

CNFE-SE: A novel approach combining complex network-based feature engineering and stacked ensemble to predict the success of Intrauterine Insemination and ranking the features

Sima Ranjbari

Tarbiat Modares University

Toktam Khatibi (✉ toktamk@gmail.com)

Tarbiat Modares University <https://orcid.org/0000-0001-5824-9798>

Ahmad Vosough Taghi Dizaj

Royan Institute

Hesamoddin Sajadi

Royan Institute

Mehdi Totonchi

Royan Institute

Firouzeh Ghaffari

Royan Institute

Research article

Keywords: IUI outcome prediction, complex networks, feature engineering, stacked ensemble classifier, feature selection

Posted Date: December 11th, 2020

DOI: <https://doi.org/10.21203/rs.3.rs-24259/v3>

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

[Read Full License](#)

Version of Record: A version of this preprint was published on January 2nd, 2021. See the published version at <https://doi.org/10.1186/s12911-020-01362-0>.

1 **CNFE-SE: A novel approach combining complex network-based feature engineering and**
2 **stacked ensemble to predict the success of Intrauterine Insemination and ranking the**
3 **features**

4 Sima Ranjbari, M.S.C.¹, Toktam Khatibi, Ph.D.^{1*}, Ahmad Vosough Taghi Dizaj, M.D.²,
5 Hesamoddin Sajadi, M.D.⁴, Mehdi Totonchi, Ph.D.^{3,4*}, Firouzeh Ghaffari⁵

- 6 1. School of Industrial and Systems Engineering, Tarbiat Modares University
- 7 2. Department of Genetics at Reproductive Biomedicine Research Center, Royan Institute for
8 Reproductive Biomedicine, ACECR, Tehran, Iran, email: vosough@royaninstitute.org.
- 9 3. Department of Reproductive Imaging, Reproductive Biomedicine Research Center, Royan
10 Institute for Reproductive Biomedicine, ACECR, Tehran, Iran.
- 11 4. Department of Andrology, Reproductive Biomedicine Research Center, Royan Institute
12 for Reproductive Biomedicine, ACECR, Tehran, Iran.
- 13 5. Department of Endocrinology and Female Infertility, Reproductive Biomedicine Research
14 Center, Royan Institute for reproductive biomedicine, ACECR, Tehran, Iran, email:
15 ghafaryf@yahoo.com

16
17 *Corresponding Authors: Toktam Khatibi
18 Email: toktam.khatibi@modares.ac.ir
19
20

1 Abstract

2
3 **Background:** Intrauterine Insemination (IUI) outcome prediction is a challenging issue which the
4 assisted reproductive technology (ART) practitioners are dealing with. Predicting the success or
5 failure of IUI based on the couples' features can assist the physicians to make the appropriate
6 decision for suggesting IUI to the couples or not and/or continuing the treatment or not for them.
7 Many previous studies have been focused on predicting the in vitro fertilization (IVF) and
8 intracytoplasmic sperm injection (ICSI) outcome using machine learning algorithms. But, to the
9 best of our knowledge, a few studies have been focused on predicting the outcome of IUI. The
10 main aim of this study is to propose an automatic classification and feature scoring method to
11 predict intrauterine insemination (IUI) outcome and ranking the most significant features.

12 **Methods:** For this purpose, a novel approach combining complex network-based feature
13 engineering and stacked ensemble (CNFE-SE) is proposed. Three complex networks are extracted
14 considering the patients' data similarities. The feature engineering step is performed on the
15 complex networks. The original feature set and/or the features engineered are fed to the proposed
16 stacked ensemble to classify and predict IUI outcome for couples per IUI treatment cycle. Our
17 study is a retrospective study of a 5-year couples' data undergoing IUI. Data is collected from
18 Reproductive Biomedicine Research Center, Royan Institute describing 11255 IUI treatment
19 cycles for 8,360 couples. Our dataset includes the couples' demographic characteristics, historical
20 data about the patients' diseases, the clinical diagnosis, the treatment plans and the prescribed drugs
21 during the cycles, semen quality, laboratory tests and the clinical pregnancy outcome

22 **Results:** Experimental results show that the proposed method outperforms the compared methods
23 with Area under receiver operating characteristics curve (AUC) of 0.84 ± 0.01 , sensitivity of 0.79
24 ± 0.01 , specificity of 0.91 ± 0.01 , and accuracy of 0.85 ± 0.01 for the prediction of IUI outcome.

25 **Conclusions:** The most important predictors for predicting IUI outcome are semen parameters
26 (sperm motility and concentration) as well as female body mass index (BMI).

27 **Key Words:** IUI outcome prediction, complex networks, feature engineering, stacked ensemble
28 classifier, feature selection

29 1 Background

30
31 Infertility is defined as the failure of the female partner to conceive after at least one year of regular
32 unprotected sexual intercourse [1]. More than 186 million people of the world's population
33 specifically people living in developing countries are suffering from infertility [2]. In most cases,
34 the causes of infertility are not clear, which complicates the treatment procedure. These problems
35 have been exacerbated for several reasons, such as lifestyle changes, infection, and genetic issues.
36 In many cases, the only way to get pregnant has been through the use of assisted reproductive
37 technology (ART), and its performance has not yet been optimized [3].

38 Every year, more than 1.5 million ART cycles are carried out all over the world [4]. ART consists
39 of three basic procedures including intrauterine insemination (IUI), in-vitro fertilization (IVF) and
40 intracytoplasmic injection (ICSI) which are generally carried out in different steps of the treatment
41 [5]. The first-line treatment, second and the third stages of ART are IUI, IVF, and ICSI,
42 respectively [6]. In comparison with other sophisticated methods of ART, IUI has been considered
43 as the easiest, minimally invasive and less expensive one. Most of the recent researches have
44 shown the efficacy of IUI [6, 7].

1 IUI outcome prediction is a challenging issue which the ART practitioners are dealing with.
2 Predicting the success or failure of IUI based on the couples' features can assist the physicians to
3 make the appropriate decision for suggesting IUI to the couples or not and/or continuing the
4 treatment or not for them [5].
5 Machine Learning approaches, as the modern scientific discipline, concentrates on how to detect
6 the hidden patterns and extract the information from data. Machine learning provides different
7 methods and algorithms to predict the output from some input predictors which can be used for
8 clinical decision making [8].
9 To the best of our knowledge, many previous studies have been focused on predicting the IVF and
10 ICSI outcome using machine learning methods as summarized in Table 1.

11 *Table 1- summarizing the previous studies of predicting ART outcome*

12
13 As illustrated by Table 1, the previous studies related to outcome prediction of ART methods are
14 listed which have analyzed data using data mining and/or statistical methods. For this purpose,
15 classifiers such as Decision Tree (DT), Logistic Regression (LR), Naïve Bayes (NB), K-Nearest
16 Neighbors (K-NN), Support Vector Machines (SVM), Random Forest (RF), and Artificial Neural
17 Networks (ANN) such as Multi-Layered Perceptron (MLP) and Radial Basis Function (RBF) have
18 been used in the previous studies for predicting the clinical pregnancy after the complete cycles of
19 different ART methods. A main drawback of the most of the considered previous studies is small
20 volume of dataset and a few number of the considered features. Small dataset increases the risk of
21 overfitting the trained models. Overfitting occurs when a model has good predictive ability for
22 training dataset but shows poor performance for test dataset. Models with high overfitting property
23 has lower generalization ability.

24 In this study, a dataset including the features of 11255 IUI treatment cycles for 8360 couples is
25 considered for IUI outcome prediction. Our dataset includes the couples' demographic
26 characteristics, historical data about the patients' diseases, the clinical diagnosis, the treatment
27 plans and the prescribed drugs during the cycles, semen quality, laboratory tests and the clinical
28 pregnancy outcome. Considering the large number of couples and their corresponding IUI
29 treatment cycles is a main advantage of this study compared to the considered previous studies.

30 On the other hand, most of the previous studies have considered the outcome prediction for IVF
31 or ISCI. To the best of our knowledge, a few studies have been focused on predicting the outcome
32 of IUI which have used clustering methods [9, 10] or regression analysis [11].

33 The previous studies which have been based on regression analysis only have considered the
34 weights of the independent features to predict the overall pregnancy probability and they have not
35 assessed the interconnection among the features [11-17]. Many previous studies have suffered
36 from the lack of statistical power due to their small dataset [17, 18]. Also, the AUC performance
37 of the previously proposed models for predicting IUI outcome have been low [12]. Therefore, it is
38 required to improve the prediction performance by proposing novel methods and considering more
39 data records.

40 Most of the considered previous studies have used single classifiers and/or RF as a simple
41 ensemble classifier. Some previous studies have illustrated that the stacked models can improve
42 the classification performance for other applications and other datasets [19-21]. Therefore, in this
43 study, a novel stacked ensemble is designed and proposed for improving the performance of IUI
44 outcome prediction.

45 The main aim of this study is to develop an automatic classification and feature scoring method to
46 predict intrauterine insemination (IUI) outcome and ranking of the most significant features, based

1 on the features describing the couples and their corresponding IUI treatment cycles. For this
2 purpose, a novel approach combining complex network-based feature engineering and stacked
3 ensemble (CNFE-SE) is proposed. Three complex networks are extracted considering the patients'
4 data similarities. The feature engineering step is performed on the complex networks. The original
5 feature set and/or the features engineered are fed to the proposed stacked ensemble to classify and
6 predict IUI outcome for couples per IUI treatment cycle. Our study is a retrospective study of a 5-
7 year couples' data undergoing IUI. Data is collected from Reproductive Biomedicine Research
8 Center, Royan Institute describing 11255 IUI treatment cycles for 8,360 couples.

