# Statistical Approaches to Establishing Bioequivalence

## Guidance for Industry

*DRAFT GUIDANCE*

**This guidance document is being distributed for comment purposes only.**

Comments and suggestions regarding this draft document should be submitted within 60 days of publication in the *Federal Register* of the notice announcing the availability of the draft guidance.  Submit electronic comments to https://www.regulations.gov.  Submit written comments to the Dockets Management Staff (HFA-305), Food and Drug Administration, 5630 Fishers Lane, Rm. 1061, Rockville, MD 20852.  All comments should be identified with the docket number listed in the notice of availability that publishes in the *Federal Register*.

For questions regarding this draft document, contact (CDER) David Coppersmith at 301-796-9193.

# Statistical Approaches to Establishing Bioequivalence
## Guidance for Industry

**U.S. Department of Health and Human Services**
**Food and Drug Administration**
**Center for Drug Evaluation and Research (CDER)**

**December 2022**
**Biopharmaceutics**

**Revision 1**

# TABLE OF CONTENTS

# Statistical Approaches to Establishing Bioequivalence Guidance for Industry[1]

This draft guidance, when finalized, will represent the current thinking of the Food and Drug Administration (FDA or Agency) on this topic.  It does not establish any rights for any person and is not binding on FDA or the public.  You can use an alternative approach if it satisfies the requirements of the applicable statutes and regulations.  To discuss an alternative approach, contact the FDA staff responsible for this guidance as listed on the title page.

## I.      INTRODUCTION

Requirements for submitting bioavailability (BA) and bioequivalence (BE) data in investigational new drugs (INDs), new drug applications (NDAs), abbreviated new drug applications (ANDAs), and supplements; the definitions of BA and BE; and the types of in vitro and in vivo studies that are appropriate to measure BA and establish BE are set forth in part 320 (21 CFR part 320).  This guidance provides recommendations on how to meet provisions of part 320 for all drug products.

In general, FDA's guidance documents do not establish legally enforceable responsibilities.  Instead, guidances describe the Agency's current thinking on a topic and should be viewed only as recommendations, unless specific regulatory or statutory requirements are cited.  The use of the word *should* in Agency guidances means that something is suggested or recommended, but not required.

### A.      Overview

This guidance provides recommendations to sponsors and applicants who intend to use equivalence criteria in analyzing in vivo or in vitro BE studies for INDs, NDAs, ANDAs, and supplements to these applications.  This guidance discusses statistical approaches for BE comparisons and focuses on how to use these approaches both generally and in specific situations.  When finalized, this guidance will replace the guidance for industry *Statistical Approaches to Establishing Bioequivalence*, which was issued in February 2001 (2001 guidance).  This guidance provides recommendations on the topics covered in the 2001 guidance as well as recommendations on additional topics, including missing data and intercurrent events, adaptive design, and specific situations, such as narrow therapeutic index drugs and highly variable drugs.

---

[1] This guidance has been prepared by the Office of Generic Drugs in the Center for Drug Evaluation and Research (CDER) in cooperation with CDER's Office of Translational Sciences and Office of Pharmaceutical Quality at the Food and Drug Administration.

41  Defined as *relative BA*, the assessment of BE involves comparison between a test (T) and
42  reference (R) drug product, where T and R can vary depending on the comparison to be
43  performed (e.g., to-be-marketed formulation versus clinical trial formulation, generic drug versus
44  reference listed drug (RLD), originally approved formulation versus postapproval formulation
45  changes).  Although BA and BE are closely related, BE comparisons normally rely on (1) a
46  criterion, (2) a confidence interval for the criterion, and (3) a predetermined BE limit.  BE
47  comparisons could also be used in certain pharmaceutical product line extensions, such as
48  additional strengths, new dosage forms (e.g., changes from immediate release to extended
49  release), and new routes of administration.[2]  In these contexts, the approaches described in this
50  guidance can be used to determine BE.  The general approaches discussed in this guidance may
51  also be useful when assessing pharmaceutical equivalence (i.e., the identical dosage form and
52  route(s) of administration that contain identical amounts of the identical active drug ingredient)
53  or performing equivalence comparisons in clinical pharmacology studies and other areas.
54
55  This guidance is intended to encourage the use of science-based approaches to making statistical
56  BE assessments.  Given the evolving nature of statistical approaches and technologies, FDA
57  encourages generic and new drug applicants to propose and discuss novel methodologies (e.g.,
58  model-based BE and novel adaptive designs for comparative clinical endpoint BE studies) with
59  the Agency through appropriate regulatory meetings, as described below.
60
61      **B.      Statistical Guidance Background**
62
63  In the July 1992 guidance on *Statistical Procedures for Bioequivalence Studies Using a Standard*
64  *Two-Treatment Crossover Design* (the 1992 guidance), the Center for Drug Evaluation and
65  Research (CDER) recommended that a standard in vivo BE study design be based on the
66  administration of either single or multiple doses of the T and R products to healthy subjects on
67  separate occasions, with random assignment to the two possible sequences of drug product
68  administration.  The 1992 guidance further recommended that statistical analysis for
69  pharmacokinetic (PK) measures, such as area under the curve (AUC) and peak concentration
70  ($C_{max}$), be based on the *two one-sided tests procedure* to determine whether the average values
71  for the PK measures determined after administration of the T and R products were comparable.
72  This approach is termed *average BE* (ABE) and involves the calculation of a 90% confidence
73  interval for the ratio of the averages (population geometric means) of the measures for the T and
74  R products.  To establish BE, the calculated confidence interval should fall within a BE limit,
75  usually 80 to 125% for the ratio of the product averages.[3]  In addition to this general approach,
76  the 1992 guidance provided specific recommendations for (1) logarithmic transformation of PK
77  data, (2) methods to evaluate sequence effects, and (3) methods to evaluate outlier data.

---

[2] For example, to submit an ANDA that is not the same as its RLD because it has a different strength, dosage form, or route of administration than that of the RLD, an applicant first must obtain permission from FDA through the citizen petition process.  See section 505(j)(2)(C) of the Federal Food, Drug and Cosmetic Act (21 U.S.C. 355(j)(2)(C)); 21 CFR 314.93(b).  Such petitions are referred to as suitability petitions.

[3] For a broad range of drugs, a BE limit of 80 to 125% for the ratio of the product averages has been adopted for use of an average BE criterion.  Generally, the BE limit of 80 to 125% is based on a clinical judgment that a test product with BA measures outside this range should be denied market access.

78
79 In addition to reiterating the key points from the 1992 guidance and replacing that guidance, the
80 2001 guidance introduced two additional approaches to assessing BE: *population BE* and
81 *individual BE*. Both of these approaches, unlike the *average BE* approach, include a comparison
82 of the variabilities of the PK metrics of the two products being compared, as well as the average
83 responses. However, the individual BE approach is not currently used in the regulatory setting
84 while the population BE approach is mainly used for certain in vitro BE studies. The 2001
85 guidance also includes discussion of *replicated crossover designs* — crossover designs in which
86 at least some of the subjects receive at least one of the products more than once. The discussion
87 of these designs in that guidance included their implications for possible carryover effects and
88 their use in screening for outliers.
89
90 This guidance provides recommendations on the topics covered by the 1992 guidance and the
91 2001 guidance, as well as recommendations on some additional topics. As noted in the
92 Overview section above, when finalized, this guidance will replace the 2001 guidance.
93
94
95 **II.    GENERAL CONSIDERATIONS**
96
97     **A.    Study Design**
98
99        *1.    Experimental Design*
100
101            a.    Nonreplicated designs
102
103 A conventional nonreplicated design, such as the standard two-formulation, two-period, two-
104 sequence crossover design, can be used to generate data when an average or population approach
105 is chosen for BE comparisons. Under certain circumstances, such as products with apparent,
106 long half-lives where crossover studies are impractical, parallel designs can be used.
107
108            b.    Replicated crossover designs
109
110 Replicated crossover designs can be used irrespective of which BE approach is selected to
111 establish BE, although they are not necessary when an average or population BE approach is
112 used. When a reference-scaled BE approach is used, replicated crossover designs are critical to
113 allow estimation of within-subject variances for the R (and T if a fully replicated study is used)
114 measures. In particular, the following four-period, two-sequence, two-formulation design is
115 recommended for fully replicated BE studies (see Appendix A for further discussion of
116 replicated crossover designs).
117

|  |  | **Period** | | | |
|---|---|---|---|---|---|
|  |  | *1* | *2* | *3* | *4* |
| *Sequence* | *1* | T | R | T | R |
|  | *2* | R | T | R | T |

118
119
120    For this design, the same lots of the T and R formulations should be used for the replicated
121    administration.  Each period should be separated by an adequate washout period.
122
123    Other fully replicated crossover designs are also possible.  For example, a three-period design, as
124    shown below, could be used.  A fully replicated design can estimate the subject-by-formulation
125    interaction variance components.
126

|          |   | Period |   |   |
|----------|---|--------|---|---|
|          |   | 1      | 2 | 3 |
| Sequence | 1 | T      | R | T |
|          | 2 | R      | T | R |

127
128    The following three-period, three-sequence, two-formulation, partially replicated design can also
129    be used for assessing reference-scaled BE, though it cannot fully estimate the subject-by-
130    formulation interaction variance component (as a fully replicated design can).
131

|          |   | Period |   |   |
|----------|---|--------|---|---|
|          |   | 1      | 2 | 3 |
| Sequence | 1 | T      | R | R |
|          | 2 | R      | T | R |
|          | 3 | R      | R | T |

132    A greater number of subjects would be needed for the three-period designs compared to the
133    recommended four-period design to achieve the same statistical power to conclude BE.
134
135        c.    Adaptive design
136
137    An adaptive design is a clinical trial design that allows for prospectively planned modifications
138    to one or more aspects of the design based on accumulating data from subjects in the trial.  An
139    adaptive design can be a group sequential design, or other design with one or more adaptive
140    features.[4]  For example, Potvin's methods (Potvin et al. 2008, Xu et al. 2016)[5] are a combination
141    of a group sequential design and an adaptive design with sample size re-estimation.
142

---

[4] See the guidance for industry *Adaptive Designs for Clinical Trials of Drugs and Biologics* (November 2019).  We update guidances periodically. For the most recent version of a guidance, check the FDA guidance web page at https://www.fda.gov/regulatory-information/search-fda-guidance-documents.

