

# Prediction of Unsuccessful Endometrial Ablation: Random Forest vs Logistic Regression

**Kelly Yvonne Roger Stevens** (✉ [kyr.stevens@gmail.com](mailto:kyr.stevens@gmail.com))

Catharina Ziekenhuis <https://orcid.org/0000-0003-3512-9181>

**Liesbet Lagaert**

Catharina Hospital: Catharina Ziekenhuis

**Tom Bakkes**

University of Technology Eindhoven: Technische Universiteit Eindhoven

**Malou Evi Gelderblom**

Catharina Hospital: Catharina Ziekenhuis

**Saskia Houterman**

Catharina Hospital: Catharina Ziekenhuis

**Tanja Gijsen**

Elkerliek Hospital: Elkerliek Ziekenhuis

**Benedictus Schoot**

Catharina Hospital: Catharina Ziekenhuis

---

## Original Article

**Keywords:** Endometrial ablation, Machine Learning, Random Forest

**Posted Date:** December 22nd, 2020

**DOI:** <https://doi.org/10.21203/rs.3.rs-124914/v1>

**License:**   This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

---

1 **Prediction of unsuccessful endometrial ablation: Random Forest vs Logistic Regression**

2 Stevens, Kelly Yvonne Roger<sup>1,5</sup>; Lagaert, Liesbet; <sup>1,5</sup> Bakkes, Tom<sup>2</sup>; Gelderblom, Malou Evi<sup>1</sup>;  
3 Houterman, Saskia<sup>3</sup>; Gijsen, Tanja<sup>4</sup>; Schoot, Benedictus<sup>1,5</sup>.

4

5 <sup>1</sup> *Department of Obstetrics and Gynaecology, Catharina Hospital, Eindhoven, the Netherlands*

6 <sup>2</sup> *Department of Department of Electrical Engineering, Biomedical Diagnostics lab, TU Eindhoven, the*  
7 *Netherlands*

8 <sup>3</sup> *Department of Education and Research, Catharina Hospital, Eindhoven, the Netherlands*

9 <sup>4</sup> *Department of Obstetrics and Gynaecology, Elkerliek Hospital, Helmond, the Netherlands*

10 <sup>5</sup> *Women's Clinic, Ghent University Hospital, Ghent, Belgium*

11

12 **Email:**

13 K.Y.R. Stevens: [Kyr.stevens@gmail.com](mailto:Kyr.stevens@gmail.com)

14 L. Lagaert: [Liesbet.lagaert@gmail.com](mailto:Liesbet.lagaert@gmail.com)

15 T Bakkes: [t.h.g.f.bakkes@tue.nl](mailto:t.h.g.f.bakkes@tue.nl)

16 M.E. Gelderblom [Malou@gelderblom.info](mailto:Malou@gelderblom.info)

17 S. Houterman: [Saskia.houterman@catharinaziekenhuis.nl](mailto:Saskia.houterman@catharinaziekenhuis.nl)

18 T. Gijsen [tgijisen@elkerliek.nl](mailto:tgijisen@elkerliek.nl)

19 B.C. Schoot: [Dick@schoot.com](mailto:Dick@schoot.com)

20

21 **Corresponding author contact information:**

22 K.Y.R. Stevens, Department of Obstetrics and Gynaecology, Catharina Hospital,

23 Present address: Michelangelolaan 2, 5623 EJ Eindhoven, the Netherlands

24 Email: kyr.stevens@gmail.com

25 Telephone number: 0031620962469

26 ORCID: 0000-0003-3512-9181

27

28 L. Lagaert, Department of Obstetrics and Gynaecology, Catharina Hospital,

29 Present address: Michelangelolaan 2, 5623 EJ Eindhoven, the Netherlands

30 Email: [Liesbet.lagaert@gmail.com](mailto:Liesbet.lagaert@gmail.com)

31 Telephone number: 0032494854176

32

33

34 Partly presented as abstract at:

35 - The 28<sup>th</sup> Annual International Congress of the European Society of Gynaecological  
36 Endoscopy, Thessaloniki, Greece, 2019, Oral (preliminary results)

37 - The 48<sup>th</sup> Global congress of the American Association of Gynaecologic  
38 Laparoscopists, Vancouver, Canada, 2019, Poster presentation (preliminary results)

39

40

41

42

43 Abstract

44 **Background** Five percent of premenopausal women experience abnormal uterine  
45 bleeding. Endometrial ablation (EA) is one of the treatment options for this common problem.  
46 However, this technique shows a decrease in patient satisfaction and treatment efficacy on the long  
47 term.

48 **Study Objective:** To develop a prediction model to predict surgical re-intervention (for example re-  
49 ablation or hysterectomy) within two years after endometrial ablation (EA) by using Machine  
50 Learning (ML). The performance of the developed prediction model was compared with a previously  
51 published multivariate logistic regression model (LR).

52 **Design** This retrospective cohort study, with a minimal follow up time of two years, included 446  
53 pre-menopausal women (18+) that underwent an EA for complaints of heavy menstrual bleeding.  
54 The performance of the ML - and the LR model was compared using the area under the Receiving  
55 Operating Characteristic (ROC) curve.

56 **Results:** We found out that the ML model (AUC of 0.65 (95% CI 0.56-0.74)) is not superior compared  
57 to the LR model (AUC of 0.71 (95% CI 0.64-0.78)) in predicting the outcome of surgical re-  
58 intervention within two years after EA.

59 **Conclusion** Although Machine Learning techniques are gaining popularity in development of clinical  
60 prediction tools, this study shows that ML is not necessarily superior to the traditional statistical LR  
61 techniques. The performance of a prediction model is influenced by the sample size, the number of  
62 features of a dataset, hyperparameter tuning and the linearity of associations. Both techniques  
63 should be considered when developing a clinical prediction model.

64 **Key words:** Endometrial ablation, Machine Learning, Random Forest

65

66 **Article**

67 **Introduction**

68 Five percent of premenopausal women has complaints of abnormal uterine bleeding (1). Endometrial  
69 ablation (EA) is one of the treatment options for this common complaint. Due to the low costs and  
70 less invasive nature of this procedure (lower intra-operative complication risks, shorter recovery  
71 time, and lower post-operative morbidity), this form of treatment seems to be a less-invasive surgical  
72 treatment for menorrhagia compared to hysterectomy (2–6). However, long-term follow up shows a  
73 decrease in patient satisfaction and treatment efficacy. Due to permanent relief, the more invasive  
74 hysterectomy remains the most effective treatment of abnormal uterine bleeding (7–14).

75 According to literature, several factors prior to endometrial ablation appear to have an influence on  
76 the success-rate of this procedure. Younger age, complaints of dysmenorrhea, multiparity, a thicker  
77 pre-procedural endometrium, a duration of menstruation above seven days, presence of an  
78 intramural leiomyoma on transvaginal sonography, a history of sterilization or caesarean section, and  
79 a longer uterine depth are some of the possible negative influencing factors (1,2,8,9,11–18).

80 To optimize the counselling of patients with abnormal uterine bleeding, a prediction model based on  
81 the combined influence of the above-mentioned predictors could provide a better insight into the  
82 individual prognosis of endometrial ablation. In times of personalised medicine this can create better  
83 individual care leading to fewer re-interventions, lower healthcare costs and more patient  
84 satisfaction. With the use of a prediction model shared decision making can be optimized (19).

85 For this reason Stevens et al. (16) developed two multivariate prediction models to help counsel  
86 patients for failure of EA and for surgical re-intervention within two years after EA. The developed  
87 prediction models have a clinically acceptable c-index of 0.68 and 0.71 respectively. In addition,  
88 Stevens et al. is performing an external validation of these models, results of these data will follow.

89 In the field of gynaecology, many prediction models are developed using statistic multivariate logistic  
90 regression as a standard approach, these are based on a combination of various predictors that are  
91 significantly related to the outcome of interest. However, this method cannot automatically estimate  
92 the interconnection between predictors and in this way can overestimate the influence of an  
93 individual predictor (20,21).

94 We were also interested in other techniques of developing a prediction model. In recent years  
95 Machine Learning (ML) methods have been increasingly used in the development of clinical  
96 prediction models. ML is a scientific discipline that focuses on models that directly and automatically  
97 learn from data without using pre-identified statistical parameters and without assumption of a  
98 preconceived relationship between predictors and outcomes (20,22). A potential advantage of  
99 Machine Learning methods compared to the traditional statistical strategies is the possibility of  
100 capturing complex, nonlinear relationships in the data (23,24). ML algorithms use training data with  
101 well-defined input and output variables. This gives the opportunity to define a model with predictors  
102 which can be used for new and similar data. Compared to statistical logistic regression models, this  
103 can be done without a priori assumption of relevant variables (25). We chose surgical re-  
104 intervention as most objective outcome measure to compare both prediction models in predicting  
105 unsuccessful endometrial ablation.