9 The main novelty of this study lies in three folds including:

- 10 - Proposing a method for feature scoring and classification based on weighted complex
11 networks and stacking ensemble classifiers
- 12 - Proposing feature engineering method based on complex networks
- 13 - Designing a novel stacked ensemble classifier for predicting IUI outcome

14 2 METHODS

15 The main steps of the proposed approach combining complex network-based feature engineering
16 and stacked ensemble (CNFE-SE) to predict the success of Intrauterine Insemination and ranking
17 the features are illustrated in Figure 1.
18

19 *Figure 1- the main steps of the proposed method (CNFE-SE) for feature scoring and classifying the patients to predict IUI*
20 *outcome*

21 The main steps of the proposed method (CNFE-SE) as depicted in Figure 1 include the modules
22 for data collection and preparation, feature scoring and classification and finally model evaluation
23 and validation. The first module consists of data collection, sampling from data, preprocessing the
24 collected data and filtering irrelevant features. In the next module, ignoring a feature, constructing
25 three complex networks from the patients, extracting features from the constructed complex
26 networks, training the classifiers based on the extracted features and finally scoring the ignored
27 feature are performed. The last module evaluates and validates the models trained in the previous
28 module. More details about the mentioned tasks are described in the following subsections.

29 2.1 Data collection

30 Our research is approved by the Institutional Review Board of the Royan Institute Research Center
31 and the Royan Ethics Committee consistent with Helsinki Declaration with the approval ID of
32 IR.ACECR.ROYAN.REC.1398.213. Anonymity and confidentiality of data were respected.

33 Dataset studied in this article is collected from Royan Institute, a public none-profitable
34 organization, affiliated to the academic center for education, culture and research (ACECR) in
35 Iran. It includes the features describing the patients having been treated by IUI method in the
36 Infertility clinic at Royan Institute between January 2011 and September 2015.

37 In this retrospective study, a completed episode is defined as a
38 sequence of treatment cycles resulting in positive clinical pregnancy or when the
39 treatment with IUI is stopped. The inclusion criteria for the couples to be treated under IUI cycles
40 were male factor, ovulatory disorders such as PCOS, hypothalamic amenorrhea, diminished
41 ovarian reserve, combined causes, and unexplained subfertility. The couples' duration of infertility
42 was at least 1 year. Male infertility was defined as the semen quality parameters lower than the
43 standards determined by WHO including sperm concentration lower than 15 million/ejaculate,
44 semen volume lower than 1.5 mL, and total motility lower than 40% [22]. The male partners with
45 donor sperms, Varicocele, and semen samples with total motile sperm count lower than 1×10^6

1 were excluded from being candidates for IUI treatment. Additionally, patients with anatomical and
2 metabolic abnormalities, severe endometriosis and/or systemic diseases were excluded from our
3 study.

4 11,255 IUI cycles related to 8,360 couples are considered in which the women age ranges from
5 16 to 47 with the average age of 29. This dataset contains 1,622 positive outcomes and 9,633
6 negative ones. Therefore, the overall pregnancy rate is 14.41% per completed cycle and 19.4% per
7 couple. Each couple is treated for 1.31 ± 0.59 (mean \pm Standard Deviation) IUI cycles which
8 ranges from 1 to 7 cycles.

9 The features describe the couples' demographic characteristics, historical data about their diseases,
10 the clinical diagnosis, the treatment plans and the prescribed drugs to the couples, male semen
11 quality, laboratory tests and the clinical pregnancy outcome. The considered demographic features
12 include age, body mass index (BMI), education level, consanguinity with spouse and some other
13 features. The information about the history of the patients' subfertility consists of the duration and
14 type of infertility, length of marriage and so on.

15 The types of feature values are numerical, binary, nominal and binominal types for 86, 152, 51
16 and 7 features, respectively. More details about the features is shown in Appendix A.

17
18 In the collected dataset, the majority of couples (almost 72%) have been treated for one cycle, 22%
19 of couples have underwent two cycles, 5% of couples have been treated for three cycles, and less
20 than 1% have been treated more than three cycles. The maximum number of cycles for treating a
21 couple is seven. Figure 2 depicts the distributions of positive and negative clinical pregnancy rates
22 for patients per treatment cycle.

23
24
25 *Figure 2 –The ratio of positive and negative clinical pregnancy per treatment cycle*

26 As illustrated by Figure 2, 63% of the couples belonging to the positive class (positive clinical
27 pregnancy after completing the cycle) have been pregnant after the first treatment cycle. 26% of
28 data records in the positive class have received positive outcome after the second cycle. Moreover,
29 74% of the couples in the negative class have been considered after the first cycle.

30 **2.2 Data sampling**

31 Data should be randomly partitioned into training and test datasets with no overlapping among
32 these two subsets. The models are trained on the training dataset and finally are evaluated by
33 applying them to the test datasets.

34 K-fold cross validation (C.V.) is a common and popular sampling strategy used for this purpose.
35 In this method, data is randomly divided into K disjoint equal-size subsets. Every time, one of
36 these K subsets is considered as the test dataset and all (K-1) remaining subsets make the training
37 one. The model is trained K times on K training datasets and applied to the corresponding test
38 datasets to evaluate the performances of the trained models.

39 Before sampling from data, the features having missing value rate higher than 20% are removed
40 from the study. Moreover, the patient records with high missing value rate (higher than 20%) are
41 excluded from the study and then, 5-fold C.V. is used for sampling from the collected dataset, in
42 this study.

43 At first, dataset is partitioned into non-overlapping subsets D_1, D_2, \dots, D_K based on K-fold Cross
44 Validation strategy. Then, the models are trained on K training datasets composed of all $D_1, \dots,$
45 D_K subsets excluding D_i for $1 \leq i \leq K$. Therefore, the i^{th} training dataset consists of all D_1, \dots, D_K
46 but D_i and the i^{th} test dataset is D_i . The i^{th} training dataset is balanced using over-sampling strategy.

Moreover, a strategy for classification structural risk assessment is used named as A-Test which will be described in the evaluation and validation subsection with more details. The number of instances of positive and negative outcomes in each folder of 5-fold is 324-325 and 1926-1927, respectively. therefore, the imbalance ratio of the training set in each of 5-folds is about 0.168.

2.3 Data preprocessing

Preprocessing of data is one of the most essential steps in the knowledge discovery tasks. A previous study have stated that 80% of total time in data mining projects is allocated for data preparation and preprocessing step [23].

In the first step, the initial collected dataset includes almost 86,000 data records describing the partners and about 1,000 features. The data records describing one couple per IUI treatment cycle are aggregated to make our dataset. Thus, the aggregated dataset includes 11,255 data records and 296 features describing a couple during an IUI treatment cycle.

The nominal features are converted to dummy binary variables. If a nominal features has m different levels or values, it will be converted to $(m-1)$ dummy binary variables. Therefore, instead of considering a nominal feature in the classification and feature ranking, its corresponding dummy binary variables are considered in the mentioned tasks.

The missing values for numeric and categorical features are imputed based on the average and the most frequent values, respectively [24]. All numerical and ordinal features are normalized using min-max normalization method and the nominal features are converted into dummy binary variables.

Outlier detection is performed in this study based on isolation forest method which has been proposed by Liu et al. [25] as an appropriate outlier detection method for high dimensional data. The hyperparameters of Isolation Forest including the number of estimators, maximum number of the samples, contamination coefficient, maximum number of the features, bootstrapping or not, and the number of jobs are tuned using grid search method. For evaluating the performance of Isolation Forest, its results are compared to other outlier detection methods such as One-class SVM with kernel of Radial Basis Function (RBF), boxplot analysis and expert's opinions. Three outliers are identified by this method and excluded from the study.

2.4 Filtering irrelevant features

Since the aggregated dataset consists of many features, the irrelevant features can be removed to reduce the computational time required for processing and analyzing data. Thus, the features having very low correlation with the output feature or very high correlation with other input features are excluded from this study. The linear correlation coefficient between pairs of the features F_p and F_q are calculated as Eq. (1):

$$Corr(F_p, F_q) = \sum_i \frac{(F_{i,p} - m_p)(F_{i,q} - m_q)}{\sqrt{\sum_j (F_{j,p} - m_p)^2} \sqrt{\sum_j (F_{j,q} - m_q)^2}} \quad (1)$$

Where $F_{x,p}$ ($F_{x,q}$) indicates the x^{th} row of the feature F_p (F_q) and m_p (m_q) denotes the average of the feature F_p (F_q), respectively.

If two features F_p and F_q have low (high) correlation, $Corr(F_p, F_q)$ tends to zero (-1 or +1).

2.5 Ignoring a feature

Breiman has proposed measuring the feature importance by mean decrease in accuracy (MDA) of random forest [26]. This study aims at ranking the features according to their predictive power for classifying the instances to positive or negative clinical pregnancy. For this purpose, all the steps 6-9 are performed by considering all the features excluding one feature each time and MDA for the trained proposed classifier is calculated on the validation dataset. MDA values show the amount of reducing the model accuracy after removing a feature. Therefore, the higher values of MDA indicate the higher predictive ability of the corresponding features.

2.6 Constructing complex networks of patients

For modeling nonlinear data, complex networks are effective method [27]. Complex network is a weighted undirected graph $G=(V, E, W)$, where V is the set of nodes, E denotes the set of edges $e(v_i, v_j)$ between the pairs of the nodes v_i and v_j and W is the weights $w(v_i, v_j)$ assigned to their corresponding edges $e(v_i, v_j)$ of E .

Three complex networks are constructed from the training datasets and one data record which should be classified independent from it belongs to training or test dataset. The first one is comprised of all the training data records and one data record which should be classified as its nodes and is called CN1. The second and the third complex networks consist of one data record which should be classified and all training data records excluding the negative and positive classes and named as CN2 and CN3, respectively. If the considered data record belongs to training dataset, its class label is excluded from its corresponding complex networks.

In other words, the nodes of CN1, CN2 and CN3 are one data record which should be classified and all the training data records, positive labeled and negative labeled training data records, respectively. Therefore, for each data record, three complex networks are constructed.

An edge between node v_i and v_j is drawn if the distance between the input features of the i^{th} and j^{th} training data records is smaller than a user-defined threshold. For calculating the pairwise distance between data records, Euclidean distance function is used and can be calculated as Eq. (2):

$$Distance(v_i, v_j) = \sqrt{\sum_{p=1}^m (F_{i,p} - F_{j,p})^2} \quad (2)$$

Where m is the number of the input features, $F_{i,p}$ and $F_{j,p}$ denote the p^{th} input feature values for data records corresponding to v_i and v_j .