[5] Potvin, D., C.E. DiLiberti, W.W. Hauck, A.F. Parr, D.J. Schuirmann, and R.A. Smith, 2008, Sequential Design Approaches for Bioequivalence Studies With Crossover Designs, Pharmaceutical Statistics: The Journal of Applied Statistics in the Pharmaceutical Industry 7, no. 4: 245-262; Xu, J., C. Audet, C.E. DiLiberti, W.W. Hauck, T.H. Montague, A.F. Parr, D. Potvin, and D.J. Schuirmann, 2016, Optimal Adaptive Sequential Designs for Crossover Bioequivalence Studies, Pharmaceutical Statistics (15) 1:15-27.

143 Adaptive design can provide ethical advantages[6] and statistical efficiency. When appropriately
144 implemented, adaptive designs can reduce resources used, decrease time to study completion,
145 and increase the chance of study success, especially when the prior information needed for the
146 study design is limited. However, use of adaptive designs can also have limitations. For
147 example, adaptive designs may call for certain statistical methods to avoid increasing the chance
148 of erroneous conclusions and introducing bias in estimates and for complex adaptive designs,
149 such methods may not be readily available.[7] The decision to use or not use an adaptive design is
150 at the applicant's discretion.

151

152 In general, the design, conduct, and analysis of a proposed adaptive study design should satisfy
153 the following recommendations:

154

155 • The details of the adaptive design should be completely specified prior to initiation of the
156 study and documented accordingly. For example, prospective planning should include
157 prespecification of the anticipated number and timing of interim analyses, the type of
158 adaptation, the statistical inference methods to be used and the specific algorithm
159 governing the adaptive decision. If a study should be stopped early (e.g., for futility or
160 for success in demonstrating BE), detailed stopping criteria should be pre-specified and
161 scientifically justified.

162

163 • The applicant should establish that estimation of treatment effect will be sufficiently
164 reliable, and the chance of erroneous conclusions will be adequately controlled. The
165 Agency will accept appropriately designed BE studies that are scientifically justified.
166 Support might include published literature in peer-reviewed journals in which the
167 applicant's proposed approach is validated or simulation results meeting desired criteria
168 (e.g., the Type I error probability of the proposed approach is controlled at a nominal
169 level of 0.05 for a BE test). Appropriate details (e.g., literature references, proofs,
170 simulation codes/results) for the methodology should be submitted.

171

172 • The applicant should ensure that study integrity will be appropriately maintained. A
173 comprehensive written data access plan defining how study integrity will be maintained
174 in the presence of the planned adaption should be included in the protocol or statistical
175 analysis plan (SAP). This applies to both adaptive comparative clinical endpoint BE
176 studies and PK BE studies, whether blinded or unblinded by design.

177

178 For details, refer to the guidance for industry *Adaptive Design for Clinical Trials of Drugs and*
179 *Biologics* (November 2019).

---

[6] See footnote 4. For example, the ability to stop a trial early if it becomes clear that the trial is unlikely to demonstrate equivalence can reduce the number of patients exposed to the unnecessary risk of an ineffective investigational treatment and allow subjects the opportunity to explore more promising therapeutic alternatives.
[7] See footnotes 4 and 5.

180  Due to the increased complexity of adaptive studies and uncertainties regarding their operating
181  characteristics, applicants are encouraged to contact the Agency early to discuss their proposed
182  adaptive study designs and statistical methods via the controlled correspondence,[8] pre-ANDA
183  meeting,[9] pre-IND meeting, or pre-NDA meeting pathway.[10]
184
185       d.       Design with sparse sampling
186
187  For certain generic products, a sparse BE design is used, where the sampling for each subject is
188  done at a single or very limited number of time points rather than the number needed to get a full
189  concentration profile.  For example, some ophthalmic products are studied using a sparse BE
190  design, where only a single sample is collected from a single eye of each subject, at one assigned
191  sampling time point for that subject.  More generally, a sparse BE study design can be a parallel
192  design where each subject should receive only one treatment, T or R, but not both.  Alternatively,
193  a crossover sparse study design can be used where each subject receives both test and reference
194  treatments (e.g., in subjects undergoing indicated cataract surgery for both eyes).
195
196  For a sparse BE study design, the mean concentration for each product at each time point of
197  measurement is calculated by using the mean concentration of the subjects measured at each time
198  point to derive the mean profile for each product.  Based on the trapezoid rule, the $AUC_{0-t}$ for
199  each product is computed as a weighted linear combination of these mean concentrations at each
200  time point through time t.  The $AUC_{0-t}$ is the area under the concentration – time curve from
201  zero to the time t.  $C_{max}$ and $T_{max}$ (time to maximum observed concentration) can be determined
202  accordingly.  The ratios of $AUC_{0-t}$ and $C_{max}$ between the test and the reference product are used
203  to assess BE.  Estimation of the standard deviation and confidence interval for the ratio of
204  $AUC_{0-t}$ may be done by bootstrap or parametric methods (e.g., Bailer's methods (Bailer 1988)[11]
205  for a parallel study design), and that for the ratio of $C_{max}$ may be done by bootstrap methods.  BE
206  is supported if the 90% confidence interval for the ratio of $AUC_t$ between the test and the
207  reference product lies within the BE margin (80.00%, 125.00%).  Model-based approaches can
208  be considered when they can reliably control the error rate of concluding BE for bio inequivalent
209  products (Type I error).[12]
210
211  For complicated issues such as other forms of sparse design or alternative statistical methods,
212  applicants are encouraged to contact the Agency early to discuss their proposed study design and
213  statistical methods via the controlled correspondence, pre-ANDA meeting, pre-IND meeting, or
214  pre-NDA meeting pathway.[13]

---

[8] See the guidance for industry *Controlled Correspondence Related to Generic Drug Development* (December 2020).

[9] See the guidance for industry *Formal Meetings Between FDA and ANDA Applicants of Complex Products Under GDUFA* (October 2022).

[10] See the draft guidance for industry *Formal Meetings Between the FDA and Sponsors or Applicants of PDUFA Products* (December 2017).  When final, this guidance will represent FDA's current thinking on this topic.

[11] Bailer, A.J., 1988, Testing for the Equality of Area Under the Curves When Using Destructive Measurement Techniques, Journal of Pharmacokinetics and Biopharmaceutics, 16(3): 303-309.

[12] Zhao, L., M.-J. Kim, L. Zhang, and R. Lionberger, 2019, Generating Model Integrated Evidence for Generic Drug Development and Assessment, Clinical Pharmacology and Therapeutics, 105(2): 338-349.

[13] See footnotes 8, 9, and 10.

215
216        2.        *Sample Size Determination*
217
218 It is an applicant's responsibility to design an adequately powered BE study for the proposed
219 study. We recommend that applicants enroll enough subjects to power the study at a level of 0.8
220 or higher, for a BE test to be carried out with a type 1 error rate of 0.05 (see section III.C.1.a for
221 more details). When determining the sample size, rates of attrition and noncompliance (e.g.,
222 protocol violation) should be taken into consideration. Enough subjects should be recruited,
223 randomized, and dosed at the beginning of the study to ensure that the desired number of
224 evaluable subjects will be available for analysis. All eligible subjects who were dosed should be
225 included in the analysis. For BE studies, add-on subjects after the pre-specified number of
226 subjects have been reached are generally not encouraged except in an adaptive study design with
227 a pre-specified adaptation to add subjects and statistical methods to control the Type I error rate
228 under the nominal level.
229
230 The number of subjects to be included in a study should be based on an appropriate sample size
231 calculation for the proposed study design.[14,15,16] For example, the standard 2×2 cross-over study
232 will use a particular calculation while studies with a different design or set of endpoints will use
233 different calculations. For sample size re-estimation in an adaptive study design, refer to Section
234 II.A.1.c. Adaptive Design.
235
236 Sample size and power calculation should be supported by established scientific practice. For
237 complex study designs with no analytical solutions for sample size calculation, simulation can be
238 used to estimate the needed sample size in order to reach a desired power. The method by which
239 the sample size is determined should be given in the protocol, together with the estimates of any
240 quantities used in the calculations (such as variances, mean values, response rates, the assumed
241 effect size). The basis for these estimates should also be given. For example, variance estimates
242 can be obtained from the biomedical literature and/or pilot studies. It is important to investigate
243 the sensitivity of the sample size calculated to a variety of deviations from the assumed
244 estimates. This may be facilitated by providing a range of sample sizes appropriate for a
245 reasonable range of deviations from the assumptions or alternative approaches supported by
246 published peer-reviewed literature.
247
248 Applicants should enter a sufficient number of subjects in the study to allow for dropouts.
249 Dropouts generally should not be replaced because replacement of subjects during the study
250 could complicate the statistical model and analysis. Applicants who wish to replace dropouts
251 during the study should indicate this intention in the protocol. The protocol should also state
252 whether samples from replacement subjects, if not used, will be assayed. If the dropout rate is
253 high and applicants wish to add more subjects, a modification of the statistical analysis may be

---

[14] Chow, S.-C. and J.-P. Liu, 2008, Design and Analysis of Bioavailability and Bioequivalence Studies, 3rd Edition, New York: Chapman and Hall/CRC.

