1-error**, or false positive) or not rejecting the null hypothesis when it is actually false (i.e., the

173  **type 2 error**, or false negative) [2]. The probability α of a type 1 error for one significant test is the
174  **comparisonwise error rate** (CER) [3].

175  Situations in individual clinical studies occur for which multiple statistical tests are performed, which
176  increases the risk of at least one false-positive test. If k independent tests are performed, the probability
177  of rejecting at least one of the k independent null hypotheses when all null hypotheses are actually true
178  is called the global **familywise error rate** (FWER, considering the **family** of k tests as one experiment
179  under the complete null hypothesis). When considering independent tests, FWER is equal to $1-(1-α)^k$
180  [3][1]. When performing multiple tests, it is not usually expected that all null hypotheses will be true
181  simultaneously. Therefore, multiple tests procedures mentioned below usually control **FWER in a**
182  **strong sense** (also called multiple level): the probability of erroneously rejecting at least one true null
183  hypothesis, irrespective of which and how many of the individual null hypotheses are true [3]. A multiple
184  test procedure that controls FWER in a strong sense also controls the global FWER (but not vice versa)
185  [4]. Most of the usual procedures described below control FWER in a strong sense [4]. Thus, for the rest
186  of the document, we use the term FWER to mean "FWER in a strong sense".

187  Three main situations for which multiple statistical tests arise for individual clinical studies are considered
188  in this guideline. First, because most diseases have more than one consequence, many clinical studies
189  are designed to estimate the effectiveness of a treatment for more than one endpoint [5,6]. Situations in
190  which, as part of the final analysis of an individual clinical study, the same consequence is assessed at
191  different time points (e.g., clinical remission at 6 months and at 12 months after inclusion) or the same
192  endpoint is assessed in different populations (e.g., intention-to-treat population and a subpopulation)
193  are considered to be related. Indeed, these lead to multiplicity issues that can be considered similar
194  from a methodological perspective.

195  Second, multiplicity issues can arise in the context of **interim analysis**. An interim analysis is any
196  analysis used to compare treatment groups with respect to effectiveness at any time before formal
197  completion of a trial [1]. Thus, for a given endpoint, multiple analyses can be performed at different
198  times. Interim analyses can be planned for making decisions on whether to stop the trial early. Reasons
199  for early stopping may include clearly established superiority of the treatment(s) of interest(s),
200  confirmation that superiority is unlikely to occur and unacceptable adverse effects. Stopping could apply
201  to the entire trial or to a subset (e.g., ending a treatment group or discontinuing accrual of a subgroup
202  of patients). The benefits of interim analyses (ethical, scientific) can be opposed by methodological
203  disadvantages such as overestimation of the treatment effectiveness, lower precision or credibility, a
204  potential increase in type 1 errors if not appropriately managed and lower power for accounting for the
205  effects of prognostic factors [7].

206  Third, clinical trials can be conducted with more than two treatment groups (this situation is called
207  "**multiple groups**" in the rest of the document). Broadly speaking, multiple-group trials can be used to
208  compare different treatments between each other ("all pairwise comparisons" situation) or different
209  treatments to a reference treatment ("many to one" situation) [8].

210  As already mentioned, FWER increases with the number of independent tests. Hence, numerous valid
211  **multiplicity procedures**, such as the Bonferroni method, multiple-step procedures (e.g., the Holm
212  procedure), parametric multiple testing procedures (e.g., the Dunnett procedure), hierarchical test
213  sequences, α allocation, gatekeeping strategies and α spending functions, were developed to control
214  FWER at an acceptable level (e.g., at a global level of 5%) [9]. These procedures differ in the way in
215  which the algorithm for decision-making (i.e., rejecting or accepting the null hypotheses) is defined, the
216  purpose of the analyses (e.g., interim analyses, analysis of multiple endpoints) and the balance they
217  achieve between FWER control and loss of power (i.e., the probability of rejecting a false null
218  hypothesis).

---

[1] If multiple tests are dependent (e.g., they assess correlated endpoints using the same data), there is no easy way to compute
the theoretical FWER as it depends on the correlation structure between the different tests, but a high FWER can be expected
anyway if many tests are performed.

219    With the **Bayesian inference** approach, decision-making does not rely on proving whether a null
220    hypothesis is false and is instead guided by estimation of the distribution of treatment effect for the
221    endpoint(s) of interest(s). This distribution, called the posterior distribution, is estimated by combining
222    data from previous knowledge about the endpoint of interest (operationalised as the prior distribution)
223    and data obtained by conducting the clinical study (operationalised as the likelihood). Decision-making,
224    or stopping rules for interim analyses (i.e., claiming that a meaningful effect does or does not occur),
225    can then be made by estimating whether the posterior probability is higher or lower than a prespecified
226    threshold, which can have multiple definitions. For example, the posterior distribution of a risk ratio can
227    be used to estimate if its real value has high probability (e.g., 97.5%) of being less than 1, or less than
228    a value that is considered clinically meaningful [10]. Thus, the relevance of concepts such as the type 1
229    error for Bayesian inference is a matter of debate [11]. Nonetheless, as RCTs are designed to answer
230    a specific research question in a binary manner (i.e., concluding if there is a true effect or not), some
231    consider that controlling for the risk of false-positive conclusions still applies even if data are analysed
232    using Bayesian inference [12]. Thus, methods for adjusting the threshold for interpreting the results of a
233    Bayesian RCT while controlling for a desired FWER level when multiple hypotheses are tested have
234    been proposed [13].

235    ### 3.2    Requirements for appropriate reporting of methods and results in a JCA

236    #### 3.2.1    Multiple outcomes

237    Prospective specification of all data analyses that are performed to test hypotheses about the
238    prespecified endpoints, including the choice of multiplicity procedures, either before initiation of a clinical
239    study, or at least before database lock (i.e., planned analyses), is considered an essential element of
240    an adequate hypothetico-deductive approach. This helps in avoiding data dredging and ultimately helps
241    in controlling the type 1 error rate.