The weight of the edge $e(v_i, v_j)$ is calculated as Eq. (3):

$$w(v_i, v_j) = \frac{distance(v_i, v_j)}{\max(distance(v_k, v_h); \forall v_k, v_h \in V)} \quad (3)$$

2.7 Feature engineering based on the complex networks

In this section, three complex networks per data record are constructed including the considered data record, all training instances as CN1 and all training instances excluding negative (positive) instances as CN2 (CN3). A simple intuitive hypothesis is that a node has more similarity with the training instances of its own class compared to the instances of the other class. Therefore, the node centrality in different complex networks CN1, CN2 and CN3 can be compared to classify the node. Features listed in Table 2 are defined based on this hypothesis.

Table 2- list of the features engineered from the complex networks in this study

Node degree is the number of its adjacent edges. Betweenness centrality for graph nodes have been introduced by Bavelas [28] and is calculated as Eq. (4). If a node lies in many shortest paths between pairs of nodes, its Betweenness centrality will be high. Nodes with high Betweenness centrality are the bridges for information flow.

$$\text{Betweenness}(v_i) = \sum_{j < k} \frac{\text{number of the shortest paths between } v_j \text{ and } v_k \text{ passing } v_i}{\text{number of the shortest paths between } v_j \text{ and } v_k} \quad (4)$$

Node closeness centrality measures the reciprocal of the sum of the length of the shortest paths between the node and all other nodes in the graph.

Node Eigen vector centrality is higher when the node is pointed to by many important nodes.

Clustering coefficient of a node is calculated as Eq. (5):

$$\text{clusteringCoefficient}(v_i) = \frac{\text{number of triangles connected to } v_i}{\text{number of triples centered around } v_i} \quad (5)$$

Since, the number of the instances are very high, the complex networks are partitioned into smaller communities to reduce the computational complexity for calculating the engineered features.

One complex network extracted from only 100 data records treated by IUI method as a sample is shown in Figure 3.

Figure 3- One complex network extracted from only 100 data records treated by IUI method

Figure 4 depicts two complex networks of the same samples of positive instances drawn by different thresholds.

Figure 4- two complex networks drawn from the positive training data samples by (a) threshold of 0.7 * average of the distance matrix, (b) threshold of 0.5 * average of the distance matrix

As shown by Figure 4, reducing the threshold for keeping the edges in the complex network even with a small value lead to the network with more sparsity and more small-sized communities.

Figure 5 illustrates three complex networks from the samples of both classes, negative and/or positive classes.

Figure 5- three complex networks extracted from the samples of (a) both classes, (b) negative class, and (c) positive class with the same threshold

As shown by Figure 5, for the same thresholds, complex network considering the instances of both classes has the most density and the complex network from only positive instances has the most sparsity and consists of several small communities.

2.8 Training the stacked ensemble classifier

Stacked ensemble classifier which is a scalable meta-modeling methodology has been first introduced by Wolpert in 1994 [29]. It has been inspired by neural networks whose classifiers have been considered as the nodes. Instead of a linear model, the stacked classifier can use any base classifier. The stacking operation has been performed by either a normal stacking or a re-stacking mode. In the normal stacking mode, the base classifiers in each layer use the output scores of the

1 previous ones as the predictors similar to a typical feedforward neural network. The formula of
2 normal stacking mode is written as Eq. (6):

$$3 \quad f_n(x,V) = V_{n,k} \left(f_{n-1}(x,V_{n-1,1}), f_{n-1}(x,V_{n-1,2}), \dots, f_{n-1}(x,V_{n-1,D_{n-1}}) \right) \quad (6)$$

5
6 Where n indicates the nth layer of the stacked ensemble, x denotes a sample of a dataset, V presents
7 a vector holding the neurons (the base classifiers), D is the number of hidden neurons through the
8 nth hidden layer and finally, k is the kth neuron in the nth layer.

9
10 Some previous studies have illustrated that the stacked models can improve the performance of
11 the classification [20, 21, 30]. Therefore, in this study, a new stacked ensemble classifier is
12 proposed and designed based on the normal stacking mode. In the beginning, some of the basic
13 classifiers are trained, and those outperforming the others are selected to be considered as the base
14 classifiers in the stacked ensemble layers. The architecture of the proposed stacked ensemble
15 classifier is shown in Figure 6.

16
17
18 *Figure 6- (a) input datasets, and (b) the architecture of the proposed Stacked Ensemble classifier*

19 As illustrated in Figure 6, input dataset consists of the features in OFS, FS-Fi, EFS and/or EFS-Fi.
20 Input dataset is fed to the base classifiers in the first layer of the proposed stacked ensemble
21 classifier.

22 Several different classifiers are trained and verified. The classifiers for using in the ensemble layers of our
23 proposed stacked ensemble classifier are chosen among different trained classifiers with different values of
24 hyperparameters based on their accuracy and diversity on the validation dataset. A previous study has
25 proposed a method to choose classifiers for ensemble learning based on accuracy and diversity which is
26 used in this study for the same purpose (Yang, 2011). The pairwise diversity of the classifiers are calculated
27 using Q statistic.

28 Logistic regression (LR) [31], support vector machines (SVM) [32], decision tree (DT) [33],
29 random forest (RF) [26], Adaboost [34] and LightGradient Boosting Machine (LightGBM) [35]
30 are the base classifiers chosen based on their accuracy and diversity in both ensemble layers.

31 LR, SVM with linear kernel and DT are appropriate classifiers for classifying linearly separable
32 data. SVM with non-linear kernels, RF, Adaboost and LightGBM are ensemble classifiers which
33 can classify nonlinearly separable data with high performance. All the mentioned classifiers can
34 be trained fast. Therefore, they are chosen as the base classifiers of the proposed stacked ensemble
35 classifier.

36 The hyperparameters of the classifiers are tuned based on grid search method and the best values
37 for hyperparameters leading to the highest accuracy for validation dataset are considered for each
38 classifier.

39 After training the base classifiers in the first layer, their outputs are considered as Meta features
40 according to the normal stacking mode. The Meta features are fed into the base classifiers of the
41 second layer for training them. Finally, the outputs of the base classifiers in the second layer are
42 aggregated by weighted voting aggregation rule.

43 The weight of each base classifier is obtained by measuring its accuracy for classifying the
44 validation dataset. The validation dataset is about 20% of the original training dataset which is
45 excluded during the base classifiers' training in both layers.

1 Mathematical calculation is performed in this study to show the performance improvement
 2 obtained by stacked ensemble compared to traditional one-layer ensemble and the individual
 3 classifiers.

4 Without loss of generality, it is assumed that each base classifier in the first layer of stacked ensemble has
 5 the error rate of $\varepsilon < 0.50$. If the aggregation of the base classifiers is performed with bagging strategy
 6 which is the simplest aggregation method and uses majority voting, the error rate of the first ensemble
 7 layer (ε_{L1}) can be calculated as Eq. (7):

$$8 \quad \varepsilon_{L1} = \sum_{i=\lfloor \frac{M}{2} + 1 \rfloor}^M \binom{M}{i} \varepsilon^i (1 - \varepsilon)^{M-i} \quad (7)$$

9 Where M is the number of the base independent classifiers in the first ensemble layer. For misclassifying
 10 a data record using bagging strategy as the aggregation method, more than half of the base classifiers
 11 should misclassify the record. If it is assumed that i is the number of the base classifiers which misclassify
 12 the data record, i should be more than M/2 for misclassifying it with the first ensemble layer. For
 13 example, if M is 25, at least 13 base classifiers should misclassify data for erroneous classifying data in
 14 ensemble of these base classifiers. Now, if ε is 0.35 for each of 25 base classifiers, ε_{L1} will be 0.04. It
 15 shows the first layer of ensemble or traditional ensemble can improve the error rate of the single
 16 independent classifiers significantly.

17 Now, it is assumed that we have one more ensemble layer such as a two-layer stacked ensemble. Bagging
 18 strategy uses simple majority voting for classifying data as Eq. (8):

$$19 \quad classLabel_{ensemble}(r_j) = \begin{cases} Positive & \text{if } \sum_{i=1}^M \delta(classLabel_i(r_j) == Positive) > \frac{M}{2} \\ Negative & \text{otherwise} \end{cases} \quad (8)$$

20 Where r_j indicates the j^{th} data record and i denotes the i^{th} base classifier. As shown in Eq. (8), a simple
 21 decision tree or SVM with linear kernel can provide rules or find hyperplanes to classify data according to
 22 Eq. (8). Therefore, it can be shown that the performance of each base classifier in the second layer will
 23 not be worse than the simple bagging aggregation strategy used in the first ensemble layer.

24 This conclusion is true because each base classifier will try to find the hyperplane or rules to discriminate
 25 the training samples of two classes. But, bagging strategy uses simple majority voting. Furthermore, the
 26 input features (the first meta feature set as shown by Fig.6) for the base classifiers of the second ensemble
 27 layer are the same as the input features fed to the bagging strategy in the first ensemble layer. These input
 28 features are the output class labels generated by the base classifiers in the first layer. Therefore, the error
 29 rate of each base classifier in the second ensemble layer would be at most ε_{L1} .

30 The aggregation rule in the first ensemble layer is majority voting in the bagging strategy. The base
 31 classifiers try to separate the instances of different classes using linear or non-linear hyperplanes or rules.
 32 The input dataset for majority voting in the first ensemble layer is the first meta feature set. Therefore, the
 33 input of the majority voting rule and the base classifiers of the second ensemble layer is the same. The
 34 majority voting rule can be stated as Eq. (9) for the first meta feature set with M columns:

$$35 \quad label_{MV}(r_j) = \begin{cases} Positive & \text{if } \sum_{i=1}^M classLabel_i(r_j) > 0 \\ Negative & \text{otherwise} \end{cases} \quad (9)$$

36 Where MV is the majority voting strategy. Majority voting strategy is similar to using a hyperplane
 37 considering the equal coefficients for all of its input features as the separator of two classes.

