# Prediction of unsuccessful endometrial ablation: Random Forest vs Logistic Regression

Kelly Stevens<sup>1</sup>, Liesbet Lagaert<sup>1</sup>, Tom Bakkes<sup>2</sup>, Malou Gelderblom<sup>3</sup>, Saskia Houterman<sup>3</sup>, Tanja Gijsen<sup>4</sup>, and Dick Schoot<sup>3</sup>

<sup>1</sup>University Hospital Ghent <sup>2</sup>Technical University Eindhoven <sup>3</sup>Catharina Hospital <sup>4</sup>Elkerliek Hospital

May 5, 2020

# Abstract

Objective: To develop a prediction model to predict surgical re-intervention within two years after endometrial ablation (EA) by using a random forest technique (RF). The performance of the developed prediction model was then compared with a previously published multivariate logistic regression model (LR) (1). Design: Retrospective cohort study. Setting: Data from two non-university teaching hospitals in the Netherlands were used. Population: 446 pre-menopausal women who have had an EA for heavy menstrual bleeding between January 2004 and April 2013. Methods: The RF model was trained in MATLAB (2018b) using the TreeBagger function in the Statistics and Machine Learning Toolbox. Main outcome measures: The performance of the two models was compared using the area under the Receiving Operating Characteristic (ROC) curve (AUROC). Measurements and Main Results: The LR model had an AUC of 0.71 (95% CI 0.64-0.78). The RF model had an AUC of 0.63 (95% CI 0.54-0.71). and an AUC of 0.65 (95% CI 0.56-0.74) after hyperparameter optimization. Conclusion: The RF model is not superior compared to the LR model in predicting the outcome of surgical re-intervention within two years after EA. Machine learning techniques are gaining popularity in development of clinical prediction tools, but they are not necessarily superior to traditional statistical logistic regression techniques. The performance of a model is influenced by the sample size and the number of features, hyperparameter tuning and the linearity of associations. Both techniques should be considered when developing a prediction model.

Prediction of unsuccessful endometrial ablation: Random Forest vs Logistic Regression

# Introduction

Abnormal uterine bleeding in premenopausal women is a common complaint in five percent of the women who experiences complaints of abnormal uterine bleeding. (2) Endometrial ablation (EA) is one of the treatment options for this common problem. Due to the less invasive nature (lower intra-operative complication risks, shorter recovery time, and lower post-operative morbidity), and low costs of this procedure, this form of treatment seems to be a less-invasive surgical treatment for menorrhagia compared to hysterectomy (3–7). However, long-term follow up shows a decrease in patient satisfaction and treatment efficacy. Due to permanent relief, the more invasive hysterectomy remains the most effective treatment of abnormal uterine bleeding (8–15).

According to literature, several factors prior to endometrial ablation appear to have an influence on the success or failure-rate of this procedure. Younger age, complaints of dysmenorrhea, parity above or equal to five, a thicker pre-procedural endometrium, a duration of menstruation above seven days, presence of

an intramural leiomyoma on transvaginal sonography, a history of sterilization or caesarean section, and a longer uterine depth are some of the possible negative influencing factors (1,2,8,9,11-18).

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

For this reason Stevens et al.(1) developed two multivariate prediction models to help counsel patients for failure of EA and for surgical re-intervention within two years after EA. The developed prediction models have a clinically acceptable c-index of 0.68 and 0.71 respectively. In addition, Stevens et al. is performing an external validation of these two prediction models, using retrospective data of similar patient groups in two non-university teaching hospitals in the Netherlands. Results of these data will follow. In the field of gynaecology, many prediction models are developed using multivariate logistic regression as a standard approach, these are based on a combination of various predictors that are significantly related to the outcome of interest. However, this method cannot automatically estimate the interconnection between predictors and in this way can overestimate the influence of an individual predictor (20,21).

We were also interested in other statistical techniques of developing a prediction model. In recent years machine learning (ML) methods have been increasingly used in the development of clinical prediction models. This method is a scientific discipline that focuses on models that directly and automatically learn from data (20,22). Potential advantage of the machine learning methods compared to the traditional statistical strategies is the possibility of capturing complex, nonlinear relationships in the data (23,24). ML computer algorithms use training data with well-defined input and output variables. This gives the opportunity to define a model with predictors which can be used for new and similar data. Compared to statistical logistic regression models, this can be done without a priori assumption of relevant variables (25).