[15] Draft guidance for industry *Bioequivalence Studies with Pharmacokinetic Endpoints for Drugs Submitted Under an ANDA* (August 2021). When final, this guidance will represent FDA's current thinking on this topic.

[16] Patterson, S.D. and B. Jones, 2017, Bioequivalence and Statistics in Clinical Pharmacology, 2nd Edition, New York: Chapman and Hall/CRC.

254    recommended. Additional subjects should not be included after data analysis unless the study
255    was designed from the beginning as an adaptive design.

257    In general, for PK BE or in vitro BE studies, sample size calculation should be based on BE
258    metrics (e.g., AUC, $C_{max}$) after log-transformation; for comparative clinical endpoint BE studies,
259    sample size calculation should be based on the un-transformed comparative clinical endpoints
260    unless otherwise noted in the relevant FDA product-specific guidance (PSG).[17] The number of
261    evaluable subjects in a PK BE study should not be less than 12. For highly variable drug
262    products, a minimum of 24 subjects are recommended for BE assessment.[18]

264    **B.    Data Preparation**

266    The drug concentration in biological fluid determined at each sampling time point should be
267    furnished on the original scale for each subject participating in the study. The PK measures of
268    systemic exposure should also be furnished on the original scale. The variables for a
269    comparative clinical endpoint BE study should also be furnished on the original scale. The
270    mean, standard deviation, and coefficient of variation for each variable should be computed and
271    tabulated in the final report.

273    *1.    Log-Transformation*

275    A general approach to assessing BE is to compare the log-transformed BA measures after
276    administration of the T and R products.

278    *a.    Logarithmic transformation for PK measures*

280    This guidance recommends that PK BE measures (e.g., AUC and $C_{max}$) be log-transformed (see
281    Appendix B). The choice of common or natural logs should be consistent and should be stated in
282    the study report. The limited sample size in a typical BE study precludes a reliable
283    determination of the distribution of the data set. Sponsors and/or applicants are not encouraged
284    to test for normality of error distribution after log-transformation, nor should they use normality
285    of error distribution as a reason for carrying out the statistical analysis on the original scale.
286    Justification should be provided if sponsors or applicants believe that their BE study data should
287    be statistically analyzed on the original rather than on the log scale.

288
289
290
291

---

[17] For the most recent version of a product-specific guidance, check the product-specific web page at
https://www.accessdata.fda.gov/scripts/cder/psg/index.cfm.
[18] Davit, B. and D. Conner, 2010, Reference-Scaled Average Bioequivalence Approach. In: I. Kanfer and L.
Shargel, editors. Generic Drug Product Development — International Regulatory Requirements for Bioequivalence,
New York, NY: Informa Healthcare, 271-272; Food and Drug Administration, Advisory Committee for
Pharmaceutical Science, October 5-6, 2006.

292                  b.        Data transformation for comparative pharmacodynamic and clinical
293                               endpoint BE study

294

295    The decision on whether and how to transform a variable for a comparative pharmacodynamic
296    (PD) or comparative clinical endpoint BE study should be specified in the protocol, especially
297    for the primary variable(s).  The basis for the variables should also be given in the protocol.  For
298    example, these variables can be obtained from the biomedical literature and/or pilot studies.
299    Similar considerations apply to other derived variables, such as the use of change from baseline,
300    percentage change from baseline, the area under the curve of repeated measures, or the ratio of
301    two different variables.  Subsequent clinical interpretation should be carefully considered.
302    Regarding comparative clinical endpoint studies, in general the log-transformation is not
303    used.  For example, in the case of the Fieller's confidence interval for the ratio of two means, the
304    raw (untransformed) data are used for the confidence interval derivation.[19]

305

306                  c.        Negative values for baseline corrected PK or PD endpoints

307

308    Because data transformation and scales might affect BE conclusions, they should be chosen
309    carefully and appropriately justified in the protocol.[20]  If a baseline correction results in a
310    negative plasma concentration value, the value should be set equal to 0 before calculating the
311    baseline-corrected AUC.

312

313        2.        *Missing Data and Intercurrent Events*

314

315    Subjects may have missing data in the study for various reasons (e.g., subject's refusal to
316    continue in the study, worsening of conditions or emergence of adverse events, subject's failure
317    to meet scheduled appointments for evaluation).  Subjects may also have intercurrent (post-
318    randomization) events that affect either the interpretation or the existence of the measurements
319    associated with the question of interest (e.g., noncompliance with the protocol for various
320    reasons, use of rescue medication due to lack of efficacy, death).  Missing data and intercurrent
321    events can introduce problems such as bias, misleading inference, loss of precision and loss of
322    power, which make it hard to interpret the trial outcome.

323

324    The ICH (Internal Council for Harmonization) E9(R1) Addendum introduces the concept of an
325    estimand, which is a precise description of the treatment effect reflecting the clinical question
326    posed by a particular study objective.[21]  The trial protocol of a BE study should include the
327    following components of an estimand: (1) the treatment of interest and alternative treatment(s) to
328    which comparison will be made: e.g., test drug compared with reference drug; (2) the analysis
329    population for BE assessment; (3) the variable (or endpoint) to be measured for each subject
330    (e.g., AUC or $C_{max}$); (4) the specification of how to account for intercurrent events in assessing
331    the scientific question of interest (for example, in a comparative clinical endpoint BE study with

---

[19] Fieller, E., Some Problems in Interval Estimation, 1954, Journal of the Royal Statistical Society, 16(2): 175-185.
[20] For example, see Sun, W., S. Grosser, and Y. Tsong, 2017, Ratio of Means vs. Difference of Means as Measures of Superiority, Noninferiority, and Average Bioequivalence, Journal Biopharmaceutical Statistics, 27(2): 338-355.
[21] Guidance for industry *E9(R1) Statistical Principles for Clinical Trials: Addendum: Estimands and Sensitivity Analysis in Clinical Trials*, Revision 1 (May 2021).

332    a binary endpoint, subjects who discontinue study treatment early  due to lack of treatment effect
333    should be included as treatment failures); and (5) the population-level summary for the variable
334    to compare between treatment conditions, e.g., the geometric mean ratio of the test to reference
335    drug in a PK BE study.
336
337    The protocol should include plans to minimize missing data.  The trial protocol should
338    prospectively define anticipated causes of missing data, the corresponding statistical assumptions
339    about reasons for the missing data, and how missing data will be treated in the statistical
340    analysis.  The treatment of missing data in the statistical analysis should be justified such that
341    valid statistical inferences can be made under the assumptions about the missing data
342    mechanism.
343
344    Statistical methods for handling missing data include complete case analysis, available case
345    analysis, weighting methods, imputation, and model-based approaches.  For example, in a two-
346    way crossover study, a complete case analysis could be a general linear model as implemented in
347    SAS PROC GLM, which removes all subjects with any missing observations for any variables
348    included in the GLM model (i.e., removes subjects missing one or both periods).  An available
349    case analysis could be done using SAS PROC MIXED, which uses all observed data (e.g., in a
350    two-way crossover study, uses all subjects with one or two complete periods of data).
351
352    Approaches for handling missing data and the statistical methods for the primary BE analysis
353    (e.g., GLM vs. MIXED) should be pre-specified in the study protocol or SAP.  Depending on the
354    nature of the assumed or likely missing data mechanism, statistical methods from any of these
355    categories may be appropriate.  The validity of a statistical approach to handle missing data
356    depends on a variety of factors, including, but not limited to, the mechanism for missingness, the
357    fraction of incomplete cases, the values that are missing, specifics of the analysis, and definition
358    of the estimand.  Sensitivity analyses using alternative approaches may also be used in the
359    statistical analysis to address missing data.  Sensitivity analyses should be pre-specified in the
360    trial protocol to evaluate the robustness of conclusions to deviations from the assumptions about
361    the missing data mechanism.  The applicant should provide detailed information about reasons
362    for missing data and any observed intercurrent events.
363
364    For a particular drug product, if the PSG recommends certain approaches to handling missing
365    data, the applicants should refer to that PSG.  Applicants may choose to contact the Agency via
366    the controlled correspondence, pre-ANDA meeting, pre-IND meeting, or pre-NDA meeting
367    pathway to discuss their proposed approach to handling missing data if such an approach is
368    different from what is recommended in the PSG or if the applicants have further questions.
369
370          *3.      Outlier Detection*
371
372    Outlier data in BE studies are defined as subject data for one or more BA measures that are
373    discordant with corresponding data for that subject and/or for the rest of the subjects in a study.
374    Because BE studies are usually carried out as crossover studies, the most important type of
375    subject outlier is the within-subject outlier, when one subject or a few subjects differ notably
376    from the rest of the subjects with respect to a within-subject T-R comparison.  The existence of a