1 The base classifiers try to find a best hyperplane for discriminating the instances of two classes.
 2 Therefore, their fitted hyperplane will not be worse than the hyperplane used with majority voting
 3 strategy. Thus, their performance will be more than or equal to the performance of the majority voting in
 4 the first ensemble layer. According to the Eq. (7), it is shown that the performance of the majority voting
 5 will be much better than the performance of the single classifiers in the first ensemble layer. Therefore,
 6 the performance of the single classifiers in the second ensemble layer will be better than the performance
 7 of the single classifiers in the first ensemble layer.

8 According to Eq. (7), if the bagging strategy is used for the second ensemble layer, the error rate of the
 9 second ensemble layer in the stacked ensemble would be ϵ_{L2} which can be calculated as Eq. (9):

$$10 \quad \epsilon_{L2} = \sum_{j=\binom{M2}{2}+1}^{M2} \binom{M2}{j} \epsilon_{b2}^j (1 - \epsilon_{b2})^{M2-j} \leq \sum_{j=\binom{M2}{2}+1}^{M2} \binom{M2}{j} \epsilon_{L1}^j (1 - \epsilon_{L1})^{M2-j} \quad (9)$$

11 Where M_2 is the number of the base classifier in the second ensemble layer of the stacked ensemble and
 12 ϵ_{b2} is the error rates of the base classifiers in the second ensemble layer. As mentioned in the previous
 13 paragraph, the error rate of each base classifier in the second layer would be at most ϵ_{L1} . Therefore, ϵ_{b2}
 14 will be not more than ϵ_{L1} .

15 According to Eq. (7) and Eq.(9), the relationship among ϵ , ϵ_{L1} and ϵ_{L2} can be shown in Eq. (10):

$$16 \quad \epsilon_{L2} \ll \epsilon_{L1} \ll \epsilon \quad (10)$$

17 A previous study have demonstrated that adding more layers to stack ensemble can improve the
 18 classification performance in terms of accuracy and AUC [1].

19 Based on the obtained results, it can be shown that adding more layers to stacked ensemble can improve its
 20 performance. Although, adding more layers has higher burden of time complexity and memory usage, too.
 21 There are a few studies considering the effect of the ensemble size or cardinality (the number of the base
 22 classifiers in the ensemble classifier) on the performance of the ensemble method [1, 2]. The previous
 23 studies have shown the ensemble size depends on the diversity of the base classifiers included in the
 24 ensemble and its aggregation rule [1, 2]. In addition, a previous study has examined different ensemble
 25 sizes including 10,20, 50 and 100 classifiers for bioinformatics applications [3]. They have shown that the
 26 best ensemble size has been 50 but the ensemble size of 10 is sufficient to achieve to highly reasonable
 27 performance [3].

29 **2.9 Scoring the ignored feature**

30 As mentioned in section 1.5, MDA score is calculated for each feature and is considered as the
 31 feature importance score.

33 **2.10 Evaluating and validating the trained models**

34 To evaluate the performances of the trained models, the performance measures for classification
 35 problems are used in this study including Accuracy, Sensitivity, Specificity and F-Score as shown
 36 in Eq. (11) -(14):

$$37 \quad Accuracy = \frac{TP + TN}{N} \quad (11)$$

$$38 \quad Sensitivity = \frac{TP}{TP + FN} \quad (12)$$

$$39 \quad Specificity = \frac{TN}{TN + FP} \quad (13)$$

$$F - Score = 2 \times \frac{Sensitivity \times Specificity}{Sensitivity + Specificity} \quad (14)$$

Where TP and FP (TN and FN) indicate the number of instances in the positive (negative) classes which are classified correctly and incorrectly, respectively.

Moreover, the area under the curve (AUC) of the receiver operating curve (ROC) is considered. In order to validate the results, the experiments are repeated 50 times, and each time the data is selected based on 5-fold C.V.

A novel method named as A-Test has been proposed in a previous study to calculate the structural risk of a classifier model as its instability with the new test data [36]. A-test calculates the misclassification error percentage $\Gamma_{\zeta,K}$ for different K values using the balanced K-fold validation. In this study, the values of $\Gamma_{\zeta,K}$ will be reported for different classifiers and different feature sets. $\Gamma_{\zeta,K}$ is calculated as Eq. (15):

$$\Gamma_{\zeta,K} = \frac{100}{N} \left(\sum_{i=1}^N \delta((predictedLabel == Negative). (realLabel == Positive)) + \sum_{i=1}^N \delta((predictedLabel == Positive). (realLabel == Negative)) \right) \quad K = 2 \dots K_{max} \quad (15)$$

Where K_{max} cannot be more than the size of the minority class. For estimating the structural risk of a classifier method, the average of the values of $\Gamma_{\zeta,K}$ is considered as Eq. (16):

$$\Gamma_{\zeta}^{\wedge} = \frac{\sum_{K=2}^{K_{max}} \Gamma_{\zeta,K}}{K_{max} - 1} \quad (16)$$

Where Γ_{ζ}^{\wedge} ranges from 0 to 100% which higher values show higher risk of classification and lower values show the higher capacity and generalization ability of the model. Therefore, the lower values of Γ_{ζ}^{\wedge} are more desired.

3 Experimental results

In this section, the features are ranked based on MDA obtained by ignoring them during the training of CNFE-SE. Then the partial dependencies between high-ranked features are discussed. Finally, the performance of the proposed model (CNFE-SE) is compared with other state-of-the-art classifiers.

Ranking the significance of features

Error! Reference source not found. Figure 7 represents top-20 important features with highest MDA score for IUI outcome prediction based on 50 repetitions of CNFE-SE training on different training samples. Post wash total motile sperm counts, female BMI, sperm motility grades a + b, total sperm motility and sperm motility grade c are high-ranked predictors of IUI outcome. Additionally, post-wash total motile sperm counts, female BMI, and total sperm counts are the features illustrated with dark blue colors in Figure 7 Error! Reference source not found., have the highest repetitions as the first informative features. Generally, the variables related to the men's semen analysis parameters are high-ranked features in this study.

Figure 7-Overview of top features ranked based on CNFE-SE

1 The Pearson correlation coefficients are calculated among the top-20 important features, and
2 Figure 8 depicts the heat map of the correlation coefficients.

3
4 *Figure 8- the pairwise correlation analysis of 20 most important features*

5 As shown by Figure 8, the male semen parameters are positively correlated to each other, the more
6 sperm concentration, the more total sperm count, and the more total motile sperm count. Also,
7 couples' duration of infertility and duration of marriage are positively correlated.
8 Figure 9 shows the exact values of MDA score for top-20 features in this study.

9
10 *Figure 9- MDA values of top-20 features in this study*

11 In addition, Table 3 lists MDA values of top-20 features.

12 *Table 3- MDA values of top-20 features*

13 14 **3.1 Partial dependency between the features**

15
16 Figure 10 depicts the partial dependency plots for the most important features. Partial dependency
17 plots show whether a feature has a positive or negative effect on the response variable when the
18 other ones are controlled. However, in order to interpret the graphs, we should note that changes
19 in the clinical pregnancy probabilities in terms of the value of the features, even the most
20 significant ones, are roughly small (the y-axis range is 0.44-0.52). Therefore, it is noteworthy that
21 none of the features could individually and significantly alter the pregnancy rates more than 0.52.
22 This finding underlines the value of the machine learning approach by determining the complicated
23 association between individual predictors to make an effective classification model.

24
25
26 *Figure 10- Partial dependency plots of nine features among the important features which the blue curves indicate locally
27 weighted smoothing. It shows pregnancy variation obtained by CNFE-SE (y-axis) as a function of a feature (x-axis) in IUI.*

28 According to the results of the partial dependency plots as shown by Figure 8, the clinical
29 pregnancy rate has raised with increased number of post-wash total motile sperm counts and after
30 processing sperm concentration. Also, when their values respectively vary upper than 100 million
31 and 30 million spermatozoa per ml, the rate of pregnancy reaches its highest rate. In addition, the
32 likelihood of IUI success increases through growing the number of total sperm counts which is
33 mentioned in the previous studies, too [37].

34 35 **3.2 Comparing the performance of CNFE-SE with other state-of-the-art classifiers**

36 Table 4 lists the performance measures for comparing CNFE-SE with other state of the art
37 classifiers.

38 *Table 4- comparing the performance of CNFE-SE with other state of the art classifiers*

39
40 Two different feature sets are considered as the input variables fed to the classifiers including all
41 296 features and only the most important features (top-20 features shown in Fig.6). Moreover,
42 CNFE-SE is trained and evaluated twice (one time without doing feature engineering (FE) and
43 another time with performing feature engineering).

1 The models are executed and trained on different random training samples up to 50 times and the
2 mean \pm standard deviation values are depicted in Table 4. The CNFE-SE outperforms the
3 compared models by AUC of 0.84 ± 0.01 , sensitivity of 0.79 ± 0.01 , specificity of 0.91 ± 0.01 , and
4 accuracy of 0.85 ± 0.01 when trains on all 296 features. Moreover, CNFE-SE has the superior
5 performance when only 20-top features are fed to it as input variables with AUC of 0.87 ± 0.01 ,
6 sensitivity of 0.82 ± 0.01 , specificity of 0.92 ± 0.01 and accuracy of 0.87 ± 0.01 . Our obtained results
7 show that feature engineering and considering only 20-top features improve the performance of
8 CNFE-SE.

9 Table 5 shows the confusion matrix of CNFE-SE for total dataset.

10
11 *Table 5- the confusion matrix of CNFE-SE for total dataset*

12
13 Figure 11 depicts ROC curve for CNFE-SE trained with all features.