106 The aim of the study was to develop a Machine Learning model to predict the chance of surgical re-  
107 intervention (for example re-ablation or hysterectomy) within two years after EA. Furthermore, we  
108 compared the performance of the ML model with the prediction by the previously published  
109 multivariate logistic regression re-intervention model of Stevens et al (16).

110

111

112 **Methods:**

113 This study used the same dataset as was used to develop the prediction models in the study from  
114 Stevens et al. , the full study protocol can be consulted there. (16)

115 This retrospective two-centred cohort study, performed in two non-university teaching hospitals in  
116 the Netherlands (Catharina Hospital, Eindhoven; Elkerliek Hospital Helmond), included 446 patients  
117 who have had an EA for complaints of abnormal uterine bleeding (16). Both hospitals used similar  
118 ablation techniques between 2004 and 2013, being Cavatherm® (Veldana Medical SA, Morges,  
119 Switzerland), Gynecare Thermachoice® (Ethicon, Sommerville, US) and Thermablate® EAS (Idoman,  
120 Ireland). Recent publications have shown that these ablation techniques were equally effective  
121 (14,26). Local medical ethical review boards approved the study. All patients gave informed consent.

122 Patients were identified in the Electronic Patient care System by using specified search terms related  
123 to endometrial ablation. Exclusion criteria were a postmenopausal status at time of EA; (suspicion of)  
124 endometrial malignancy or uterine cavity deformations (adenomyosis; anomalies; fibroids; or a  
125 polyp). Follow-up period after treatment was at least two years. This time-interval was chosen  
126 because previous literature stated that most re-interventions were done within two years. Follow-up  
127 ended on the day of hysterectomy, in case of death or on April 15, 2015 (9,17,18,26–28).

128 Data were extracted from individual patient files by two researchers (K.S. & D.M. (16)). Next, patients  
129 were asked to fill in a questionnaire regarding follow-up information. In case of non-response,  
130 patients were contacted by letter and ultimately by telephone by the authors of Stevens et al (16).  
131 The used questionnaire contained questions based on significant variables predicting surgical re-  
132 intervention after EA that were previously published (2,5,8,11–17, 31,32).

133 The entire dataset consists of 446 patients with different categorical and continuous variables. For  
134 the Machine Learning algorithms all features were extracted from the original dataset of Stevens et  
135 al. (16) A total of five pre-operative variables was used to develop the Machine Learning model. This  
136 were the preoperative variables that were significant predictors in the final multivariate re-

137 intervention model of Stevens et al. (age, duration of menstruation, dysmenorrhea, parity and  
138 previous caesarean section) (16). The continuous data were not discretized into categories as was  
139 done in the development of the previously published logistic regression model(16).

140

#### 141 *Development of the Logistic regression model*

142 Statistical analysis of the data was performed by using SPSS 21.0 for Windows (IBM Corp., Armonk,  
143 NY, USA).

144 To determine which variables were significant, univariable logistic regression analysis was used.  
145 The variables with a p-value  $<.10$  were used in the multivariable analysis. This was followed by a  
146 backward stepwise manual selection process, progressively excluding the variable with the highest p-  
147 value (16).

148 As described by Steyerberg et al., the p-value of 0.10 was used to prevent a potential incorrect  
149 exclusion of a predictive factor. This would be far more detrimental for the test than missing a  
150 potential discriminating factor (31,32).

151 Multicollinearity and interaction between the significant variables in the model was tested. Bootstrap  
152 resampling was used for internal validation (n=5000) (32,33). To correct for over-optimism of the  
153 model, regression coefficients were multiplied by the calculated shrinkage factor. A detailed  
154 description of the development of the LR-model van be found in the study of Stevens et al. (16).

155

#### 156 *Development of the Machine Learning model (Random Forest model)*

157 For the development of the Machine Learning model, we used a Random Forest (RF) technique. This  
158 is a Machine Learning method used for classification and regression, which operates by constructing  
159 a large ensemble of decision trees on training data (22,23,34). Each tree in the Random Forest is built



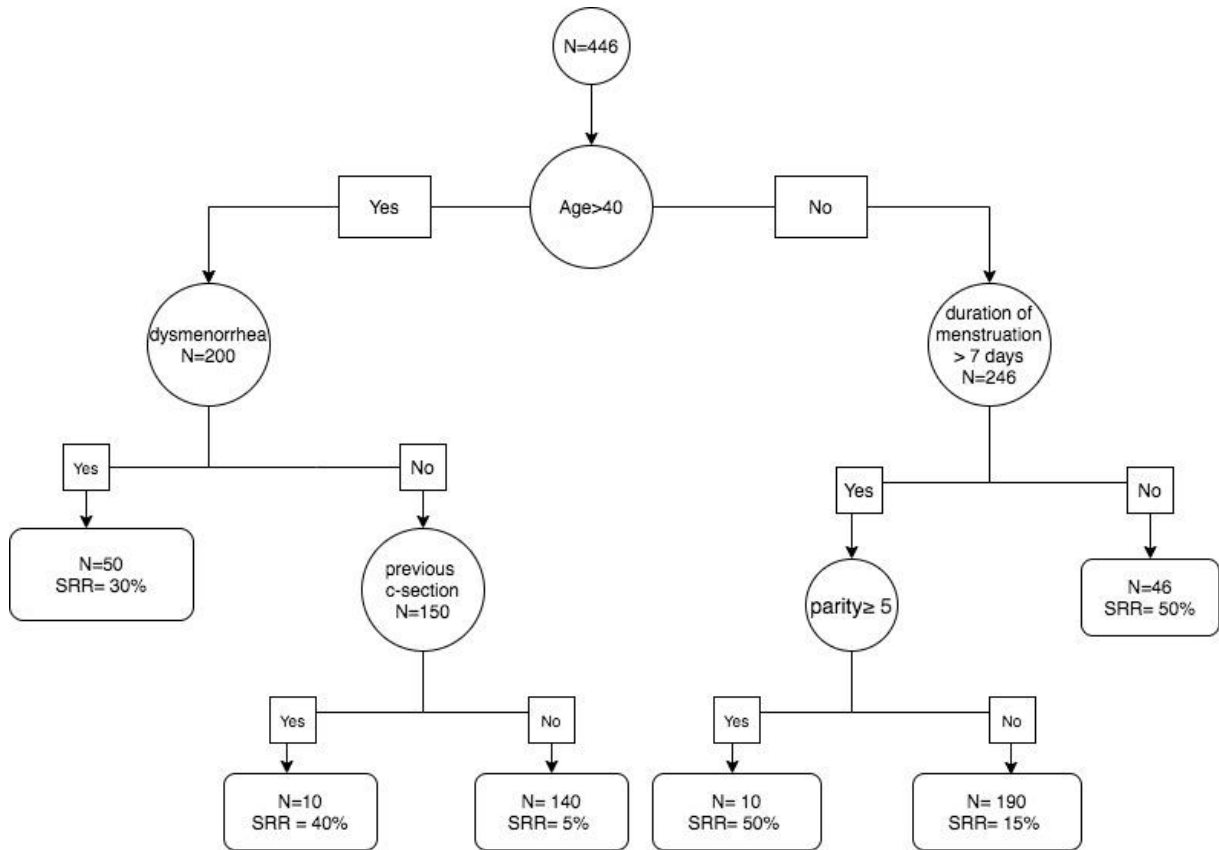
160 using a bootstrap sample randomly drawn from a training dataset. This results in a reduction of  
161 variance and corrects for a single decision trees ability to overfit to a training set. Each tree in the  
162 forest gives an individual prediction on the outcome measure. For a classification problem (in this  
163 case, surgical re-intervention or no surgical re-intervention after EA) the final Random Forest model  
164 averages the prediction of all the trees in the forest (21,23,34,35).

165 Making the model, we first trained a RF model using the five following pre-operative predictors: age,  
166 duration of menstruation, dysmenorrhea, parity and previous caesarean section. These factors were  
167 associated with a higher probability of surgical re-intervention within two years after EA in the  
168 previously published multivariate logistic regression model (16).

169 As described above, a RF model is an ensemble of many decision tree models. When building  
170 decision trees, each tree in the forest uses random samples (patients) from the training set (“tree  
171 bagging”). Figure 1 shows an example of an individual decision tree in the Random Forest. The  
172 decision tree is a flowchart-like binary branch structure. At each ‘node split’ in the tree the data are  
173 divided in two, based on the value of variable of the decision node. If no more splits are possible a  
174 prediction will be calculated for the cases in the final leaf node (23,34,36).

