# Predicting the cumulative chance of live birth over multiple complete cycles of in vitro fertilisation: an external validation study

1    **TITLE PAGE**

2

3    **Title: Predicting the cumulative chance of live birth over multiple complete cycles of in vitro**

4    **fertilisation: an external validation study**

5

6    **Running title:** Validation of an IVF model predicting live birth

7

8    **Authors:**

9    J.A. Leijdekkers[1,*], M.J.C. Eijkemans[2], T.C. van Tilborg[1], S.C. Oudshoorn[1], D.J. McLernon[3], S.

10   Bhattacharya[4], B.W.J. Mol[5], F.J.M. Broekmans[1], H.L. Torrance[1], on behalf of the OPTIMIST group.

11

12

13   [1]Department of Reproductive Medicine and Gynaecology, University Medical Centre Utrecht, Utrecht

14   University, PO box 85500, 3508 GA Utrecht, The Netherlands. [2]Julius Centre for Health Sciences and

15   Primary Care, University Medical Centre Utrecht, Utrecht University, PO box 85500, 3508GA Utrecht,

16   The Netherlands. [3]Institute of Applied Health Sciences, Medical Statistics Team, University of Aberdeen,

17   Aberdeen AB25 2ZD,  UK. [4]School of Medicine, College of Biomedical and Life Sciences, Cardiff

18   University School of Medicine, Cardiff CF14 4XN, UK [5]Department of Obstetrics and Gynaecology,

19   Monash University, VIC 3800 Clayton, Australia.

20

21

22   *Correspondence address:** J.A. Leijdekkers, Department of Reproductive Medicine and Gynaecology,

23   University Medical Centre Utrecht, Utrecht University, PO box 85500, 3508 GA Utrecht, The

24   Netherlands. E-mail: j.a.leijdekkers@umcutrecht.nl

25

26    **ABSTRACT**

27

28    *Study question*

29    Are the published pre-treatment and post-treatment McLernon models, predicting cumulative live birth

30    rates (LBR) over multiple complete IVF cycles, valid in a different context?

31

32    *Summary answer*

33    With minor recalibration of the pre-treatment model, both McLernon models accurately predict

34    cumulative LBR in a different geographical context and a more recent time period.

35

36    *What is known already*

37    Previous IVF prediction models have estimated the chance of a live birth after a single fresh embryo

38    transfer, thereby excluding the important contribution of embryo cryopreservation and subsequent IVF

39    cycles to cumulative LBR. In contrast, the recently developed McLernon models predict the cumulative

40    chance of a live birth over multiple complete IVF cycles at two certain time points: a) before initiating

41    treatment using baseline characteristics (pre-treatment model) and b) after the first IVF cycle adding

42    treatment related information to update predictions (post-treatment model). Before implementation of

43    these models in clinical practice, their predictive performance needs to be validated in an independent

44    cohort.

45

46    *Study design, size, duration*

47    External validation study in an independent prospective cohort of 1515 Dutch women who participated in

48    the OPTIMIST study (NTR2657) and underwent their first IVF treatment between 2011 and 2014.

49    Participants underwent a total of 2881 complete treatment cycles, with a complete cycle defined as all

50    fresh and frozen thawed embryo transfers resulting from one episode of ovarian stimulation. The follow

51    up duration was 18 months after inclusion, and the primary outcome was ongoing pregnancy leading to

52    live birth.

53

54    *Participants/materials, setting, methods*

55    Model performance was externally validated up to three complete treatment cycles, using the linear

56    predictor as described by McLernon et al. to calculate the probability of live birth. Discrimination was

57    expressed by the c-statistic and calibration was depicted graphically in a calibration plot. In contrast to the

58    original model development cohort, anti-Müllerian hormone (AMH), antral follicle count (AFC) and body

59    weight were available in the OPTIMIST cohort, and evaluated as potential additional predictors for model

60    improvement.

61

62    *Main results and the role of chance*

63    Applying the McLernon models to the OPTIMIST cohort, the c-statistic of the *pre-treatment model* was

64    0.62 (95% confidence interval (CI) 0.59-0.64) and of the *post-treatment model* 0.71 (95% CI 0.69-0.74).

65    The calibration plot of the *pre-treatment model* indicated slight overestimation of the cumulative LBR. To

66    improve calibration, the *pre-treatment model* was recalibrated by subtracting 0.35 from the intercept. The

67    *post-treatment model* calibration plot revealed accurate cumulative LBR predictions. After addition of

68    AMH, AFC and body weight to the McLernon models, the c-statistic of the *updated pre-treatment model*

69    improved slightly to 0.66 (95% CI 0.64-0.68), and of the *updated post-treatment model* remained at the

70    previous level of 0.71 (95% CI 0.69-0.73).

71    Using the *recalibrated pre-treatment model*, a woman aged 30 years with two years of primary infertility

72    who starts ICSI treatment for male factor infertility has a chance of 40% of a live birth from the first

73    complete cycle, increasing to 72% over three complete cycles. If this woman weighs 70 kilograms, has an

74    AMH of 1.5 ng/mL and an AFC of 10 measured at the beginning of her treatment, the *updated pre-*

75    *treatment model* revises the estimated chance of a live birth to 30% in the first complete cycle and 59%

76    over three complete cycles. If this woman then has 5 retrieved oocytes, no embryos cryopreserved and a

77    single fresh cleavage stage embryo transfer in her first ICSI cycle, the *post-treatment model* estimates the

78    chances of a live birth at 28% and 58%, respectively.

79

80    *Limitations, reasons for caution*

81    Two randomised controlled trials (RCT) evaluating the effectiveness of gonadotropin dose

82    individualisation on basis of the AFC were nested within the OPTIMIST study. The strict dosing

83    regimens, the RCT in- and exclusion criteria and the limited follow up time of 18 months might have

84    influenced model performance in this independent cohort. Also, consistent with the original model

85    development study, external validation was performed using the optimistic assumption that the

86    cumulative LBR in couples who discontinue treatment without a live birth would have been equal to that

87    of those who continue treatment.

88

89    *Wider implications of the findings*

90    After national recalibration to account for geographical differences in IVF/ICSI treatment, the McLernon

91    prediction models can be introduced as new counselling tools in clinical practice to inform patients and to

92    complement clinical reasoning. These models are the first to offer an objective and personalised estimate

93    of the cumulative probability of live birth over multiple complete IVF cycles.

94

95    *Study funding/competing interest(s):*

96    No external funds were obtained for this study. M.J.C.E., D.J.M. and S.B. have nothing to disclose. J.A.L,

97    S.C.O, T.C.v.T. and H.LT. received an unrestricted personal grant from Merck BV. B.W.M. is supported

98    by a NHMRC Practitioner Fellowship (GNT1082548)  and  reports consultancy for ObsEva, Merck and

99    Guerbet. F.J.M.B. receives monetary compensation as a member of the external advisory board for Merck

100   BV (the Netherlands) and Ferring pharmaceutics BV (the Netherlands), for consultancy work for Gedeon

101   Richter (Belgium) and  Roche Diagnostics on automated AMH assay development, and for a research

102   cooperation with Ansh Labs (USA).

103

104     *Trial registration number*

105     Not applicable

106

107     **KEYWORDS**

108     Prediction model, external validation, live birth, IVF/ICSI, infertility, cumulative live birth, personalised,

109     counselling, prognostic research

110   **Introduction**

111   Infertility is defined as the failure to conceive within 12 months of regular unprotected intercourse, and

112   affects approximately one in six couples (Oakley *et al.*, 2008; Zegers-Hochschild *et al.*, 2017). The

113   majority of infertile couples seek fertility care, and many of those with prolonged unresolved infertility

114   will be treated with ART regardless of cause (Boivin *et al.*, 2007; Datta *et al.*, 2016). IVF and ICSI are

115   both widely used techniques for couples with infertility. Globally more than 1.6 million annual cycles of

116   IVF/ICSI are performed and while success rates have increased over time (Dyer *et al.*, 2016; McLernon *et*

117   *al.*, 2016), this treatment is still not effective for all infertile couples, with live birth rates (LBR) at around

118   25-30% per treatment cycle (Malizia *et al.*, 2009; McLernon *et al.*, 2016; de Neubourg *et al.*, 2016). Since

119   IVF/ICSI is expensive and carries several risks, the probability of a live born child should be weighed

120   against the risks and costs of this treatment.

121   Several prognostic models have been developed to objectively estimate the probability of a live birth after

122   IVF/ICSI treatment (Leushuis *et al.*, 2009; van Loendersloot *et al.*, 2014). It is known that prediction

123   models often perform optimistically in their development sample, even after correction by internal

124   validation. This is caused by overfitting, which occurs when the model corresponds too closely to the

125   development data due to the inclusion of too many predictors (Moons, Kengne, Woodward, *et al.*, 2012).

126   External validation in an independent cohort of women is thus essential to examine the performance and

127   generalisability of the prediction model (Altman *et al.*, 2009; Harrell *et al.*, 1996). Unfortunately, most of

128   the currently available models that predict the chance of a live birth after IVF/ICSI treatment have never

129   been externally validated (Leushuis *et al.*, 2009; van Loendersloot *et al.*, 2014). Also, the majority of

130   these models predict the probability of a live birth after a single fresh embryo transfer, excluding the

131   important contribution of embryo cryopreservation and subsequent treatment cycles to LBR. This limits

132   their potential as counselling tools for couples and clinicians, especially considering the increased use and

133   improved techniques of embryo cryopreservation and frozen thawed embryo transfer cycles in recent

134   years (Wong *et al.*, 2014).

135    Three of the largest model development studies for prediction of live birth after IVF and/or ICSI

136    treatment used data from the Human Fertilisation and Embryology Authority (HFEA) database in the UK

137    (McLernon *et al.*, 2016; Nelson and Lawlor, 2011; Templeton *et al.*, 1996). Treatment and outcome data

138    from all licenced fertility clinics within the UK have been recorded in this database since 1992. The two

139    models developed by Templeton et al. and Nelson et al. were both externally validated, and their

140    predictive performance was compared to one another in several studies (Arvis *et al.*, 2012; van

141    Loendersloot *et al.*, 2011; Smeenk *et al.*, 2000; Smith *et al.*, 2015; te Velde *et al.*, 2014). Although these

142    models have been recommended in previous studies and used internationally to predict live birth after

143    IVF and ICSI (Leushuis *et al.*, 2009; Smith *et al.*, 2015; te Velde *et al.*, 2014), neither model predicts

144    cumulative LBR over multiple IVF/ICSI treatment cycles including frozen thawed embryo transfer

145    cycles.

146    Recently, a new model was developed by McLernon et al. using the HFEA database (McLernon *et al.*,

147    2016). This model is the first to provide an individualised estimate of the cumulative chance of a live

148    birth over multiple complete cycles of IVF/ICSI, with a complete cycle defined as all fresh and frozen

149    thawed embryo transfers resulting from one episode of ovarian stimulation. For model development, data

150    from 113 873 women and 184 269 complete cycles between 1999 and 2009 were used. Internal validation

151    of the model showed promising results, however evaluation of the predictive performance of the model in

152    a different geographical context using more contemporary data has yet to be performed. Additionally, a

153    number of potential key predictors, such as measures for ovarian reserve and female body weight, were

154    unavailable in the HFEA database and could not be included in the original model (McLernon *et al.*,

155    2016).