14
15
16 *Figure 11- ROC curve for CNFE-SE trained with all features*

17 As shown by Figure 11, AUC of CNFE-SE trained on all features is 0.84 ± 0.01 . As illustrated by
18 Table 4, the compared single classifiers show almost weak performances. The main reason is that
19 the patients treated with IUI do not have complicated conditions and the leading cause of their
20 infertility is idiopathic. Therefore, the data of the two classes have high similarity with each other,
21 and their differentiation using single classifier is not an easy task. However, among these models,
22 Light-GBM as one of state-of-the-art machine learning algorithms has the second best
23 performance because it is a gradient boosting framework that uses tree-based learning algorithms
24 and not only covers multi hyper-parameters but also has more focus on the accuracy of the results
25 [35].

26
27 When the classes are imbalanced, Precision-Recall curve is a useful instrument for the
28 presentation of prediction success. A great area under this curve shows both high precision, which
29 is related to low false-positive rate, and high recall, refers to low false-negative rate. Figure 12
30 indicates the precision-recall curves for CNFE-SE trained using top-20 features.

31
32 *Figure 12 -Precision-recall curves for CNFE-SE trained using top-20 features*

33 As shown in Figure 12, CNFE-SE predicts both classes with highly reasonable performance.
34 Moreover, the results of A-test method for structural risk calculation for different combinations of
35 feature sets and classifiers are shown in Table 6.

36 *Table 6- Results of the A-Test: The values of $\Gamma^{\wedge}\zeta$ and the minimum value of $\Gamma\zeta$*

37
38 Lower values of $\Gamma^{\wedge}\zeta$ and $\Gamma\zeta$ shows lower risk of the classifier for classifying previously unseen
39 records and the higher capacity and generalization ability of the model. Therefore, the feature set
40 and classifier achieving the lower values of $\Gamma^{\wedge}\zeta$ and $\Gamma\zeta$ is more desired. As shown by Table 6,
41 CNFE-SE trained using top-20 features has the superior performance based on A-Test results.

4 DISCUSSION

In the current study, among the various features that significantly affect the IUI outcome, the most potential predictors are female BMI and semen quality parameters. Semen data such as sperm count and motility are illustrated as the most prognostic factors in pregnancies, conceived by IUI and their association with IUI outcome have demonstrated in some previous studies [38]. Moreover, some previous studies have confirmed that semen descriptors, after the swim-up procedure have been more important than the ones before sperm washing process [39, 40]. Similarly, the percentage of motile sperm and its progression in the ejaculate have been known as significant predictors in IUI outcome prediction in the literature [41, 42]. Sperm motility grades a + b (progressive motility) and grade d (immotile sperms) are also determined in this study as potential predictive factors for a successful IUI [43]. Thus, if their corresponding values are more than 20% and less than 15%, respectively, the IUI success rate is higher.

Furthermore, the results of this study indicates that the IUI success rate is almost low when the female BMI is abnormal (BMI is lower than 20 or larger than 30). If female BMI is about 25 as the normal BMI value, the probability of pregnancy increases. This finding is mentioned in the previous studies, too [44].

Previous studies have shown that pregnancy rate could be reduced by increase in the female age [42, 45]. The present study identifies that the women older than 38 have a lower chance of successful IUI. However, Edrem et al. have not found the female age to be a prognostic factor in the prediction of IUI outcome [46].

As shown in Figure 8, the duration of infertility inversely affects the fertility rate, and the decline in fecundity is acclaimed by some previous works, as well. Also, the previous studies have shown that when the couples' duration of infertility is less than six years, the pregnancy success rate is higher [47].

The total dose of gonadotropins is taken into account in this study as an important feature. Moreover, its significance has been considered recently, too [11]. This study identifies that the total dose of gonadotropin is positively correlated with the pregnancy rate. Moreover, other factors contributing to failure or success of IUI outcome according to this study's findings include semen volume, male age, sperm normal and amorphous morphology, duration of the marriage, and endometrial thickness which some of them have been demonstrated as the influential attributes in some previous studies [48-50].

Eventually, the CNFE-SE is trained using the 20 most important features and it yields surprisingly good performances (AUC =0.87, 95% CI 0.86-0.88). It shows that the model carried out by these features, demonstrates a highly reasonable performance.

Some studies consider different patients' cycles as independent of each other, which may lead to a biased result. For example, they have considered the first cycle information [16, 51]. Our reanalysis of the primary cycle data revealed that the AUC performances of Light-GBM and CNFE-SE are 0.62 ± 0.01 and 0.84 ± 0.01 , respectively, which does not change significantly when all the cycles are taken into account. Moreover, as shown in the materials and methods section, increasing the number of cycles augment the clinical pregnancy rate which are in line with the importance of this feature in subsequent IUI outcome [52, 53]. On the contrary, the variable cycle number has not identified as an important feature according to CNFE-SE feature scores. This

1 finding may be due to the high number of data in the first cycle compared to the second, third and
2 more cycles, which approximately 74% of the data belongs to the first cycle of IUI treatment.

3
4 Finally, our study has some restrictions. Some of the female hormonal tests including FSH, TSH,
5 LH, and AMH have not been measured in all the patients before beginning IUI cycle, and therefore
6 they are eliminated from the analysis due to their high missing value rate. At the Royan center, the
7 patients who are entering the IUI treatment cycles are those who do not have complicated
8 conditions, and the women's hormonal tests are usually normal. Moreover, the male BMI is
9 excluded because of its high rate of missing values. The features describing the geographic
10 information of couples' habitats are removed from the study due to their low quality data entry.

11
12 Currently machine learning algorithms has been increasingly employed in different medical fields
13 [8]. Therefore, through using machine learning methods, we are able to predict the success or
14 failure of the IUI cycle treatment outcome for each couple, based on their demographic
15 characteristics and cycle information. In other words, our proposed CNFE-SE model shows
16 superior performance among the compared state of the art classifiers. A decision support system
17 (DSS) can be designed and implemented based on CNFE-SE. This DSS can help the physicians to
18 choose other treatment plans for the couples and reduce patients' costs if their IUI cycle success
19 rate is low. The schematic of this medical assistance system is shown in Figure 13.

20
21
22 *Figure 13- Schematic of our proposed medical decision support system for IUI outcome prediction*

23 The proposed DSS is trained on the training dataset by CNFE-SE after preprocessing the collected
24 dataset. After completing the training of CNFE-SE, every time a new data record is registered in
25 the DSS, it can be classified by CNFE-SE into positive or negative outcome. The predicted
26 outcome for the new data record can assist the physicians to decide to treat the couple with IUI
27 method or not.

28 **Conclusion:**

29 In conclusion, the use of machine learning methods to predict the success or failure rate of the IUI
30 could effectively improve the evaluation performances in comparison with other classical
31 prediction models such as regression analysis. Furthermore, our proposed CNFE-SE model
32 outperforms the compared methods with highly reasonable accuracy. CNFE-SE can be used as
33 clinical decision-making assistance for the physicians to choose a beneficial treatment plan with
34 regards to their patients' therapy options, which would reduce the patients' costs as well.

35 The experimental results in this study show that the most important features for predicting IUI
36 outcome are semen parameters (sperm motility and concentration) as well as female BMI.

37 Some features which have been identified as good discriminative features for IUI outcome
38 prediction in the previous studies are excluded from this study because of their high missing value
39 rate. For example, some of the female hormonal tests including FSH, TSH, LH, and AMH are not
40 routinely measured in all the patients before IUI and they are excluded from the study. It is
41 proposed to augment dataset with data records without missing value in the mentioned features
42 and consider the excluded features to CNFE-SE, and then try to rank the augmented feature set
43 and evaluate the performance of the classifier.

44 On the other hand, some data records have noisy information which can reduce the performance
45 of the classifiers. As future work, it is suggested that improving the robustness of CNFE-SE against
46 the noisy data by including vote-boosting and other previously proposed methods for increasing

1 the noise robustness of the classifiers. Moreover, the data is highly imbalanced which can have
2 negative effect on the classifiers' performance. As another research opportunity, it is suggested
3 that reducing the influence of data distribution per class by incorporating the advanced balanced
4 sampling strategies.
5 Determining the optimal ensemble size is a challenging issue, yet. It is suggested that the impact
6 of the ensemble size on the overall performance of stacked ensemble is studied in the future studies
7 on different tasks and different datasets.

8 **Abbreviations:**

9 ACECR: academic center for education, culture and research
10 AMH: Anti-müllerian hormone
11 ANN: Artificial neural networks
12 ART: assisted reproductive technology
13 AUC: Area under curve
14 BMI: Body mass index
15 CN1: Complex network which is comprised of all the training data records as its nodes
16 CN2: Complex network which includes all training data excluding negative class
17 CN3: Complex network which includes all training data excluding positive class
18 CNFE-SE: Complex network-based feature engineering and stacked ensemble
19 C.V.: Cross validation
20 DSS: Decision support system
21 DT: Decision tree
22 FN: False negative
23 FP: False Positive
24 FSH: Follicle-stimulating hormone
25 HCA: Hierarchical clustering analysis
26 ICSI: Intracytoplasmic injection
27 IUI: Intrauterine Insemination
28 IVF: In-vitro fertilization (IVF)
29 K-NN: K-nearest neighbors
30 LH: Luteinizing Hormone
31 LR: Logistic regression
32 MDA: Mean decrease of accuracy
33 MLP: Multi-layer perceptron
34 N: Negative
35 NB: Naïve Bayes
36 P: positive
37 PCA: Principal component analysis
38 RF: Random forest
39 RBF: Radial basis function
40 SVM: Support vector machines
41 SD: Standard deviation
42 TN: True negative
43 TP: True positive
44 TSH: Thyroid-stimulating hormone
45

1 **Declarations**

2 **Ethics approval and consent to participate:**

3 This study is approved by the institutional review board of the ROYAN Institute
4 (IR.ACECR.ROYAN.REC.1398.213). The informed consent requirement for this study was waived
5 because this was a retrospective study with little patients' sensitive or personal information, and all data
6 were anonymized.

7 The full name of the ethics committee who approved this study is IR.ACECR.ROYAN.REC which
8 ROYAN Institute belongs to. The committee's reference number is
9 IR.ACECR.ROYAN.REC.1398.213.