---

**Requirements for JCA reporting**

o    Accurate and unambiguous endpoint definitions (time points, measurement, method of assessment).

o    Null and alternative hypotheses that are tested.

o    How the endpoints were tested (statistical methods), including, if relevant, the multiplicity procedure that was used, the order of testing (if a hierarchical test procedure was used), the desired FWER level and which FWER was controlled (global level or multiple level).

o    The α level used to determine if the study was a success.

o    The CER level for each statistical test (i.e., the significance level required for each test).

o    The results, with appropriate statistics (position and dispersion indices in the intervention and comparator groups, appropriate effect measure, p value for the corresponding test and appropriate measures of statistical precision).

o    For the results for a given test, whether the test was appropriately controlled for multiplicity and if it was a planned analysis or not.

---

242    #### 3.2.2    Interim analyses

243    As already mentioned, the design of planned analyses, including interim analyses, is considered an
244    essential element of an adequate hypothetico-deductive approach.

245    Interim analyses are conducted with a definite cutoff date, usually defined as the occurrence of a specific
246    number of events of interest (e.g., number of deaths for overall survival analyses). Statistical analyses
247    of the corresponding database cannot be initiated before timely quality control and data management.
248    Data and corresponding interim analyses become available with a time lag between the cutoff date and
249    the clinical study report date (i.e., date of validation of the report of statistical analyses). However, the
250    appropriate date to report when assessing an interim analysis should be its appropriate cutoff date and

251  not the clinical study report date. Several interim analyses can be reported in the same clinical study
252  report.

253  Interim analyses can lead to early stopping of a clinical study. Results for interim analysis and the
254  decisions regarding clinical study continuation should be reported.

---

**Requirements for JCA reporting**

- o  Accurate and unambiguous endpoint definitions (time points, measurement, method of assessment).

- o  Schedule of all interim analyses.

- o  Respective cutoff date, with corresponding follow-up (the clinical study report date should not be used when reporting interim analyses in a JCA report).

- o  How the endpoints were tested (statistical methods), including the method chosen for controlling for multiplicity, for example, a triangular design, a group sequential design and its boundaries (e.g., Pocock, O'Brien-Fleming), α spending designs and their boundaries, and the desired FWER level.

- o  The α level used to determine if the study was a success.

- o  The CER level for each statistical test (i.e., the significance level required for each test).

- o  Implications (or not) and recommendations from an independent committee (e.g., data and safety monitoring board, data and safety monitoring committee).

- o  Any consequence of interim analyses for the conduct of the clinical trial (modification of study protocol, continuing or early stopping, no change, data release).

- o  The results, with appropriate statistics (position and dispersion indices in the intervention and comparator group, appropriate effect measure, p-value of the corresponding test and appropriate measures of statistical precision).

- o  For the results for a given test, if it was appropriately controlled for multiplicity and if it was a planned analysis or not.

- o  When unplanned interim analyses were conducted, why they were deemed necessary and by whom (sponsor or regulatory body).

---

255  **3.2.3   More than two treatment groups**

256  As already mentioned, planned analyses, including multiple-group comparison, are considered an
257  essential element of an adequate hypothetico-deductive approach.

258

---

**Requirements for JCA reporting**

- o Accurate and unambiguous endpoint definitions (time points, measurement, method of assessment).

- o Null and alternative hypotheses that are tested, with an unambiguous definition of the comparator.

- o How the endpoints were tested (statistical methods), including, if relevant, the multiplicity procedure that was used, the desired FWER level, and which FWER was controlled (global or multiple level).

- o The α level used to determine if the study was a success.

- o The CER level for each statistical test (i.e., the significance level required for each test).

- o The results, with appropriate statistics (position and dispersion indices in the intervention and comparator groups, appropriate effect measure, p value for the corresponding test and appropriate measures of statistical precision).

- o For the results for a given test, if it was appropriately controlled or not for multiplicity and if it was a planned analysis or not.

---

## 4  MULTIPLE STATISTICAL HYPOTHESIS TESTING IN EVIDENCE SYNTHESIS STUDIES

### *4.1  Purposes, definitions and general methodological considerations*

When conducting evidence synthesis studies such as pairwise meta-analysis (i.e., synthesis of direct evidence (multiple head-to-head RCTs) for when exactly two treatments are compared) or analysis of more complex evidence networks (see EUnetHTA 21 Methodological Guideline D4.3.2 *Direct and indirect comparisons*) via network meta-analysis (NMA), multiplicity issues can arise in a multiple of ways. These issues can be similar to those encountered when dealing with individual clinical studies or they can be specific to the design of an evidence synthesis study [14]. However, the possibilities and necessities to deal with multiplicity in evidence synthesis studies are limited because the data are already observed. Therefore, it is not possible to plan for multiplicity adjustments in a strong confirmatory sense.

### *4.2  Requirements for appropriate reporting of methods and results in a JCA*

#### 4.2.1  Multiple outcomes

As evidence synthesis can concern a wide range of outcomes, it is usually recommended to prespecify before data extraction the outcome(s) that will be of interest to collect information about these outcomes only [14,15]. However, because users of evidence synthesis studies have heterogeneous interests in the consequence of a medical condition, evidence synthesis studies are frequently performed with the inclusion of all outcomes that are likely to be of importance. Moreover, decisions on including outcomes encountered during the data extraction process are frequently made. Thus, an unambiguous definition of what constitutes a family of tests in the context of evidence synthesis studies is difficult to achieve. Therefore, while the procedures mentioned in Section 3.1 for dealing with multiple statistical hypothesis testing can theoretically be used (some require access to individual patient data, while those based on p values can be performed using aggregated data only), in practice this is almost never the case [14].