Random forest is a machine learning method used for classification and regression that operates by constructing a large ensemble of decision trees on training data (22,23,26). Each tree in the random forest is built using a bootstrap sample randomly drawn from the training dataset. This results in a reduction of variance and corrects for a single decision trees ability to overfit to a training set. Each tree in the forest gives an individual prediction on the outcome measure. For a classification problem (in this case, surgical re-intervention or no surgical re-intervention after EA) the final random forest model averages the prediction of all the trees in the forest (21,23,27).

The aim of the study was to develop a random forest prediction model to predict the chance of surgical re-intervention within two years after EA. Furthermore, it was our aim to compare the performance of the random forest model with the prediction by previously published the multivariate logistic regression model (1). In both models the surgical re-intervention within two years after EA is used as primary outcome measure.

## Methods:

This retrospective two-centred cohort study, performed in two non-university teaching hospitals in the Netherlands (Catharina Hospital, Eindhoven; Elkerliek Hospital Helmond), included 446 patients who have had an EA for complaints of abnormal uterine bleeding (1). Both hospitals used similar ablation techniques between 2004 and 2013, being Cavatherm® (Veldana Medical SA, Morges, Switzerland), Gynecare Thermachoice® (Ethicon, Sommerville, US) and Thermablate® EAS (Idoman, Ireland). Recent publications have shown that these ablation techniques were equally effective (15,28). Local medical ethical review boards approved the study. All patients gave informed consent.

Patients were identified in the Electronic Patient care System by using specified search terms related to endometrial ablation. Exclusion criteria were a postmenopausal status at time of EA; (suspicion of) endometrial malignancy or uterine cavity deformations (adenomyosis; anomalies; fibroids; or a polyp). Follow-up period after treatment was at least two years. This time-interval was chosen because previous literature stated that most re-interventions were done within two years. Follow-up ended on the day of hysterectomy, in case of death or on April 15, 2015 (10,16,18,28–30).

Data were extracted from individual patient files by two researchers. Next, patients were asked to fill in a questionnaire regarding follow-up information. In case of non-response, patients were contacted by letter and ultimately by telephone. The questionnaire contained questions based on significant variables predicting surgical re-intervention after EA that were previously published (2,5,8,11–16,18,31,32).

The entire dataset consists of 446 patients with different categorical and continuous variables. For the machine learning algorithms all features were extracted from the original dataset and a total of five pre-operative variables were used to develop the machine learning model. This were the preoperative variables that were significant predictors in the final multivariate prediction model of Stevens et al. (age, duration of menstruation, dysmenorrhea, parity and previous caesarean section) (1). The continuous data were not discretized into categories as was done in the development of the previously published logistic regression model(1).

# Development of the Logistic regression model

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

To determine which variables were significant, univariable logistic regression analysis was used. The variables with a p-value <.10 were used in the multivariable analysis. This was followed by a backward stepwise manual selection process, progressively excluding the variable with the highest p-value (1).

As described by Steyerberg et al., the p-value of 0.10 was used to prevent a potential incorrect exclusion of a predictive factor. This would be far more detrimental for the test than missing a potential discriminating factor (33,34).

Interaction terms were used to test possible interaction between the significant variables in the model. Furthermore, multicollinearity was tested. Bootstrap resampling was used for internal validation (n=5000). (34,35) To correct for over-optimism of the model, regression coefficients were multiplied by the calculated shrinkage factor(1).

## Development of the Random forest model

We first trained a RF model using the five following pre-operative predictors: age, duration of menstruation, dysmenorrhea, parity and previous caesarean section. These factors were associated with a higher probability of surgical re-intervention within two years after EA in the previously published multivariate logistic regression model (1).

As described above, a RF model is an ensemble of many decision tree models. When building decision trees, each tree in the forest uses random samples (patients) from the training set ("tree bagging"). Figure 1 shows an example of an individual decision tree in the random forest. A decision tree is a flowchart-like binary branch structure. At each 'node split' in the tree the data are divided in two, based on the value of variable of the decision node. If no more splits are possible a prediction will be calculated for the cases in the final leaf node (23,26,36).