175 At each node split a random subset of features (such as duration of menstruation and parity) is  
176 considered (“feature bagging”), this to avoid over-selection of strong predictive features, leading to  
177 similar splits in the trees. This finally leads to a robust model and prevents model overfitting  
178 (21,23,34–37).

179



180

181 *Figure 1. An illustration of a decision tree in the Random Forest model. The decision tree directs each case from the root node to the leaf*  
 182 *nodes, resulting in a prediction.*

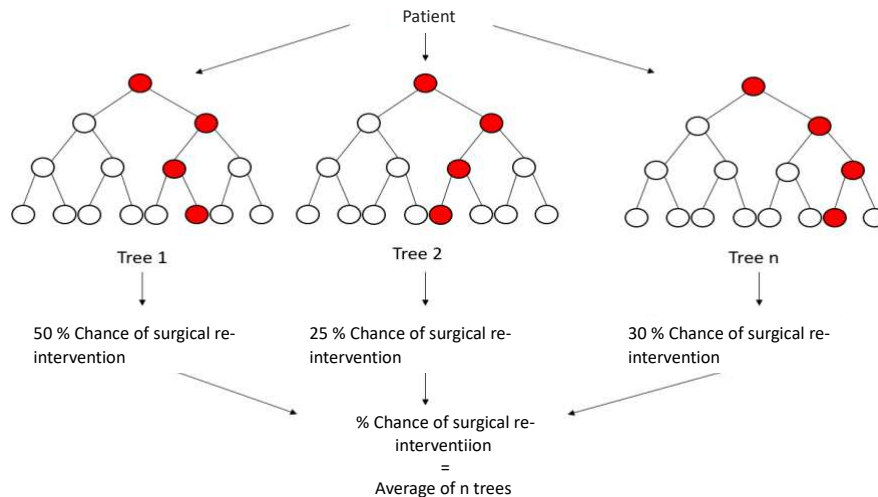
183 *N= Number, SRR = Surgical Re-intervention Rate.*

184

185 Following this process, the classification result of a RF model is produced by computing a large  
 186 ensemble of those trees and averaging the prediction of each single decision tree on surgical re-  
 187 intervention. Figure 2 shows a simplified example of the RF model. In practice, the decision trees and  
 188 the resulting prediction model contain a large number of leaf nodes(34,38).

189

190



191

192 *Figure 2: A simplified Random Forest model for the prediction of the surgical re-intervention .*

193

194 The RF was trained in MATLAB (2018b) using the TreeBagger function in the Statistics and Machine  
 195 Learning Toolbox.

196 We began running the RF module with default parameter values before starting to improve the RF's  
 197 performance by hyperparameter optimization. Default parameters are pre-set values for the  
 198 hyperparameters on which the construction of the decision trees is based, for example 500 for ntree  
 199 ( number of trees in the forest) (34,35). Hyper-parameter optimization refers to the automatic  
 200 optimization of the hyper-parameters of a ML model. Hyper-parameters are all the parameters of a  
 201 model that are used to configure the model (e.g. minimum leaf size, number of splits, ntree and  
 202 mtry, which are the number of features randomly selected as candidate feature at each split "feature  
 203 bagging").

204 To predict the chance of surgical re-intervention within two years after EA, the model was initially  
 205 trained and internally validated on the 446 cases. To make a good comparison between de RF and LR  
 206 the same validation technique was used. Therefore, a bootstrap resampling of 5000 was used to

207 make training bags and test bags. The performance measure Area Under The Receiver Operating  
208 Curve (AUROC) was calculated on the test sets and averaged for the 5000 bootstrap samples.

209

#### 210 *Comparison of the prediction models*

211 The performance of the models was tested and compared using the AUROC. Accuracy was not used  
212 as performance measure, since the database is unbalanced (ratio between re-intervention and no re-  
213 intervention 1:8 (53:446)) (43). It was chosen to use the performance measures (AUC) as used in the  
214 previous study of Stevens et al (16). In this way a good comparison can be made.

215

#### 216 Predictors of surgical re-intervention: Variable importance measure (VIM)

217 To identify important predictors of surgical re-intervention we used two methods for analysis.  
218 First, a statistical univariate logistic regression analysis was applied to assess the importance of each  
219 variable. For each variable an odds ratio (OR) with a 95% confidence interval (CI) was calculated.  
220 Secondly, a permutation-based variable importance was used. This VIM is based on AUC statistic of  
221 the RF model. The AUC statistic is computed by randomly permutating the values of predictor x, and  
222 comparing the resulting AUC to the not permuted AUC. Leaving out an important feature will result  
223 in a lower AUC of the RF model, while leaving out an unimportant feature will not change the AUC  
224 significantly (23,38,41).

225

#### 226 **Results**

227 Seven hundred sixty-two patients were identified retrospectively. Thirty-three patients were  
228 excluded, thirty did not meet the inclusion criteria and three underwent an incomplete endometrium

229 ablation. The remaining 729 patients were contacted, resulting in a response-rate of 61% (N = 446).

230 A total amount of 446 patients was available for analysis (16).

231

232 Fifty-three (11.9%) of these patients required a surgical re-intervention within two years after EA.

233 Patients mean age during their EA was 43.8 years (SD  $\pm$ 5.5, range 20-55, missing values 0). The mean

234 number of parity was 2.2 (SD  $\pm$  1.0, missing values 0). Sixty-one (13.7%) of the patients underwent a

235 caesarean section. The mean number of previous caesarean section was 0.2 (SD  $\pm$  0.6, missing values

236 0)

237 Hundred sixty-nine (39.4%) of the patients had a menstruation period longer than seven days, the

238 mean number of menstrual days was 9.4 (SD  $\pm$  6.0, missing values 17). Two hundred fifty-six (57.4%)

239 of the patients had complaints of dysmenorrhea and four hundred thirty-four (97.3%) of the patients

240 had complaints of abnormal uterine bleeding (16).

241

242 Prediction models:

243 Logistic regression model

244 Univariate analysis showed six significant predictors, multivariate analyses resulted in a logistic

245 regression model consisting of five significant predictors: age (OR 0.95, 95% CI 0.90 – 1.00), duration

246 of menstruation >7 days (OR 2.05, 95% CI 1.10 – 3.82), dysmenorrhea (OR 2.48, 95% CI 1.21 – 5.07),

247 parity  $\geq$ 5 (OR 7.63, 95% CI 1.51 – 38.46), and previous caesarean section (OR 2.21, 95% CI 1.05 –

248 4.64). The AUC of the final prediction model after correcting by the shrinkage factor was 0.71 (95% CI

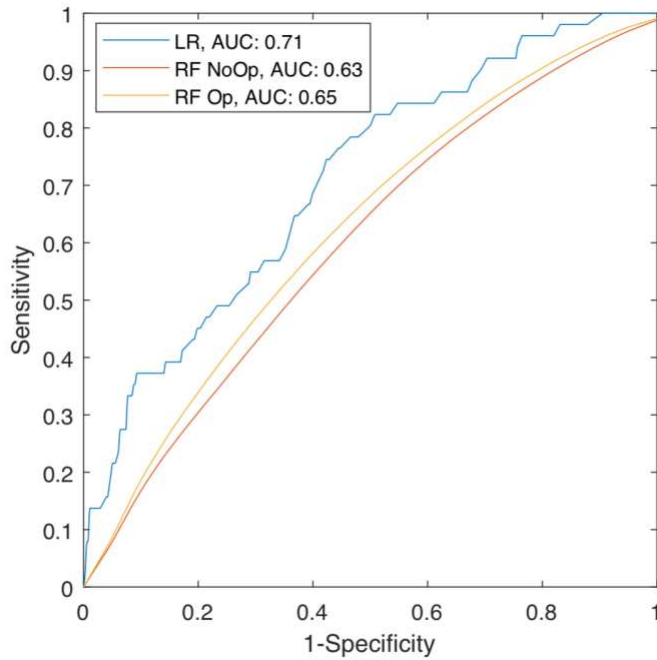
249 0.64-0.78) (Figure 4).

250 The final model is described in the article of Stevens et al (16).

251

252 Random forest model

253 The Random Forest method resulted in a model which predicts the chance of re-intervention within  
254 two years after EA with an AUC of 0.63 (95% CI 0.54-0.71). An AUC of 0.65 (95% CI 0.56-0.74) was  
255 achieved after optimization of this model (Figure 4).