**Requirements for JCA reporting**

- o For evidence synthesis studies, control for multiplicity for dealing with multiple outcomes is a possibility but should not be expected.

- o Accurate and unambiguous endpoint definitions (time point, measurement, method of assessment).

- o Null and alternative hypotheses that were tested.

- o How the endpoints were tested (statistical methods).

- o If control for multiplicity was performed, how it was performed and the desired FWER level, and which FWER was controlled (global or multiple level).

- o The CER level for each statistical test.

- o The results (overall estimation at the evidence synthesis level) with appropriate statistics (position and dispersion indices in the intervention and comparator groups, appropriate effect measure, p value for the corresponding test and appropriate measures of statistical precision).

- o For the results for a given test, whether the outcome was prespecified before data extraction or not.

- o If control for multiplicity was performed, if it was appropriately conducted or not.

### 4.2.2    More than two treatment groups

Evidence synthesis comparing the relative effectiveness of more than two interventions using aggregated data are performed under the general framework for NMA. This framework can allow, if an appropriate network of evidence can be constructed, simultaneous estimation of the relative effectiveness of all pairwise comparisons of treatments included in the network, or at least for multiple two-by-two comparisons. Thus, an issue with multiple hypothesis testing may arise in NMA as several groups are observed in such a framework. Methodologically, how the potential issue of multiple hypothesis testing should be addressed when this type of evidence synthesis is performed is currently a matter of debate, and usually no multiplicity procedures are applied when estimating the relative effectiveness of the different treatments compared using an NMA [16]. Nonetheless, in the context of JCA, the assessment scope, expressed by the PICO question(s), defines the relevant comparator(s) for relative effectiveness assessment. Therefore, in a JCA, multiplicity due to multiple groups in evidence synthesis studies is not the main problem if the relevant comparison is the new intervention versus one control. When only one effect estimate of an NMA is of interest (and all the other effect estimates are only reported for completeness), multiplicity problems are therefore not an issue.

**Requirements for JCA reporting**

- o Accurate and unambiguous endpoint definitions (time point, measurement, method of assessment).

- o Null and alternative hypotheses that were tested.

- o How the endpoints were tested (statistical methods).

- o The CER level for each statistical test.

- o The results (overall estimation at the evidence synthesis level) with appropriate statistics (position and dispersion indices in the intervention and comparator groups, appropriate effect measure, p value for the corresponding test and appropriate measures of statistical precision).

### 4.2.3    Multiple time points

When performing an evidence synthesis study, it is possible that outcomes are measured at different time points depending on the time points defined or the follow-up duration in the original studies that are pooled. It is also possible to perform multiple evidence syntheses for the multiple time points available. Thus, multiplicity issues due to multiple time points for analysis can arise in evidence syntheses.

304 A solution to limit the issue of multiple time points can be to choose a single time point for the analysis.
305 However, this is only feasible if comparable time points are available from the studies included. Whether
306 time points are comparable strongly depends on the clinical indication and the treatment assessed in
307 the evidence synthesis. A different solution to the problem of different time points is to use a summary
308 effect measure over time, such as repeated-measures analysis of variance for continuous outcomes or
309 Cox regression for time-to-event data in the single studies [14]. The evidence synthesis can then be
310 performed by using the estimates of the summary effect measure (e.g., the hazard ratio), avoiding
311 multiplicity due to multiple time points. If individual patient data are available, methods for dealing with
312 multiple time points can be used directly in the evidence synthesis, such as NMA for survival data with
313 fractional polynomials to estimate, for example, the difference in restricted mean survival time for a
314 selected time point (see EUnetHTA 21 Practical Guideline D.4.3.1 *Direct and indirect comparisons* and
315 EUnetHTA 21 Methodological Guideline D.4.3.2 *Direct and indirect comparisons*).

---

**Requirements for JCA reporting**

- o Accurate and unambiguous endpoint definitions (time point, measurement, method of assessment).

- o If one common time point has been chosen for analysis for an endpoint, whether it was prespecified or not before data extraction and if choice of this common time point was justified (with its justification).

- O Null and alternative hypotheses that were tested.

- O How the endpoints were tested (statistical methods), especially methods that were used to estimate summary effect measures based on multiple time points.

- O The CER level for each statistical test.

- o The results (overall estimation at the evidence synthesis level) with appropriate statistics (position and dispersion indices in the intervention and comparator groups, appropriate effect measure, p value for the corresponding test and appropriate measures of statistical precision).

---

316 **4.2.4 Multiple operationalisations and multiple effect measures**

317 Multiplicity might arise in evidence synthesis because of the various ways available for analysing an
318 endpoint and operationalising an endpoint. In evidence synthesis studies, different methods for
319 analysing the same endpoint may be used to assess the robustness of a result, for example, via
320 sensitivity analyses (see Section 8). However, in such cases the effect measure or operationalisation
321 primarily used for the assessment should be prespecified [14].

322 For continuous data, a common effect measure is the mean difference. If the studies included in the
323 evidence synthesis measure the same endpoint using a different operationalisation (e.g., level of
324 depression using different depression scales), the effect measure can be standardised on a common
325 metric. If individual patient data are available, another option could be to dichotomise continuous data,
326 for example, as responders versus non-responders. However, the threshold for assignment as a
327 responder or non-responder must be scrutinised regarding whether it was prespecified before data
328 extraction and if it corresponds to validated and consensus cutoff values. For binary data, common effect
329 measures are the risk ratio, odds ratio and absolute risk difference. If analysis of data with different effect
330 measures leads to different results and the conclusion from the evidence synthesis would differ
331 depending on the effect measure used, multiplicity becomes problematic. Therefore, the effect measure
332 should be prespecified before data extraction and appropriate for the type of data being analysed [14].
333 Nevertheless, as member states may have different requirements regarding the effect measure for their
334 national assessment of the health technology, it cannot be ruled out that multiple effect measures for an
335 outcome need to be reported in the JCA.