At each node split a random subset of features (such as duration of menstruation and parity) is considered ("feature bagging"), this to avoid over-selection of strong predictive features, leading to similar splits in the trees. This finally leads to a robust model and prevents model overfitting (21,23,26,27,36,37).

Following this process, the classification result of a RF model is produced by computing a large ensemble of those trees and averaging the prediction of each single decision tree on surgical re-intervention. Figure 2 shows a simplified example of the RF model. In practice, the decision trees and the resulting prediction model contain a large number of leaf nodes(26,38).

A RF method has its own hyperparameters: ntree, mtry, minimum leaf size and maximum node splits. Ntree is the number of trees in the forest. It should at the one hand be as large as possible, so that each feature (variable) can have enough opportunities to be picked, but not too large to reduce unnecessary calculation time. A default value of ntree = 500 was used (39). Mtry is the number of features randomly selected as candidate feature at each split ("feature bagging"). This was set as square root of the number of variables ([?]n). Minimum leaf size is the minimum number of cases that is required to produce another node split. Maximum node splits is the maximum amount of splits. Neither a minimum leaf size nor maximum node splits was set (23,40).

We began running the RF module with default parameter values before starting to improve the RF's performance by hyperparameter optimization. Default parameters are pre-set values for the hyperparameters on which the construction of the decision trees is based, for example 500 for ntree (26,27).

To predict the chance of surgical re-intervention within two years after EA, the model was initially trained and internally validated on the 446 cases. To make a good comparison between de RF and LR the same validation technique had to be used. Therefore, a bootstrap resampling of 5000 was used to make training bags and test bags (Out Of Bag (OOB) samples). The cases that were not selected by the bootstrap resampling form the test bag which was used as a validation sample to assess the performance of the trained model on new observations. (Figure 3) The performance measure Area Under The Receiver Operating Curve (AUROC) was calculated on the test sets (the OOB samples) and averaged for the 5000 bootstrap samples. These two bags must not be confused with the "tree bags" and the "feature bags" used to construct the decision trees in the random forest (21,23,26,36).

The RF was trained in MATLAB (2018b) using the TreeBagger function in the Statistics and Machine Learning Toolbox. The curvature test was used for split-predictor selection to get an unbiased selection between the continuous and categorical variables. The Gini-impurity index was used to evaluate the accuracy of a split and to predict the variable importance. A perfect separation results in a Gini score of zero (all observations belonging to one label, in this case surgical re-intervention or no surgical re-intervention), whereas the worst case split results in 50/50 classes (23).

The parameter optimization was performed by a random grid-search of the minimum leaf size and the maximum number of splits. The minimum leaf size can take a value between 1 and half the sample size (N/2 = 223). The maximum number of splits can take a value between 1 and the sample size minus one (N-1 = 445). A random search was chosen since it has been shown to have a similar performance to a full grid search, but has a reduced computation time (38,41).

For each random combination of minimum leaf size and maximum number of splits, a RF was trained on the training bag. The combination was scored using OOB prediction of the tree bags. This was repeated for 20 random combinations, the combination with the highest area under the curve (AUC) on the OOB-predictions was used to train a RF which was tested on the validation test set (42).

## Comparison of the prediction models

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

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

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

## Results

Seven hundred sixty-two patients were identified retrospectively. Thirty-three patients were excluded, thirty did not meet the inclusion criteria and three underwent an incomplete endometrium ablation. The remaining 729 patients were contacted, resulting in a response-rate of 61% (N = 446). A total amount of 446 patients was available for analysis (1).

Fifty-three (11.9%) of these patients required a surgical re-intervention within two years after EA. Patients mean age during their EA was 43.8 years (SD +-5.5, range 20-55, missing values 0). The mean number of parity was 2.2 (SD +- 1.0, missing values 0). Sixty-one (13.7%) of the patients underwent a caesarean section. The mean number of previous caesarean section was 0.2 (SD +- 0.6, missing values 0)

Hundred sixty-nine (39.4%) of the patients had a menstruation period longer than seven days, the mean number of menstrual days was 9.4 (SD +- 6.0, missing values 17). Two hundred fifty-six (57.4%) of the patients had complaints of dysmenorrhea and four hundred thirty-four (97.3%) of the patients had complaints of abnormal uterine bleeding (1).