156    The main objective of the current study was therefore to perform geographical and temporal validation of

157    the new HFEA model by using recent data from a different country. We also wanted to determine whether

158    inclusion of additional parameters, such as female body weight and ovarian reserve test results i.e. antral

159    follicle count (AFC) and anti-Müllerian hormone (AMH), could improve the predictive performance of

160    the model.

161    **Materials and methods**

162    *Data sources*

163    External validation was performed on data from the OPTIMIST study (van Tilborg, Oudshoorn, *et al.*,

164    2017). This multicentre prospective cohort study included 1515 women from 25 infertility centres in the

165    Netherlands between May 2011 and May 2014. Participants were younger than 44 years of age, had

166    regular menstrual cycles and no significant uterine or ovarian abnormalities on transvaginal ultrasound.

167    Women with polycystic ovarian syndrome, metabolic or endocrine abnormalities or undergoing oocyte

168    donation were excluded. All participants were included before their first IVF/ICSI cycle, or the first cycle

169    after a previous live birth. The primary outcome was ongoing pregnancy, achieved within 18 months of

170    follow up, and resulting in live birth. Ethical approval for the OPTIMIST study was obtained from the

171    Institutional Review Board of the University Medical Centre Utrecht (MEC 10-273), and all participants

172    provided written informed consent. A more detailed description of study procedures and results were

173    reported previously (Oudshoorn *et al.*, 2017; van Tilborg *et al.*, 2012; van Tilborg, Oudshoorn, *et al.*,

174    2017; van Tilborg, Torrance, *et al.*, 2017).

175    *McLernon model*

176    The McLernon model consists of two clinical prediction models to estimate the individualised cumulative

177    chance of a live birth over a maximum of six complete treatment cycles. **Before initiating treatment**, the

178    *pre-treatment model* predicts the probability of a live birth from both fresh and frozen thawed embryo

179    transfers based on couple characteristics and the use of IVF or ICSI. Included predictors are: female age

180    (years), duration of infertility (years), previous pregnancy, causes of infertility (tubal factor, anovulation,

181    male factor, unexplained infertility), type of treatment (IVF or ICSI) and treatment year (see

182    Supplementary Text 1).

183    **After the first fresh treatment cycle,** treatment specific characteristics from this cycle are added in the

184    *post-treatment model* to update the predicted probability. Added predictors are: number of oocytes,

185    cryopreservation of embryos, and the number and stage of embryos at the first fresh embryo transfer

186    (single, double or triple embryo transfer; blastocyst or cleavage stage). All causes of infertility are

187    excluded as predictors in the post-treatment model, except for tubal factor (see Supplementary Text 2).

188    For women with zero oocytes collected in the first cycle, a separate post-treatment model is available.

189    To predict the probability of a live birth in the *i*th cycle, assuming no live birth occurred in the previous

190    cycle(s), complete cycle number is included in both models as a discrete time variable. A complete cycle

191    includes all fresh and frozen thawed embryo transfers resulting from one episode of ovarian stimulation.

192    With the predicted probability of a live birth per subsequent complete cycle, the cumulative probability of

193    a live birth can be calculated up to six complete cycles (see Supplementary Text 1 and 2).

194    *Statistical analysis*

195    Nine predictor variables had missing values (Table I). The proportion of missing values was low (<

196    2.5%), except for AMH (11.2%). During the OPTIMIST study, blood sampling was performed on the day

197    of randomisation. Logistic issues prevented blood sampling in some cases, thus compromising the ability

198    to undertake post-hoc measurements of AMH in the total population. As the reasons for missing values

199    were considered to be unrelated to the AMH value itself or the measurement, these were defined as

200    missing (completely) at random.

201    Multiple imputation was applied for predictors with missing values in the OPTIMIST database (Sterne *et*

202    *al.*, 2009). In this process 10 imputed datasets were created using a multivariate imputation by chained

203    equations (MICE) algorithm (van Buuren and Groothuis-Oudshoorn, 2011). Predicted probabilities for a

204    live birth were calculated on each imputed dataset, using the predictors and parameter-estimates of both

205    the pre-treatment model as well as the post-treatment model as described by McLernon et al 2016

206    (McLernon *et al.*, 2016). In accordance with the original models, the variables female age, treatment year

207    and number of oocytes were treated with restricted cubic splines in the validation process. The separate

208    post-treatment model for women with zero oocytes collected in the first treatment cycle was not validated

209    in this study, as the number of women for this analysis was too low in the OPTIMIST database.

210    Cumulative probabilities were calculated up to three complete IVF/ICSI cycles, as most couples in the

211    Netherlands only have three treatment cycles due to the current reimbursement policy. Also, the

212    OPTIMIST follow up period was 18 months, reducing the number of women with more than three

213    treatment cycles. The validation process was performed ten times on each of the imputed datasets and

214    separate results were pooled using Rubin's rules (Rubin, 2004).

215    The predictive performance of the McLernon models was evaluated in terms of discrimination and

216    calibration. Discrimination quantifies the ability of a model to correctly differentiate between subjects

217    with an event and subjects without an event (Moons, Kengne, Woodward, *et al.*, 2012). In the context of

218    fertility treatment, it is the ability of the models to distinguish between women with a live birth and

219    women without a live birth after IVF/ICSI treatment. It is expressed by the c-statistic or the area under the

220    receiver operating curve (AUROC), which ranges between 0.5 and 1. A c-statistic of 1 indicates perfect

221    discrimination, whereas a c-statistic of 0.5 represents a model with no discrimination at all. In this study,

222    the c-statistic (and 95% CI) was calculated using the method suggested by Harrell et al. (Harrell *et al.*,

223    1996).

224    Calibration describes the degree of agreement between predicted probabilities and observed outcomes

225    (Moons, Kengne, Woodward, *et al.*, 2012), in this context the predicted probability of a live birth and the

226    observed LBR. Calibration can be assessed graphically by forming subgroups of patients determined by

227    ranges of predicted probabilities, and then plotting the observed proportion of events against the mean

228    predicted probability within these subgroups. When perfect calibration is present, the plot shows a

229    diagonal line with a slope of one and an intercept of zero. In the current study, five equal subgroups of

230    patients were formed. This was based on the sample size of the OPTIMIST cohort and the related

231    precision of the point estimates in the calibration plot. Within these subgroups, the Kaplan Meier

232    estimates of the observed cumulative LBR over three complete treatment cycles were plotted against the

233    mean predicted probability of cumulative live birth. A smoothed line was then added in this plot using the

234    proportional hazard regression approach described by Harrell et al (Harrell *et al.*, 1996). In addition to

235    this, a systematic difference in the predicted and observed LBR was assessed by using calibration-in-the-

236    large (Steyerberg, 2009), and the intercept of the prediction models was adjusted in case a systematic

237    over- or underestimation was present.

238    *Updating the models*

239    Following the external validation of the models, the additional value of updating the McLernon models

240    with pre-specified new biomarkers was evaluated. AMH (ng/mL), AFC (2-10 mm) and body weight (kg)

241    were added to the pre-treatment and post-treatment model in a multivariable logistic regression analysis,

242    in which the linear predictor of the McLernon model was entered as a fixed variable. The final model was

243    established using a manual backward selection process. Predictors were eliminated from the model

244    according to the Akaike Information Criterion (AIC) (Akaike, 1974).

245    The predictive performance of the new updated models was evaluated by calculating the c-statistic (and

246    95% CI). To assess for overfitting, internal validation was performed by bootstrapping (Steyerberg,

247    2009). Two hundred bootstrap samples, all of which were of the same size as the original validation

248    sample, were created by random sampling with replacement (Harrell, 2001; Steyerberg, 2009). In each

249    bootstrap sample, a new model was fitted with the same predictors as the updated models. The c-statistic

250    was calculated for each of the 200 sample derived models, in both the bootstrap sample as well as the

251    original validation cohort. The difference between these two c-statistics was calculated for each of the 200

252    sample derived models, and averaged to give the optimism estimate. This was subtracted from the

253    original c-statistic to obtain the optimism corrected c-statistic for the updated models.

254    All statistical analyses were performed using R for Windows (version 3.3.2; R Foundation for Statistical

255    Computing, Vienna, Austria).

256    **Results**

257    Of the 1515 women included in the OPTIMIST study, four were excluded in the current study as they

258    never started IVF/ICSI treatment. A total of 2881 IVF/ICSI cycles were performed over a period of 18

259    months of follow up. Table I shows the patient and first cycle treatment characteristics of the OPTIMIST

260    cohort (validation sample) and the HFEA cohort (development sample). Women included in the

261    validation sample were about the same age as women in the development sample, but had a shorter

262    average duration of infertility. The causes of infertility showed a similar distribution across both samples,

263    with the exception of anovulation which rendered women ineligible for the OPTIMIST study. The

264    treatment characteristics showed that embryo cryopreservation was more frequently performed after the

265    first IVF/ICSI cycle in the validation sample and that these women most often had a cleavage stage single

266    embryo transfer in the first fresh cycle, whereas women in the development sample most often had a

267    cleavage stage double embryo transfer. No formal assessment was performed for the differences and

268    similarities between the cohorts, as a description rather than a p-value is considered to be useful for

269    interpretation of the models' performance in this external validation study.

270    The flowchart in Figure 1 shows the number of women in the OPTIMIST and HFEA cohorts who started

271    a treatment cycle, had a live birth or discontinued treatment without having a live birth. The LBR per

272    cycle was similar in both cohorts for the first, second and fourth treatment cycle. In the third cycle the

273    LBR was slightly higher in the OPTIMIST cohort compared to the HFEA cohort. As few women in the

274    OPTIMIST cohort received a fifth or sixth cycle, LBR in these cycles could not be compared. The

275    proportion of women without a live birth that continued treatment was higher after the first and second

276    cycle in the OPTIMIST cohort as compared to the HFEA cohort. After the third cycle, the proportion

277    continuing treatment in the OPTIMIST cohort decreased, while it remained constant in the HFEA cohort.

278    At the end of follow up, 52% of the women in the OPTIMIST study had a treatment related live birth. The

279    overall LBR of the HFEA cohort was 43% over six complete IVF/ICSI cycles.

280    As mentioned previously, external validation of the McLernon models was performed up to three

281    complete treatment cycles, and therefore the fourth, fifth and sixth complete treatment cycle in the

282    OPTIMIST dataset (n=102 complete treatment cycles, n= 15 live births) were excluded from further

283    analysis. Also, for the post-treatment model validation, women with zero oocytes collected in the first

284    treatment cycle were excluded (n= 226 women, n = 526 complete treatment cycles, n= 82 live births) as a

285    separate model was developed for this group of women by McLernon et al (McLernon *et al.*, 2016). Due

286    to the small numbers, this separate model could not be validated in this study.

287    *Discrimination and calibration*

288    In the validation sample, the pooled c-statistic for the pre-treatment model was 0.62 (95% CI 0.59-0.64)

289    and for the post-treatment model 0.71 (95% CI 0.69-0.74). Figure 2a and 3 show the calibration plots for

290    both original models, depicting the correlation between the observed and predicted cumulative LBR. The

291    pre-treatment calibration plot had an intercept of -0.23 (95% CI -0.36- -0.10) and a slope of 0.98 (95% CI

292    0.69-1.27), and the post-treatment calibration plot had an intercept of -0.01 (95% CI -0.12-0.11) and a

293    slope of 0.97 (95% CI 0.77-1.19).