377   subject outlier with no protocol violations and for which there are not bioanalytical errors could
378   indicate one of the following situations:
379
380           a.        Product failure
381
382   Product failure could occur, for example, when a subject exhibits an unusually high or low
383   response to one or the other of the products because of a problem with the specific dosage unit
384   administered.  This could occur, for example, with a sustained and/or delayed-release dosage
385   form exhibiting dose dumping or a dosage unit with a coating that inhibits dissolution.
386
387           b.        Subject-by-formulation interaction
388
389   A subject-by-formulation interaction could occur when an individual is representative of subjects
390   present in the general population in low numbers, for whom the relative BA of the two products
391   is markedly different from that for most of the population, and for whom the two products are
392   not bioequivalent, even though they might be bioequivalent in most of the population.  In the
393   case of product failure, the unusual response could be present for either the T or R product.
394   However, in the case of a subpopulation, even if the unusual response is observed on the R
395   product, there could still be concern about lack of bioequivalence of the two products.  For these
396   reasons, applicants should not remove data from the statistical analysis of BE studies solely
397   because those data are identified as statistical outliers.
398
399   In general, outlier data (whether due to product failure, subject-by-formulation interaction, or
400   another cause) may only be removed from the BE statistical analysis if there is real-time
401   documentation demonstrating a protocol violation during the clinical and/or
402   analytical/experimental phase of the BE study.  Applicants should include a prospective plan in
403   the BE study protocol for handling subjects (experimental outliers) in the BE statistical analysis.
404   Data from redosing studies are not considered valid evidence to support removal of outlier data
405   from the statistical analysis.  All subject data should be submitted, with potential outliers flagged
406   with appropriate documentation as part of the submission.  However, for a replicated PK BE
407   study, if reference-scaled average BE is used, the applicant should ensure that the calculated
408   intra-subject variability is not inflated due to extreme values or situations.
409
410   To characterize aberrant observations for exploratory or quality control purposes, the choice of
411   the appropriate technique depends on whether there are outlying subjects or outlying
412   observations, as well as on the study design.
413
414       **C.    Statistical Models**
415
416       *1.    General Statistical Criteria for Bioequivalence*
417
418   The general structure of a BE criterion is that a function ($\Theta$) of population measures should be
419   demonstrated to be no greater than a specified value ($\theta$).  Using the terminology of statistical
420   hypothesis testing, this is accomplished by testing the hypothesis $H_0$: $\Theta \geq \theta$ versus $H_a$: $\Theta < \theta$ at a

421     desired level of significance, often 5%. Rejection of the null hypothesis $H_0$ (i.e., demonstrating
422     that the estimate of $\Theta$ is statistically significantly less than $\theta$) results in a conclusion of BE.
423
424                  a.       Use of confidence intervals to do two one-sided tests
425
426     In BE assessment we are frequently interested in testing whether a parameter (for example, the
427     difference of means for a T and R product for a specific endpoint) is contained within a defined
428     interval, call it $[\theta_1, \theta_2]$. The recommended method for doing such a test is the *Two One-Sided*
429     *Tests Procedure*.[22] A one-sided statistical test is carried out to determine whether the parameter
430     is $\geq \theta_1$, and a second one-sided test is carried out to determine whether the parameter is $\leq \theta_2$;
431     both tests are carried out at a level of significance $\alpha$, which is usually 0.05. If both tests are
432     successful (that is, we reject the null hypothesis in both cases), we conclude that the parameter is
433     contained in $[\theta_1, \theta_2]$.
434
435     These two one-sided tests are sometimes carried out by calculating a $100 (1-2\alpha)$ % confidence
436     interval for the parameter and determining whether this confidence interval is completely
437     contained in the interval $[\theta_1, \theta_2]$. For this confidence interval method of carrying out the tests to
438     be valid, the confidence interval should be an *equal tails* confidence interval. If the lower and
439     upper confidence limits of the $100 (1-2\alpha)$ % confidence interval are $L_1$ and $L_2$, respectively, then
440     the confidence interval is *equal tails* if $L_1$, by itself, is at least a $100 (1-\alpha)$ % lower confidence
441     bound for the parameter and $L_2$, by itself, is at least a $100 (1-\alpha)$ % upper confidence bound for
442     the parameter.
443
444     In some cases, there may not be general agreement as to the best choice of a particular statistical
445     testing methodology for carrying out the two one-sided tests (for example, if the parameter of
446     interest is the difference between the success probabilities for a T and R product for a binary
447     endpoint). In such cases, careful consideration should be given to the choice of statistical
448     methods for doing the two one-sided tests, which may or may not correspond to a confidence
449     interval method.
450
451          2.       *Statistical Information and Implementation of Criteria for PK Measures (AUC$_{0-t}$,*
452               *AUC$_{0-\infty}$, and C$_{max}$)*
453
454     We recommend that applicants provide the following statistical information for $AUC_{0-t}$,
455     $AUC_{0-\infty}$, and $C_{max}$:
456
457        •    Geometric means for the formulations tested
458        •    Arithmetic means for the formulations tested
459        •    Geometric mean ratios of Test vs. Reference and their corresponding 90% confidence
460             intervals or 95% upper confidence bounds (e.g., for highly variable drugs or narrow
461             therapeutic index drugs)

---

[22] Schuirmann, D. J., 1987, A Comparison of the Two One-Sided Tests Procedure and the Power Approach for
Assessing the Equivalence of Average Bioavailability, Journal of Pharmacokinetics and Biopharmaceutics, 15(6):
657-680.

462
463 Recommended statistical information for other types of outcome measures is discussed in section
464 III: Specific Situations.
465
466 To facilitate BE comparisons, for crossover studies, the measures for each individual should be
467 displayed in parallel for the formulations tested. For each BE measure, the ratio of the individual
468 geometric mean of the T product to the individual geometric mean of the R product should be
469 tabulated side by side. The summary tables should indicate in which sequence each subject
470 received the product.
471
472 Statistical analyses of BE data are typically based on a statistical model for the logarithm of the
473 BA measures (e.g., AUC and $C_{max}$). The model is a mixed-effects or two-stage linear model.
474 Each subject, j, theoretically provides a mean for the log-transformed BA measure for each
475 formulation, $\mu_{Tj}$ and $\mu_{Rj}$ for the T and R formulations, respectively. The model assumes that
476 these subject-specific means come from a distribution with population means $\mu_T$ and $\mu_R$, and
477 between-subject variances $\sigma_{BT}^2$ and $\sigma_{BR}^2$, respectively. The model allows for a correlation, $\rho$,
478 between $\mu_{Tj}$ and $\mu_{Rj}$. The subject-by-formulation interaction variance component, $\sigma_D^2$, is related
479 to these parameters as follows:
480
481 $$\sigma_D^2 = \text{variance of } (\mu_{Tj} - \mu_{Rj})$$
482
483 $$= (\sigma_{BT} - \sigma_{BR})^2 + 2 (1-\rho)\sigma_{BT}\sigma_{BR} \text{ [23]}$$
484
485 For a given subject, the observed data for the log-transformed BA measure are assumed to be
486 independent observations from distributions with means $\mu_{Tj}$ and $\mu_{Rj}$, and within-subject variances
487 $\sigma_{WT}^2$ and $\sigma_{WR}^2$. The total variances for each formulation are defined as the sum of the within-
488 and between-subject components (i.e., $\sigma_{TT}^2 = \sigma_{WT}^2 + \sigma_{BT}^2$ and $\sigma_{TR}^2 = \sigma_{WR}^2 + \sigma_{BR}^2$). For analysis
489 of crossover studies, the means are given additional structure by the inclusion of period and
490 sequence effect terms.
491
492 The applicant may also consider prespecifying inclusion of important demographic and baseline
493 prognostic covariates in the statistical model for parallel studies. This sort of adjustment can
494 increase the precision and power of the statistical analysis and compensate for any lack of
495 balance between treatment groups with no inflation of Type 1 error.
496
497
498
499

---

[23] Schall, R., and H. G. Luus, 1993, On Population and Individual Bioequivalence, Statistics in Medicine, 12(12): 1109-1124.

500 **III.    SPECIFIC SITUATIONS[24]**

501

502        **A.       In Vitro Bioequivalence and Population Bioequivalence**

503

504    This section discusses statistical methods for assessment of in vitro BE, including population BE
505    (PBE), a similarity index ($f_2$), statistical approaches respectively for in vitro release tests (IVRT),
506    in vitro permeation tests (IVPT) and in vitro abuse-deterrent formulations (ADF) comparative
507    studies, and a profile comparison approach based on Earth Mover's Distance (EMD).

508

509        *1.       Population Bioequivalence*

510

511    One of the recommended statistical approaches for evaluating in vitro BE is population BE
512    (PBE). To test for PBE, the null and alternative hypotheses are given as follows:

513 $$H_0: \theta \geq \theta_P \quad \text{vs.} \quad H_a: \theta < \theta_P$$

514    where $\theta = \frac{(\mu_T - \mu_R)^2 + \sigma_T^2 - \sigma_R^2}{\sigma_R^2}$ if the estimated $\sigma_R > \sigma_0$ or $\theta = \frac{(\mu_T - \mu_R)^2 + \sigma_T^2 - \sigma_R^2}{\sigma_0^2}$ if the estimated

515    $\sigma_R \leq \sigma_0$.
516    Here, $\mu_T$ and $\mu_R$ are the population means, $\sigma_T^2$ and $\sigma_R^2$ are the population variances of the log-
517    transformed measure for T and R products, respectively; $\sigma_0^2$ is a regulatory constant for variance;
518    and $\theta_P$ is the PBE limit. The concept of PBE is to compare the difference of the T and R
519    products with that of the reference versus reference itself. This comparison can be denoted in
520    terms of the population difference ratio as follows:

521 $$\sqrt{\frac{E(Y_T - Y_R)^2}{E(Y_R - Y_R')^2}} = \sqrt{\frac{(\mu_T - \mu_R)^2 + \sigma_R^2 + \sigma_T^2}{2\sigma_R^2}} = \sqrt{\frac{\theta}{2} + 1}.$$

522    The regulatory constant variance, $\sigma_0^2$, is set based on the following considerations. Due to the
523    low variability of in vitro measurements, this guidance recommends that the ratio of geometric
524    means should fall within 0.90 and 1.11. As a result, an upper BE limit of 1.11 is recommended
525    for the average BE limit for in vitro data. Assuming $\sigma_R^2 = \sigma_T^2 = \sigma_0^2$, $\mu_T - \mu_R = \ln 1.11$ and the
526    maximum allowable limit for population difference ratio is 1.25, this leads to the recommended
527    choice of $\sigma_0^2 = 0.01$.