## Prediction models:

## Logistic regression model

Univariate analysis showed six significant predictors, multivariate analyses resulted in a logistic regression model consisting of five significant predictors: age (OR 0.95, 95% CI 0.90 – 1.00), duration of menstruation >7 days (OR 2.05, 95% CI 1.10 – 3.82), dysmenorrhea (OR 2.48, 95% CI 1.21 – 5.07), parity [?]5 (OR 7.63, 95% CI 1.51 – 38.46), and previous caesarean section (OR 2.21, 95% CI 1.05 – 4.64). The AUC of the final prediction model after correcting by the shrinkage factor was 0.71 (95% CI 0.64-0.78) (Figure 4).

The final model is described in the article of Stevens et al (1).

# Random forest model

The random forest method resulted in a model which predicts the chance of re-intervention within 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 achieved after optimization of this model (Figure 4).

# Predictors of surgical re-intervention: Variable importance

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

#### Discussion

# Main findings

In this study, a RF model was made to predict surgical re-intervention within two years after EA. Comparison of the predictive value of a RF model with the existing logistic regression model of Stevens et al. was made (1).

The existing logistic regression model has a C-index of 0.71 (95% CI 0.64-0.78) (1). The RF model, developed in this study, shows a C-index of 0.65 (95% CI 0.56-0.74) after hyperparameter optimization. This shows that the LR prediction model developed by Stevens et al. probably performs better in predicting surgical

re-intervention within two years after EA than the newly developed RF model. However, this difference in performance is not statistically significant when we look at the confidence intervals. Significant predictors of the model are age, duration of menstruation >7 days, dysmenorrhea, parity [?]5 and previous caesarean section(1).

In our database, high parity ([?]5) is a predictive variable for surgical re-intervention. This can be related to the larger uterine cavity of grand multiparous women. However, when considering our RF model, parity has no large impact on the AUC. This is in line with previously reported studies that show no significant increased risk of treatment failure with increasing parity (2,17).

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

In our cohort, pre-operative dysmenorrhea is associated with a higher risk of surgical re-intervention. There is evidence that gynaecologic pathology causing this dysmenorrhea (adenomyosis and endometriosis) reduces the success of endometrial ablation (9,18,32,45,46). This can be explained by the fact that EA is not an appropriate treatment for these diseases due to the superficial effect of energy to the uterine wall of ablation. It could help to diagnose these diseases before performance of EA. However, sensitivity and specificity of the diagnostic tools for determining these diseases in the pre-operative setting are still low (47).

In line with previous studies, we found that younger age was associated with a higher risk of surgical reintervention (8,10-14,31).

The duration of menstruction > 7 days is also a negative predictive factor for surgical re-intervention after EA. This may be caused by a thicker endometrium which is more difficult to completely remove by the device (8,11).

# Interpretation in light of other evidence

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

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

Secondly, some studies demonstrate that ML is performing better when only a small set of pre-specified predictors are used in the prediction model. There seems to be an influence of the number of predictors (p) and the ratio of p:n (sample size). RF tends to perform better for increasing p and p:n. (20,24,49,50) In our study, to limit potential bias, the five identical predictors as published before (1) were considered for the LR and RF algorithms. We did this to allow a fair comparison between the two models, probably in disadvantage of the RF model (20,24,49,50).

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

Finally, we can also consider that for this clinical problem a logistic approach is better than a RF model for modelling the relationship between surgical re-intervention and the explanatory variables. Probably the previously mentioned complex, nonlinear relationships that a ML approach can better capture are not present in this dataset.

## Strengths and limitations

The predictors obtained by univariate and multivariate logistic regression are in accordance with the existing literature (51). However, when we compare the variable importance between the OR (LR) and the difference in AUC (ML) of each variable, we identify a different ranking in variable importance.

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

Dysmenorrhea (OR 2.48) and a parity>5 (OR 7.63) have the highest odds ratio in the multivariate analysis, while for the difference in area under the curve the duration of menstruation and dysmenorrhea are the most important variables. We consider two possible reason for the difference in importance. The first reason is that for the LR model all continuous variables (except age) were discretized, while for the RF model continuous variables were handled. A second reason is that in the LR the predictors have different units, and these were not standardized. This means that a subjective assessment of variable importance cannot easily be made by simply comparing the raw sizes of the OR (2,8,13-18,31,32,44). This can be seen as a strength of our study since the difference in AUC for each predictor (permuted vs. not permuted) reflects the variable importance in a standardized way.