256

257 *Fig4. ROC-curve of the logistic regression and Random Forest model. LR AUC 0.71 (95% CI 0.64-0.78), NoOp AUC 0.63 (0.54-0.71), Op AUC:*  
258 *(0.56 – 0.74)*

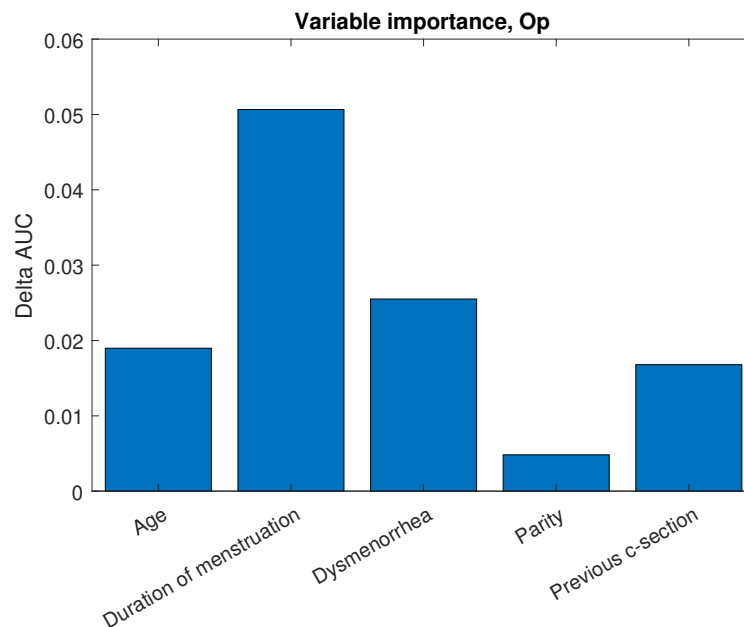
259 *LR= logistic regression, RF= Random Forest, Op= after hyperparameter optimization, NoOp= before hyperparameter optimization*

260

261 Predictors of surgical re-intervention: Variable importance

262 The AUC was used to quantify the importance of the predictor. For each RF model, the AUC was  
263 calculated on the test set. Then the same was done after permuting each predictive variable. By  
264 calculating the difference between the permuted and non-permuted AUC, the importance of each  
265 individual predictor can be quantified (Figure 5). The difference in AUC for the different predictors in  
266 the optimized model were in ascending order of importance: 0.005 for parity, 0.017 for previous  
267 caesarean section, 0.019 for age, 0.026 for dysmenorrhea and 0.051 for duration of menstruation.

268 This means dysmenorrhea and duration of menstruation have the highest impact on the AUC of the  
269 RF model. (Figure 5)



270

271 *Figure. 5. Contribution of predictors of surgical re-intervention within 2 years after endometrial ablation, after hyperparameter optimization*

272

273

## 274 **Discussion**

### 275 Main findings

276 In this study, a ML model was made using Random Forest technique to predict surgical re-  
277 intervention within two years after EA. Comparison of the predictive performance of the RF model  
278 with the existing logistic regression model of Stevens et al. was made (16).

279 The existing logistic regression model has a C-index of 0.71 (95% CI 0.64-0.78) (16). The ML model,  
280 developed in this study, shows a C-index of 0.65 (95% CI 0.56-0.74) after hyperparameter  
281 optimization. This shows that the LR prediction model developed by Stevens et al. (16) probably  
282 performs better in predicting surgical re-intervention within two years after EA than the newly

283 developed RF model. However, this difference in performance is not statistically significant when we  
284 look at the confidence intervals.

285 In the LR model, high parity ( $\geq 5$ ) is a predictive variable for surgical re-intervention. This can be  
286 related to the larger uterine cavity of grand multiparous women. However, when considering our RF  
287 model, parity has no large impact on the AUC. This is in line with previously reported studies that  
288 show no significant increased risk of treatment failure with increasing parity (1,15).

289 Previous caesarean section is also related to higher rates of surgical re-intervention which can be  
290 explained by irregularity of the uterine wall caused by the uterine scar (44). This can inhibit complete  
291 contact of the ablation device with the uterine wall, leading to residual active endometrium.

292 In our cohort, pre-operative dysmenorrhea is associated with a higher risk of surgical re-intervention.  
293 There is evidence that gynaecologic pathology causing this dysmenorrhea (adenomyosis and  
294 endometriosis) reduces the success of endometrial ablation (8,17,30,45,46). This can be explained by  
295 the fact that EA is not an appropriate treatment for these diseases due to the superficial effect of  
296 energy to the uterine wall of ablation. It could help to diagnose these diseases before performance of  
297 EA. However, sensitivity and specificity of the diagnostic tools for determining these diseases in the  
298 pre-operative setting are still low (47).

299 In line with previous studies, we found that younger age was associated with a higher risk of surgical  
300 re-intervention (7,9–13,29).

301 The duration of menstruation  $> 7$  days is also a negative predictive factor for surgical re-intervention  
302 after EA. This may be caused by a thicker endometrium which is more difficult to completely remove  
303 by the device (7,10).

304

305 Interpretation in light of other evidence



306 There are several possible reasons to understand why the LR model probably performs better  
307 compared to the ML model.

308 Firstly, ML tends to work better for variables with strong predictive power (20,48). We observed that  
309 most of the candidate predictors in this model have low predictive power. The variables parity, age  
310 and previous c-section show low predictive power. The difference in area under the curve for these  
311 predictors that was produced using the permutation based variable importance was  $<0.02$ . There are  
312 different reasons to explain that this specific dataset, and its separate and combined predictors  
313 appeared to have a low predictive power. On one hand, the outcome can be unpredictable, meaning  
314 these candidate predictors have little influence on the outcome measure. On the other hand, the  
315 dataset can be too small to identify the predictive power of a candidate predictor. A larger dataset  
316 could possibly identify more predictors (20,48).

317 Secondly, some studies demonstrate that ML is performing better when a larger set of potential  
318 predictors are used in the prediction model. There seems to be an influence of the number of  
319 predictors ( $p$ ) and the ratio of  $p:n$  (sample size). RF tends to perform better for increasing  $p$  and  $p:n$ .  
320 (20,24,49,50) In our study, to limit potential bias, the five identical predictors as published before  
321 (16) were considered for the LR and RF algorithms. We did this to allow a fair comparison between  
322 the two models, probably in disadvantage of the RF model (20,24,49,50).

323 Another possible reason for a lower AUC of the RF model is the necessity of big datasets to reach an  
324 optimal performance. A dataset with 446 participants might be too small for ML to make robust  
325 conclusions. For LR however, this number of patients can be enough to develop a prediction model.

326 Besides that, we didn't discretize the continuous variables in the RF model, we found some literature  
327 suggesting discretization of variables can improve the classification performance (51).

328 Finally, we can also consider that for this clinical problem a logistic approach is better than a RF  
329 model for modelling the relationship between surgical re-intervention and the explanatory variables.

330 Probably the previously mentioned complex, nonlinear relationships that a ML approach can better  
331 capture are not present in this dataset.

332

### 333 Strengths and limitations

334 The predictors obtained by univariate and multivariate logistic regression are in accordance with the  
335 existing literature (1,8,10–15,17,51). However, when we compare the variable importance between  
336 the OR (LR) and the difference in AUC (ML) of each variable, we identify a different ranking in  
337 variable importance.

338 The difference in ranking of variable importance is a limitation of the study because there is no  
339 proper way to compare the importance of each predictor on surgical re-intervention between the RF  
340 and LR model. For the LR model the OR is defined for each predictor X as the odds of a surgical re-  
341 intervention in participants having predictor X over participants not having predictor X (Beta). While  
342 for the RF model the variable importance is defined as the difference in AUC when predictor X is not  
343 permuted.

344 Dysmenorrhea (OR 2.48) and a parity>5 (OR 7.63) have the highest odds ratio in the multivariate LR  
345 analysis, while for the RF model the duration of menstruation and dysmenorrhea are the most  
346 important variables. We consider two possible reason for the difference in importance. The first  
347 reason is that for the LR model all continuous variables (except age) were discretized, while for the  
348 RF model continuous variables were handled. A second reason is that in the LR the predictors have  
349 different units, and these were not standardized. This means that a subjective assessment of variable  
350 importance cannot easily be made by simply comparing the raw sizes of the OR (21,23,34,48). This  
351 can be seen as a strength of our study since the difference in AUC for each predictor (permuted vs.  
352 not permuted) reflects the variable importance in a standardized way.