10

11 **Consent for publication:**

12 Not applicable.

13

14 **Availability of data and materials**

15 Our study is a retrospective study of a 5-year couples' data undergoing IUI. Data is collected from
16 Reproductive Biomedicine Research Center, Royan Institute for 8,360 couples who underwent
17 11,255 IUI cycles were included. But, we are not allowed to share the original dataset because of
18 the privacy and security issues.

19 **Competing interests:**

20 The authors declare that there are no conflicts of interest.

21

22 **Funding:**

23 This study was not funded by any organization.

24

25 **Authors' Contributions**

26 Conceptualization: SR, TK and MT

27 Data curation: SR, TK and MT

28 Formal analysis: SR, TK and MT

29 Funding acquisition: there is no funding.

30 Investigation: SR, TK and MT

31 Methodology: SR and TK

32 Project administration: TK

33 Software: SR and TK

34 Supervision: TK and MT

35 Validation: SR, TK, AVTD, HS, MT, FG

36 Visualization: SR, TK and MT

37 Writing – original draft: SR and TK

38 Writing – review & editing: SR, TK, AVTD, HS, MT, FG

39 All authors have read and approved the manuscript.

40

41 **Acknowledgments:**

42 The authors acknowledge the Royan institute staffs, especially the informatics department for their
43 valuable contributions. There is no conflict of interest in this study.

1 REFERENCES

- 2 1. Medicine, P.C.o.t.A.S.f.R., *Definitions of infertility and recurrent pregnancy*
3 *loss: a committee opinion*. Fertil Steril, 2013. **99**(1): p. 63.
- 4 2. Borgh, M. and C. Wyns, *Fertility and infertility: Definition and*
5 *epidemiology*. Clinical Biochemistry, 2018. **62**: p. 2-10.
- 6 3. Milewska, A.J., et al., *Prediction of infertility treatment outcomes using*
7 *classification trees*. Studies in Logic, Grammar Rhetoric, 2016. **47**(1): p. 7-
8 19.
- 9 4. Blank, C., et al., *Prediction of implantation after blastocyst transfer in in vitro*
10 *fertilization: a machine-learning perspective*. Fertil Steril, 2019. **111**(2): p.
11 318-326.
- 12 5. Patil, A.S., *A Review of Soft Computing Used in Assisted Reproductive*
13 *Techniques (ART)*. International Journal of Engineering Trends and
14 Applications (IJETA), 2015. **2**(3): p. 88-93.
- 15 6. Bahadur, G., et al., *First line fertility treatment strategies regarding IUI and*
16 *IVF require clinical evidence*. Human Reproduction, 2016. **31**(6): p. 1141-
17 1146.
- 18 7. Ombelet, W., P. Puttemans, and E. Bosmans, *Intrauterine insemination: a*
19 *first-step procedure in the algorithm of male subfertility treatment*. Human
20 Reproduction, 1995. **10**(suppl_1): p. 90-102.
- 21 8. Deo, R.C., *Machine learning in medicine*. Circulation, 2015. **132**(20): p.
22 1920-1930.
- 23 9. Milewska, A.J., et al., *Analyzing Outcomes of Intrauterine Insemination*
24 *Treatment by Application of Cluster Analysis or Kohonen Neural Networks*.
25 Studies in Logic, Grammar Rhetoric, 2013. **35**(1): p. 7-25.
- 26 10. Kooptiwoot, S. and M.A. Salam, *IUI mining: human expert guidance of*
27 *information theoretic network approach*. Soft Computing, 2006. **10**(4): p.
28 369-373.
- 29 11. Ghaffari, F., et al., *Evaluating the effective factors in pregnancy after*
30 *intrauterine insemination: a retrospective study*. International journal of
31 fertility and sterility, 2015. **9**(3): p. 300.
- 32 12. Steures, P., et al., *Prediction of an ongoing pregnancy after intrauterine*
33 *insemination*. Fertil Steril, 2004. **82**(1): p. 45-51.
- 34 13. Goldman, R.H., et al., *Patient-specific predictions of outcome after*
35 *gonadotropin ovulation induction/intrauterine insemination*. Fertil Steril,
36 2014. **101**(6): p. 1649-1655. e2.
- 37 14. Marshburn, P.B., et al., *Spermatozoal characteristics from fresh and frozen*
38 *donor semen and their correlation with fertility outcome after intrauterine*
39 *insemination*. Fertil Steril, 1992. **58**(1): p. 179-186.

- 1 15. Moro, F., et al., *Anti-Müllerian hormone concentrations and antral follicle*
2 *counts for the prediction of pregnancy outcomes after intrauterine*
3 *insemination*. International Journal of Gynecology and Obstetrics, 2016.
4 **133**(1): p. 64-68.
- 5 16. Lemmens, L., et al., *Predictive value of sperm morphology and progressively*
6 *motile sperm count for pregnancy outcomes in intrauterine insemination*.
7 Fertil Steril, 2016. **105**(6): p. 1462-1468.
- 8 17. Arslan, M., et al., *Predictive value of the hemizona assay for pregnancy*
9 *outcome in patients undergoing controlled ovarian hyperstimulation with*
10 *intrauterine insemination*. Fertil Steril, 2006. **85**(6): p. 1697-1707.
- 11 18. Florio, P., et al., *Evaluation of endometrial activin A secretion for prediction*
12 *of pregnancy after intrauterine insemination*. Fertil Steril, 2010. **93**(7): p.
13 2316-2320.
- 14 19. Shah, S. and A. Kusiak, *Cancer gene search with data-mining and genetic*
15 *algorithms*. Computers in Biology and Medicine, 2007. **37**(2): p. 251-261.
- 16 20. Kaya, A., *Cascaded classifiers and stacking methods for classification of*
17 *pulmonary nodule characteristics*. Computer Methods and Programs in
18 Biomedicine, 2018. **166**: p. 77-89.
- 19 21. Wang, S.Q., J. Yang, and K.C. Chou, *Using stacked generalization to predict*
20 *membrane protein types based on pseudo-amino acid composition*. Journal of
21 theoretical biology, 2006. **242**(4): p. 941-946.
- 22 22. Tocci, A. and C. Lucchini, *WHO reference values for human semen*. Human
23 reproduction update, 2010. **16**(5): p. 559-559.
- 24 23. Zhang, S., C. Zhang, and Q. Yang, *Data preparation for data mining*. Applied
25 artificial intelligence, 2003. **17**(5-6): p. 375-381.
- 26 24. Han, J., J. Pei, and M. Kamber, *Data mining: concepts and techniques*. 2011:
27 Elsevier.
- 28 25. Liu, F.T., K.M. Ting, and Z.H. Zhou, *Isolation Forest*, in *2008 Eighth IEEE*
29 *International Conference on Data Mining*. 2008, IEEE. p. 413-422.
- 30 26. Breiman, L., *Random Forests*. Machine Learning, 2001. **45**: p. 5-32.
- 31 27. Diykh, M., Y. Li, and S. Abdulla, *EEG Sleep Stages Identification Based on*
32 *Weighted Undirected Complex Networks*. Computer Methods and Programs
33 in Biomedicine, 2020. **184**: p. 105116.
- 34 28. Bavelas, A., *A mathematical model for group structure, human organization*.
35 Appl. Anthropol., 1948. **7**(3): p. 16-30.
- 36 29. Wolpert, D.H., *Stacked generalization*. Neural networks, 1992. **5**(2): p. 241-
37 259.
- 38 30. Güneş, F., R. Wolfinger, and P.Y. Tan. *Stacked ensemble models for improved*
39 *prediction accuracy*. in *Static Anal. Symp*. 2017.

- 1 31. Sperandei, S., *Understanding logistic regression analysis*. Biochem med,
2 2014. **24**(1): p. 12-18.
- 3 32. Cortes, C. and V. Vapnik, *Support-vector network*. Machine Learning, 1995.
4 **20**: p. 1-25.
- 5 33. Quinlan, J.R., *Induction of Decision Trees*. Machine Learning, 1986. **1**: p. 81-
6 106.
- 7 34. Zhu, J., et al., *Multi-class AdaBoost*. Statistics and its interfere, 2009. **2**: p.
8 349-360.
- 9 35. Ke, G., et al. *Lightgbm: A highly efficient gradient boosting decision tree*. in
10 *Advances in Neural Information Processing Systems*. 2017.
- 11 36. Gharehbaghi, A. and M. Linden, *A Deep Machine Learning Method for*
12 *Classifying Cyclic Time Series of Biological Signals Using Time-Growing*
13 *Neural Network*. IEEE Transactions on Neural Networks and Learning
14 Systems, 2018. **29**(9): p. 4102-4115.
- 15 37. Campana, A., et al., *Intrauterine insemination: evaluation of the results*
16 *according to the woman's age, sperm quality, total sperm count per*
17 *insemination and life table analysis*. Human Reproduction, 1996. **11**(4): p.
18 732-736.
- 19 38. Kuriya, A., C. Agbo, and M.H. Dahan, *Do pregnancy rates differ with intra-*
20 *uterine insemination when different combinations of semen analysis*
21 *parameters are abnormal?* Journal of the Turkish German Gynecological
22 Association, 2018. **19**(2): p. 57.
- 23 39. Zhang, E., et al., *Effect of sperm count on success of intrauterine insemination*
24 *in couples diagnosed with male factor infertility*. Materia socio-medica, 2014.
25 **26**(5): p. 321.
- 26 40. Ombelet, W., et al., *Semen quality and intrauterine insemination*.
27 Reproductive BioMedicine Online, 2003. **7**(4): p. 485-492.
- 28 41. Dickey, R.P., et al., *Comparison of the sperm quality necessary for successful*
29 *intrauterine insemination with World Health Organization threshold values*
30 *for normal sperm*. Fertil Steril, 1999. **71**(4): p. 684-689.
- 31 42. Duran, H.E., et al., *Sperm DNA quality predicts intrauterine insemination*
32 *outcome: a prospective cohort study*. Human Reproduction, 2002. **17**(12): p.
33 3122-8.
- 34 43. Muriel, L., et al., *Value of the sperm chromatin dispersion test in predicting*
35 *pregnancy outcome in intrauterine insemination: a blind prospective study*.
36 Human Reproduction, 2006. **21**(3): p. 738-744.
- 37 44. Thijssen, A., et al., *Predictive factors influencing pregnancy rates after*
38 *intrauterine insemination with frozen donor semen: a prospective cohort*
39 *study*. Reproductive biomedicine online, 2017. **34**(6): p. 590-597.