294    The pre-treatment model systematically overestimated the cumulative LBR over three complete cycles for

295    women in the validation sample. This is shown by a calibration curve with most of the confidence

296    intervals under the reference line (Figure 2a), indicating significantly higher predicted probabilities than

297    observed LBR. The calibration-in-the-large analysis confirmed this systematic overestimation with an

298    intercept of  -0.35. To improve calibration, the pre-treatment model was thus adjusted by subtracting 0.35

299    from the intercept of the original linear predictor, which decreased the predicted odds of a live birth by a

300    factor of 1.42 (see Supplementary Text 3). The calibration plot of the recalibrated pre-treatment model

301    showed improved accuracy of the predictions, with all confidence intervals overlapping the reference line

302    (Figure 2b). In contrast to the pre-treatment model, the post-treatment model correctly estimated the

303    cumulative LBR in the validation sample, as is shown by a calibration plot with confidence intervals

304    overlapping the reference line indicating no significant over- or underestimation (Figure 3).

305    *Updating of the models*

306    Addition of the biomarkers AMH, AFC and body weight to the pre-treatment and post-treatment model in

307    a multivariable regression analysis resulted in two new updated models. The updated pre-treatment model

308    included all three biomarkers as additional predictors for live birth. Since the relationship between both

309    AMH and AFC with the probability of live birth was non-linear, these predictors were included using

310    restricted cubic splines (see Supplementary Figure 1). The updated post-treatment model included only

311    AFC and AMH as additional predictors for live birth, of which AFC was modelled by using restricted

312    cubic splines (see Supplementary Figure 2). After internal validation of the updated models by

313    bootstrapping, the updated pre-treatment model had a corrected c-statistic of 0.66 (95% CI 0.64-0.68) and

314    the updated post-treatment model had a corrected c-statistic of 0.71 (95% CI 0.69-0.73). The addition of

315    AFC, AMH and body weight thus resulted in a slight improvement of the discriminatory capacity of the

316    pre-treatment model, while addition of AFC and AMH had no beneficial effect on the discriminative

317    performance of the post-treatment model.

318    *Examples of model predictions*

319    Figures 4, 5 and 6 show examples of model predictions as illustration for clinical application. Figure 4

320    presents predictions of the *recalibrated pre-treatment model* for couples with primary infertility caused

321    by a male factor. Cumulative probabilities of live birth are calculated up to three complete ICSI cycles,

322    and are differentiated by female age (30 or 40 years) and duration of infertility (2 years or 5 years). As is

323    shown in figure 4, age is the most important predictor in the pre-treatment model. A 30-year-old woman

324    with 2 years of infertility has a predicted probability of a live birth of 0.40 in the first ICSI cycle,

325    increasing to 0.72 over three complete cycles. For a 40-year-old woman with 2 years of infertility, these

326    probabilities are 0.15 and 0.32 respectively.

327     Figure 5 shows predictions of the *updated pre-treatment model*, with AMH, AFC and body weight as new

328     predictors in the model. Predictions are presented for couples with two years of primary infertility caused

329     by a male factor, and differentiation is based on female age (30 or 40 years), AMH (2.0 or 0.5 ng/mL) and

330     AFC (15 or 7). In all scenarios the female body weight is 70 kilograms. A 30-year-old woman with an

331     average ovarian reserve at the start of her first treatment – indicated by an AMH of 2.0 ng/mL and an

332     AFC of 15 –  has a predicted probability of a live birth of 0.37 in the first cycle and 0.69 over three cycles

333     (0.17 and 0.37 for a 40-year-old woman). If this woman has a reduced ovarian reserve – indicated by an

334     AMH of 0.5 ng/mL and an AFC of 7 – the predicted probabilities decrease to 0.19 and 0.42, respectively

335     (0.08 and 0.18 for a 40-year-old woman).

336     Figure 6 shows predictions of the *post-treatment model,* which revises the predicted probabilities of the

337     pre-treatment models by adding information of the first treatment cycle. Predictions are calculated for

338     women with two years of primary, non-tubal infertility and are differentiated by female age (30 or 40

339     years), number of oocytes (10 or 5) and embryo cryopreservation (yes or no). In all scenarios the woman

340     received a cleavage stage single embryo transfer. The predicted probabilities of a live birth for women

341     with a favourable prognosis – aged 30-years, 10 oocytes retrieved and cryopreserved embryos – is 0.49 in

342     the first ICSI cycle, increasing to 0.83 over 3 complete cycles. In contrast, for women with a poorer

343     prognosis – aged 40 years, 5 oocytes retrieved and no embryos cryopreserved – the predicted probabilities

344     are 0.11 and 0.26, respectively.

345    **Discussion**

346    *Main findings*

347    This external validation study of the McLernon pre-treatment and post-treatment model found that, after

348    minor recalibration of the intercept of the pre-treatment model, both models accurately predict the

349    cumulative probability of live birth up to three complete IVF/ICSI cycles in a more contemporary cohort

350    in another country. The discriminatory capacity of the pre-treatment model in an external cohort was

351    limited, whereas the post-treatment model had a fair ability to discriminate between couples with and

352    without a live birth after treatment.

353    *Strengths*

354    This study focuses on the external validation of an IVF prediction model, which is an essential but

355    frequently overlooked step before implementation of a prediction model in clinical practice (Altman *et*

356    *al.*, 2009). In contrast to redeveloping new models for the same outcome, external validation and updating

357    of existing models prevents the loss of scientific information by combining the information captured in

358    the original model with information of a new patient cohort (Moons, Kengne, Grobbee, *et al.*, 2012).

359    Embryo cryopreservation has become an important part of IVF/ICSI treatment, and most couples have

360    more than just one complete treatment cycle (Wong *et al.*, 2014). Unlike previous prediction models

361    (Leushuis *et al.*, 2009; van Loendersloot *et al.*, 2014), the McLernon models provide a more useful

362    estimate of cumulative treatment success. As such, the validation of these models represents a significant

363    step forward in creating a clinically useful tool to manage expectations and to inform decision making

364    around IVF.

365    This study benefits from the prospective design of the OPTIMIST study, which has ensured reliable data

366    collection, with relatively low numbers of missing values and a low risk of selection bias. The multicentre

367    design resulted in a highly representable cohort for Dutch fertility care. And although it is known that the

368     IVF/ICSI success rates vary between fertility centres, the inclusion of multiple centres will increase the

369     generalisability and applicability of the external validation of the McLernon models within the

370     Netherlands.

371     Furthermore, the external validation was performed on data collected in a recent time period (2011-2014).

372     Due to changing patient populations, new treatment protocols, improving technologies and increasing

373     success rates over time, prediction models in reproduction medicine have no static form and should be

374     regularly updated to optimally reflect the latest circumstances in which they are used (Altman *et al.*,

375     2009). As the McLernon models were developed on data collected between 1999 and 2009, data of the

376     more recently performed OPTIMIST study were helpful to investigate if model performance was still

377     accurate in current practice.

378     *Weaknesses*

379     This study has a number of limitations. First, the external validation involved data from a prospective

380     cohort study within which two randomised controlled trials were embedded evaluating the effectiveness

381     of individualised doses of gonadotropins based on AFC. Strict dosing regimens might have affected some

382     treatment outcomes, such as cancellation rates and number of oocytes, thus influencing the predictive

383     capacity of the models in the validation sample. However, as the OPTIMIST study found no difference

384     between the dosing regimens on cumulative live birth rates, the impact on model performance is likely to

385     be minimal.

386     Second, the OPTIMIST study used strict eligibility criteria. Therefore, the validation sample does not

387     fully represent the diversity of the patient population initiating IVF/ICSI treatment in the Netherlands. As

388     none of the women in the validation sample were anovulatory, external validation of the models was only

389     performed for an ovulatory population. This limits the generalisability of the models to some extent, as

390     the original McLernon models were developed in a population which also included anovulatory women.

391     Also, it could have had some impact on model performance. However, since anovulation had only a small

392    predictive value in the pre-treatment model, and the majority of couples underwent IVF/ICSI for other

393    indications, a large impact on model performance is unlikely.

394    Third, the OPTIMIST study had a follow up period of 18 months, leading to small numbers of women

395    with more than three complete treatment cycles. Model performance could therefore only be reliably

396    validated up to three complete cycles. However, most couples in the Netherlands complete a maximum of

397    three treatment cycles which is partly due to the national reimbursement policy, but also by the high rates

398    of embryo cryopreservation, increasing the number of embryo transfers and LBR per cycle. Therefore,

399    model validation up to three complete cycles has particular clinical relevance for current Dutch fertility

400    care.

401    Last, the original McLernon prediction models were developed on linked cycle data, which were then

402    used to estimate cumulative pregnancy chances. Therefore, these models used the optimistic assumption

403    that the cumulative LBR in couples who discontinue IVF treatment without a live birth would have been

404    equal to that of couples who continue further treatment cycles, after correction of predictor effects. This

405    assumption tends to lead to overestimation of the cumulative LBR, as women with a low prognosis of

406    achieving a live birth are generally more likely to discontinue treatment (Brandes *et al.*, 2009; Olivius *et*

407    *al.*, 2004). Since the reasons for treatment withdrawal were unknown in the current external validation

408    study, a similar method was used that probably resulted in some degree of overestimation of the

409    cumulative LBR in the validation cohort. However, as the original McLernon models were developed

410    with this approach, and the predictions for cumulative LBR over multiple complete cycles were

411    considered to be clinically more relevant than per cycle predictions, we feel that the current method is the

412    best option for the external validation of the McLernon models.

413    *Explanation of findings*

414    The discriminatory capacity of the pre-treatment model was markedly lower in the validation sample than

415    in the development sample. In the development study, a c-statistic of 0.73 (95% CI 0.72-0.74) was

416   reported, whereas the present study found a c-statistic of 0.62 (95% CI 0.59-0.64). For the post-treatment

417   model, the discriminatory performance in the validation sample was comparable to that in the

418   development sample, with a c-statistic of 0.71 (95% CI 0.69-0.74) and 0.72 (95% CI 0.71-0.73)

419   respectively (McLernon *et al.*, 2016). As it is known that prediction models tend to perform too

420   optimistically in the development dataset due to overfitting, some reduction in model performance is to be

421   expected during external validation due to the differences between samples (Altman *et al.*, 2009; Moons,

422   Kengne, Woodward, *et al.*, 2012). This, to some extent, also explains the lower overall performance of

423   the pre-treatment model. The comparable performance of the post-treatment model in both samples

424   indicates that the treatment related variables that were added to this model (number of oocytes,

425   cryopreservation of embryos, and the number and stage of embryos) are important predictors for live birth

426   after treatment.