528

529    The determination of PBE limit, $\theta_P$, is based on the consideration of average BE criterion and
530    the addition of variance terms to PBE criterion as the following form:

531 $$\frac{(\mu_T - \mu_R)^2 + \sigma_T^2 - \sigma_R^2}{\max\{\sigma_0^2, \sigma_R^2\}} = \frac{\text{Average BE limit} + \text{Variance term}}{\text{Scaled variance term}}.$$

532

533    The FDA recommended allowance for the variance term is 0.01. This value may be adjusted
534    depending on the average BE limit for in vitro data based on further communication with the
535    Agency. Accordingly, the PBE limit, $\theta_P$, is recommended as follows:

---

[24] Some specific situations are addressed in the following subsections with specified choices of BE criteria. Further
discussion regarding these specified choices can be found in the guidances cited in those subsections.

536
$$\theta_P = \frac{(\ln 1.11)^2 + 0.01}{0.01} = 2.089$$

537
538 A linearized form is recommended to use to test $H_0: \theta \geq \theta_P$. That is, testing $H_0: \theta \geq \theta_P$ is
539 equivalent to testing $H_0: \gamma \geq 0$ where $\gamma = (\mu_T - \mu_R)^2 + (\sigma_T^2 - \sigma_R^2) - \theta_P \sigma_R^2$ if the estimated
540 $\sigma_R > \sigma_0$ or $\gamma = (\mu_T - \mu_R)^2 + (\sigma_T^2 - \sigma_R^2) - \theta_P \sigma_0^2$ if the estimated $\sigma_R \leq \sigma_0$. Here, $\gamma_1 =$
541 $(\mu_T - \mu_R)^2$, $\gamma_2 = \sigma_T^2$ and $\gamma_3 = \sigma_R^2 + \theta_P \sigma_R^2$ if the estimated $\sigma_R > \sigma_0$ or $\gamma_3 = \sigma_R^2 + \theta_P \sigma_0^2$ if the
542 estimated $\sigma_R \leq \sigma_0$.
543 Suppose $\hat{\gamma}_U$ is a 95% upper confidence bound for $\gamma$. Then, PBE is supported if and only if $\hat{\gamma}_U \leq$
544 0. Based on the work of Howe (1974)[25] and Ting et al. (1990)[26], an approximate 95% upper
545 confidence bound for $\gamma$ is given as follows:

546
$$\hat{\gamma}_U = \hat{\gamma}_1 + \hat{\gamma}_2 - \hat{\gamma}_3 + \sqrt{(\tilde{\gamma}_1 - \hat{\gamma}_1)^2 + (\tilde{\gamma}_2 - \hat{\gamma}_2)^2 + (\tilde{\gamma}_3 - \hat{\gamma}_3)^2}$$

547
548 where $\hat{\gamma}_1$, $\hat{\gamma}_2$, and $\hat{\gamma}_3$ are point estimators of $\gamma_1$, $\gamma_2$, and $\gamma_3$, respectively; $\tilde{\gamma}_1$ and $\tilde{\gamma}_2$ are 95%
549 upper confidence bounds for $\gamma_1$ and $\gamma_2$ and $\tilde{\gamma}_3$ is a 95% lower confidence bound for $\gamma_3$. For
550 further detail, see, e.g., the draft PSGs for Budesonide suspension (September 2012) and
551 Fluticasone Propionate metered spray (June 2020).[27]

552
553     *2.    Similarity Index (f₂)*
554
555 For a comparison of dissolution profiles, similarity is assessed using the similarity index, $f_2$
556 (Shah et al., 1998),[28] as described in detail in the guidance for industry *Immediate Release Solid*
557 *Oral Dosage Forms Scale-Up and Postapproval Changes: Chemistry, Manufacturing, and*
558 *Controls, In Vitro Dissolution Testing, and In Vivo Bioequivalence Documentation* (November
559 1995). In particular, given that all profiles are conducted on a minimum of 12 individual dosage
560 units, 2 profiles are similar if the value of their similarity factor $f_2$ is between 50 and 100.

561
562     *3.    In-Vitro Release Test*
563
564 When an in-vitro release test (IVRT) is used to support a demonstration of BE for topical
565 dermatological drug products as part of an in vitro characterization-based BE approach, a two-
566 stage, nonparametric statistical approach is recommended, and described in the draft guidance
567 for industry *In Vitro Release Test Studies for Topical Drug Products Submitted in ANDAs*
568 (October 2022).[29] The statistical approach is the same as that used to assess the equivalence of
569 drug release rates for non-sterile semisolid dosage forms evaluated by a comparative IVRT study
570 in the context of certain postapproval changes; this is shown in detail in the guidance for industry

---

[25] Howe, W.G., 1974, Approximate Confidence Limits of the Mean of X+Y Where X and Y are Two Tabled Independent Random Variables, Journal of the American Statistical Association, 69:789-794.
[26] Ting, N., R.K. Burdick, F. Graybill, S. Jeyaratnam, and T.F.C. Lu, 1990, Confidence Intervals on Linear Combinations of Variance Components That Are Unrestricted in Sign, Journal of Statistical Computation and Simulation, 35:135-143.
[27] When final, these guidances will represent FDA's current thinking on these topics.
[28] Shah, V.P., Y. Tsong, P. Sathe, and J.P. Liu, 1998, In Vitro Dissolution Profile Comparison—Statistics and Analysis of the Similarity Factor, f2, Pharmaceutical Research, 15(6):889-896.
[29] When final, this guidance will represent FDA's current thinking on this topic.

571 *Nonsterile Semisolid Dosage Forms — Scale-Up and Postapproval Changes: Chemistry,*
572 *Manufacturing, and Controls; In Vitro Release Testing and In Vivo Bioequivalence*
573 *Documentation* (May 1997).

574
575 The assessment of equivalence by an IVRT involves a comparison of the median in vitro drug
576 release rates of two formulations using a non-parametric statistical test which is resistant to
577 outliers that are expected to occur under the particular testing conditions.
578
579     *4.      In-Vitro Permeation Test*
580
581 When an in-vitro permeation test (IVPT) is used to support a demonstration of BE for topical
582 dermatological drug products as part of an in vitro characterization-based BE approach, a mixed
583 scaled criterion is recommended, and described in detail in the draft guidance for industry *In*
584 *Vitro Permeation Test Studies for Topical Drug Products Submitted in ANDAs* (October 2022).[30]
585 According to that methodology, a confidence interval is calculated for each of the endpoints, log-
586 transformed maximum flux ($J_{max}$) and log-transformed total (cumulative) amount (AMT)
587 permeated.  The permeation test is performed with excised skin sections from patients
588 undergoing a surgical procedure or from cadaver donors and the statistical test uses the within-
589 reference standard deviation, $S_{WR}$, as the threshold that prompts use of either the unscaled or
590 scaled confidence interval.
591
592 The mixed-scaled criterion uses the within-reference standard deviation as a threshold,
593 independently, for each endpoint.  Specifically, for $J_{max}$ or log-transformed total (cumulative)
594 amount permeated, the reference-scaled average BE approach is used for the endpoint only if it
595 has a $S_{WR} > 0.294$. The regular ABE approach (refer to Schuirmann, 1987)[31] is used for the
596 endpoint with $S_{WR} \leq 0.294$.
597
598 In the reference-scaled average BE approach, the hypotheses to be tested are:
599

600 $$H_0: \frac{(\mu_T - \mu_R)^2}{\sigma_{WR}^2} \geq \theta$$

601 $$H_a: \frac{(\mu_T - \mu_R)^2}{\sigma_{WR}^2} < \theta$$

602 Here we determine the 100(1-α)% upper confidence bound for $(\mu_T - \mu_R)^2 - \theta\sigma_{WR}^2$
603 where:
604     -   $\mu_T - \mu_R$ = mean difference of T and R products
605     -   $\sigma_{WR}^2$ = within-subject variance of R product
606     -   $\theta = \frac{(\ln(m))^2}{(\sigma_{W0})^2}$, $m = 1.25$, and $\sigma_{W0} = 0.25$ (regulatory constant)

607 For the T product to be bioequivalent to the R product, both of the following conditions must be
608 satisfied for each endpoint tested:

---

[30] When final, this guidance will represent FDA's current thinking on this topic.
[31] See footnote 22.

609
610         a.   The 95% upper confidence bound for $(\mu_T - \mu_R)^2 - \theta\sigma_{WR}^2$ must be less than
611             or equal to zero (numbers should be kept to a minimum of four significant
612             figures for comparison).
613
614         b.   The point estimate of the T/R geometric mean ratio must fall within the pre-
615             specified limits $\left[\frac{1}{m}, \; m\right]$, where $m = 1.25$.
616
617 In the case of the non-scaled approach, we calculate the $100(1\text{-}2\alpha)\%$ confidence interval for
618 $\mu_T - \mu_R$ as
619

$$\bar{I} \pm t_{(1-\alpha),(n-1)} * \sqrt{\frac{S_I^2}{n}}$$

620

621
622 where:
623     -   $\bar{I}$ is the point estimate for the mean difference of T and R products
624     -   $S_I^2$ stimate of inter-donor variability

625     -   $t_{(1-\alpha),(n-1)}$ is the $100(1-\alpha)$ percentile of the student's t-distribution with $(n-1)$
626       degrees of freedom

627     -   n is the number of donors

628     -   the value of α is usually set at 0.05
629
630 For the T product to be bioequivalent to the R product, the $100(1\text{-}2\alpha)\%$ confidence interval for
631 $\mu_T - \mu_R$ must be contained within the limits $\left[\frac{1}{m}, \; m\right]$ in the original scale for each endpoint
632 tested, where $m = 1.25$.
633
634         5.      *Abuse-Deterrent Formulation Comparative Studies*
635
636 An ADF is a formulation that has abuse-deterrent properties, which are defined as drug product
637 properties that are expected to meaningfully deter certain types of abuse, even if they do not fully
638 prevent abuse.[32] The general BE statistical considerations for in vitro ADF comparative studies
639 presented in this guidance align with the guidance for industry – *Abuse-Deterrent Opioids —*
640 *Evaluation and Labeling*[33] and the guidance for industry – *General Principles for Evaluating the*
641 *Abuse Deterrence of Generic Solid Oral Opioid Drug Products* (November 2017). The potential
642 route of abuse (i.e., ingestion (oral route), injection (parenteral route), insufflation (nasal route), or
643 smoking (inhalation route)) and its relevance to ADF design feature(s) will determine how an
644 applicant should evaluate the abuse deterrence of the product utilizing a tier-based approach. To
645 support in vitro ADF comparative studies, the Agency recommends applicants provide

---

[32] See the guidance for industry *Abuse-Deterrent Opioids - Evaluation and Labeling* (April 2015).
[33] Ibid.