353 We used bootstrap resampling for internal validation (n=5000) in the LR and RF model. Using the  
354 same validation method limits potential bias.  
355 Furthermore, the same predictors were considered for the LR and ML algorithms. This limits  
356 potential bias, but will limit the potential power of a RF technique as well.  
357 Another important strength of this study is the use of all participants in evaluating the performance  
358 of the RF model. By using the test sets, there is no need for an independent validation dataset.  
359 It could be seen as a limitation of this study that we did not perform an external validation in another  
360 cohort. However, we did not expect it to be significantly better in performance, since the internal  
361 validation of the RF did not perform better than the logistic regression model. In addition, an  
362 external validation for the logistic regression model is being performed at the time of this study.  
363 Finally, we can state that ML models are in our experience not easily implemented in the clinical  
364 practice; since these are often not available in commonly used software packages in clinical practice.  
365 However, future structured data-registration is increasing, which makes it easier to create big  
366 datasets available for ML-programs.

367

#### 368 **Conclusion:**

369 In conclusion we can state that for the prediction of surgical re-intervention within two years after  
370 EA, the logistic regression model gives a better prediction compared to the Machine Learning model.  
371 However, Machine Learning algorithms should always be considered as candidate prediction tool for  
372 classification or regression problems because of the possible advantages. So far there is no evidence  
373 for one single algorithm that outperforms the other in general use. Further research is needed for  
374 the evaluation of Machine Learning based predictive modelling.

375

376 **Declarations**

377 **Disclosure of interests:** There are no conflicts of interest to disclose.

378 **Contribution to authorship**

379 K.Y.R. Stevens: Project development, Data collection/management, Data analysis,  
380 Manuscript writing/editing

381 L. Lagaert: Project development, Data collection/management, Data analysis,  
382 Manuscript writing/editing

383 T. Bakkes: Development of Random Forest Model (Machine learning)

384 M. Gelderblom: Manuscript editing

385 S. Houterman: Manuscript editing

386 T. Gijzen: Data collection, Manuscript editing

387 B.C. Schoot: Project development, Data collection, Manuscript editing

388

389 **Details of ethics approval**

390 All methods were carried out in accordance with relevant guidelines and regulations. The data  
391 collection was done in the first study (development of LR model) performed by Stevens et al (16).

392 This study was approved by the local medical ethical review board of Catharina hospital and Elkerliek  
393 hospital. All patients gave informed consent. For this second study (using the same data), the ethical  
394 board in the Catharina hospital and in the Elkerliek hospital concluded this ethics approval was valid.

395 **Availability of data and material:**

396 The datasets generated and analysed during the current study are not publicly available due to  
397 privacy, but they are available from the corresponding author on a reasonable request.

398 **Funding:** None

399

400 **Acknowledgements:**

401 The authors want to thank the patients for completing the questionnaires and for  
402 consenting to participate in our study.

403 **Consent for publication:** Not applicable

404 **Availability of data and materials:**

405 The datasets generated and analyzed during the current study are not publicly available due to  
406 privacy, but they are available from the corresponding author on a reasonable request.

407

408

409 **Literature references**

410

411 1. Peeters JAH, Penninx JPM, Mol BW, Bongers MY. Prognostic factors for the success of  
412 endometrial ablation in the treatment of menorrhagia with special reference to previous  
413 cesarean section. *Eur J Obstet Gynecol Reprod Biol* [Internet]. 2013 Mar [cited 2018 Dec  
414 3];167(1):100–3. Available from:  
415 <https://linkinghub.elsevier.com/retrieve/pii/S0301211512005301>

416 2. Waddell G, Pelletier J, Desindes S, Anku-Bertholet C, Blouin S, Thibodeau D. Effect of  
417 Endometrial Ablation on Premenstrual Symptoms. *J Minim Invasive Gynecol* [Internet]. 2015  
418 May [cited 2018 Dec 3];22(4):631–6. Available from:  
419 <https://linkinghub.elsevier.com/retrieve/pii/S1553465015000886>

420 3. Laberge P, Leyland N, Murji A, Fortin C, Martyn P, Vilos G, et al. Endometrial Ablation in the  
421 Management of Abnormal Uterine Bleeding. *J Obstet Gynaecol Canada*. 2015;

422 4. Bouzari Z, Yazdani S, Azimi S, Delavar MA. Thermal balloon endometrial ablation in the  
423 treatment of heavy menstrual bleeding. *Med Arch (Sarajevo, Bosnia Herzegovina)*.  
424 2014;68(6):411–3.

425 5. Miller J, Troeger KA, Lenhart GM, Bonafede M, Basinski CM, Lukes AS. Cost effectiveness of  
426 endometrial ablation with the NovaSure® system versus other global ablation  
427 modalities and hysterectomy for treatment of abnormal uterine bleeding: US commercial and  
428 Medicaid payer perspectives. *Int J Womens Health* [Internet]. 2015 Jan [cited 2018 Dec 3];59.  
429 Available from: [http://www.dovepress.com/cost-effectiveness-of-endometrial-ablation-with-](http://www.dovepress.com/cost-effectiveness-of-endometrial-ablation-with-the-novasurereg-system-peer-reviewed-article-IJWH)  
430 [the-novasurereg-system-peer-reviewed-article-IJWH](http://www.dovepress.com/cost-effectiveness-of-endometrial-ablation-with-the-novasurereg-system-peer-reviewed-article-IJWH)

431 6. Angioni S, Pontis A, Nappi L, Sedda F, Sorrentino F, Litta P, et al. Endometrial ablation: First-vs.  
432 second-generation techniques. *Minerva Ginecologica*. 2016.

- 433 7. El-Nashar SA, Hopkins MR, Creedon DJ, St Sauver JL, Weaver AL, McGree ME, et al. Prediction  
434 of treatment outcomes after global endometrial ablation. *Obstet Gynecol* [Internet]. 2009 Jan  
435 [cited 2018 Dec 3];113(1):97–106. Available from:  
436 <https://insights.ovid.com/crossref?an=00006250-200901000-00016>
- 437 8. Wishall KM, Price J, Pereira N, Butts SM, Della Badia CR. Postablation risk factors for pain and  
438 subsequent hysterectomy. *Obstet Gynecol* [Internet]. 2014 Nov [cited 2018 Dec  
439 3];124(5):904–10. Available from: [https://insights.ovid.com/crossref?an=00006250-](https://insights.ovid.com/crossref?an=00006250-201411000-00007)  
440 [201411000-00007](https://insights.ovid.com/crossref?an=00006250-201411000-00007)
- 441 9. Thomasee MS, Curlin H, Yunker A, Anderson TL. Predicting pelvic pain after endometrial  
442 ablation: which preoperative patient characteristics are associated? *J Minim Invasive Gynecol*  
443 [Internet]. 2013 Sep [cited 2018 Dec 3];20(5):642–7. Available from:  
444 <https://linkinghub.elsevier.com/retrieve/pii/S1553465013001957>
- 445 10. Bongers MY, Mol BWJ, Brölmann HAM. Prognostic factors for the success of thermal balloon  
446 ablation in the treatment of menorrhagia. *Obstet Gynecol* [Internet]. 2002 Jun [cited 2018  
447 Dec 3];99(6):1060–6. Available from: <http://www.ncbi.nlm.nih.gov/pubmed/12052600>
- 448 11. Longinotti MK, Jacobson GF, Hung Y-Y, Learman LA. Probability of hysterectomy after  
449 endometrial ablation. *Obstet Gynecol* [Internet]. 2008 Dec [cited 2018 Dec 3];112(6):1214–20.  
450 Available from:  
451 [http://content.wkhealth.com/linkback/openurl?sid=WKPTLP:landingpage&an=00006250-](http://content.wkhealth.com/linkback/openurl?sid=WKPTLP:landingpage&an=00006250-200812000-00006)  
452 [200812000-00006](http://content.wkhealth.com/linkback/openurl?sid=WKPTLP:landingpage&an=00006250-200812000-00006)
- 453 12. Shaamash AH, Sayed EH. Prediction of successful menorrhagia treatment after thermal  
454 balloon endometrial ablation. *J Obstet Gynaecol Res* [Internet]. 2004 Jun [cited 2018 Dec  
455 3];30(3):210–6. Available from: <http://doi.wiley.com/10.1111/j.1447-0756.2004.00189.x>
- 456 13. Klebanoff J, Makai GE, Patel NR, Hoffman MK. Incidence and predictors of failed second-