336 Likewise, the analysis of several operationalisations for an endpoint (e.g., number or patients with event,
337 time-to-event, event rate) may lead to issues with multiple hypothesis testing. Similar to the situation for
338 multiple effect measures, member states may have different requirements regarding the
339 operationalisation for their national assessment and several operationalisations for an endpoint may
340 need to be reported in the JCA report.

---

**Requirements for JCA reporting**

- o Accurate and unambiguous endpoint definitions (time points, measurement, method of assessment).

- o If one operationalisation and/or effect measure was chosen for a specific endpoint, whether it was prespecified before data extraction or not and if it was justified (along with its justification).

- o If an endpoint that was continuous in the original individual studies has been dichotomised for analysis in an evidence synthesis study, whether the threshold for dichotomisation was planned before data extraction or not, along with the justification for choosing this threshold.

- o Null and alternative hypotheses that were tested.

- o How the endpoints were tested (statistical methods) and, if performed, the multiplicity procedure that was used and the desired FWER level.

- o The CER level for each statistical test.

- o The results (overall estimation at the evidence synthesis level), with appropriate statistics (position and dispersion indices in the intervention and comparator groups, appropriate effect measure, p value for the corresponding test and appropriate measures of statistical precision).

---

## 341  5   SUBGROUP ANALYSES IN INDIVIDUAL CLINICAL STUDIES

### 342  *5.1   Purposes, definitions and general methodological considerations*

343  Patients may respond differently to treatments because of demographic factors, disease characteristics,
344  comorbidities, environmental aspects, or characteristics related to other treatments, such as pre-
345  treatment or concomitant treatment. To examine whether an estimated overall effect in a single study is
346  driven by a specific patient group, subgroup analyses are conducted.

347  The term **subgroup** refers to a subset of the clinical trial population defined by one or more specific
348  patient characteristics measured at baseline. Postbaseline factors are not appropriate for defining
349  subgroups for the investigation of a treatment effect as they may be affected by the treatment itself
350  received by patients during the study. The term subgroup is not to be confounded with the term
351  **subpopulation**, which is defined as a subset of the patient population targeted as described in the
352  therapeutic indication. Subpopulations of interest may be specified during the assessment scope (see
353  EUnetHTA Practical Guideline D4.2.1 *Scoping process*) and are analysed as separate PICOs.

354  An **effect modifier** is a variable that modifies a treatment effect, that is, a variable that alters the relative
355  effectiveness between two treatments. Effect modifiers may be patient characteristics as listed above,
356  for example. Variables that represent methodological characteristics of a study are not regarded as
357  potential effect modifiers and therefore their potential impact on estimating treatment effectiveness
358  should be analysed in sensitivity analyses. Subgroup analyses refer to the comparison of treatment
359  effects in the (disjoint) subgroups of a potential effect modifier. In statistical terms, an evident effect
360  modification is referred to as an **interaction** between a treatment and the relevant variable.

361  A priori planning of subgroup analyses

362  Prespecification of subgroups is being encouraged in the planning of individual clinical studies as it can
363  lend credibility to positive or negative subgroup findings. However, a priori planned subgroup analyses
364  are often limited to the primary endpoint. From the perspective of assessment of an individual clinical
365  study, all other subgroup analyses, such as analyses of subgroups or subgroup analyses for further
366  endpoints not prespecified in the SAP, are unplanned analyses. These are not controlled for multiple
367  hypothesis testing and lack statistical robustness.

368  Nevertheless, member states may require further subgroup analyses than those planned at the single
369  study level for assessment at a national level (see EUnetHTA Practical Guideline D4.2.1 *Scoping*

370   *process*). Precise investigation of these subgroups depends on the use of results from unplanned
371   analyses at the single study level.

372   A particular point of attention is also the choice of cutoff value(s) for performing subgroup analyses when
373   the characteristic that defines the subgroup is initially a continuous variable. Indeed, to comply with an
374   adequate hypothetico-deductive approach, cutoff value(s) should be prespecified and the choice of the
375   value should be justified with an adequate rationale. In the case of subgroup analyses performed
376   because of the assessment scope, justification for the choice of cutoff value(s) pertains to the member
377   state(s) that require specific subgroup analyses.

378   <u>An interaction test is a requirement</u>

379   When interpreting subgroup analyses, it should be considered that a statistically significant effect in one
380   subgroup and a lack of effect or the reverse effect in another subgroup cannot be interpreted as the
381   existence of different treatment effects between subgroups on its own. Instead, demonstration of
382   different effects between different subgroups should be conducted using an appropriate interaction test
383   (e.g., adequate regression or analysis-of-variance model). Within an individual clinical study, interaction
384   can be tested on the basis of individual patient data. Different homogeneity and interaction tests have
385   been discussed in the literature [17–20]. For this guideline, the term "interaction test" refers to all of
386   these tests.

387   It should be kept in mind that owing to potential small sample sizes for subgroups, the power of
388   interaction tests for detecting heterogeneity can be low. Furthermore, in very small sample sizes,
389   prognostic variables (i.e., a patient characteristic that affects the outcome of interest irrespective of
390   which treatment is received) may be unbalanced within subgroups between treatment groups if
391   randomisation is not stratified according to the subgroup characteristic analysed [21,22]. In such cases,
392   the unbalanced prognostic variable is therefore a confounder (i.e., a variable that affects both the
393   treatment received and the outcome). Thus, the effect estimates within the subgroups may be biased
394   due to confounding, and this bias can lead to different results in the different subgroups. Therefore, in
395   the case of very small sample sizes, it cannot be ruled out that any differences detected between
396   subgroups are caused by systematic errors such as confounding.