646   justification for the sample size, statistical test, and number of batches to assess the abuse-deterrent
647   properties and demonstrate consistency of abuse-deterrent performance throughout the drug
648   product shelf-life and lifecycle (i.e., postapproval changes).   Applicants should consider a
649   standardized accept/reject criterion based on delta or confidence interval relevant to the abuse-
650   deterrent outcome.  The Agency recommends the use of relevant statistics (e.g., sampling plans)
651   to support evaluation of abuse-deterrent properties.

653   For ANDA submissions, a non-inferiority approach should be taken when comparing T product
654   with R product to conclude that T product is no less abuse deterrent than R product.[34]  The Agency
655   recommends inferential analyses to evaluate the abuse deterrence of T product versus R product.
656   In the analyses, a hierarchical set of null hypotheses serves as a gatekeeper for subsequent null
657   hypotheses, evaluating the abuse deterrence of T and R products under progressively more
658   challenging conditions.  A hierarchical inferential approach is used to maintain a fixed family-wise
659   experiment Type I error rate.  Typically, the acceptable Type I error probability ($\alpha$) will be set at
660   5%.

662              6.         *Earth Mover's Distance* Based Profile Comparison Approach

664   EMD is a statistical metric that measures the discrepancy (distance) between distributions
665   without a prior assumption of the distribution.[35]  The EMD has been recommended in a profile
666   comparison approach to assess equivalence of particle size distribution profile,[36] where the
667   profile exhibits complex distribution (i.e., multiple peaks) that cannot be accurately described by
668   some conventional descriptors (e.g., the D50 and SPAN).  The EMD-based profile comparison
669   approach is briefly described as follows. To assess equivalence between the T and R product
670   formulations in the particle size distribution shape, an average profile of all R product samples
671   (i.e., R center) is calculated and serves as the reference profile to compute the distance between
672   an R or a T product sample to the R center using the EMD algorithm.  After obtaining the profile
673   distances between each R product sample and the R product average (R – R center distance), and
674   the profile distances between each T product sample and the R product average (T – 'R center'
675   distance), a statistical equivalence method, e.g., the PBE, is then applied to the two groups of
676   distances to indicate whether the T and R products are statistically equivalent in the particle size
677   distribution shape.  For details, refer to Rubner et al. (2000).[37]

679   Importantly, considering the increasingly emerging technologies and methods for in vitro BE
680   studies, applicants are encouraged to contact the Agency early to discuss their proposed study
681   designs and statistical methods via the controlled correspondence, pre-ANDA meeting, pre-IND
682   meeting, or pre-NDA meeting pathway.[38]

---

[34] Guidance for Industry *Evaluating the Abuse Deterrence of Generic Solid Oral Opioid Drug Products* (November 2017).

[35] Rubner, Y., C. Tomasi, and L.J. Guibas, 2000, The Earth Mover's Distance as a Metric for Image Retrieval, International Journal of Computer Vision, 40(2):99-121.

[36] Draft PSG for industry on Cyclosporine emulsion (October 2016). When final, this guidance will represent the FDA's current thinking on this topic.

[37] See footnote 35.

[38] See footnotes 8, 9, and 10.

684    **B.    Statistical Methods for Narrow Therapeutic Index and Highly Variable Drug**
685         **Products**
686
687        *1.    Statistical Method for Narrow Therapeutic Index Drugs*
688
689    If a drug is a narrow therapeutic index drug, a fully replicated cross-over design should be used.
690    The statistical analysis should be carried out using both the ABE and the reference-scaled
691    average BE tests for both AUC and $C_{max}$.
692
693    The reference-scaled average BE is evaluated by testing the null hypothesis:
694        $$H_0 : \quad \frac{(\mu_T - \mu_R)^2}{\sigma_{WR}^2} \geq \theta$$
695    versus the alternative hypothesis:
696        $$H_a : \quad \frac{(\mu_T - \mu_R)^2}{\sigma_{WR}^2} < \theta$$
697
698    where:
699        – $\mu_T$ is the population average response of the log-transformed measure for the Test
700            formulation.

701        – $\mu_R$ is the population average response of the log-transformed measure for the
702            Reference formulation.

703        – $\sigma_{WR}^2$ is the population within subject variance of the Reference formulation.

704        – $\theta = \frac{[\ln(\Delta)]^2}{\sigma_{W0}^2}$ is the BE limit.

705        – $\Delta$ and $\sigma_{W0}^2$ are predetermined constants.  Refer to the draft guidance for industry
706            *Bioequivalence Studies With Pharmacokinetic Endpoints for Drugs Submitted*
707            *Under an ANDA* (August 2021) for the values of $\Delta$ and $\sigma_{W0}^2$.[39]

708    Testing is usually done at $\alpha$=0.05 and that rejection of the null hypothesis supports the
709    conclusion of bioequivalence.
710
711    Narrow therapeutic index BE studies should pass both the reference-scaled approach and the
712    unscaled average BE limits of 80.00 to 125.00%.
713
714    In addition, the test/reference ratio of the within-subject standard deviation should be evaluated.
715    The within-subject variability comparison of the T and R drug products is carried out by a one-
716    sided F test.  The null hypothesis for this test is the following.
717
718    $H_0 : \frac{\sigma_{WT}}{\sigma_{WR}} \geq \delta$

---

[39] When final, this guidance will represent FDA's current thinking on this topic.

719
720     And the alternative hypothesis is:
721
722     $H_a : \frac{\sigma_{WT}}{\sigma_{WR}} < \delta$
723
724     where $\sigma_{WT}$ is the within-subject standard deviation for the test product, $\sigma_{WR}$ is the within-subject
725     standard deviation for the reference product and $\delta$ is the limit to declare the within-subject
726     variability of the test product is not greater than that of the reference product (refer to the draft
727     guidance for industry *Bioequivalence Studies With Pharmacokinetic Endpoints for Drugs*
728     *Submitted Under an ANDA* (August 2021) where $\delta$ was set to 2.5).[40]
729
730     • The 100(1-α)% CI for $\sigma_{WT}/\sigma_{WR}$ is given by
731     • $\left( \frac{s_{wt}/s_{wR}}{\sqrt{F_{\frac{\alpha}{2}}(v_1, v_2)}}, \frac{s_{wt}/s_{wR}}{\sqrt{F_{1-\frac{\alpha}{2}}(v_1, v_2)}} \right)$

732         Here, α=0.1, $F_{\frac{\alpha}{2}}(v_1, v_2)$ and $F_{1-\frac{\alpha}{2}}(v_1, v_2)$ are the values of the F-distribution with $v_1$
733         (numerator) and $v_2$ (denominator) degrees of freedom that has probability of α/2 and 1-
734         α/2 to its right, respectively.
735
736         2.      *Statistical Method for Highly Variable Drugs*
737
738     If a drug is a high variable drug, a partial or fully replicated cross-over design should be used.
739     The statistical analysis should be carried out using the mixed scaling approach below for both
740     AUC and $C_{max}$.
741
742     The mixed scaling approach:
743
744     If the estimated within-subject standard deviation of the RLD is < 0.294, the two one-sided test
745     procedure should be used to determine BE for the individual PK parameter.  Otherwise, the
746     reference-scaled procedure should be used to determine BE for the individual PK parameter
747     together with a point estimate constraint for the estimated test/reference geometric mean ratio.
748
749     For the reference-scaled approach the upper BE limit for Test/Reference ratio of geometric
750     means is $\Delta = \frac{1}{0.8}$ , the regulatory constant is $\sigma_{w0} = 0.25$ and the point estimate constraint is
751     80.00 to 125.00%.
752
753     Refer to the draft guidance for industry *Bioequivalence Studies With Pharmacokinetic Endpoints*
754     *for Drugs Submitted Under an ANDA* (August 2021) for further details.[41]
755

---

[40] When final, this guidance will represent FDA's current thinking on this topic.
[41] When final, this guidance will represent FDA's current thinking on this topic.