- 457 generation endometrial ablation. *Gynecol Surg* [Internet]. 2017 Dec 15 [cited 2018 Dec  
458 3];14(1):26. Available from: [https://gynecolsurg.springeropen.com/articles/10.1186/s10397-](https://gynecolsurg.springeropen.com/articles/10.1186/s10397-017-1030-4)  
459 017-1030-4
- 460 14. Louie M, Wright K, Siedhoff MT. The case against endometrial ablation for treatment of heavy  
461 menstrual bleeding. *Curr Opin Obstet Gynecol* [Internet]. 2018 Aug [cited 2018 Dec  
462 3];30(4):287–92. Available from: [http://insights.ovid.com/crossref?an=00001703-900000000-](http://insights.ovid.com/crossref?an=00001703-900000000-99318)  
463 99318
- 464 15. Lybol C, van der Coelen S, Hamelink A, Bartelink LR, Nieboer TE. Predictors of Long-Term  
465 NovaSure Endometrial Ablation Failure. *J Minim Invasive Gynecol*. 2018;
- 466 16. Stevens KYR, Meulenbroeks D, Houterman S, Gijzen T, Weyers S, Schoot BC. Prediction of  
467 unsuccessful endometrial ablation: a retrospective study. *Gynecol Surg* [Internet].  
468 2019;16(1):7. Available from: <https://doi.org/10.1186/s10397-019-1060-1>
- 469 17. Shavell VI, Diamond MP, Senter JP, Kruger ML, Johns DA. Hysterectomy subsequent to  
470 endometrial ablation. *J Minim Invasive Gynecol* [Internet]. 2012 Jul [cited 2018 Dec  
471 3];19(4):459–64. Available from:  
472 <https://linkinghub.elsevier.com/retrieve/pii/S1553465012001161>
- 473 18. Kreider SE, Starcher R, Hoppe J, Nelson K, Salas N. Endometrial ablation: is tubal ligation a risk  
474 factor for hysterectomy. *J Minim Invasive Gynecol* [Internet]. 2013 Sep [cited 2018 Dec  
475 3];20(5):616–9. Available from:  
476 <https://linkinghub.elsevier.com/retrieve/pii/S1553465013001581>
- 477 19. van Montfort P, Smits LJM, van Dooren IMA, Lemmens SMP, Zelis M, Zwaan IM, et al.  
478 Implementing a Preeclampsia Prediction Model in Obstetrics: Cutoff Determination and  
479 Health Care Professionals' Adherence. *Med Decis Making*. 2020 Jan;40(1):81–9.
- 480 20. Evangelia christodoulou, Jie MA, Collins GS, Steyerberg EW, Verbakel JY, van Calster B. A



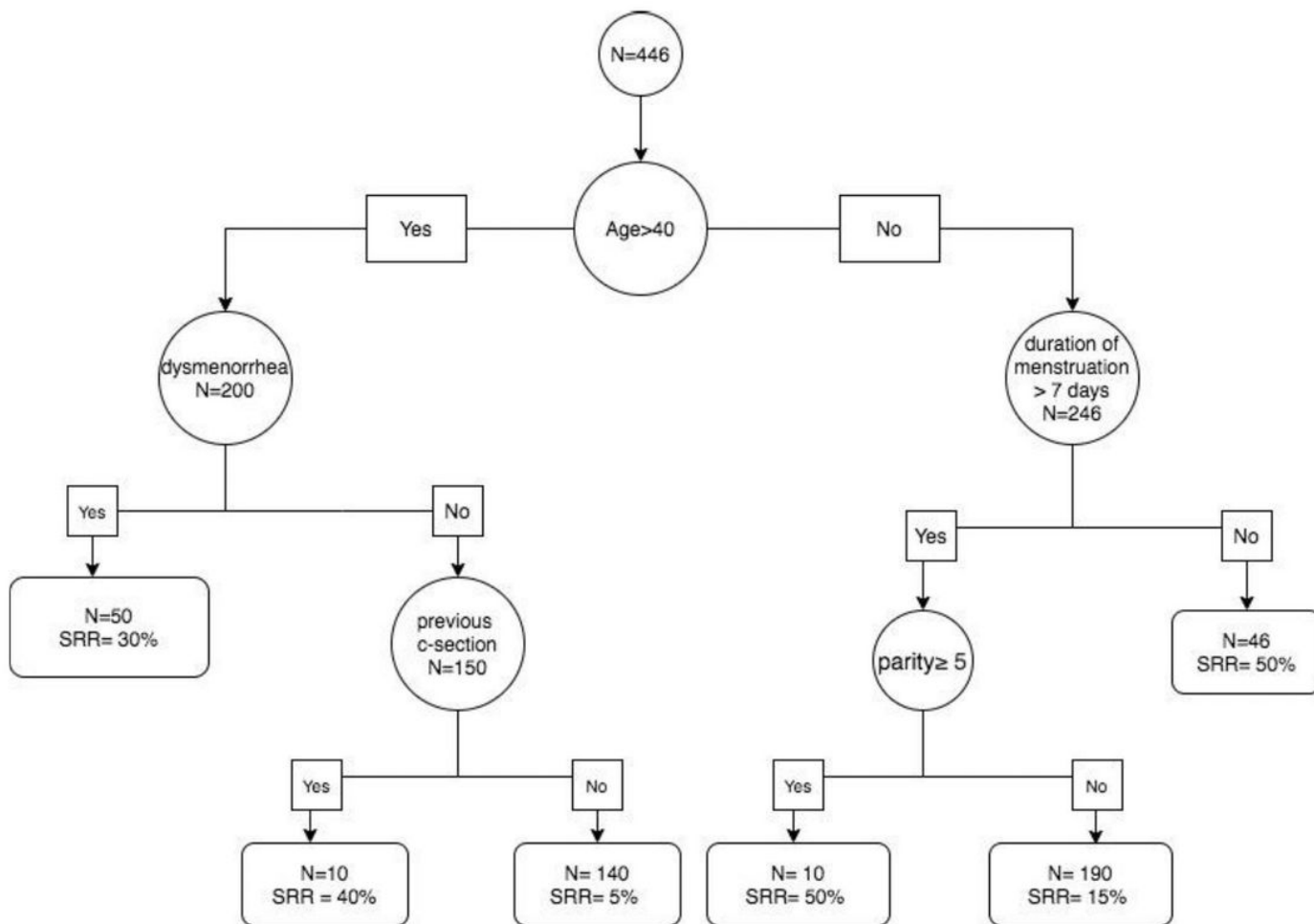
- 481 systematic review shows no performance benefit of machine learning over logistic regression  
482 for clinical prediction models. *J Clin Epidemiol.* 2019;
- 483 21. Breiman L. *Statistical Modeling: The Two Cultures.* *Stat Sci.* 2001;
- 484 22. Deo RC. *Machine learning in medicine.* *Circulation.* 2015;
- 485 23. Couronné R, Probst P, Boulesteix AL. *Random forest versus logistic regression: A large-scale*  
486 *benchmark experiment.* *BMC Bioinformatics.* 2018;
- 487 24. Chen JH, Asch SM. *Machine Learning and Prediction in Medicine — Beyond the Peak of*  
488 *Inflated Expectations.* *N Engl J Med.* 2017;
- 489 25. Panesar SS, D’Souza RN, Yeh FC, Fernandez-Miranda JC. *Machine Learning Versus Logistic*  
490 *Regression Methods for 2-Year Mortality Prognostication in a Small, Heterogeneous Glioma*  
491 *Database.* *World Neurosurg X.* 2019;
- 492 26. Sambrook AM, Bain C, Parkin DE, Cooper KG. *A randomised comparison of microwave*  
493 *endometrial ablation with transcervical resection of the endometrium: Follow up at a*  
494 *minimum of 10 years.* *BJOG An Int J Obstet Gynaecol.* 2009;
- 495 27. Herman MC, Penninx JPM, Mol BW, Bongers MY. *Ten-year follow-up of a randomized*  
496 *controlled trial comparing bipolar endometrial ablation with balloon ablation for heavy*  
497 *menstrual bleeding.* *Obstetrical and Gynecological Survey.* 2014.
- 498 28. Penninx JPM, Herman MC, Mol BW, Bongers MY. *Five-year follow-up after comparing bipolar*  
499 *endometrial ablation with hydrothermablation for menorrhagia.* *Obstet Gynecol [Internet].*  
500 *2011 Dec [cited 2018 Dec 3];118(6):1287–92. Available from:*  
501 *<http://insights.ovid.com/crossref?an=00006250-201112000-00012>*
- 502 29. Bansi-Matharu L, Gurol-Urganci I, Mahmood T, Templeton A, van der Meulen J, Cromwell D.  
503 *Rates of subsequent surgery following endometrial ablation among English women with*