397   If for one outcome there is a difference, for example, between two age groups as well as between men
398   and women, separate analyses would theoretically be required for each age group and for men and
399   women (i.e., analyses of four subgroups) to interpret the results. However, such analyses are rarely
400   available and may result in subgroups with rather small sample sizes.

401   ## *5.2   Requirements for appropriate reporting of methods and results in a JCA*

402   In the JCA report, information regarding a priori planning of subgroup analyses, consideration of
403   multiplicity and definitions of subgroups in the protocol and SAP of the clinical studies assessed must
404   be provided.

405

---

**Requirements for JCA reporting**

- o Accurate and unambiguous definitions of subgroups and endpoints.

- o Null and alternative hypotheses that are tested.

- o How the subgroup analysis was performed (statistical methods), including the multiplicity procedure that was used, if performed, and the desired FWER level.

- o When cutoff value(s) for a continuous variable were chosen for defining subgroups, whether they were prespecified and how the choice of these values was justified.

- o The CER level for each subgroup analysis.

- o The results (p values) of an appropriate interaction test for all subgroup analyses conducted.

- o The results (appropriate estimates and effect measures for each subgroup, with a corresponding measure of statistical precision and p values for the effect in each subgroup).

- o Whether each statistical test for subgroup analysis was appropriately controlled for multiplicity or not, and if it was a planned analysis or not.

- o Visual presentation of the results using a forest plot is strongly encouraged.

---

# 6 SUBGROUP ANALYSES IN EVIDENCE SYNTHESIS STUDIES

## 6.1 Purposes, definitions and general methodological considerations

The purpose of subgroup analyses and the definition of an effect modification (an interaction) described for subgroup analyses in individual clinical studies also apply to evidence synthesis studies.

The term subgroup should refer to a subset of the patient population included in the evidence synthesis defined by one or more specific patient characteristics measured at baseline. In an evidence synthesis study, the individual studies included may represent subgroups if they included patients with specific characteristics of the respective subgroup.

A priori planning of subgroup analyses

In an evidence synthesis, subgroup analyses that were planned in each of the studies included are mostly not available. Nonetheless, subgroup analyses can be planned within the study protocol and/or SAP of a specific evidence synthesis study. In addition, within the assessment scope, subgroups may be defined together with the PICO framework.

An interaction test is a requirement

Within an evidence synthesis, the results from several studies can be summarised via meta-analyses. To investigate whether an estimated overall effect in meta-analysis is driven by a specific patient group, common tests for heterogeneity (in this case, heterogeneity between subgroups rather than studies) or meta-regression may be considered. In the case of evidence synthesis performed using individual patient data, for example, an adequate regression or analysis-of-variance model with a corresponding interaction term can be used. When only aggregated data are available, a Q test in a meta-analysis and an F test in a meta-regression are examples of appropriate tests for interaction. As for subgroup analyses in single studies, statistical tests for interaction may have low power and may not be sufficient to exclude the possibility of meaningful subgroup interactions.

In a meta-regression, the statistical association between the effect sizes in individual studies and the study characteristics is investigated, so that study characteristics can possibly be identified that explain the different effect sizes, that is, the heterogeneity. However, it is important that the limitations of such analyses are considered when interpreting any results. Meta-regressions that attempt to show an association between the different effect sizes and the average patient characteristics in individual studies are subject to the same limitations as the results from ecological studies in epidemiology [23].

435 The high risk of bias in such analyses based on aggregated data cannot be balanced by adjustment. An
436 alternative approach is therefore the use of individual patient data, as meta-analyses that include
437 individual patient data generally provide greater certainty of results, that is, more precise results not
438 affected by ecological bias [24,25]. If heterogeneity is plausible, it can threaten the certainty of results
439 associated with an evidence synthesis study. More information on this issue can be found in EUnetHTA
440 21 Practical Guideline D4.3.1 *Direct and indirect comparisons*.

441 ## *6.2 Requirements for appropriate reporting of methods and results in a JCA*

442 The requirements for appropriate reporting of methods and results from evidence syntheses are similar
443 to those for single studies. In general, a priori planned subgroup analyses should not be replaced by
444 unplanned analyses. All analyses should be reported.

---

**Requirements for JCA reporting**

- o Accurate and unambiguous definitions of the subgroups and endpoints are required, as well as the null and alternative hypotheses that were tested.

- o How the subgroup analysis was performed (statistical methods), including the multiplicity procedure that was used, if performed, and the desired FWER level.

- o When cutoff value(s) for a continuous variable were chosen for defining subgroups, whether they were prespecified and how the choice of these values is justified.

- o The CER level against which each subgroup analysis was performed.

- o The results (p values) of an appropriate interaction test for all subgroup analyses conducted.

- o The results (appropriate estimates and effect measures in each subgroup, with a corresponding measure of statistical precision and p values for the effect in each subgroup).

- o Whether each statistical test for subgroup analysis was appropriately controlled for multiplicity or not, and if it was a planned analysis or not.

- o Visual presentation of the results using a forest plot is strongly encouraged.

---

445 # 7   SENSITIVITY ANALYSES IN INDIVIDUAL STUDIES

446 ## *7.1 Purposes, definitions, and general methodological considerations*

447 **Sensitivity analyses** are an integral part of the reporting of clinical study results and are essential in
448 investigating the robustness of the effect observed in the clinical study to variations in the assumptions
449 and their impact.