756        **C.**       <mark>**Comparative Clinical Endpoint**</mark> **Bioequivalence Studies**
757
758    For some products, the <mark>PSG may recommend</mark> an appropriately designed comparative clinical
759    endpoint BE study. In particular, a comparative clinical endpoint BE study is an option to be
760    considered for measuring BA or demonstrating BE of dosage forms intended to deliver the active
761    moiety locally, e.g., topical preparations for the skin, eye, and mucous membranes; oral dosage
762    forms not intended to be systemically absorbed, e.g., an antacid; bronchodilators administered by
763    oral inhalation.
764
765    In general, these studies will have a <mark>randomized, parallel group design, with three arms: test,</mark>
766    <mark>reference, and placebo/vehicle.</mark>
767
768       •    A placebo/vehicle arm is recommended to demonstrate that the T product and R product
769           are active and to establish that the study is sufficiently sensitive to detect differences
770           between products at the lower end of the dose/response curve.
771
772    To establish BE, it is recommended that the following <mark>compound hypotheses (continuous</mark>
773    <mark>endpoint or dichotomous endpoint)</mark> be tested. Rejection of the null hypothesis supports the
774    conclusion of equivalence of the two products.
775
776    For a continuous endpoint:
777    The null hypothesis for this test is:
778
779    $H_0$: $\mu_T / \mu_R \leq \theta_1$ or $\mu_T / \mu_R \geq \theta_2$
780
781    versus the alternative hypothesis:
782    $H_a$: $\theta_1 < \mu_T / \mu_R < \theta_2$
783
784    where:
785          –    $\mu_T$ = mean of the primary endpoint for the test group, and
786          –    $\mu_R$ = mean of the primary endpoint for the reference group.
787
788    The null hypothesis, $H_0$, is rejected with a Type I error ($\alpha$) of 0.05 (two one-sided tests) if the
789    90% confidence interval for the ratio of the means between T and R products ($\mu_T / \mu_R$) is
790    contained within the interval [$\theta_1$, $\theta_2$].
791
792    For a dichotomous endpoint:
793    The null hypothesis for this test is:
794
795    $H_0$: $\pi_T - \pi_R \leq \Delta_1$ or $\pi_T - \pi_R \geq \Delta_2$
796
797    versus the alternative hypothesis:
798    $H_a$: $\Delta_1 < \pi_T - \pi_R < \Delta_2$
799

800    where:

801        –   $\pi_T$= the success rate of the primary endpoint for the treatment group, and $\pi_R$= the

802           success rate of the primary endpoint for the reference group.

803

804    The null hypothesis, $H_0$, is rejected with a Type I error ($\alpha$) of 0.05 (two one-sided tests) if the

805    estimated 90% confidence interval for the difference of the success rates between T and R

806    products ($\pi_T$–$\pi_R$) is contained within the interval [$\Delta_1$, $\Delta_2$].

807

808      •   For continuous and binary endpoints, in order to demonstrate adequate study sensitivity,

809         the test product and reference product should both be statistically superior to placebo

810         ($p<0.05$) with regard to the primary endpoint.

811

812      •   Refer to PSGs for comparative clinical endpoint BE study designs, definitions of study

813         populations, regulatory constant (e.g., equivalence interval limit), and analyses specific to

814         a given product.

815

816      **D.**      **Studies in Multiple Groups**

817

818    There can be multiple sources of group[42] effects in BE studies.  Sometimes, groups reflect

819    factors arising from study design and conduct. For example, a PK BE study can be carried out in

820    two or more clinical centers and the study may be considered a multi-group BE study.  The

821    combination of multiple factors may complicate the designation of group.  Therefore, sponsors

822    should minimize the group effect in a PK BE study as recommended below:

823

824      (1) Dose all groups at the same clinic unless multiple clinics are needed to enroll a

825        sufficient number of subjects.

826

827      (2) Recruit subjects from the same enrollment pool to achieve similar demographics

828        among groups.

829

830      (3) Recruit all subjects, and randomly assign them to group and treatment arm, at study

831        outset.

832

833      (4) Follow the same protocol criteria and procedures for all groups.

834

835      (5) When feasible (e.g., when healthy volunteers are enrolled), assign an equal sample

836        size to each group.

837

838    Bioequivalence should be determined based on the overall treatment effect in the whole study

839    population.  In general, the assessment of BE in the whole study population should be done

840    without including the treatment and group interaction(s) term in the model, but applicants may

841    also use other pre-specified models, as appropriate (Fleiss 1986, Permutt 2003, Tsiatis et al.

---

[42] In literature, the term *group* is sometimes referred to as *subgroup*.

842    2008).[43]  The assessment of interaction between the treatment and group(s) is important,
843    especially if any of the first four study design criteria recommended above are not met and the
844    PK BE data are considered pivotal information for drug approval.  If the interaction term of
845    group and treatment is significant (Alosh et al. 2015, Grizzle 1965),[44] heterogeneity of treatment
846    effect across groups should be carefully examined and interpreted with care.  If the observed
847    treatment effect of the products varies greatly among the groups, vigorous attempts should be
848    made to find an explanation for the heterogeneity in terms of other features of trial management
849    or subject characteristics, which may suggest appropriate further analysis and interpretation.
850
851    It is important that statistical methods and models for the primary BE analysis are fully pre-
852    specified in the protocol or SAP (e.g., in an ANDA study, the applicant should pre-specify
853    detailed statistical criteria and models to be used if the interaction term of group and treatment is
854    applicable).  In addition, the statistical model should reflect the multigroup nature of the study.
855    For example, if subjects are dosed in two groups in a crossover BE study, the model should
856    reflect the fact that the periods for the first group are different from the periods for the second
857    group, i.e., the period effect should be nested within the group effect.
858
859    When there are multiple centers with very few subjects in some centers and sponsors want to
860    combine centers in the analysis, any rules for combination should be pre-specified in the protocol
861    or SAP and a sensitivity analysis is recommended.  More complicated scenarios may be
862    discussed with the appropriate CDER review division before submission.
863
864         **E.        Bioequivalence Statistics for Adhesion and Irritation Studies**
865
866    In terms of the statistical method used in irritation, sensitization or/and adhesion studies for
867    Transdermal and Topical Delivery Systems, refer to the Statistical Consideration section in the
868    draft guidance for industry Assessing *the Irritation and Sensitization Potential of Transdermal*
869    *and Topical Delivery Systems for ANDAs* (October 2018) and the Considerations for Statistical
870    Analysis section in the draft guidance for industry *Assessing Adhesion With Transdermal and*
871    *Topical Delivery Systems for ANDAs* (October 2018).[45]
872
873
874

---

[43] Fleiss, J.L., 1986, Analysis of Data from Multiclinic Trials, Controlled Clinical Trials, 7(4):267-275;
Permutt, T., 2003, Probability Models and Computational Models for ANOVA in Multicenter Clinical Trials,
Journal of Biopharmaceutical Statistics, 13(3):495-505; Tsiatis, A.A., M. Davidian, M. Zhang, and X. Lu, 2008,
Covariate Adjustment for Two-Sample Treatment Comparisons in Randomized Clinical Trials: A Principled Yet
Flexible Approach, Statistics in Medicine, 27(23):4658-4677.

[44] Alosh, M., K. Fritsch, M. Huque, K. Mahjoob, G. Pennello, M. Rothmann, E. Russek-Cohen, F. Smith, S. Wilson,
and L. Yue, 2015, Statistical Considerations on Subgroup Analysis in Clinical Trials, Statistics in Biopharmaceutical
Research, 7(4):286-303; Grizzle, J.E., 1965, The Two-Period Change-Over Design and Its Use in Clinical Trials,
Biometrics, 21(2):467-480.

[45] See also the draft guidance for industry *Assessment of Adhesion for Topical and Transdermal Systems Submitted in
New Drug Applications* (July 2021).  When final, these guidances will represent FDA's current thinking on these
topics.

875     **F.      Dose Scale for Bioequivalence Assessment**
876
877     In this method, the BE assessment is based on relative bioavailability of the test and reference
878     formulations at the site(s) of action.  The relative bioavailability, F, is the ratio of the doses of
879     test and reference formulations that produce an equivalent PD response.
880
881     Generally, the F is estimated by fitting an Emax model that describes the within-study dose-
882     response relationship.  Among available statistical methods for Emax model fitting, nonlinear
883     mixed effect (NLME) modeling is recommended, because the NLME modeling is capable of
884     characterizing between-subject variability and residual unexplained variability, and less sensitive
885     to aberrant observation and missing values.
886
887     For model fitting details, refer to the PSG on Orlistat oral capsule.[46]
888
889     To determine BE, the 90% confidence interval for F can be estimated by a bootstrap procedure.
890     Each bootstrap estimation includes the calculation of F by fitting the selected model to a sample
891     dose-response data set, which is generated by resampling with replacement.  To maintain the
892     correlation of observations within subject, resampling by subject (remaining observations from
893     all T and R treatment arms) is recommended rather than resampling by observations.  The
894     Agency has also recommended using Efron's bias corrected and accelerated method to compute a
895     90% confidence interval for F.[47]  Alternatively, the 90% confidence interval for F can be
896     estimated without a bootstrap procedure, directly from the point estimate of logF and its standard
897     error calculated using NLME modeling.
898
899     Given the complexity of dose scale analysis for comparative PD BE studies, applicants are
900     encouraged to contact the Agency early to discuss their proposed study designs and statistical
901     methods (e.g., alternative modeling approaches, impact of the missing data and the handling
902     strategy) via the controlled correspondence, pre-ANDA meeting, pre-IND meeting, or pre-NDA
903     meeting pathway.[48]
904
905     **G.      Bioequivalence Studies Using Multiple References**
906
907     In BE studies with more than two reference treatment arms (e.g., a three-period study including
908     two references, one from the European Union (EU) and another from the United States, or a
909     four-period study including test and reference in fed and fasted states), the BE determination
910     should be based on the comparison between the relevant test and reference products, using only
911     the data from those products.  The BE analysis for this comparison should be conducted
912     excluding the data from the non-relevant treatment(s) — for example, in a BE study with a T
913     product, an EU reference product, and a U.S. reference product, the comparison of the T product
914     to the U.S. reference product should be based on an analysis excluding the data from the EU
915     reference.  However, full data from the BE studies, including data comparing the T product that

---

[46] Draft PSG for industry on Orlistat oral capsule (August 2021).  When final, this guidance will represent FDA's
current thinking on this topic.
[47] Ibid.
[48] See footnotes 8, 9, and 10.