We used bootstrap resampling for internal validation (n=5000) in the LR and RF model. Using the same validation method limits potential bias. Furthermore, the same predictors were considered for the LR and ML algorithms. This limits potential bias, but will limit the potential power of a RF technique as well. Another important strength of this study is the use of all participants in evaluating the performance of the RF model. By using the test sets, there is no need for an independent validation dataset.

It could be seen as a limitation of this study that we did not perform an external validation in another cohort. However, we did not expect it to be significantly better in performance, since the internal validation of the RF did not perform better than the logistic regression model. In addition, an external validation for the logistic regression model is being performed at the time of this study (52).

Finally, we can state that ML models are in our experience not easily implemented in the clinical practice; since these are often not available in commonly used software packages in clinical practice. However, future structured data-registration is increasing, which makes it easier to create big datasets available for ML-programs.

## **Conclusion:**

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

#### Declarations

# Disclosure of interests

There are no conflicts of interest to disclose.

# Contribution to authorship

KS, LL and BS were the chief investigators and were responsible for the project development, data management, data analysis and manuscript writing. TB was the trial statistician and developed the RF model. SH and MG contributed to the manuscript editing. TG contributed to the data collection and manuscript editing.

# Details of ethics approval

The ethical board in the Catharina hospital and in the Elkerliek hospital concluded that ethics approval was not necessary for this study.

Funding

None

## Acknowledgements

The authors want to thank the patients for completing the questionnaires and for

consenting to participate in our study.

## Literature references

1. Stevens KYR, Meulenbroeks D, Houterman S, Gijsen T, Weyers S, Schoot BC. Prediction of unsuccessful endometrial ablation: a retrospective study. Gynecol Surg. 2019;

2. Peeters JAH, Penninx JPM, Mol BW, Bongers MY. Prognostic factors for the success of endometrial ablation in the treatment of menorrhagia with special reference to previous cesarean section. Eur J Obstet Gynecol Reprod Biol. 2013 Mar;167(1):100–3.

3. Waddell G, Pelletier J, Desindes S, Anku-Bertholet C, Blouin S, Thibodeau D. Effect of endometrial ablation on premenstrual symptoms. J Minim Invasive Gynecol. 2015;

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

5. Bouzari Z, Yazdani S, Azimi S, Delavar M. Thermal Balloon Endometrial Ablation in the Treatment of Heavy Menstrual Bleeding. Med Arch. 2014;

6. Miller JD, Lenhart GM, Bonafede MM, Basinski CM, Lukes AS, Troeger KA. Cost effectiveness of endometrial ablation with the NovaSure(r) system versus other global ablation modalities and hysterectomy for treatment of abnormal uterine bleeding: US commercial and medicaid payer perspectives. Int J Womens Health. 2015;

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

8. El-Nashar SA, Hopkins MR, Creedon DJ, St Sauver JL, Weaver AL, McGree ME, et al. Prediction of treatment outcomes after global endometrial ablation. Obstet Gynecol. 2009 Jan;113(1):97–106.

9. Wishall KM, Price J, Pereira N, Butts SM, Badia CRD. Postablation risk factors for pain and subsequent hysterectomy. In: Obstetrics and Gynecology. 2014.

10. Thomassee MS, Curlin H, Yunker A, Anderson TL. Predicting Pelvic Pain After Endometrial Ablation: Which Preoperative Patient Characteristics Are Associated? J Minim Invasive Gynecol. 2013;

11. Bongers MY, Mol BWJ, Brolmann HAM. Prognostic factors for the success of thermal balloon ablation in the treatment of menorrhagia. Obstet Gynecol. 2002 Jun;99(6):1060–6.

12. Longinotti MK, Jacobson GF, Hung Y-Y, Learman LA. Probability of hysterectomy after endometrial ablation. Obstet Gynecol. 2008 Dec;112(6):1214–20.

13. Shaamash AH, Sayed EH. Prediction of successful menorrhagia treatment after thermal balloon endometrial ablation. J Obstet Gynaecol Res. 2004 Jun;30(3):210–6.