- 1 45. Merviel, P., et al., *Predictive factors for pregnancy after intrauterine*
2 *insemination (IUI): An analysis of 1038 cycles and a review of the literature.*
3 *Fertility and Sterility*, 2010. **93**(1): p. 79-88.
- 4 46. Erdem, A., et al., *Factors affecting live birth rate in intrauterine insemination*
5 *cycles with recombinant gonadotrophin stimulation.* *Reproductive*
6 *biomedicine online*, 2008. **17**(2): p. 199-206.
- 7 47. Kamath, M.S., et al., *Predictive factors for pregnancy after intrauterine*
8 *insemination: A prospective study of factors affecting outcome.* *Human*
9 *reproductive sciences*, 2010. **3**(3): p. 129.
- 10 48. Licht, R.S., L. Handel, and M. Sigman, *Site of semen collection and its effect*
11 *on semen analysis parameters.* *Fertil Steril*, 2008. **89**(2): p. 395-397.
- 12 49. Francavilla, F., et al., *Effect of sperm morphology and motile sperm count on*
13 *outcome of intrauterine insemination in oligozoospermia and/or*
14 *asthenozoospermia.* *Fertil Steril*, 1990. **53**(5): p. 892-897.
- 15 50. Luco, S.M., et al., *The evaluation of pre and post processing semen analysis*
16 *parameters at the time of intrauterine insemination in couples diagnosed with*
17 *male factor infertility and pregnancy rates based on stimulation agent. A*
18 *retrospective cohort study.* *European Journal of Obstetrics Gynecology:*
19 *Reproductive Biology Endocrinology*, 2014. **179**: p. 159-162.
- 20 51. Blank, C., et al., *Prediction of implantation after blastocyst transfer in in vitro*
21 *fertilization: a machine-learning perspective.* 2019. **111**(2): p. 318-326.
- 22 52. Nuojua-Huttunen, S., et al., *Intrauterine insemination treatment in*
23 *subfertility: an analysis of factors affecting outcome.* *Human Reproduction*,
24 1999. **14**(3): p. 698-703.
- 25 53. Liu, W., et al., *Comparing the pregnancy rates of one versus two intrauterine*
26 *inseminations (IUIs) in male factor and idiopathic infertility.* *Journal of*
27 *assisted reproduction genetics*, 2006. **23**(2): p. 75-79.

28 **Figure legend:**

29 Figure 1- the main steps of the proposed method (CNFE-SE) for feature scoring and classifying
30 the patients to predict IUI outcome

31 Figure 2 –The ratio of positive and negative clinical pregnancy per treatment cycle

32 Figure 3- One complex network extracted from only 100 data records treated by IUI method

33 Figure 4- two complex networks drawn from the positive training data samples by (a) threshold of
34 $0.7 * \text{average of the distance Matrix}$, (b) threshold of $0.5 * \text{average of the distance matrix}$

35 Figure 5- three complex networks extracted from the samples of (a) both classes, (b) negative class,
36 and (c) positive class with the same threshold

37 Figure 6- (a) input datasets, and (b) the architecture of the proposed Stacked Ensemble classifier

38 Figure 7-Overview of top features ranked based on CNFE-SE

39 Figure 8- Partial dependency plots of nine features among the important features which the blue
40 curves indicate locally weighted smoothing. It shows pregnancy variation obtaining from CNFE-
41 SE (y-axis) as a function of a feature (x-axis) in IUI.

- 1 Figure 9- ROC curve for CNFE-SE trained with all features
- 2 Figure 10-Precision-recall curves for CNFE-SE
- 3 Figure 11- Schematic of our proposed medical decision support system for IUI outcome prediction

Figures

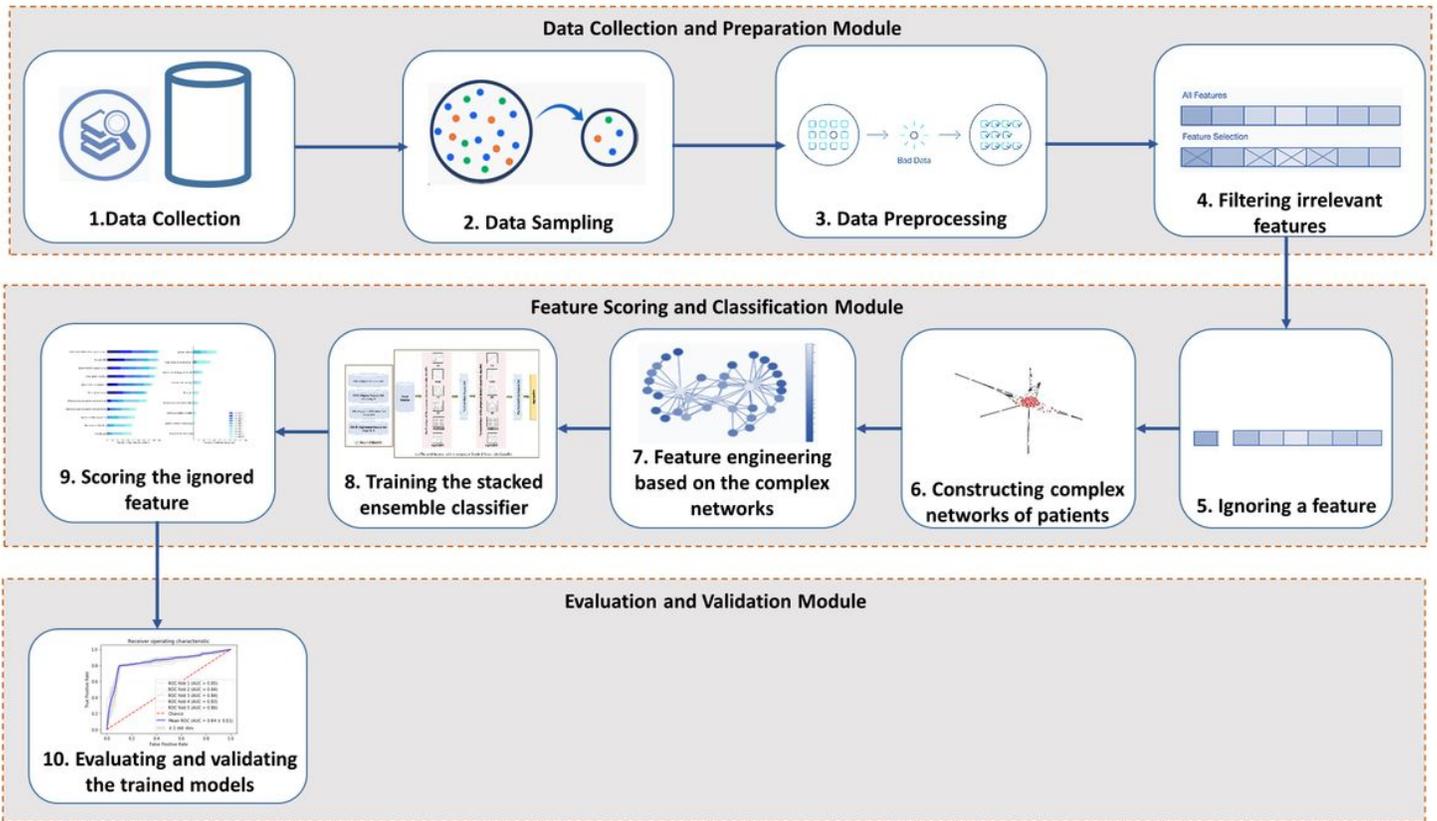


Figure 1

the main steps of the proposed method (CNFE-SE) for feature scoring and classifying the patients to predict IUI outcome