427   Other than the influence of overfitting, some key differences between the Dutch and UK healthcare

428   systems may also have affected the models' performance in this external validation study. An important

429   factor is the reimbursement policy for fertility treatment. All Dutch infertile couples are insured for a

430   minimum of three complete IVF/ICSI cycles. In contrast, most couples in the UK receive no standard

431   funding for ART (Berg Brigham *et al.*, 2013). Since IVF/ICSI treatment is expensive, this induces

432   discrepancies in the patient population initiating and continuing treatment between the two study samples

433   (Rajkhowa *et al.*, 2006). As can be seen in the baseline table (Table I) and flowchart (Figure 1), couples

434   in the UK had a longer average duration of infertility before starting treatment and were more likely to

435   discontinue treatment after the first and second cycles than couples in the Netherlands. Also, the decrease

436   in LBR is more evident in the UK than in the Netherlands over the first three cycles, which suggests that

437   differences exist in both reasons for discontinuation as well as prognostic profiles of women

438   discontinuing treatment in the two countries. These phenomena are, in part, financially driven, and could

439   partially explain the difference in predictive ability of the UK models in the Dutch cohort.

440    Furthermore, despite the fact that the infertility guidelines of both countries include similar approaches

441    for treatment of infertile couples, there are important variations in treatment characteristics between the

442    two study samples (Dutch Society of Obstetrics and Gynaecology (NVOG), 2010; National Institute for

443    Health and Care Excellence (NICE), 2013). Some of these differences are mainly due to changes in

444    clinical practice over time. As is shown by the baseline table (Table 1), women in the more recent Dutch

445    cohort (2011-2014) generally had a single embryo transfer in their first fresh treatment cycle, whereas

446    women in the earlier UK cohort (1999-2009) most often had a double embryo transfer. Also, embryo

447    cryopreservation was performed in over half of the Dutch women as compared to only a quarter of the

448    women in the UK. Other differences are explained by variation in treatment protocols between

449    geographic locations. For one, no blastocyst stage embryos transfers were performed in the Netherlands in

450    contrast to the proportion of blastocyst stage embryo transfers in the UK of more than 10%. Also, Dutch

451    women more frequently had no embryo available for transfer after their first treatment cycle, which is

452    most likely caused by strict cancellation criteria particularly for hyper response. These differences in

453    treatment characteristics suggest that the development sample does not fully reflect clinical practice in a

454    more recent time period and in a different geographic context. As cumulative LBR are substantially

455    affected by the variation in treatment characteristics (Glujovsky *et al.*, 2016; Pandian *et al.*, 2013; Wong

456    *et al.*, 2014), this could explain part of the different performance of the pre-treatment model in the

457    validation sample . The stable performance of the post-treatment model, which includes embryo stage and

458    embryo cryopreservation as important predictors, seems to confirm the impact of the variation in these

459    variables on model performance.

460    The addition of measures of ovarian reserve, i.e. AMH and AFC, and body weight to the McLernon

461    prediction models revealed only a marginal improvement of model performance in the OPTIMIST

462    dataset. The additional value of these tests can therefore be questioned, especially in view of the extra

463    costs and physical burden on the patient. Female age is one of the most important predictors in the

464    McLernon models (McLernon *et al.*, 2016). As female age is correlated with the ovarian reserve, adding

465    AMH and AFC provides limited new information to the prediction models. This is in line with previous

466    studies that showed that ovarian reserve tests have no added value to the use of female age alone in the

467    prediction of ongoing pregnancy after treatment (Broer *et al.*, 2013). Other potential predictors for live

468    birth, such as ethnicity, smoking status and alcohol intake, were not included in this update of the

469    McLernon model (Dhillon *et al.*, 2015; Rossi *et al.*, 2011; Waylen *et al.*, 2009). The additional value of

470    these variables for model performance was considered uncertain, as the reporting is remarkably subjective

471    and/or often incomplete (Liber and Warner, 2018; Stockwell *et al.*, 2016).

472    *Clinical implications*

473    Discrimination and calibration have been recognized as measures to evaluate the performance of

474    prediction models (Altman *et al.*, 2009; Steyerberg, 2009). However, the discriminative ability at the

475    binary level of most prediction models in reproductive medicine, as expressed by the c-statistic, is

476    considerably low (Leushuis *et al.*, 2009). As at the moment of prediction the outcome of pregnancy has

477    not yet occurred, the c-statistic is determined using the calculated probability of pregnancy. The

478    maximum value of the c-statistic depends on the variability of these calculated probabilities in the

479    infertile population. Since infertility is a complex and multifactorial health problem and due to the

480    absence of strong predictors for live birth – particularly pre-treatment – , the probability distribution in

481    infertile couples that have a live birth has a considerable overlap with the distribution of those without a

482    live birth. Therefore the maximum c-statistic can be expected to be low (Cook, 2007; Coppus *et al.*,

483    2009), as is seen in the external validation of the pre-treatment model. However, this does not necessarily

484    imply that such prediction models have limited use in clinical practice. Models with reliable predictions

485    and a clinically useful distribution of probabilities for achieving a live birth, as assessed by calibration,

486    can still support patients and clinicians in clinical decision making around infertility treatment (Coppus *et*

487    *al.*, 2009).

488      As the calibration plots of both the recalibrated pre-treatment model and the post-treatment model

489      indicate accurate predictions with a useful range of prognoses, these models can be used within the

490      Netherlands as counselling tools to complement clinical reasoning at two certain time points. Before

491      initiating treatment, the recalibrated pre-treatment model offers couples and clinicians a personalised and

492      objective estimate of success over multiple complete treatment cycles. And after the first fresh embryo

493      transfer, the post-treatment model provides a revised estimate using treatment related information to

494      personalize the predictions even more. Despite the applicability of the models as counselling tools to

495      inform patients about their prognosis, the McLernon models should not yet be used for decisions on

496      whether or not to withhold fertility treatment. The impact of such model-based decisions on cost-benefit

497      outcomes should be investigated first and proven to be beneficial. To implement the McLernon models as

498      counselling tools in other countries as well, national recalibration is recommended to account for

499      geographical differences in IVF/ICSI treatment.

500      The McLernon models were converted into an online calculator to facilitate the use of the models in

501      clinical practice (https://w3.abdn.ac.uk/clsm/opis). As the original pre-treatment model overestimates

502      cumulative LBR for couples in the Netherlands, conversion of the recalibrated pre-treatment model into a

503      new online calculator is needed for implementation in Dutch clinical practice. This tailored online

504      calculator can then provide accurate and up to date predictions for couples and clinicians in the

505      Netherlands. Ultimately, the online calculator will be offered for implementation on the websites of the

506      Dutch Patient Association for people with fertility problems 'Freya' and the Dutch Association of

507      Obstetrics and Gynaecology (NVOG) to increase the accessibility of the models.

508      *Research implications*

509      Following this external validation study, future studies could focus on the impact of introducing the

510      McLernon prediction models in clinical practice, and assess changes in patient and clinicians' behaviour

511      and its effects on LBR and cost-effectiveness.

512    In conclusion, after minor recalibration of the pre-treatment model, the McLernon models have proven to

513    be valid in predicting the chance of cumulative live birth after multiple complete treatment cycles in

514    another geographical context and in a more recent time period. Updating the models with AMH, AFC and

515    body weight revealed only a marginal improvement of predictive performance. Following national

516    recalibration, implementation of the McLernon models as counselling tools in clinical practice will

517    provide infertile couples and clinicians with objective and personalized estimates of success over multiple

518    complete IVF/ICSI cycles.

519

523     **Authors' roles**

524     T.C.v.T. and S.C.O and all other members from the OPTIMIST study group collected the data. D.J.M.,

525     S.B., F.J.M.B. and H.L.T were involved in study conception and study design. J.A.L. and M. J. C. E.

526     performed the statistical analysis. J.A.L. drafted the manuscript. J.A.L., M.J.C.E., F.J.M.B. B.W.M.,

527     H.L.T interpreted the data. All authors participated to the discussion of the findings and revised the

528     manuscript.

531     **Conflict of interest**

532     M.J.C.E., D.J.M. and S.B. have nothing to disclose. J.A.L, S.C.O, T.C.v.T. and H.LT. received an

533     unrestricted personal grant from Merck BV. B.W.M. is supported by a NHMRC Practitioner Fellowship

534     (GNT1082548)  and  reports consultancy for ObsEva, Merck and Guerbet. F.J.M.B. receives monetary

535     compensation as a member of the external advisory board for Merck Serono (the Netherlands) and

536     Ferring pharmaceutics BV (the Netherlands), for consultancy work for Gedeon Richter (Belgium) and

537     Roche Diagnostics on automated AMH assay development, and for a research cooperation with Ansh

538     Labs (USA).

539    **References**

540    Akaike H. A new look at the statistical model identification. *IEEE Trans Automat Contr* 1974;**19**:716–

541         723.

542    Altman DG, Vergouwe Y, Royston P, Moons KGM. Prognosis and prognostic research: validating a

543         prognostic model. *BMJ* 2009;**338**:b605.

544    Arvis P, Lehert P, Guivarc'h-Leveque A. Simple adaptations to the Templeton model for IVF outcome

545         prediction make it current and clinically useful. *Hum Reprod* 2012;**27**:2971–2978.

546    Berg Brigham K, Cadier B, Chevreul K. The diversity of regulation and public financing of IVF in

547         Europe and its impact on utilization. *Hum Reprod* 2013;**28**:666–75.

548    Boivin J, Bunting L, Collins JA, Nygren KG. International estimates of infertility prevalence and

549         treatment-seeking: potential need and demand for infertility medical care. *Hum Reprod*

550         2007;**22**:1506–12.

551    Brandes M, van der Steen JOM, Bokdam SB, Hamilton CJCM, de Bruin JP, Nelen WLDM, Kremer

552         JAM. When and why do subfertile couples discontinue their fertility care? A longitudinal cohort

553         study in a secondary care subfertility population. *Hum Reprod* 2009;**24**:3127–35.

554    Broer SL, van Disseldorp J, Broeze KA, Dolleman M, Opmeer BC, Bossuyt P, Eijkemans MJC, Mol B-

555         WJ, Broekmans FJM, Broer SL, *et al.* Added value of ovarian reserve testing on patient

556         characteristics in the prediction of ovarian response and ongoing pregnancy: an individual patient

557         data approach. *Hum Reprod Update* 2013;**19**:26–36.

558    Cook NR. Statistical Evaluation of Prognostic versus Diagnostic Models: Beyond the ROC Curve. *Clin*

559         *Chem* 2007;**54**:17–23.

560    Coppus SFPJ, van der Veen F, Opmeer BC, Mol BWJ, Bossuyt PMM. Evaluating prediction models in

561      reproductive medicine. *Hum Reprod* 2009;**24**:1774–1778.

562      Datta J, Palmer MJ, Tanton C, Gibson LJ, Jones KG, Macdowall W, Glasier A, Sonnenberg P, Field N,

563          Mercer CH, *et al.* Prevalence of infertility and help seeking among 15 000 women and men. *Hum*

564          *Reprod* 2016;**31**:2108–2118.

565      de Neubourg D, Bogaerts K, Blockeel C, Coetsier T, Delvigne A, Devreker F, Dubois M, Gillain N,

566          Gordts S, Wyns C. How do cumulative live birth rates and cumulative multiple live birth rates over

567          complete courses of assisted reproductive technology treatment per woman compare among

568          registries? *Hum Reprod* 2016;**31**:93–99.

