# Modelling and simulation SIG update: Best Practice document

Michael O'Kelly, member of Modelling and Simulation
SIG

# Summary

- Work on Best Practice since 2011.
  - > European Federation of Pharmaceutical Industries and Associations (EFPIA) MID3 working group
  - > EFSPI Special Interest Group (SIG)
- PSI Board agrees to adopt SIG Best Practice proposal pending publication of the proposal in *Pharmaceutical Statistics.*
- EFSPI Modelling and Simulation SIG is working with MID3 to gain agreement among practitioners for Best Practice in modelling and simulation globally.

**QUINTILES**

# EMA: "Best practice" depends on importance of project



EUROPEAN MEDICINES AGENCY
SCIENCE MEDICINES HEALTH

EMA-EFPIA Modelling
and Simulation Workshop

Good practices and next steps

Robert Hemmings, EMA

## M&S good practices

- Different standards for different exercises (L,M,H)
- Standard should be high!
  - Assumptions (not only mathematical)
  - Model building rationale
  - Model testing
  - Inference
  - Sensitivity analyses / Challenge assumptions
  - Reporting
- Detail of regulatory response might be vary according to impact

# EMA: "Best practice" depends on importance of project



EUROPEAN MEDICINES AGENCY
SCIENCE MEDICINES HEALTH

EMA-EFPIA Modelling and Simulation Workshop

Good practices and next steps

Robert Hemmings, EMA

## M&S good practices

- Different standards for different exercises (L,M,H)
- Standard should be high!
  - Assumptions (not only mathematical)
  - Model building rationale
  - Model testing
  - Inference
  - Sensitivity analyses / Challenge assumptions
  - Reporting
- Detail of regulatory response might be vary according to impact

# MID3 paper on Good Practices

- MID3 paper "Good Practices in Model-Informed Drug Discovery and Development (MID3): practice, application and documentation.
    › Wide industry representation: Pfizer, Bayer, Roche, AstraZeneca, GSK, J&J, Merck, BI, Novartis, Novo Nordisk.

- Paper plus supplementary spreadsheet of 103 MID3 example applications, published January 2016.

# MID3 headings for Good Practice

| Components of Good Practice plans | | |
| --- | --- | --- |
| Analysis plan | Simulation plan | Report |
| <ul><li>Introduction</li><li>Objectives</li><li>Data plan</li><li>Data exploration</li><li>Methods<ul><li>Model building</li><li>Selection+evaluation</li><li>Qualification</li></ul></li><li>Assumptions</li><li>Results</li></ul> | <ul><li>Introduction</li><li>Objectives</li><li>Additional data</li><li>Methods<ul><li>Identify model</li><li>Limitations</li><li>Qualification</li></ul></li><li>Assumptions</li><li>Results</li></ul> | <ul><li>Synopsis</li><li>Introduction</li><li>Objectives</li><li>Data</li><li>Methods<ul><li>Identify model</li><li>Limitations</li><li>Qualification</li></ul></li><li>Assumptions</li><li>Results</li><li>Applications/simulations</li><li>Discussion</li><li>Conclusion</li><li>Appendices</li></ul> |

# MID3 headings for Good Practice

| Components of Good Practice plans | | |
|---|---|---|
| Analysis plan | Simulation plan | Report |
| • Introduction<br>• Objectives<br>• Data plan<br>• Data exploration<br>• Methods<br>   • Model building<br>   • Selection+evaluation<br>   • Qualification<br>• Assumptions<br>• Results | • Introduction<br>• Objectives<br>• Additional data<br><br>• Methods<br>   • Identify model<br>   • Limitations<br>   • Qualification<br>• Assumptions<br>• Results | • Synopsis<br>• Introduction<br>• Objectives<br>• Data<br><br>• Methods<br>   • Identify model<br>   • Limitations<br>   • Qualification<br>• Assumptions<br>• Results<br>• Applications/simulations<br>• Discussion<br>• Conclusion<br>• Appendices |