916    is the subject of the application with non-U.S. reference products, should be submitted in the
917    application for completeness.  The applicant may discuss the study design and statistical
918    approach with the appropriate CDER review division before study conduct.
919
920

921 **V.  APPENDICES**
922
923    A.    **Choice of Specific Replicated Crossover Designs**
924
925    Appendix A describes why FDA prefers replicated crossover designs with only two sequences,
926    and why the Agency recommends the specific designs described in section II.A.1.b of this
927    guidance.
928
929        1.    *Reasons Unrelated to Carryover Effects*
930
931    Each unique combination of sequence and period in a replicated crossover design can be called a
932    cell of the design.  For example, the two-sequence, four-period design recommended in section
933    II.A.1.b has eight cells.  The four-sequence, four-period design below has 16 cells.
934
935                                    Period
936
937                        1     2     3     4
938
939              1     **T**   **R**   **R**   **T**
940
941              2     **R**   **T**   **T**   **R**
942    Sequence
943              3     **T**   **T**   **R**   **R**
944
945              4     **R**   **R**   **T**   **T**
946
947    The total number of degrees-of-freedom attributable to comparisons among the cells is just the
948    number of cells minus one (unless there are cells with no observations).
949
950    The fixed effects that are usually included in the statistical analysis are sequence, period, and
951    treatment (i.e., formulation).  The number of degrees-of-freedom attributable to each fixed effect
952    is generally equal to the number of levels of the effect, minus one.  Thus, in the case of the two-
953    sequence, four-period design recommended in section V.A.1, there would be 2-1=1 degree-of-
954    freedom due to sequence, 4-1=3 degrees-of-freedom due to period, and 2-1=1 degree-of-freedom
955    due to treatment, for a total of 1+3+1=5 degrees-of-freedom due to the three fixed
956    effects.  Because these 5 degrees-of-freedom do not account for all 7 degrees-of-freedom
957    attributable to the eight cells of the design, the fixed-effects model is not saturated.  There could
958    be some controversy as to whether a fixed-effects model that accounts for more or all of the
959    degrees-of-freedom due to cells (i.e., a more saturated fixed-effects model) should be used. For
960    example, a sequence-by-period-by-treatment interaction effect might be included, which would
961    fully saturate the fixed-effects model.
962
963    If the replicated crossover design has only two sequences, use of only the three main effects
964    (sequence, period, and treatment) in the fixed-effects model or use of a more saturated model
965    makes little difference to the results of the analysis, provided there are no missing observations,

966   and the study is carried out in one group of subjects.  The least squares point estimate of $\mu_T - \mu_R$
967   will be the same for the main-effects model and for the saturated model.
968
969   If the replicated crossover design has more than two sequences, these advantages are no longer
970   present.  Main-effects models will generally produce different point estimates of $\mu_T - \mu_R$ than
971   saturated models (unless the number of subjects in each sequence is equal), and there is no well-
972   accepted basis for choosing between these different estimates (though $\mu_T - \mu_R$ from the
973   saturated model was determined to be appropriate for use in the reference-scaled average BE
974   assessment).  Thus, use of designs with only two sequences minimizes or avoids certain
975   ambiguities due to specific choices of fixed effects to be included in the statistical model.
976
977          2.        *Reasons Related to Carryover Effects*
978
979   One of the reasons to use the four-sequence, four-period design described above is that it is
980   thought to be optimal if carryover effects are included in the model.
981
982   Similarly, the two-sequence, three-period design is thought to be optimal among three-period
983   replicated crossover designs.  Both of these designs are strongly balanced for carryover effects,
984   meaning that each treatment is preceded by each other treatment and itself an equal number of
985   times.
986

                                     Period

                                1      2      3

                         1      **T**    **R**    **R**
              Sequence
                         2      **R**    **T**    **T**

995   With these designs, no efficiency is lost by including simple first-order carryover effects in the
996   statistical model.  However, if the possibility of carryover effects is to be considered in the
997   statistical analysis of BE studies, the possibility of direct-by-carryover interaction should also be
998   considered.  If direct-by-carryover interaction is present in the statistical model, these favored
999   designs are no longer optimal.  Indeed, the TRR/RTT design does not permit an unbiased within-
1000  subject estimate of $\mu_T - \mu_R$ in the presence of general direct-by-carryover interaction.
1001
1002  The issue of whether a purely main-effects model or a more saturated model should be specified,
1003  as described in the previous section, also is affected by possible carryover effects.  If carryover
1004  effects, including direct-by-carryover interaction, are included in the statistical model, these
1005  effects will be partially confounded with sequence-by-treatment interaction in four-sequence or
1006  six-sequence replicated crossover designs, but not in two-sequence designs.
1007
1008  In the case of the four-period and three-period designs recommended in section II.A.1.b, the
1009  estimate of $\mu_T - \mu_R$, adjusted for first-order carryover effects, including direct-by-carryover

1010 interaction, is as efficient or more efficient than for any other two-treatment replicated crossover
1011 designs.
1012
1013        *3.*       *Two-Period Replicated Crossover Designs*
1014
1015 For most drug products, two-period replicated crossover designs such as the Balaam design
1016 (which uses the sequences TR, RT, TT, and RR) should be avoided. However, the modified
1017 Balaam design (TR, RT, RR) may be useful for particular drug products (e.g., a long half-life
1018 drug for which a two-period study would be feasible, but a three-or-more-period study would
1019 not) when reference-scaled average BE is needed.
1020
1021     **B.**     **Rationale for <mark>Logarithmic Transformation</mark> of Pharmacokinetic Data**
1022
1023        *1.*       *Clinical Rationale*
1024
1025 The FDA Generic Drugs Advisory Committee recommended in 1991 that the primary comparison of
1026 interest in a BE study is the ratio, rather than the difference, between average PK parameter data from
1027 the T and R formulations. Using logarithmic transformation, the general linear statistical model
1028 employed in the analysis of BE data allows inferences about the difference between the two means on
1029 the log scale, which can then be retransformed into inferences about the ratio of the two averages
1030 (geometric means) on the original scale. Logarithmic transformation thus achieves a general
1031 comparison based on the ratio rather than the differences.
1032
1033        *2.*       *Pharmacokinetic Rationale*
1034
1035 Westlake observed that a multiplicative model is postulated for PK measures in BA/BE studies (i.e.,
1036 AUC and $C_{max}$, but not $T_{max}$) (Westlake 1973 and 1988).[49,50] Assuming that elimination of the drug is
1037 first order and only occurs from the central compartment, the following equation holds after an
1038 extravascular route of administration:
1039
1040       $AUC_{0-\infty} = F*D/CL$
1041
1042           $= F*D/(V*Ke)$
1043
1044 where F is the fraction absorbed, D is the administered dose, and F*D is the amount of drug absorbed.
1045 CL is the clearance of a given subject that is the product of the apparent volume of distribution (V) and
1046 the elimination rate constant (Ke). The use of AUC as a measure of the amount of drug absorbed
1047 involves a multiplicative term (CL) that might be regarded as a function of the subject. For this reason,

---

[49] Westlake, W. J., 1973, The Design and Analysis of Comparative Blood-Level Trials, J. Swarbick, editor, Current Concepts in the Pharmaceutical Sciences, Dosage Form Design and Bioavailability, Philadelphia: Lea and Febiger, 149-179.

[50] Westlake, W. J., 1988, Bioavailability and Bioequivalence of Pharmaceutical Formulations, Biopharmaceutical Statistics for Drug Development, 329-352.

1048 Westlake contends that the subject effect is not additive if the data are analyzed on the original scale of
1049 measurement.
1050
1051 Logarithmic transformation of the AUC data will bring the CL (i.e., V*Ke) term into the following
1052 equation in an additive fashion:
1053
1054 $\qquad$ $\ln AUC_{0-\infty} = \ln F + \ln D - \ln V - \ln Ke$
1055
1056 Similar arguments were given for $C_{max}$. The following equation applies for a drug exhibiting one
1057 compartmental characteristic:
1058
1059 $\qquad$ $C_{max} = (F*D/V) * \exp(-Ke*T_{max})$
1060
1061 where again F, D and V are introduced into the model in a multiplicative manner. However, after
1062 logarithmic transformation, the equation becomes:
1063
1064 $\qquad$ $\ln C_{max} = \ln F + \ln D - \ln V - Ke*T_{max}$
1065
1066 Thus, log transformation of the $C_{max}$ data also results in the additive treatment of the V term.
1067
1068 **C.** ==**SAS Program Statements for Average Bioequivalence Analysis** of Replicated==
1069 **Crossover Studies**
1070
1071 The following illustrates an example of program statements to run the unscaled average BE
1072 analysis using PROC MIXED in SAS version 9, with SEQ, SUBJ, PER, and TRT identifying
1073 sequence, subject, period, and treatment variables, respectively, and Y denoting the response
1074 measure (e.g., log (AUC), log ($C_{max}$)) being analyzed:
1075

```
1076          PROC MIXED;
1077          CLASSES SEQ SUBJ PER TRT;
1078          MODEL Y = SEQ PER TRT/ DDFM=SATTERTH;
1079          RANDOM TRT/TYPE=FA0(2) SUB=SUBJ G;
1080          REPEATED/GRP=TRT SUB=SUBJ;
1081          ESTIMATE 'T vs. R' TRT 1 -1/CL ALPHA=0.1;
```

1082
1083 The *Estimate* statement assumes that the code for the test formulation precedes the code for the
1084 reference formulation in sort order (this would be the case, for example, if T were coded as 1 and
1085 R were coded as 2). If the R code precedes the T code in sort order, the coefficients in the
1086 Estimate statement would be changed to -1 1.
1087
1088 In the *Random* statement, TYPE=FA0(2) could possibly be replaced by TYPE=CSH or UNR.
1089
1090 In the *Model* statement, DDFM=SATTERTH could possibly be replaced by DDFM=KR2.
1091 However, the detailed model specification should be pre-specified in the protocol or SAP and
1092 data driven post hoc selection of the model is not allowed.

1093

1094    Additions and modifications to these statements can be made if the study is carried out in more
1095    than one group of subjects or other complicated scenarios.  Alternative software could also be
1096    used if same results are generated as in PROC MIXED in SAS.