450 In any clinical trial, the primary **estimand** should be defined according to the principles outlined in ICH
451 E9 and its addendum (E9(R1)) [1,26]. The aim of the estimand framework is to define "*a precise
452 description of the treatment effect reflecting the clinical question posed by the trial objective*". It
453 summarises at a population level what the outcomes would be in the same patients under different
454 treatment conditions being compared. The statistical analyses should be aligned to the estimand (not
455 vice versa) and sensitivity analyses should be planned in the study protocol to "*explore the robustness
456 of inference from the main estimators to deviations from its underlying modelling assumptions and
457 limitations of the data*" [26].

458 The estimand is defined by its five attributes: (1) population, (2) treatment, (3) variable (endpoint), (4)
459 **intercurrent events** (ICEs) and (5) the summary measure. ICEs (events occurring after treatment
460 initiation that affect either the interpretation or the existence of the measurements associated with the
461 clinical question of interest) should be addressed when describing the clinical question of interest to
462 precisely define the treatment effect that is to be estimated. ICEs are context-dependent; the same event
463 can be defined as **missing data** in one setting and as an ICE in another. For examples see the ICH

464 E9(R1). Sensitivity analyses (i.e., a series of analyses conducted with the intent of exploring the
465 robustness of inferences from the main estimator to deviations from its underlying modelling
466 assumptions and limitations in the data) should be used to explore the impact that changes to the
467 assumptions for any or all of these elements might have on the primary outcome of a study [26]. In
468 addition, the ICH E9 Addendum differentiates between sensitivity analyses and supplementary analyses
469 (i.e., a general description for analyses that are conducted in addition to the main and sensitivity analysis
470 with the intent of providing additional insights into understanding the treatment effect) [26]. The JCA
471 should indicate clearly which analyses are primary, sensitivity or supplementary analyses.

472 Focus should be given to the difference between missing data (i.e., data that would be meaningful for
473 analysis of a given estimand but were not collected) and ICEs and their handling in analyses. Indeed,
474 missing data should be distinguished from data that do not exist or data that are not considered
475 meaningful because of an ICE. Guidelines on handling of missing data are available [26,27] and
476 describe appropriate sensitivity analysis strategies (see the definitions above). Handling of missing data
477 (e.g., missing laboratory assessments) is considered a statistical problem that needs to be addressed
478 via appropriate statistical analyses with the aim of explore the impact of the level of missing data on the
479 basis of certain assumptions [27]. Avoiding missing data is considered of utmost importance, although
480 some degree of missing data should be anticipated in any clinical study. Therefore, appropriate handling
481 of this issue should be predefined. The methodologies chosen should be reported in the JCA.

482 Continued collection of data even after ICEs (e.g., treatment discontinuation or initiation of a rescue
483 medication) to support assessment of their impact on the clinical questions is highly supported and
484 different strategies to do so are described in the ICH E9 addendum [26]. Again, the JCA should report
485 on the strategy for the primary and any additional estimands and the strategies chosen for handling of
486 ICEs.


487 ### 7.2    Requirements for appropriate reporting of methods and results in a JCA

488 The study protocol and SAP should always be submitted to allow assessment of the estimand strategy.
489 Results should be presented according to the prespecified analyses based on the estimand framework
490 in the study protocol as well as the strategies for handling missing data and accompanying analyses,
491 and this should be reflected in the JCA.

492 There is no rule for the amount of missing data that is considered acceptable. Therefore, reports should
493 highlight the uncertainty with respect to the amount as well as the handling of missing data. The
494 acceptability of missing data is subject to member state differences in interpretation of their relevance
495 within their respective decision-making process.

496 The primary estimand describes the objective of the study via definition of the five attributes. The
497 objectives of studies might therefore be aligned with a certain PICO or not. If a secondary estimand
498 better aligns with other relevant PICO(s), this should be highlighted in the report and the JCA should be
499 clear in distinguishing the different estimands. Because estimands describe the treatment in the context
500 of the attributes, it is possible that different HTAs could also prefer different estimands. This situation
501 might be rare but is addressed in the PICO scoping process and should then be reflected in the report.
502 If secondary estimands have less statistical rigour (because they are based on outcomes not included
503 in the inferential testing strategy), this should be clearly highlighted in the report.

504 The acceptability of sensitivity analyses is subject to member state differences in interpretation of their
505 relevance within their respective decision-making process.

506

---

**Requirements for JCA reporting**

- o The JCA should present the relevance of the chosen estimand with respect to the original trial protocol as well as relevant PICO(s).

- o There should be a detailed description of the chosen estimand(s), with a focus on the five attributes as well as the ICE strategy.

- o Strategies for handling of ICEs are distinct from strategies to handle missing data and these differences should be clearly conveyed.

- o Sensitivity and supplementary analyses should be distinguished from the primary estimand(s) and its analyses.

- o There should be a clear definition of the purpose and the underlying assumption for each sensitivity analysis.

- o All sensitivity analyses should be presented in the report, preferably as a summary table. Such table(s) should include the attribute(s) that the sensitivity analyses address as well as the analysis method used for each individual analysis and the results.

- o When the results of a sensitivity analysis are not of the same directionality as for the results of the primary analysis, this should be highlighted.

---

## 8   SENSITIVITY ANALYSES IN EVIDENCE SYNTHESIS STUDIES

### *8.1   Purposes, definitions, and general methodological considerations*

Evidence synthesis results are sensitive to the inclusion/exclusion of individual studies and sensitivity analyses can help to explore the impact of individual studies on the overall conclusions, assess the robustness of the analyses in general and confirm assumptions underlying the evidence synthesis.

Sensitivity analyses are a set of analyses estimating the same effect but with different methodology to assess the impact of different decisions compared to the primary assumptions on the analysis. These alternative decisions can pertain to the inclusion of studies (size, population and outcomes, among others), certain groups of the patient population (in range/out of range), the risk of bias or the use of fixed-effect versus random-effect models.