# MID3 headings for Good Practice – includes recommendations for each heading

| Components of Good Practice plans | | |
|---|---|---|
| Analysis plan | Simulation plan | Report |
| • Introduction<br>• Objectives<br>• Data plan<br>• Data exploration<br>• Methods<br>   • Model building<br>   • Selection+evaluation<br>   • Qualification<br>• Assumptions<br>• Results | • Introduction<br>• Objectives<br>• Additional data<br><br>• Methods<br>   • Identify model<br>   • Limitations<br>   • Qualification<br>• Assumptions<br>• Results | • Synopsis<br>• Introduction<br>• Objectives<br>• Data<br><br>• Methods<br>   • Identify model<br>   • Limitations<br>   • Qualification<br><br>• Assumptions<br>• Results<br>• Applications/simulations<br>• Discussion<br>• Conclusion<br>• Appendices |

QUINTILES®

# EFSPI proposed Best Practice document

- Authored by volunteers from the SIG
  - › SIG members from variety of pharmaceutical companies.
  - › All authors of the Best Practice document from contract research organizations.
  - › Agreed to be adopted by Board of PSI.

- PSI Board agrees to adopt SIG Best Practice proposal pending its publication in *Pharmaceutical Statistics*.
  - › Best Practice paper submitted to *Pharmaceutical Statistics* November 2016.
  - › Paper is currently in review (second round of referees' comments).

# EFSPI SIG Best Practice recommendations – during, after, and next time around…

- Template for specification
  - › describe listed key elements or justify why not
  - › justify level of detail of the pre-specification.

- Quality control – level of QC should be appropriate -
  - › from simple review of specification (low-impact project)
  - › to independent programming of project (some high-impact projects)

- Presentation of results
  - › may vary depending on audience – plan in advance outputs for each audience
  - › Statistical clarity: use of confidence intervals; operating characteristics; measure of stochastic variability in simulations….

- Changes to specification
  - › Specification should be auditable, e.g.,
    - » revision history
    - » formal amendment (as in protocol amendment)
    - » include old versions as appendices

QUINTILES

# Principle: do what is necessary for Best Practice, but not more

- SIG document allows the flexibility necessary for Best Practice in this area where the regulatory and scientific importance of the projects varies widely.

QUINTILES

# Example best-practice specification for low-impact work

## Planning and Reporting for Projects that Involve Modelling and Simulation
### Best Practice Document

**Appendix B: example specification with a low level of detail**

**Using simulated data to verify an estimate of probability of success**

Specification of simulations

### B.1 Introduction
Given five treatment development programs with known probability of success, it is desired to know the probability of zero successes and of four and five successes. These probabilities have been calculated analytically. It is requested that a simulation be run to verify that the analysis is correct.

Since this is a one-off query on whose evidence alone no decision will be made, this is judged a project of low importance. Therefore the clinical background is not described; nor are metrics and criteria for decisions appropriate.

### B.2 Simulation and analysis/design
As noted, this project is of low importance and no decision will be made by it alone. Therefore the description of the elements of the simulation and analysis will not be detailed and some elements are not applicable.

#### B.2.1 Scenarios assumed and assumptions made
Probabilities of 0.1, 0.2, 0.2, 0.05 and 0.4 were given for programs 1-5, respectively. Since the objective was simple verification of an existing calculation, no justification is given here of these probabilities. Since the question answered is theoretical, just one given scenario is used.

#### B.2.1.1 Sensitivity analyses
This project is not required to assess assumptions, so sensitivity to assumptions is not planned to be analysed via sensitivity analyses

#### B.2.2 Data sets generated
Temporary sets of binary outcomes will be generated. Data will not be bootstrapped because a simple verification is sufficient. Three million binary outcomes are simulated for each program.

Page 1

#### B.2.3 Statistical analysis
The number of instances of zero, four and five successes was calculated for each of the 3 million simulations, and the probability of zero, four and five successes in a simulated instance was calculated and plotted.

#### B.2.4 Operating characteristics
Given that this modelling and simulation task is to be a sanity check, the number of simulations required to achieve a given accuracy with 5% confidence will be approximated. The probability of five successes is small (<1/100) so a precision of 0.001 is desired. Using the formula of Burton *et al.* (2006), with alpha=0.05, and approximating the variance of the probabilities as $5*p(1-p)$ where p==0.2, 3 million simulations will provide precision of approximately 0.001.