14. Klebanoff J, Makai GE, Patel NR, Hoffman MK. Incidence and predictors of failed second-generation endometrial ablation. Gynecol Surg. 2017 Dec;14(1):26.

15. Louie M, Wright K, Siedhoff MT. The case against endometrial ablation for treatment of heavy menstrual bleeding. Curr Opin Obstet Gynecol. 2018 Aug;30(4):287–92.

16. Kreider SE, Starcher R, Hoppe J, Nelson K, Salas N. Endometrial ablation: is tubal ligation a risk factor for hysterectomy. J Minim Invasive Gynecol. 2013 Sep;20(5):616–9.

17. Lybol C, van der Coelen S, Hamelink A, Bartelink LR, Nieboer TE. Predictors of Long-Term NovaSure Endometrial Ablation Failure. J Minim Invasive Gynecol. 2018;

18. Shavell VI, Diamond MP, Senter JP, Kruger ML, Johns DA. Hysterectomy subsequent to endometrial ablation. J Minim Invasive Gynecol. 2012 Jul;19(4):459–64.

19. van Montfort P, Smits LJM, van Dooren IMA, Lemmens SMP, Zelis M, Zwaan IM, et al. Implementing a Preeclampsia Prediction Model in Obstetrics: Cutoff Determination and Health Care Professionals' Adherence. Med Decis Making. 2020 Jan;40(1):81–9.

20. Evangelia christodoulou, Jie MA, Collins GS, Steyerberg EW, Verbakel JY, van Calster B. A systematic review shows no performance benefit of machine learning over logistic regression for clinical prediction models. J Clin Epidemiol. 2019;

21. Breiman L. Statistical Modeling: The Two Cultures. Stat Sci. 2001;

22. Deo RC. Machine learning in medicine. Circulation. 2015;

23. Couronne R, Probst P, Boulesteix AL. Random forest versus logistic regression: A large-scale benchmark experiment. BMC Bioinformatics. 2018;

24. Chen JH, Asch SM. Machine Learning and Prediction in Medicine — Beyond the Peak of Inflated Expectations. N Engl J Med. 2017;

25. Panesar SS, D'Souza RN, Yeh FC, Fernandez-Miranda JC. Machine Learning Versus Logistic Regression Methods for 2-Year Mortality Prognostication in a Small, Heterogeneous Glioma Database. World Neurosurg X. 2019;

26. Breiman L. Randomforest2001. Mach Learn. 2001;

27. Liu Y, Zhang Y, Liu D, Tan X, Tang X, Zhang F, et al. Prediction of ESRD in IgA nephropathy patients from an asian cohort: A random forest model. Kidney Blood Press Res. 2018;

28. Sambrook AM, Bain C, Parkin DE, Cooper KG. A randomised comparison of microwave endometrial ablation with transcervical resection of the endometrium: Follow up at a minimum of 10 years. BJOG An Int J Obstet Gynaecol. 2009;

29. Herman MC, Penninx JPM, Mol BW, Bongers MY. Ten-year follow-up of a randomized controlled trial comparing bipolar endometrial ablation with balloon ablation for heavy menstrual bleeding. Obstetrical and Gynecological Survey. 2014.

30. Penninx JPM, Herman MC, Mol BW, Bongers MY. Five-year follow-up after comparing bipolar endometrial ablation with hydrothermablation for menorrhagia. Obstet Gynecol. 2011;

31. Bansi-Matharu L, Gurol-Urganci I, Mahmood T, Templeton A, van der Meulen J, Cromwell D. Rates of subsequent surgery following endometrial ablation among English women with menorrhagia: population-based cohort study. BJOG An Int J Obstet Gynaecol. 2013 Nov;120(12):1500–7.

32. Cramer MS, Klebanoff JS, Hoffman MK. Pain is an Independent Risk Factor for Failed Global Endometrial Ablation. J Minim Invasive Gynecol. 2018 Sep;25(6):1018–23.

33. Steyerberg EW, Harrell FE, Borsboom GJJM, Eijkemans MJC, Vergouwe Y, Habbema JDF. Internal validation of predictive models: Efficiency of some procedures for logistic regression analysis. J Clin Epidemiol. 2001; 34. Steyerberg EW, Eijkemans MJ, Habbema JD. Stepwise selection in small data sets: a simulation study of bias in logistic regression analysis. J Clin Epidemiol. 1999 Oct;52(10):935–42.