### *8.2   Requirements for appropriate reporting of methods and results in a JCA*

**Requirements for JCA reporting**

- o It should be stated whether the analysis was prespecified in the study protocol and/or SAP, was identified during the assessment process or is the result of the PICO process.

- o There should be a clear definition of the purpose and underlying assumption for each sensitivity analysis.

- o All sensitivity analyses should be presented in the report, preferably as a summary table. Such table(s) should include the elements the sensitivity analyses address, such as the evidence included, the eligibility criteria, the data used with the underlying assumptions and the analysis method used for each individual analysis, and the results. Sensitivity analyses can be conducted not only for single factors but also for multifactorial situations, so the report should be clear on what type of analysis has been performed.

- o When the results of a sensitivity analysis are not of the same directionality as for the results of the primary analysis, this should be highlighted.

518 **9   POST HOC ANALYSES IN INDIVIDUAL CLINICAL STUDIES**

519 *9.1   Purposes, definitions, and general methodological considerations*

520 The term post hoc is derived from the Latin phrase *post hoc ergo propter hoc*, meaning "after this,
521 therefore because of this". Thus, in the strictest sense, post hoc analyses are all analyses that are
522 performed because of the results of a previous analysis. Therefore, there can exist post hoc analyses
523 that can be planned, such as a statistical hypothesis test performed for a particular outcome because of
524 the results of the previous test when hierarchical test sequence procedures for controlling multiplicity
525 issues are used. However, as mentioned earlier, the scope of this document mainly addresses
526 unplanned post hoc analyses, as these are the ones that can be considered problematic in terms of
527 deviation from an adequate hypothetico-deductive approach. Indeed, while planned analyses are
528 acceptable when appropriate measures are taken regarding emerging multiplicity issues, unplanned
529 post hoc analyses violate the principles of inferential hypothesis testing. Both the power of a study and
530 the certainty for correctly rejecting the null hypothesis are built on the principle of defining the parameters
531 of the hypothesis to be tested before the real data are observed.

532 However, during a HTA it might be desirable to obtain data for a patient subset that, for example, reflects
533 a PICO more closely than the strategy pursued by the applicant. In principle, post hoc analyses can
534 address all elements of the trial and not just subgroups of the population, as well as different outcome
535 measures or statistical methods.

536 In such situations an explorative investigation based on post hoc–defined subgroups might be
537 considered, reflective of the known methodological caveats. Post hoc analyses should be clearly
538 identified as such to distinguish them from the primary analyses in the JCA.

539 *9.2   Requirements for appropriate reporting of methods and results in a JCA*

540 Reporting of post hoc analyses should follow the principles outlined in the European Medicines Agency
541 guideline on the investigation of subgroups in confirmatory trials [28]. Subgroup analyses need to reflect
542 on the heterogeneity of the overall population versus any subgroups, the consistency of results across
543 the subgroups and the credibility of any subgroup, which is directly linked to the biological plausibility
544 and support for the findings from external sources.

545 If analyses derived from unplanned post hoc assessment of data are presented, they should preferably
546 be reported using descriptive statistics with clear identification that they have not been generated within
547 the inferential framework of the trial (p values must be clearly marked as nominal, i.e., as unplanned
548 analyses and not controlled for multiplicity).

549 HTDs have to provide all information available on the characteristics of the subgroups, substantiate any
550 claims regarding balance in terms of randomisation, provide evidence that no interactions with other
551 prognostic or predictive factors might be the underlying cause of any differences observed and provide
552 a strong biological rationale if a specific subgroup performs better or worse than the overall trial
553 population.

---

**Requirements for JCA reporting**

○ Planned post hoc analyses, such as a procedure for control of multiplicity issues, are to be
   reported according to the requirements described in the corresponding sections of the
   document (e.g., Section 3 if these analyses deal with controlling for multiplicity).

○ Unplanned post hoc analyses, such as those requested by a HTA as a consequence of the
   PICO process, should be clearly flagged as unplanned.

---

554  # 10  POST HOC ANALYSES IN EVIDENCE SYNTHESIS STUDIES

555  ## 10.1  Purposes, definitions and general methodological considerations

556  Evidence generation should follow a planned protocol to reduce the likelihood of drawing biased
557  conclusions. Full prespecification is difficult and often not possible for systematic reviews because
558  knowledge is already available for the underlying studies. Therefore, if an important aspect was not
559  addressed in the planning stage (PICO scoping) but proves to be of importance for the assessment,
560  additional post hoc analyses might be required.

561  ## 10.2  Requirements for appropriate reporting of methods and results in a JCA

562  The JCA should report post hoc analyses but highlight them to distinguish them from other planned
563  analyses.

---

**Requirements for JCA reporting**

- o  A report of the protocol-defined analyses and their relevance to the PICO.

- o  The report should clearly distinguish between planned analyses and unplanned post hoc analyses.

---

564  # 11  RELATED EUNETHTA DOCUMENTS

565  - **EUnetHTA 21 Practical Guideline D4.3.1: Direct and indirect comparisons**

566  A practical guideline for assessors and co-assessors that describes possible approaches and specific
567  instructions for reporting and assessing the evidence from evidence synthesis studies (pairwise meta-
568  analyses and indirect comparisons).

569  - **EUnetHTA 21 Methodological Guideline D4.3.2: Direct and indirect comparisons**

570  A methodological guideline for assessors and co-assessors that provides an understanding of the basic
571  methodological principles behind conducting evidence synthesis studies (pairwise meta-analyses and
572  indirect comparisons).

573  - **EUnetHTA 21 Practical Guideline D4.2.1: Scoping process**

574  A practical guideline for assessors and co-assessors that describes the methods and principal steps for
575  the scoping process.