- 504           menorrhagia: population-based cohort study. *BJOG An Int J Obstet Gynaecol* [Internet]. 2013  
505           Nov [cited 2018 Dec 3];120(12):1500–7. Available from:  
506           <http://www.ncbi.nlm.nih.gov/pubmed/23786246>
- 507   30.   Cramer MS, Klebanoff JS, Hoffman MK. Pain is an Independent Risk Factor for Failed Global  
508           Endometrial Ablation. *J Minim Invasive Gynecol* [Internet]. 2018 Sep [cited 2018 Dec  
509           3];25(6):1018–23. Available from:  
510           <https://linkinghub.elsevier.com/retrieve/pii/S1553465018300591>
- 511   31.   Steyerberg EW, Harrell FE, Borsboom GJJM, Eijkemans MJC, Vergouwe Y, Habbema JDF.  
512           Internal validation of predictive models: Efficiency of some procedures for logistic regression  
513           analysis. *J Clin Epidemiol*. 2001;
- 514   32.   Steyerberg EW, Eijkemans MJ, Habbema JD. Stepwise selection in small data sets: a simulation  
515           study of bias in logistic regression analysis. *J Clin Epidemiol* [Internet]. 1999 Oct [cited 2018  
516           Dec 3];52(10):935–42. Available from: <http://www.ncbi.nlm.nih.gov/pubmed/10513756>
- 517   33.   Steyerberg EW, Eijkemans MJ, Harrell FE, Habbema JD. Prognostic modelling with logistic  
518           regression analysis: a comparison of selection and estimation methods in small data sets. *Stat  
519           Med* [Internet]. 2000 Apr 30 [cited 2018 Dec 3];19(8):1059–79. Available from:  
520           <http://www.ncbi.nlm.nih.gov/pubmed/10790680>
- 521   34.   Breiman L. *Randomforest2001*. *Mach Learn*. 2001;
- 522   35.   Liu Y, Zhang Y, Liu D, Tan X, Tang X, Zhang F, et al. Prediction of ESRD in IgA nephropathy  
523           patients from an asian cohort: A random forest model. *Kidney Blood Press Res*. 2018;
- 524   36.   Fawagreh K, Gaber MM, Elyan E. Random forests: From early developments to recent  
525           advancements. *Syst Sci Control Eng*. 2014;
- 526   37.   Kaitlin ;, Smith T;, Sadler B. *Random Forest vs Logistic Regression: Binary Classification for  
527           Heterogeneous Datasets*. *Recommended Citation Kirasich*. 2018.

- 528 38. Gareth J, Daniela W, Trevor H, Rober T. An Introduction to Statistical Learning with  
529 Applications in R. Current medicinal chemistry. 2000.
- 530 39. Loh WY. Regression trees with unbiased variable selection and interaction detection . Stat Sin.  
531 2002;
- 532 40. James G, Witten D, Hastie T, Tibshirani R. An introduction to Statistical Learning. Current  
533 medicinal chemistry. 2000.
- 534 41. Hastie TT. The Elements of Statistical Learning Second Edition. Math Intell. 2017;
- 535 42. Bergstra JAMESBERGSTRA J, Yoshua Bengio YOSHUA BENGIO U. Random Search for  
536 HyperParameter Optimization. J Mach Learn Res. 2012;
- 537 43. Jeni LA, Cohn JF, De La Torre F. Facing imbalanced data - Recommendations for the use of  
538 performance metrics. In: Proceedings - 2013 Humaine Association Conference on Affective  
539 Computing and Intelligent Interaction, ACII 2013. 2013.
- 540 44. Bouzari Z, Yazdani S, Naeimi Rad M, Bijani A. Is thermal balloon ablation in women with  
541 previous cesarean delivery successful? TURKISH J Med Sci [Internet]. 2018 Apr 30 [cited 2018  
542 Dec 3];48(2):266–70. Available from: <http://www.ncbi.nlm.nih.gov/pubmed/29714438>
- 543 45. Riley KA, Davies MF, Harkins GJ. Characteristics of patients undergoing hysterectomy for failed  
544 endometrial ablation. J Soc Laparoendosc Surg. 2013;
- 545 46. Kalish GM, Patel MD, Gunn MLD, Dubinsky TJ. Computed Tomographic and Magnetic  
546 Resonance Features of Gynecologic Abnormalities in Women Presenting With Acute or  
547 Chronic Abdominal Pain. Ultrasound Q [Internet]. 2007 Sep [cited 2018 Dec 3];23(3):167–75.  
548 Available from: <http://www.ncbi.nlm.nih.gov/pubmed/17805165>
- 549 47. Gordts S, Grimbizis G, Campo R. Symptoms and classification of uterine adenomyosis,  
550 including the place of hysteroscopy in diagnosis. Fertility and Sterility. 2018.

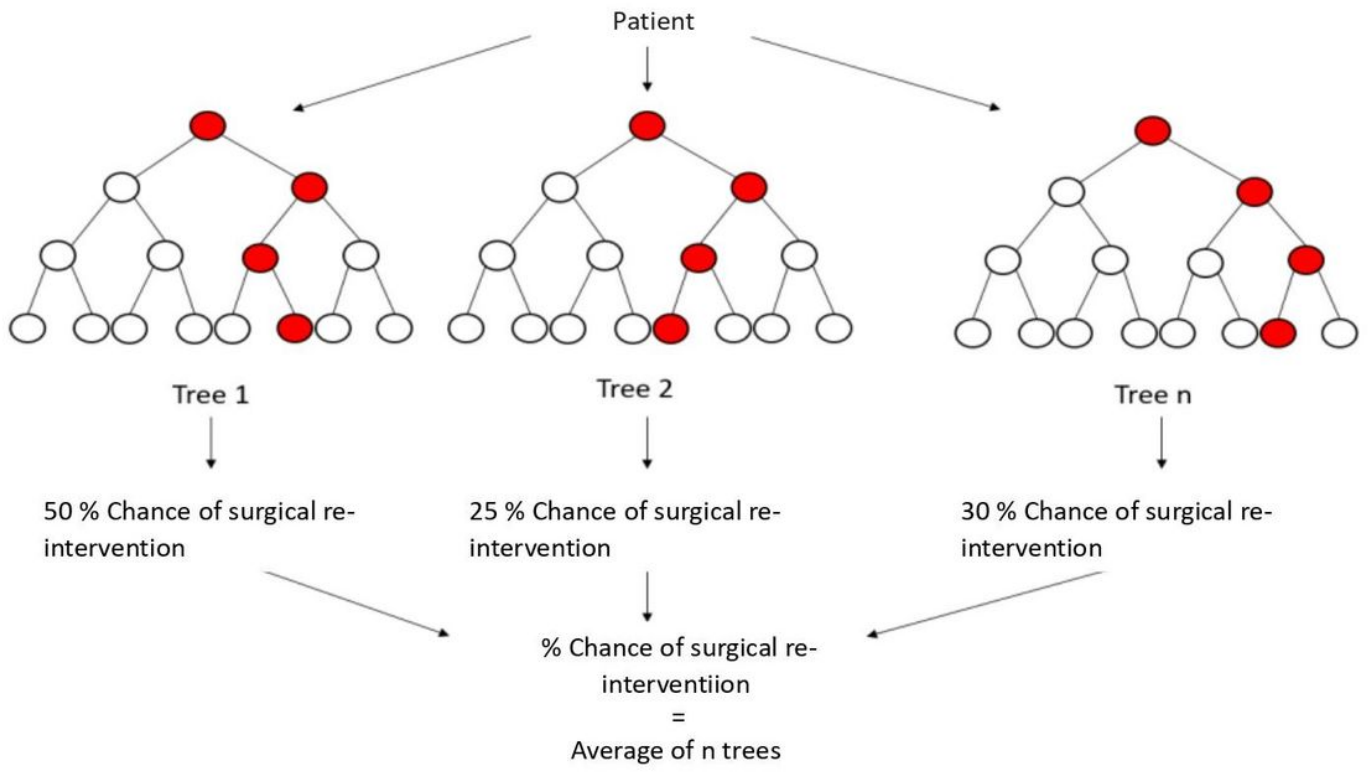
- 551 48. Ennis M, Hinton G, Naylor D, Revow M, Tibshirani R. A comparison of statistical learning  
552 methods on the GUSTO database. Stat Med. 1998;
- 553 49. Rajkomar A, Dean J, Kohane I. Machine learning in medicine. New England Journal of  
554 Medicine. 2019.
- 555 50. Kononenko I. Machine learning for medical diagnosis: History, state of the art and  
556 perspective. Artif Intell Med. 2001;
- 557 51. Lustgarten JL, Gopalakrishnan V, Grover H, Visweswaran S. Improving classification  
558 performance with discretization on biomedical datasets. AMIA . Annu Symp proceedings  
559 AMIA Symp. 2008;
- 560