#### B.2.5 Logistics
The R language package mvtBinaryEP will be used to simulate the binary outcomes. The package allows for correlations between the outcomes, but this was not required for the primary objective. R version 3.0.1 will be used. See Appendix for the R code used. The seed used was 1.

### B.3 Quality control
Given that this modelling and simulation task is of low importance and will not of itself lead to a decision, the specification will be submitted to the requestor of the calculation, but no further QC of the production of results is planned. The output will be checked against the requestor's calculations.

### B.4 Presentation of results
A table will be presented of the probability of zero, four and five successes among the five programs, calculated as the proportion of instances of zero, four and five successes in 3 million simulations of the five programs. The number of instances will also be plotted in a histogram with one stack for each level of successes, a stack for zero successes, 1 success, and so on up to five successes. These outputs are judged sufficient to act as a check of the analytic estimates, which is the objective of this project.

The precision of the result (standard error) will be presented in a footnote to the plot. Given the inclusion of precision, no confidence intervals will be presented. Given the theoretical nature of the problem and the corresponding simplicity of the simulation, no bias is to be expected in the simulation-based estimates.

The results of the modelling and simulation will not be stored. The R code will be stored in [location]. A note of the contents of the table output will be included as a comment in the R code.

Page 2

QUINTILES

12

# Example best-practice specification, high-impact work

# Best Practice for modelling and simulation, the work of the two groups, MID3 and EFSPI

- Agreement all aspects of the process

- Some differences in emphasis

- The two groups are working together to promote good practice

- The two groups participated at session on Best Practice at 2016 annual PSI conference, Berlin.

- The two groups will participate along with FDA presenters in 2016 September ASA-Biopharmaceutical Section workshop in Washington.

# Summary

EFSPI Best Practice document can be used as a tool or template to implement Best Practice as described by MID3 and/or EFSPI.

MID3 and EFSPI SIG share vision of good practice harmonised across the uses of modelling and simulation.

# Questions?

# Back-up slides

QUINTILES®

# Best Practice in modelling and simulation

## MID3

- When to use simulation

- Key elements for a good plan

- Quality control

- Iterative nature of the MID3 process

## EFSPI

- When to use simulation

- Key elements for a good plan

- Quality control

- Iterative nature of modelling and simulation

**QUINTILES**

# Best Practice in modelling and simulation

## MID3

- Agreed across 10 companies.

- Emphasis on integrating MID3 into the general pharmaceutical development process

- Three planning documents.

- Lists key recommended elements.

- Report: specifies sections, with potentially different audiences.

## EFSPI

- Authored by SIG, to be adopted by PSI.

- Emphasis on providing a tool for Best Practice for the working statistician

- One specification for a project.

- Emphasis on flexibility – specification should include key elements or justify their absence.

**QUINTILES**

# Best Practice in modelling and simulation

## MID3

- Emphasis on hypothesis **generation** rather than hypothesis **testing**.

- Tends not to go into detail on technical requirements.

## EFSPI

- Allows for possibility of hypothesis testing.

- Suggests including "less likely" scenarios in simulations.

- Considers technical detail, e.g., operating characteristics; use of confidence intervals; measure of stochastic variability in simulations; randomisation seed; software version.

# Best Practice in modelling and simulation

## MID3

- Emphasis on hypothesis **generation** rather than hypothesis **testing**.

- Tends not to go into detail on technical requirements.

## EFSPI

- Allows for possibility of hypothesis testing.

- Suggests including "less likely" scenarios in simulations.

- Considers technical detail, e.g., operating characteristics; use of confidence intervals; measure of stochastic variability in simulations; randomisation seed; software version.

EFSPI Best Practice document can be used as a tool or template to implement Best Practice as described by MID3 and/or EFSPI.
MID3 and EFSPI SIG share vision of good practice harmonised across the uses of modelling and simulation.

QUINTILES