569      Dhillon RK, Smith PP, Malhas R, Harb HM, Gallos ID, Dowell K, Fishel S, Deeks JJ, Coomarasamy A.

570          Investigating the effect of ethnicity on IVF outcome. *Reprod Biomed Online* 2015;**31**:356–363.

571      Dutch Society of Obstetrics and Gynaecology (NVOG). Landelijke Netwerkrichtlijn Subfertiliteit. 2010.

572          Available at: http://nvog-documenten.nl.

573      Dyer S, Chambers GM, de Mouzon J, Nygren KG, Zegers-Hochschild F, Mansour R, Ishihara O, Banker

574          M, Adamson GD. International Committee for Monitoring Assisted Reproductive Technologies

575          world report: Assisted Reproductive Technology 2008, 2009 and 2010. *Hum Reprod* 2016;**31**:1588–

576          1609.

577      Glujovsky D, Farquhar C, Quinteiro Retamar AM, Alvarez Sedo CR, Blake D. Cleavage stage versus

578          blastocyst stage embryo transfer in assisted reproductive technology. *Cochrane Database Syst Rev*

579          2016:CD002118.

580      Harrell FE. *Regression Modeling Strategies: With Applications to Linear Models, Logistic Regression,*

581          *and Survival Analysis*. New York: Springer-Verlag , 2001.

582      Harrell FE, Lee KL, Mark DB. Multivariable prognostic models: issues in developing models, evaluating

583     assumptions and adequacy, and measuring and reducing errors. *Stat Med* 1996;**15**:361–387.

584     Leushuis E, van der Steeg JW, Steures P, Bossuyt PMM, Eijkemans MJC, van der Veen F, Mol BWJ,

585     Hompes PGA. Prediction models in reproductive medicine: a critical appraisal†. *Hum Reprod*

586     *Update* 2009;**15**:537–552.

587     Liber AC, Warner KE. Has Underreporting of Cigarette Consumption Changed Over Time? Estimates

588     Derived From US National Health Surveillance Systems Between 1965 and 2015. *Am J Epidemiol*

589     2018;**187**:113–119.

590     Malizia BA, Hacker MR, Penzias AS. Cumulative Live-Birth Rates after In Vitro Fertilization. *N Engl J*

591     *Med* 2009;**360**:236–243.

592     McLernon DJ, Steyerberg EW, te Velde ER, Lee AJ, Bhattacharya S. Predicting the chances of a live

593     birth after one or more complete cycles of in vitro fertilisation: population based study of linked

594     cycle data from 113 873 women. *BMJ* 2016;**355**:i5735.

595     Moons KGM, Kengne AP, Grobbee DE, Royston P, Vergouwe Y, Altman DG, Woodward M. Risk

596     prediction models: II. External validation, model updating, and impact assessment. *Heart*

597     2012;**98**:691–8.

598     Moons KGM, Kengne AP, Woodward M, Royston P, Vergouwe Y, Altman DG, Grobbee DE. Risk

599     prediction models: I. Development, internal validation, and assessing the incremental value of a new

600     (bio)marker. *Heart* 2012;**98**:683–90.

601     National Institute for Health and Care Excellence (NICE). Fertility problems: assessment and treatment.

602     Clinical guideline. 2013. Available at: https://www.nice.org.uk.

603     Nelson SM, Lawlor DA. Predicting live birth, preterm delivery, and low birth weight in infants born from

604     in vitro fertilisation: A prospective study of 144,018 treatment cycles. *PLoS Med* 2011;**8**:e1000386.

605     Oakley L, Doyle P, Maconochie N. Lifetime prevalence of infertility and infertility treatment in the UK:

606          results from a population-based survey of reproduction. *Hum Reprod* 2008;**23**:447–450.

607     Olivius C, Friden B, Borg G, Bergh C. Why do couples discontinue in vitro fertilization treatment? A

608          cohort study. *Fertil Steril* 2004;**81**:258–61.

609     Oudshoorn SC, van Tilborg TC, Eijkemans MJC, Oosterhuis GJE, Friederich J, van Hooff MHA, van

610          Santbrink EJP, Brinkhuis EA, Smeenk JMJ, Kwee J, *et al.* Individualized versus standard FSH

611          dosing in women starting IVF/ICSI: an RCT. Part 2: The predicted hyper responder. *Hum Reprod*

612          2017;**32**:2506–2514.

613     Pandian Z, Marjoribanks J, Ozturk O, Serour G, Bhattacharya S. Number of embryos for transfer

614          following in vitro fertilisation or intra-cytoplasmic sperm injection. *Cochrane database Syst Rev*

615          2013;**7**:CD003416.

616     Rajkhowa M, Mcconnell A, Thomas GE. Reasons for discontinuation of IVF treatment: a questionnaire

617          study. *Hum Reprod* 2006;**21**:358–363.

618     Rossi B V, Berry KF, Hornstein MD, Cramer DW, Ehrlich S, Missmer SA. Effect of Alcohol

619          Consumption on In Vitro Fertilization. *Obstet Gynecol* 2011;**117**:136–142.

620     Rubin DB. Multiple Imputation for Nonresponse in Surveys. In: John Wiley & Sons, 2004.

621     Smeenk JM, Stolwijk AM, Kremer JA, Braat DD. External validation of the templeton model for

622          predicting success after IVF. *Hum Reprod* 2000;**15**:1065–8.

623     Smith ADAC, Tilling K, Lawlor DA, Nelson SM. External Validation and Calibration of IVFpredict: A

624          National Prospective Cohort Study of 130,960 In Vitro Fertilisation Cycles. Sun Q-Y (ed). *PLoS*

625          *One* 2015;**10**:e0121357.

626     Sterne JAC, White IR, Carlin JB, Spratt M, Royston P, Kenward MG, Wood AM, Carpenter JR. Multiple

627    imputation for missing data in epidemiological and clinical research: potential and pitfalls. *BMJ*

628    2009;**338**:b2393.

629    Steyerberg EW. *Clinical Prediction Models: A Practical Approach to Development, Validation, and*

630    *Updating*. New York, NY: Springer New York, 2009.

631    Stockwell T, Zhao J, Greenfield T, Li J, Livingston M, Meng Y. Estimating under- and over-reporting of

632    drinking in national surveys of alcohol consumption: identification of consistent biases across four

633    English-speaking countries. *Addiction* 2016;**111**:1203–1213.

634    Templeton A, Morris JK, Parslow W. Factors that affect outcome of in-vitro fertilisation treatment.

635    *Lancet* 1996;**348**:1402–1406.

636    te Velde ER, Nieboer D, Lintsen AM, Braat DDM, Eijkemans MJC, Habbema JDF, Vergouwe Y.

637    Comparison of two models predicting IVF success; the effect of time trends on model performance.

638    *Hum Reprod* 2014;**29**:57–64.

639    van Buuren S, Groothuis-Oudshoorn K. MICE : Multivariate Imputation by Chained Equations in R. *J*

640    *Stat Softw* 2011;**45**:1–67.

641    van Loendersloot L, Repping S, Bossuyt PMM, van der Veen F, van Wely M. Prediction models in in

642    vitro fertilization; where are we? A mini review. *J Adv Res* 2014;**5**:295–301.

643    van Loendersloot LL, van Wely M, Repping S, van der Veen F, Bossuyt PMM. Templeton prediction

644    model underestimates IVF success in an external validation. *Reprod Biomed Online* 2011;**22**:597–

645    602.

646    van Tilborg TC, Eijkemans MJ, Laven JS, Koks CA, de Bruin JP, Scheffer GJ, van Golde RJ, Fleischer

647    K, Hoek A, Nap AW, *et al.* The OPTIMIST study: optimisation of cost effectiveness through

648    individualised FSH stimulation dosages for IVF treatment. A randomised controlled trial. *BMC*

649       *Womens Health* 2012;**12**:29.

650   van Tilborg TC, Oudshoorn SC, Eijkemans MJC, Mochtar MH, van Golde RJT, Hoek A, Kuchenbecker

651       WKH, Fleischer K, de Bruin JP, Groen H, *et al.* Individualized FSH dosing based on ovarian reserve

652       testing in women starting IVF/ICSI: a multicentre trial and cost-effectiveness analysis. *Hum Reprod*

653       2017;**32**:2485–2495.

654   van Tilborg TC, Torrance HL, Oudshoorn SC, Eijkemans MJC, Koks CAM, Verhoeve HR, Nap AW,

655       Scheffer GJ, Manger AP, Schoot BC, *et al.* Individualized versus standard FSH dosing in women

656       starting IVF/ICSI: an RCT. Part 1: The predicted poor responder. *Hum Reprod* 2017;**32**:2496–2505.

657   Waylen AL, Metwally M, Jones GL, Wilkinson AJ, Ledger WL. Effects of cigarette smoking upon

658       clinical outcomes of assisted reproduction: a meta-analysis. *Hum Reprod Update* 2009;**15**:31–44.

659   Wong KM, Mastenbroek S, Repping S. Cryopreservation of human embryos and its contribution to

660       in vitro fertilization success rates. *Fertil Steril* 2014;**102**:19–26.