576  - **EUnetHTA 21 Practical Guideline D4.4.1: Endpoints**

577  A practical guideline for assessors and co-assessors that describes how to deal with several issues
578  encountered around the assessment of endpoints, and guides member states on defining relevant
579  endpoints during the scoping process.

580  - **EUnetHTA 21 Practical Guideline D4.6.1: Validity of clinical studies**

581  A practical guideline for assessors and co-assessors that defines the main aspects of the validity of
582  individual clinical studies, defines and classifies the different types of clinical studies that can be
583  conducted, and describes how to report and assess the validity of individual clinical studies whether
584  they are RCTs or not.

## 12 REFERENCES

1.  International Conference on Harmonisation of Technical Requirements for Registration of Pharmaceuticals for Human Use. *ICH Harmonised Tripartite Guideline. Statistical Principles for Clinical Trials E9*. Geneva: ICH; 1998.

2.  Neyman J. "Inductive behavior" as a basic concept of philosophy of science. *Rev Int Stat Inst* 1957;25(1/3):7.

3.  Bender R, Lange S. Adjusting for multiple testing—when and how? *J Clin Epidemiol* 2001;54(4):343–9.

4.  Bauer P. Multiple testing in clinical trials. *Stat Med* 1991;10(6):871–90.

5.  Pocock SJ. Clinical trials with multiple outcomes: a statistical perspective on their design, analysis, and interpretation. *Control Clin Trials* 1997;18(6):530–45.

6.  Zhang J, Quan H, Ng J, et al. Some statistical methods for multiple endpoints in clinical trials. *Control Clin Trials* 1997;18(3):204–21.

7.  Armitage P, McPherson CK, Rowe BC. Repeated significance tests on accumulating data. *J R Stat Soc Ser Gen* 1969;132(2):235–44.

8.  Bretz F, Hothorn T, Westfall P. *Multiple Comparisons Using R*. New York, NY: Chapman and Hall/CRC; 2016.

9.  Dmitrienko A, Tamhane AC, Bretz F, editors. *Multiple Testing Problems in Pharmaceutical Statistics*. Boca Raton, FL: Chapman & Hall/CRC; 2010.

10. Lesaffre E, Lawson A. *Bayesian Biostatistics*. Chichester: Wiley; 2012.

11. Ryan EG, Brock K, Gates S, et al. Do we need to adjust for interim analyses in a Bayesian adaptive trial design? *BMC Med Res Methodol* 2020;20(1):150.

12. Kapur J, Elm J, Chamberlain JM, et al. Randomized trial of three anticonvulsant medications for status epilepticus. *N Engl J Med* 2019;381(22):2103–13.

13. Guo M, Heitjan DF. Multiplicity-calibrated Bayesian hypothesis tests. *Biostatistics* 2010;11(3):473–83.

14. Bender R, Bunce C, Clarke M, Gates S, Lange S, Pace NL, et al. Attention should be given to multiplicity issues in systematic reviews. J Clin Epidemiol. 2008;61(9):857–65.

15. Li T, Higgins JPT, Deeks JJ. Collecting data. In: Higgins JPT, Thomas J, Chandler J, et al., editors. *Cochrane Handbook for Systematic Reviews of Interventions version 6.3*. London: Cochrane Collaboration; 2022. Chapter 5.

16. Efthimiou O, White IR. The dark side of the force: multiplicity issues in network meta-analysis and how to address them. *Res Synth Methods* 2020;11(1):105–22.

17. Christensen R, Bours MJL, Nielsen SM. Effect modifiers and statistical tests for interaction in randomized trials. *J Clin Epidemiol* 2021;134:174–7.

18. Tanniou J, van der Tweel I, Teerenstra S, et al. Subgroup analyses in confirmatory clinical trials: time to be specific about their purposes. *BMC Med Res Methodol* 2016;16:20.

19. Dmitrienko A, Muysers C, Fritsch A, et al. General guidance on exploratory and confirmatory subgroup analysis in late-stage clinical trials. *J Biopharm Stat* 2016;26(1):71–98.

624   20.   Alosh M, Huque MF, Bretz F, et al. Tutorial on statistical considerations on subgroup analysis in
625         confirmatory clinical trials. *Stat Med* 2017;36(8):1334–60.

626   21.   Cui L, Hung HM, Wang SJ, et al. Issues related to subgroup analysis in clinical trials. *J Biopharm*
627         *Stat* 2002;12(3):347–58.

628   22.   Sun X, Briel M, Walter SD, Guyatt GH. Is a subgroup effect believable? Updating criteria to
629         evaluate the credibility of subgroup analyses. *BMJ* 2010;340:c117.

630   23.   Greenland S, Morgenstern H. Ecological bias, confounding, and effect modification. *Int J*
631         *Epidemiol* 1989;18(1):269–74.

632   24.   Simmons LA. Self-perceived burden in cancer patients: validation of the Self-perceived Burden
633         Scale. *Cancer Nurs* 2007;30(5):405–11.

634   25.   Berlin JA, Santanna J, Schmid CH, et al, Anti-Lymphocyte Antibody Induction Therapy Study G.
635         Individual patient- versus group-level data meta-regressions for the investigation of treatment
636         effect modifiers: ecological bias rears its ugly head. *Stat Med* 2002;21(3):371–87.

637   26.   International Conference on Harmonisation of technical requirements for registration of
638         pharmaceuticals for human use. *Addendum on Estimands and Sensitivity Analysis in Clinical*
639         *Trials to the Guideline on Statistical Principles for Clinical Trials*. E9(R1). Geneva: ICH; 2019.

640   27.   European Medicines Agency. *Guideline on Missing Data in Confirmatory Clinical Trials*. London:
641         EMA; 2010.

642   28.   European Medicines Agency. *Guideline on the Investigation of Subgroups in Confirmatory*
643         *Clinical Trials*. London: EMA; 2019.