# Figures



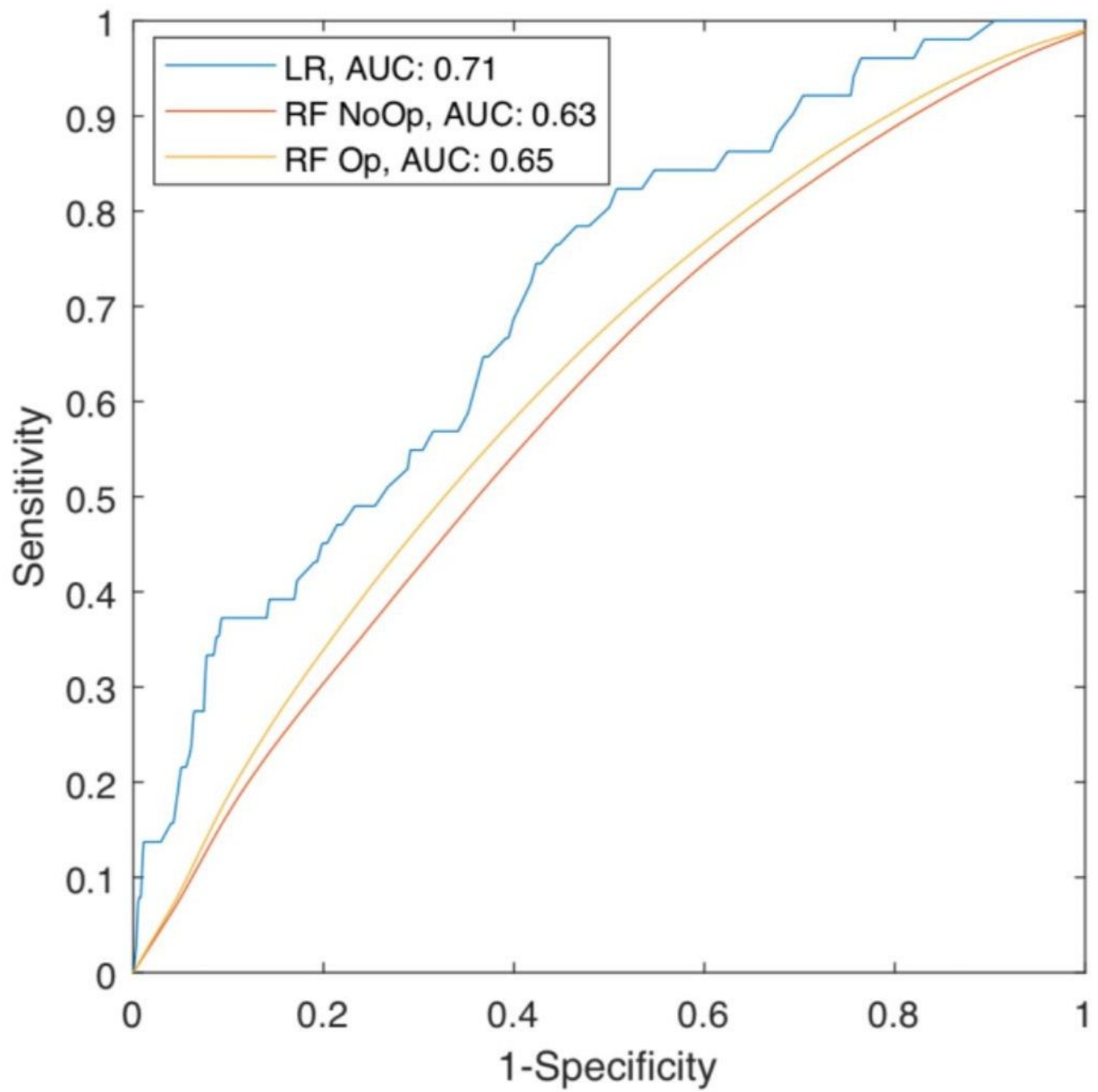
**Figure 1**

An illustration of a decision tree in the Random Forest model. The decision tree directs each case from the root node to the leaf nodes, resulting in a prediction. N= Number, SRR = Surgical Re-intervention Rate.



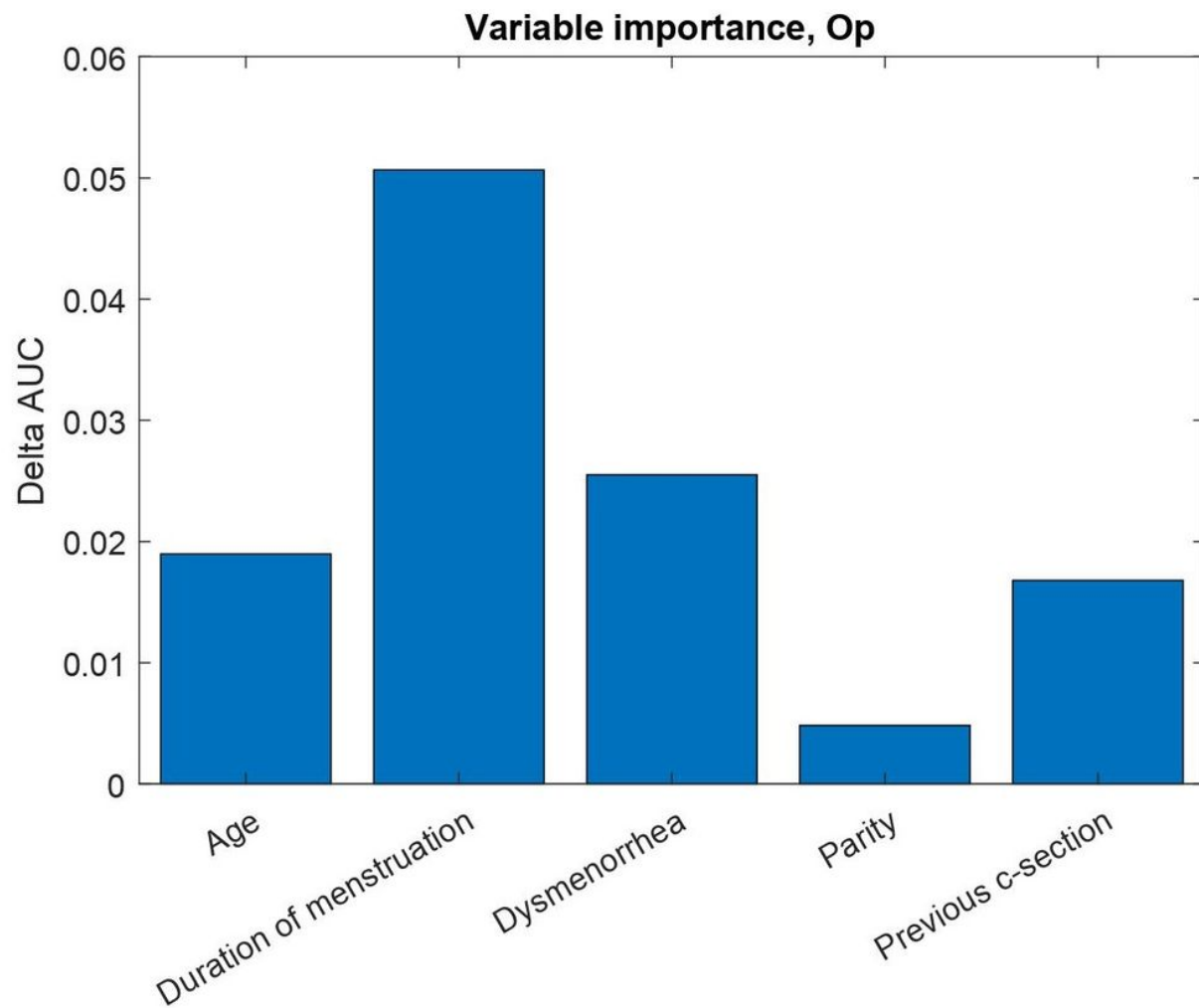
**Figure 2**

A simplified Random Forest model for the prediction of the surgical re-intervention .



**Figure 3**

ROC-curve of the logistic regression and Random Forest model. LR AUC 0.71 (95% CI 0.64-0.78). , NoOp AUC 0.63 (0.54-0.71), Op AUC: (0.56 – 0.74) LR= logistic regression, RF= Random Forest, Op= after hyperparameter optimization, NoOp= before hyperparameter optimization



**Figure 4**

Contribution of predictors of surgical re-intervention within 2 years after endometrial ablation, after hyperparameter optimization