661   Zegers-Hochschild F, Adamson GD, Dyer S, Racowsky C, de Mouzon J, Sokol R, Rienzi L, Sunde A,

662       Schmidt L, Cooke ID, *et al.* The International Glossary on Infertility and Fertility Care, 2017†‡§.

663       *Hum Reprod* 2017;**32**:1786–1801.
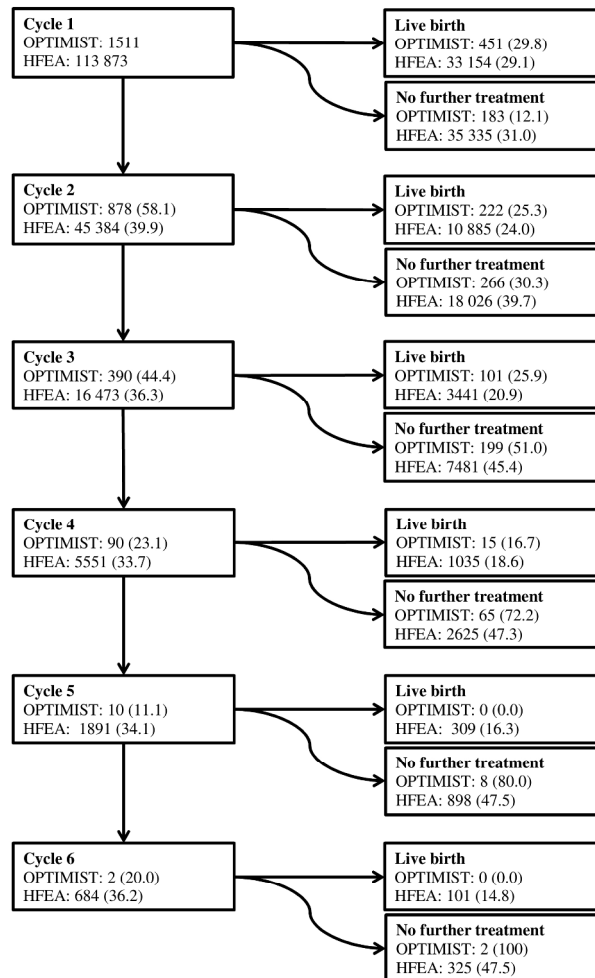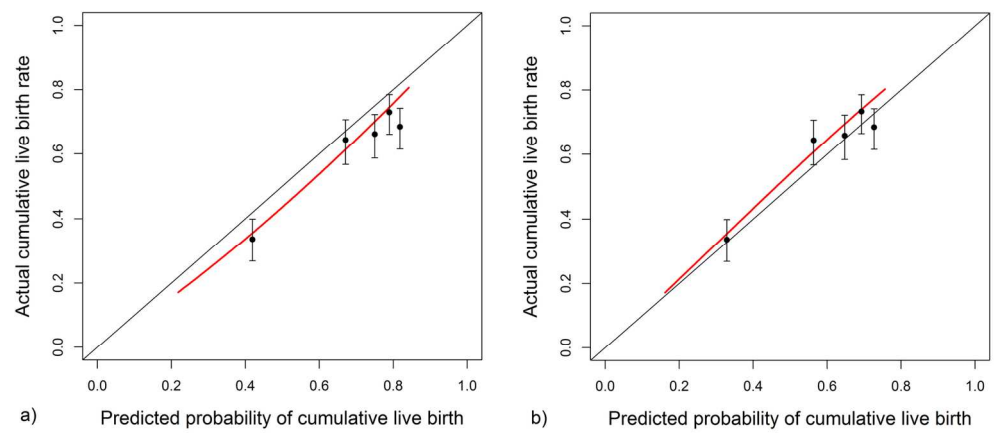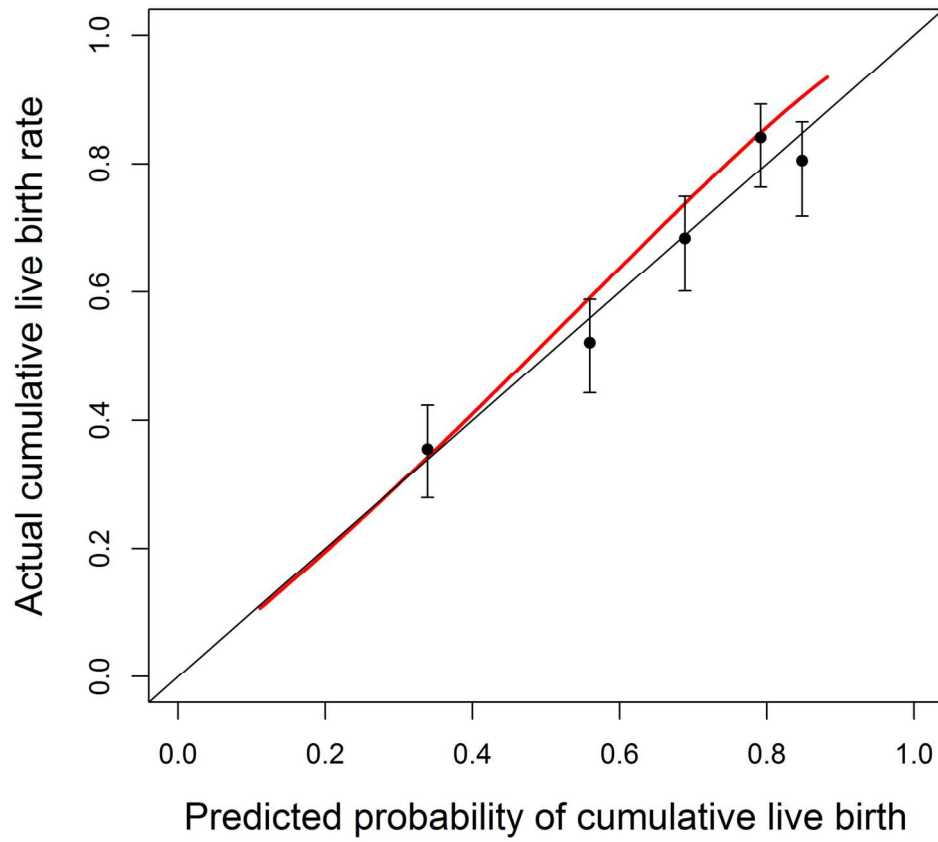
664

665

666　**Figure legends**

667　**Figure 1:** Flow chart presenting the numbers (%) of live birth, treatment continuation and discontinuation

668　over six complete cycles in the OPTIMIST and HFEA databases (McLernon *et al.*, 2016).

669

670　**Figure 2:** Calibration plots showing the association between the calculated and observed cumulative live

671　birth rates over 3 complete IVF/ICSI cycles in the OPTIMIST cohort for **a)** the *original pre-treatment*

672　*model* as described by McLernon et al (McLernon *et al.*, 2016) **b)** *recalibrated pre-treatment model* with

673　adjustment of the intercept.

674

675　**Figure 3:** Calibration plot showing the association between the calculated and observed cumulative live

676　birth rates over 3 complete IVF/ICSI cycles in the OPTIMIST cohort for the *original post-treatment*

677　*model* as described by McLernon (McLernon *et al.*, 2016).

678

679　**Figure 4:** Example of the *recalibrated pre-treatment model* predicting the cumulative probability of a

680　live birth up to three complete ICSI cycles for a woman with primary infertility caused by a male factor,

681　aged 30 or 40 years with an infertility duration of two or five years.

682

683　**Figure 5:** Example of the with AMH, AFC and body weight *updated pre-treatment model* predicting the

684　cumulative probability of a live birth up to three complete ICSI cycles for a woman with two years of

685　primary infertility caused by a male factor, aged 30 or 40 years, a total body weight of 70 kilograms, with

686　an AMH of 2.0 or 0.5 ng/mL and an AFC of 15 or 7.

687

688    **Figure 6:** Example of the *post-treatment model* predicting the cumulative probability of a live birth up to

689    three complete ICSI cycles for a woman with two years of primary infertility caused by a male factor,

690    aged 30 or 40 years, with 5 or 10 oocytes retrieved, a cleavage stage single embryo transfer, with or

691    without embryo cryopreservation.

692

693    **Supplementary Figure 1.** Plots showing the adjusted relation between the predictors included in the

694    *updated McLernon pre-treatment model* and the probability of a live birth after IVF/ICSI treatment.

695    Predictor; linear predictor (XB) of the original pre-treatment model as described by McLernon

696    (McLernon et al. 2016), Weight; female body weight in kg, AFC; antral follicle count (2-10mm), AMH;

697    anti-Müllerian hormone (ng/mL)

698

699    **Supplementary Figure 2.** Plots showing the adjusted relation between the predictors in the *updated*

700    *McLernon post-treatment model* and the probability of a live birth after IVF/ICSI treatment.

701    Predictor: linear predictor (XB) of the original post-treatment model as described by McLernon

702    (McLernon et al 2016); AFC; antral follicle count (2-10mm), AMH; anti-Müllerian hormone (ng/mL)

703

704

705

Figure 1: Flow chart presenting the numbers (%) of live birth, treatment continuation and discontinuation over six complete treatment cycles in the OPTIMIST and HFEA databases (McLernon et al., 2016).