35. Steyerberg EW, Eijkemans MJ, Harrell FE, Habbema JD. Prognostic modelling with logistic regression analysis: a comparison of selection and estimation methods in small data sets. Stat Med. 2000 Apr;19(8):1059–79.

36. Fawagreh K, Gaber MM, Elyan E. Random forests: From early developments to recent advancements. Syst Sci Control Eng. 2014;

37. Kaitlin ;, Smith T;, Sadler B. Random Forest vs Logistic Regression: Binary Classification for Heterogeneous Datasets. Recommended Citation Kirasich. 2018.

38. Gareth J, Daniela W, Trevor H, Rober T. An Introduction to Statistical Learning with Applications in R. Current medicinal chemistry. 2000.

39. Loh WY. Regression trees with unbiased variable selection and interaction detection. Stat Sin. 2002;

40. James G, Witten D, Hastie T, Tibshirani R. An introduction to Statistical Learning. Current medicinal chemistry. 2000.

41. Hastie TT. The Elements of Statistical Learning Second Edition. Math Intell. 2017;

42. Bergstra JAMESBERGSTRA J, Yoshua Bengio YOSHUABENGIO U. Random Search for HyperParameter Optimization. J Mach Learn Res. 2012;

43. Jeni LA, Cohn JF, De La Torre F. Facing imbalanced data - Recommendations for the use of performance metrics. In: Proceedings - 2013 Humaine Association Conference on Affective Computing and Intelligent Interaction, ACII 2013. 2013.

44. Bouzari Z, Yazdani S, Naeimi Rad M, Bijani A. Is thermal balloon ablation in women with previous cesarean delivery successful? TURKISH J Med Sci. 2018 Apr;48(2):266–70.

45. Riley KA, Davies MF, Harkins GJ. Characteristics of patients undergoing hysterectomy for failed endometrial ablation. J Soc Laparoendosc Surg. 2013;

46. Kalish GM, Patel MD, Gunn MLD, Dubinsky TJ. Computed Tomographic and Magnetic Resonance Features of Gynecologic Abnormalities in Women Presenting With Acute or Chronic Abdominal Pain. Ultrasound Q. 2007 Sep;23(3):167–75.

47. Gordts S, Grimbizis G, Campo R. Symptoms and classification of uterine adenomyosis, including the place of hysteroscopy in diagnosis. Fertility and Sterility. 2018.

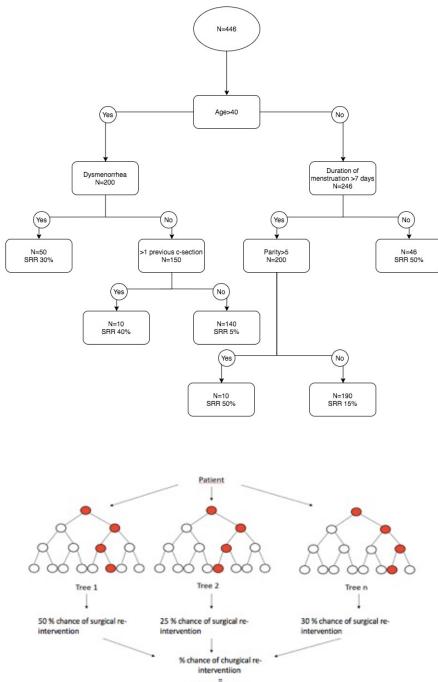
48. Ennis M, Hinton G, Naylor D, Revow M, Tibshirani R. A comparison of statistical learning methods on the GUSTO database. Stat Med. 1998;

49. Rajkomar A, Dean J, Kohane I. Machine learning in medicine. New England Journal of Medicine. 2019.

50. Kononenko I. Machine learning for medical diagnosis: History, state of the art and perspective. Artif Intell Med. 2001;

51. Lustgarten JL, Gopalakrishnan V, Grover H, Visweswaran S. Improving classification performance with discretization on biomedical datasets. AMIA . Annu Symp proceedings AMIA Symp. 2008;

52. Muller I, Houterman S, Stevens KY, Schoot D, Weyers S. Models to Predict Unsuccessful Endometrial Ablation: External Validation. J Minim Invasive Gynecol. 2018;



average of n trees