209x297mm (300 x 300 DPI)
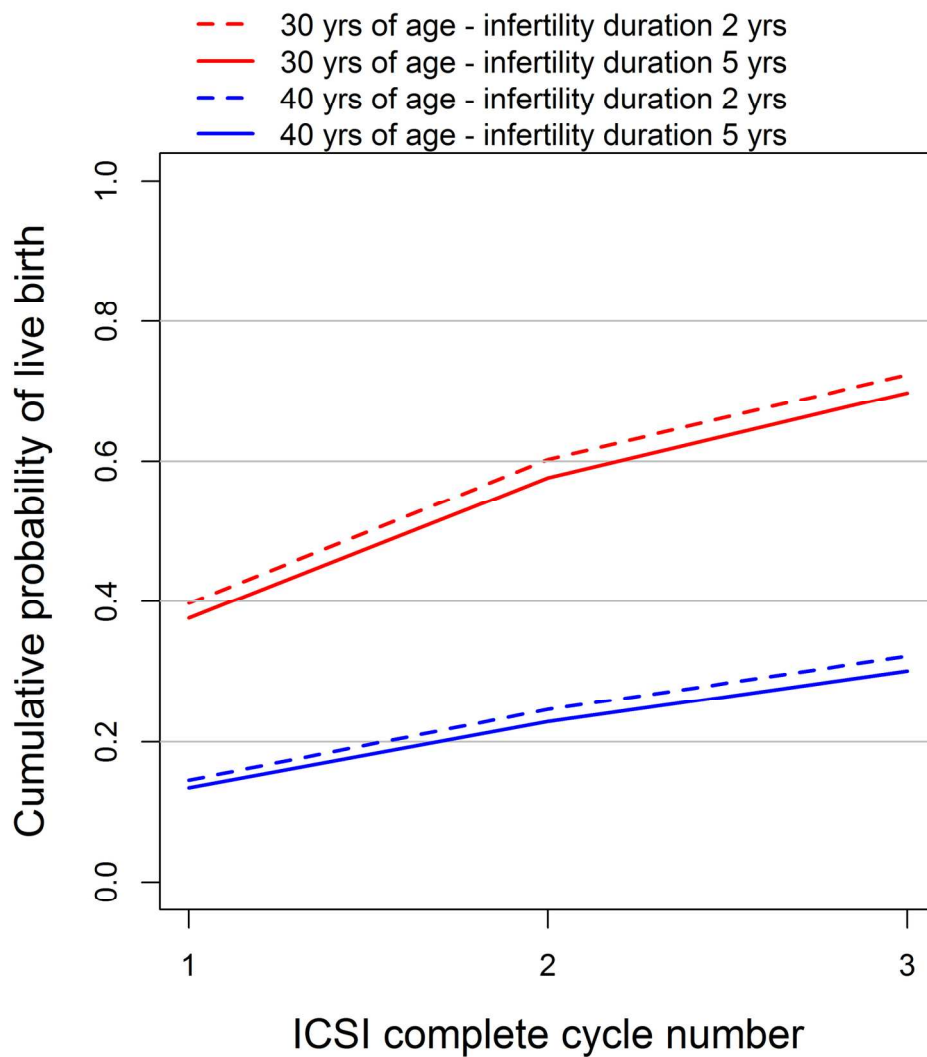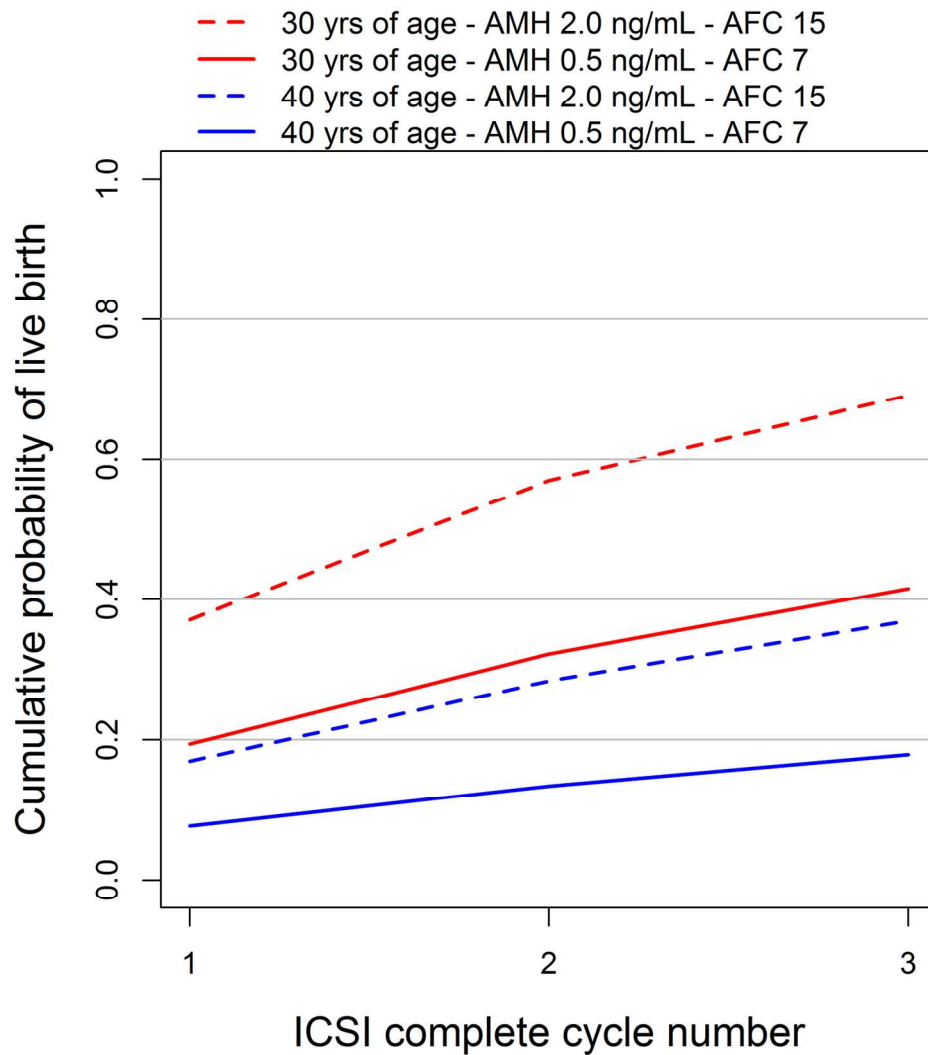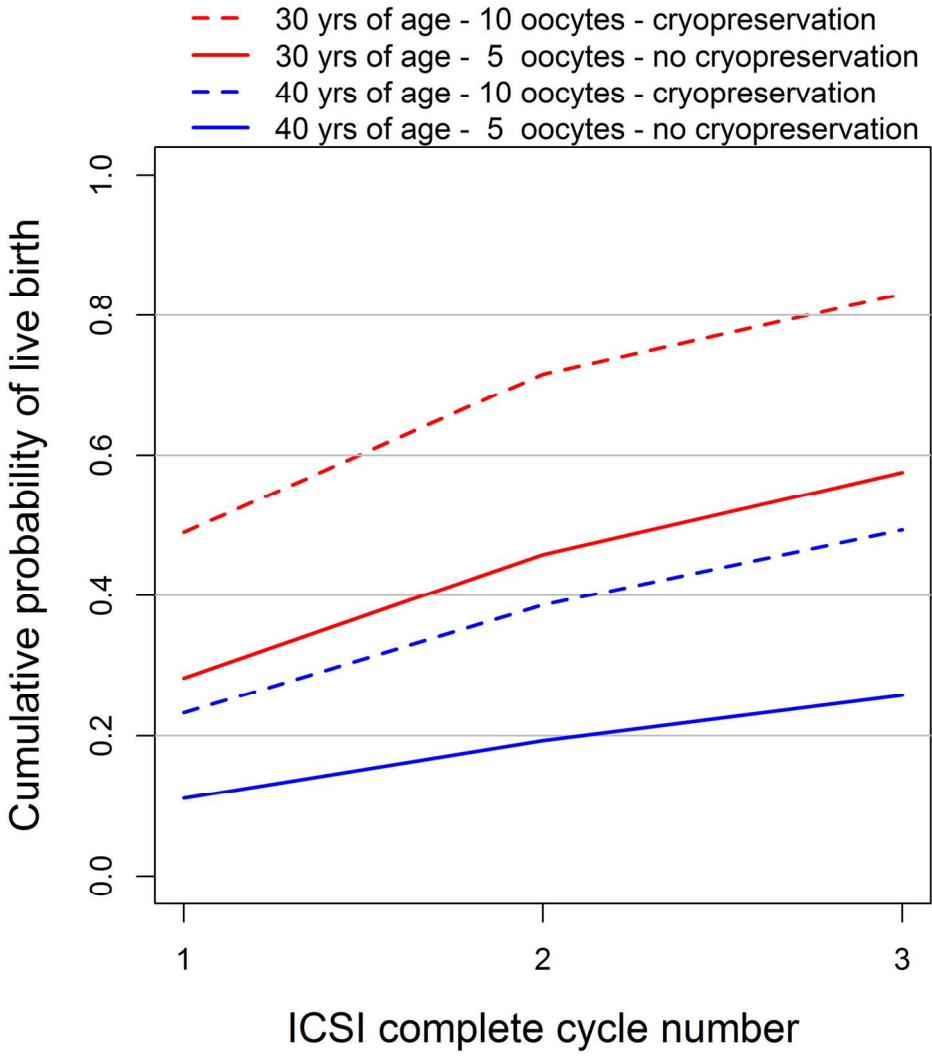
152x76mm (300 x 300 DPI)

152x152mm (300 x 300 DPI)

Figure 4:  Example of the recalibrated pre-treatment model predicting the cumulative probability of live birth up to three complete ICSI cycles for a woman with primary infertility caused by a male factor, aged 30 or 40 years with an infertility duration of two or five years.

152x166mm (300 x 300 DPI)

Figure 5: Example of the with AMH, AFC and body weight updated pre-treatment model predicting the cumulative probability of live birth up to three complete ICSI cycles for a woman with two years of primary infertility caused by a male factor, aged 30 or 40 years, a total body weight of 70 kilograms, with an AMH of 2.0 or 0.5ng/mL and an AFC of 15 or 7.

152x166mm (300 x 300 DPI)

152x166mm (300 x 300 DPI)

**Tables**

**Table I** Characteristics of patient and treatment variables included as predictors in the development sample (HFEA cohort) and the validation sample (OPTIMIST cohort) (McLernon *et al.*, 2016). ~~Unless stated otherwise data are n (%).~~
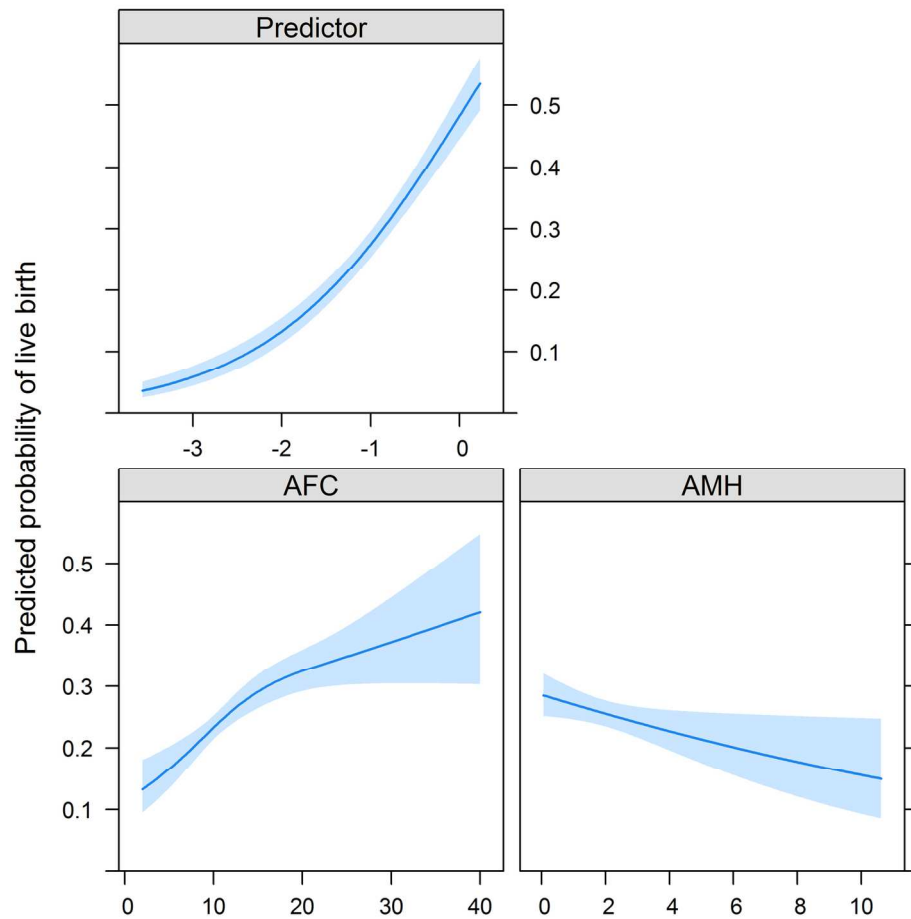
| Characteristics | HFEA cohort | OPTIMIST cohort | Missing values in OPTIMIST cohort (%) |
|---|---|---|---|
| No of women | 113 873 | 1 511 | |
| No of complete cycles | 184 269 | 2 881 | |
| **Patient characteristics** | | | |
| Age (years), mean (SD) | 34.1 (5) | 33.5 (5) | 2 (0.1) |
| Duration of infertility (years), median (IQR) | 4 (3-6) | 2 (2-3) | 18 (1.2) |
| No previous pregnancy in couple~~,~~ | 75 541 (66) | 917 (61) | 2 (0.1) |
| Cause of infertility: | | | |
| - Tubal factor~~)~~ | 26 545 (23) | 158 (11) | |
| - Male factor | 49 753 (44) | 839 (56) | |
| - Anovulatory | 15 942 (14) | NA by protocol | |
| - Endometriosis | 7 590 (7) | 60 (4) | |
| - Unexplained | 32 693 (29) | 521 (35) | |
| Body weight (kg), mean (SD) | NA | 69.5 (13) | 36 (2.4) |
| Anti-Müllerian hormone (ng/mL), median (IQR) | NA | 1.9 (1-3) | 169 (11.2) |
| Antral follicle count (2-10mm), median (IQR) | NA | 13 (9-18) | |
| **Treatment characteristics of first completed cycle** | | | |
| IVF | 67 511 (59) | 830 (55) | |
| ICSI | 46 362 (41) | 681 (45) | |
| No of oocytes retrieved, median (IQR) | 8 (5-13) | 8 (5-13)[a] | 1 (0.1) |
| No of embryos created, median (IQR) | 5 (2-8) | 4 (2-7)[a] | 4 (0.3) |
| No of embryos frozen, median (IQR) | 0 (0-1) | 1 (0-3)[a] | 6 (0.5) |
| Cryopreservation of embryos | 28 950 (25) | 726 (48) | |
| Fresh embryo transfer: stage and no. of transferred embryos: | | | 24 (1.6) |
| - Cleavage stage SET | 9 248 (8) | 1 004 (66) | |
| - Cleavage stage DET | 75 701 (66) | 125 (8) | |
| - Cleavage stage TET | 8 649 (8) | 4 (0.3) | |
| - Blastocyst stage SET | 662 (1) | NA | |
| - Blastocyst stage DET | 2 960 (3) | NA | |
| - Blastocyst stage TET~~)~~ | 130 (0.1) | NA | |
| - No transfer | 15 501 (14) | 354 (23) | |

Data are presented as number (%) unless Unless stated otherwise specified. data are n (%).IQR; interquartile range, NA; not available, SET; single embryo transfer, DET; double embryo transfer, TET; triple embryo transfer. a) Median is calculated over 1293 women who had an ovarian follicle aspiration.

Supplementary Figure 1. Plots showing the adjusted relation between the predictors included in the updated McLernon pre-treatment model and the probability of a live birth after IVF/ICSI treatment. ‖ + ‖ + Predictor; linear predictor (XB) of the original pre-treatment model as described by McLernon (McLernon et al. 2016), Weight; female body weight in kg, AFC; antral follicle count (2-10mm), AMH; anti-Müllerian hormone (ng/mL).

152x152mm (300 x 300 DPI)

Supplementary Figure 2. Plots showing the adjusted relation between the predictors in the updated McLernon post-treatment model and the probability of a live birth after IVF/ICSI treatment.‖ ┼ ‖ ┼ Predictor: linear predictor (XB) of the original post-treatment model as described by McLernon (McLernon et al 2016); AFC; antral follicle count (2-10mm), AMH; anti-Müllerian hormone (ng/mL).

152x152mm (300 x 300 DPI)

**Supplementary text 1.** McLernon pre-treatment model.

**Table showing the predictors in the original McLernon pre-treatment model** (McLernon *et al.*, 2016)**.**

| Name predictor | Description | Range of possible values |
|---|---|---|
| Age | Female age | 18 to 50 years |
| Duration | How long have you been trying to conceive? | 0 to 21 years |
| Previous | Have you been pregnant before? | 1 = No;  0 = Yes |
| Tubal | Do you have a problem with your tubes? | 1 = Yes; 0 = No |
| Anovulation | Do you have an ovulation problem? | 1 = Yes; 0 = No |
| MaleFactor | Do you have a male factor fertility problem? | 1 = Yes; 0 = No |
| Unexplained | Do you have an unexplained fertility problem? | 1 = Yes; 0 = No |
| Treatment | Which fertility treatment are you planning on having? | 1 = ICSI; 0 = IVF |

**Original pre-treatment model formulas as described by McLernon et al.** (McLernon *et al.*, 2016)**.**

1. For the non-linear relation between Age and the probability of live birth, the following Age1, Age2 and Age3 equations are first calculated and then used in the XB equation below (point 3).

   - Age1 = max((Age-26)/k,0)\*\*3+(11\*max((Age-41)/k,0)\*\*3-(15)\*max((Age-37)/k,0)\*\*3)/4;
   - Age2 = max((Age -31)/k,0)\*\*3+(6\*max((Age -41)/k,0)\*\*3-(10)\*max((Age -37)/k,0)\*\*3)/4;
   - Age3 = max((Age -34)/k,0)\*\*3+(3\*max((Age -41)/k,0)\*\*3-(7)\*max((Age -37)/k,0)\*\*3)/4; *k=15\*\*(2/3); \*\*means 'to the power of'*

2. For the non-linear relation between Year and the probability live birth, the following Year1 and Year2 equations are first calculated and then used in the XB equation below (point 3). The value Year= 0 is used for the most up to date predictions.

   - Year1 = max((Year+9)/k,0)\*\*3+((6)\*max((Year)/k,0)\*\*3-(9)\*max((Year+3)/k,0)\*\*3)/(3);
   - Year2 = max((Year+6)/k,0)\*\*3+((3)\*max((Year)/k,0)\*\*3-(6)\*max((Year+3)/k,0)\*\*3)/(3); *k= 9\*\*(2/3).*

3. Calculate XB.

   XB =   -0.9948 + 0.0362[a] + (0.0275\*Age) + (-0.1805\*Age1) + (0.4553\*Age2) + (-1.1990\*Age3) + (-0.0295\*Duration) + (-0.0772\*Previous) + (-0.0957\*Tubal) + (0.0492\*Anovulation) + (-0.1005\*MaleFactor) + (0.0602\*Unexplained) + (0.2155\*Treatment) + (0.0334\*Year) + (-0.0370\*Year1) + (0.2173\*Year2).

   a) To inflate predictions to 2013 an additional 0.0362 is added.

4. Calculate the predicted probability of live-birth after the first, second, …., sixth IVF cycle.

   PCycle1 = exp(XB)/(1+exp(XB))
   PCycle2 = exp(XB - 0.2394)/(1+exp(XB - 0.2394))
   PCycle3 = exp(XB - 0.4110)/(1+exp(XB - 0.4110))
   PCycle4 = exp(XB - 0.5628)/(1+exp(XB - 0.5628))

2

$\text{PCycle5} = \exp(XB - 0.7189)/(1+\exp(XB - 0.7189))$
$\text{PCycle6} = \exp(XB - 0.8138)/(1+\exp(XB - 0.8138))$

5.  Calculate the predicted *cumulative* probability of a live-birth after 1, 2, 3,…., 6 completed IVF cycles:

$\text{CumPCycle1} = 1-(1-p1)$
$\text{CumPCycle2} = 1-((1-p1)*(1-p2))$
$\text{CumPCycle3} = 1-((1-p1)*(1-p2)*(1-p3))$
$\text{CumPCycle4} = 1-((1-p1)*(1-p2)*(1-p3)*(1-p4))$
$\text{CumPCycle5} = 1-((1-p1)*(1-p2)*(1-p3)*(1-p4)*(1-p5))$
$\text{CumPCycle6} = 1-((1-p1)*(1-p2)*(1-p3)*(1-p4)*(1-p5)*(1-p6))$

**Supplementary text 2.** McLernon post-treatment model.


**Table showing the predictors in the original McLernon post-treatment model** (McLernon *et al.*, 2016)**.**

| Name predictor | Description | Range of possible values |
|---|---|---|
| Age | Female age | 18 to 50 years |
| Duration | How long have you been trying to conceive? | 0 to 21 years |
| Previous | Have you been pregnant before? | 1 = No; 0 = Yes |
| Tubal | Do you have a problem with your tubes? | 1 = Yes; 0 = No |
| Eggs | How many eggs were collected on your first IVF cycle? | (1 to 28) |
| Treat | Was your first cycle IVF or ICSI? | (1 = ICSI; 0 = IVF) |
| Cryo | In your first cycle did you have embryos frozen? | (1 = Yes; 0 = No) |
| Stage | What type of embryo transfer did you have in your first fresh embryo transfer? | (No embryos transferred; Single cleavage stage; Single blastocyst stage; Double cleavage stage; Double blastocyst stage; Triple cleavage stage; Triple blastocyst stage) |


**Original post-treatment model formulas as described by McLernon et al.** (McLernon *et al.*, 2016)**:**

1. For the non-linear relation between Age and the probability of live birth, the following Age1, Age2 and Age3 equations are first calculated and then used in the XB equation below (point 4):

   - Age1 = max((Age-26)/k,0)**3+(11*max((Age-41)/k,0)**3-(15)*max((Age-37)/k,0)**3)/4;
   - Age2 = max((Age -31)/k,0)**3+(6*max((Age -41)/k,0)**3-(10)*max((Age -37)/k,0)**3)/4;
   - Age3 = max((Age -34)/k,0)**3+(3*max((Age -41)/k,0)**3-(7)*max((Age -37)/k,0)**3)/4; *k=15**(2/3), **means 'to the power of'*

2. For the non-linear relation between Year and the probability of live birth, the following Year1 equation is first calculated and then used in the XB equation below (point 4). The value Year = 0 is used for the most up to date predictions.

   - Year1 = max((Year+8)/k,0)**3+((4)*max((Year+1)/k,0)**3-(7)*max((Year+4)/k,0)**3)/(3); *k= 7**(2/3).*

3. For the non-linear relation between Eggs and the probability of live birth, the following Eggs1 equation is first calculated and then used in the XB equation below (point 4):

   - Eggs1=max((Eggs-3)/k,0)**3+((6)*max((Eggs-18)/k,0)**3-(15)*max((Eggs-9)/k,0)**3)/(9); *k= 15**(2/3).*

4. Calculate XB

   XB =   -1.7564 + 0.0362[a] + (0.0272*Age) + (-0.1556*Age1) + (0.3812*Age2) + (-1.0184*Age3) + (-0.0208*Duration) + (-0.0504*Previous) + (-0.2207*Tubal) + (0.0018*Year) +

2

$(0.0619*Year1) + (0.0630*Eggs) + (-0.0479*Eggs1) + (-0.0968*Treat) + (0.6490*Cryo) + Stage^b$

a) To inflate predictions to 2013 an additional 0.0362 is added.
b) Stage equals the following values depending on group chosen:
        If Double cleavage stage then Stage=0;
        If No embryos transferred then Stage= -1.0842;
        If Single cleavage stage then Stage= -0.5675;
        If Single blastocyst stage then Stage= 0.0684;
        If Double blastocyst stage then Stage= 0.5802;
        If Triple cleavage stage then Stage= 0.0218;
        If Triple blastocyst stage then Stage= 0.4547.

1. Calculate the predicted probability of live-birth after the first, second, …., sixth IVF cycle:

$PCycle1 = exp(XB)/(1+exp(XB))$
$PCycle2 = exp(XB - 0.1933)/(1+exp(XB - 0.1933))$
$PCycle3 = exp(XB - 0.3537)/(1+exp(XB - 0.3537))$
$PCycle4 = exp(XB - 0.5122)/(1+exp(XB - 0.5122))$
$PCycle5 = exp(XB - 0.6788)/(1+exp(XB - 0.6788))$
$PCycle6 = exp(XB - 0.7666)/(1+exp(XB - 0.7666))$

2. Calculate the predicted *cumulative* probability of a live-birth after 1, 2, 3,…., 6 complete IVF cycles:

$CumPCycle1 = 1-(1-p1)$
$CumPCycle2 = 1-((1-p1)*(1-p2))$
$CumPCycle3 = 1-((1-p1)*(1-p2)*(1-p3))$
$CumPCycle4 = 1-((1-p1)*(1-p2)*(1-p3)*(1-p4))$
$CumPCycle5 = 1-((1-p1)*(1-p2)*(1-p3)*(1-p4)*(1-p5))$
$CumPCycle6 = 1-((1-p1)*(1-p2)*(1-p3)*(1-p4)*(1-p5)*(1-p6))$

1

**Supplementary Text 3.** Recalibrated pre-treatment model.

**Recalibrated pre-treatment model formula**

*The included predictors and formulas 1, 2, 4 and 5 are unchanged to the original McLernon pre-treatment model (see Supplementary Text 1)*

3. Calculate XB.

$$XB = -0.3474^a - 0.9948 + 0.0362^b + (0.0275*Age) + (-0.1805*Age1) + (0.4553*Age2) + (-1.1990*Age3) + (-0.0295*Duration) + (-0.0772*Previous) + (-0.0957*Tubal) + (0.0492*Anovulation) + (-0.1005*MaleFactor) + (0.0602*Unexplained) + (0.2155*Treatment) + (0.0334*Year) + (-0.0370*Year1) + (0.2173*Year2).$$

a) To recalibrate the pre-treatment model, 0.3474 is subtracted from the intercept.
b) To inflate predictions to 2013 an additional 0.0362 is added.