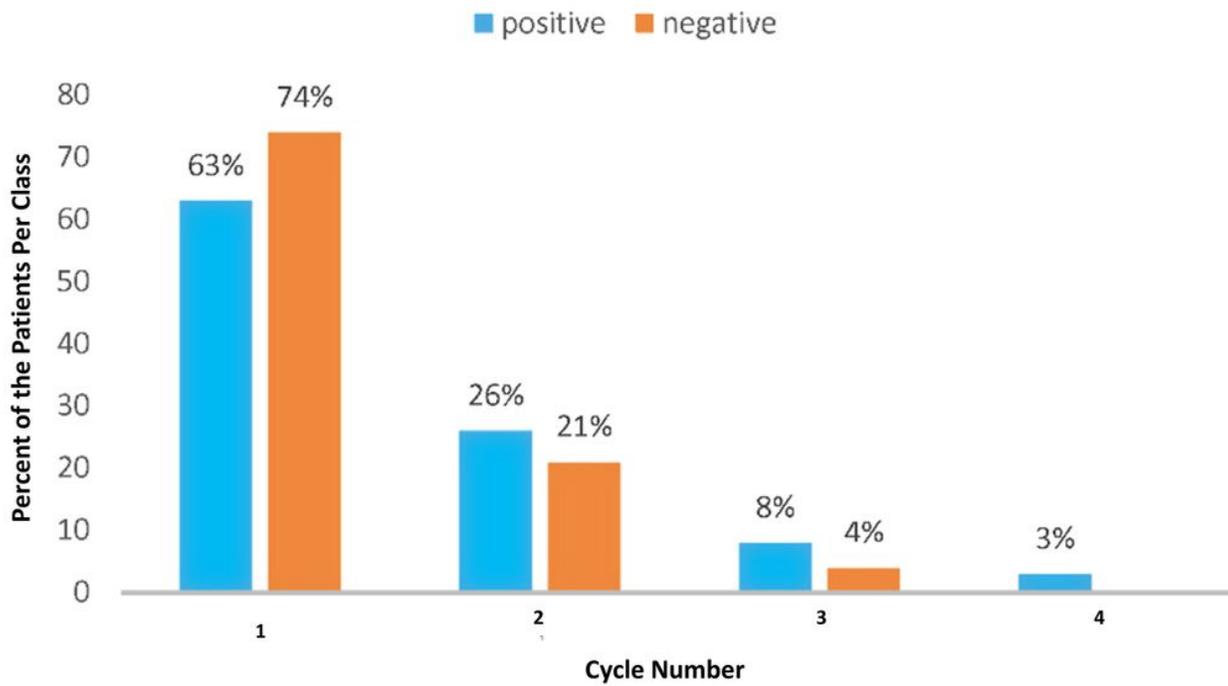


Figure 2

The ratio of positive and negative clinical pregnancy per treatment cycle

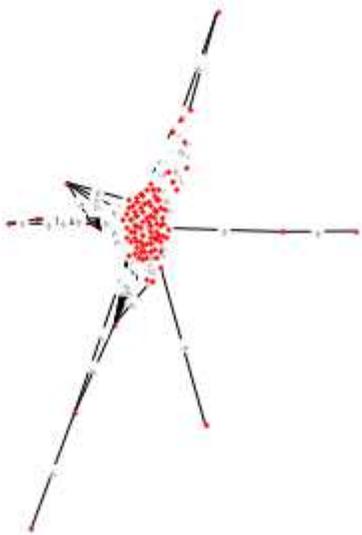


Figure 3

One complex network extracted from only 100 data records treated by IUI method

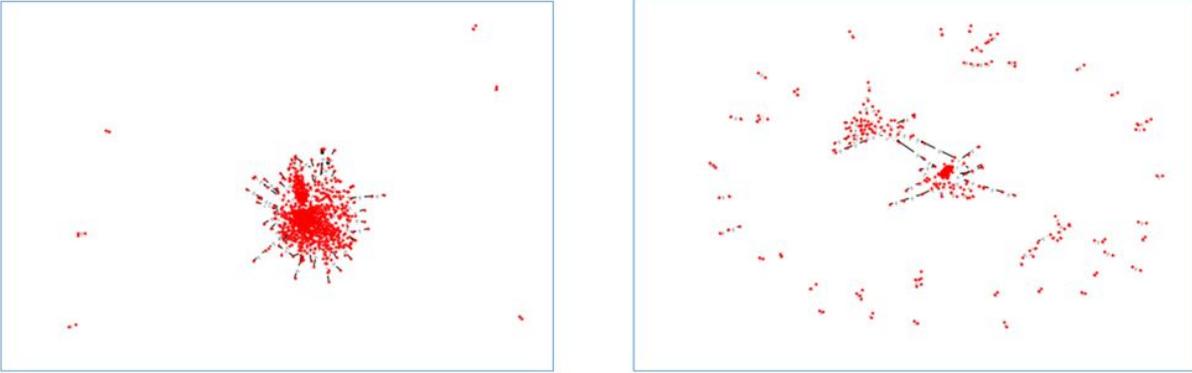


Figure 4

two complex networks drawn from the positive training data samples by (a) threshold of $0.7 * \text{average of the distance Matrix}$, (b) threshold of $0.5 * \text{average of the distance matrix}$

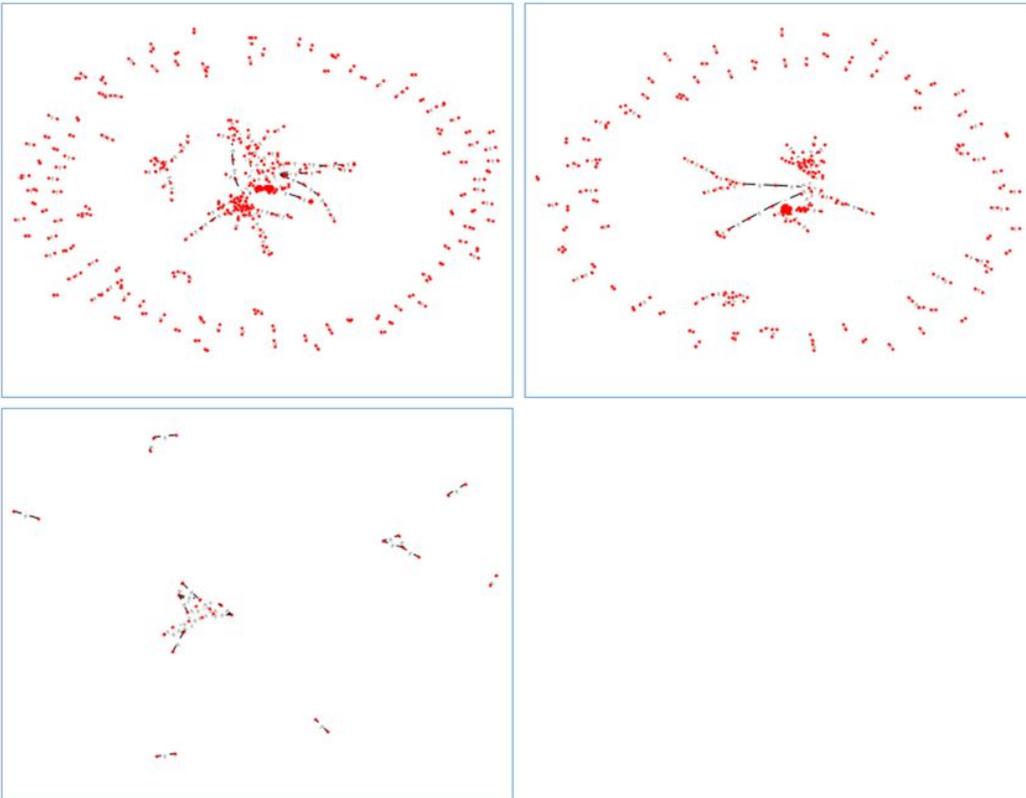


Figure 5

three complex networks extracted from the samples of (a) both classes, (b) negative class, and (c) positive class with the same threshold

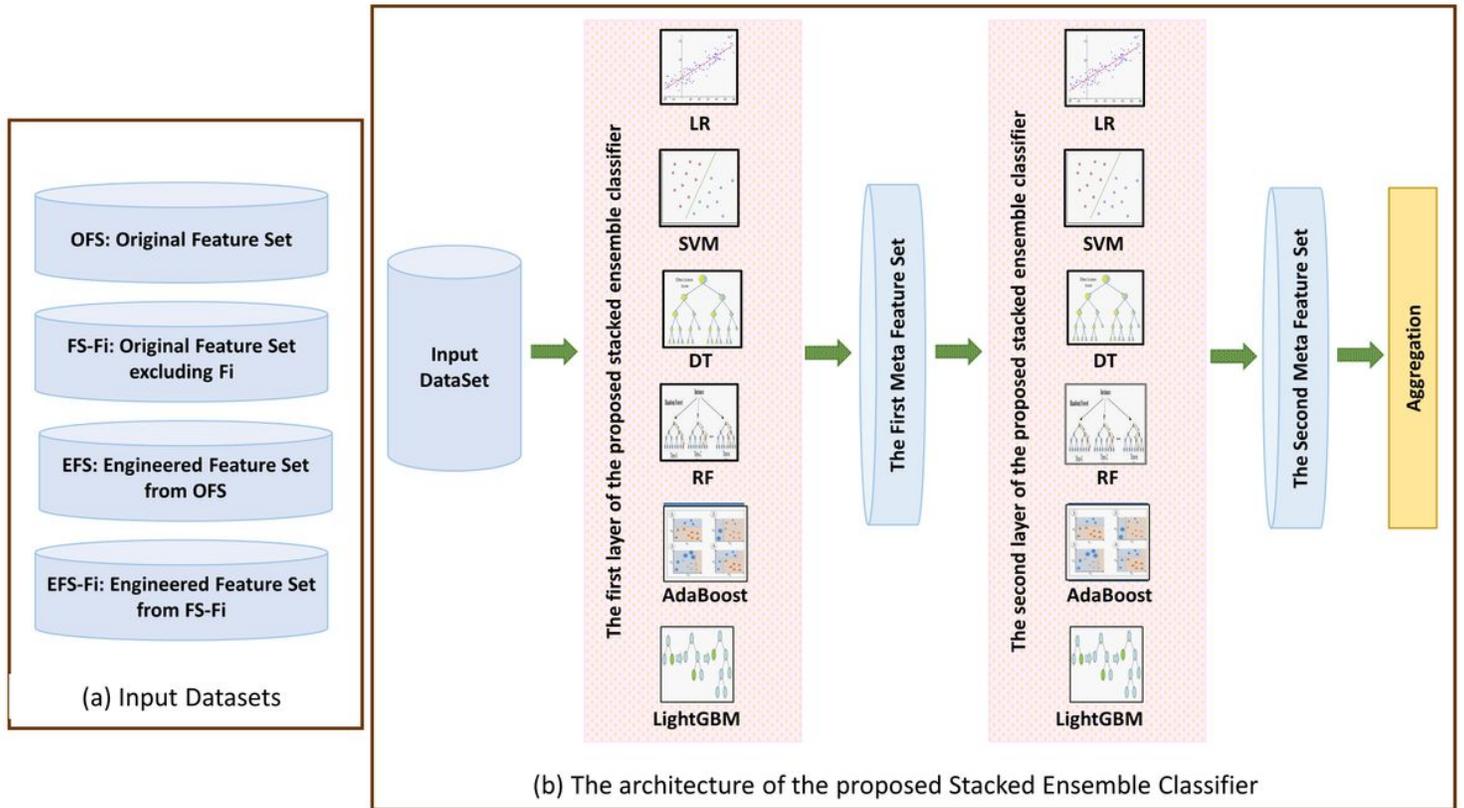


Figure 6

(a) input datasets, and (b) the architecture of the proposed Stacked Ensemble classifier

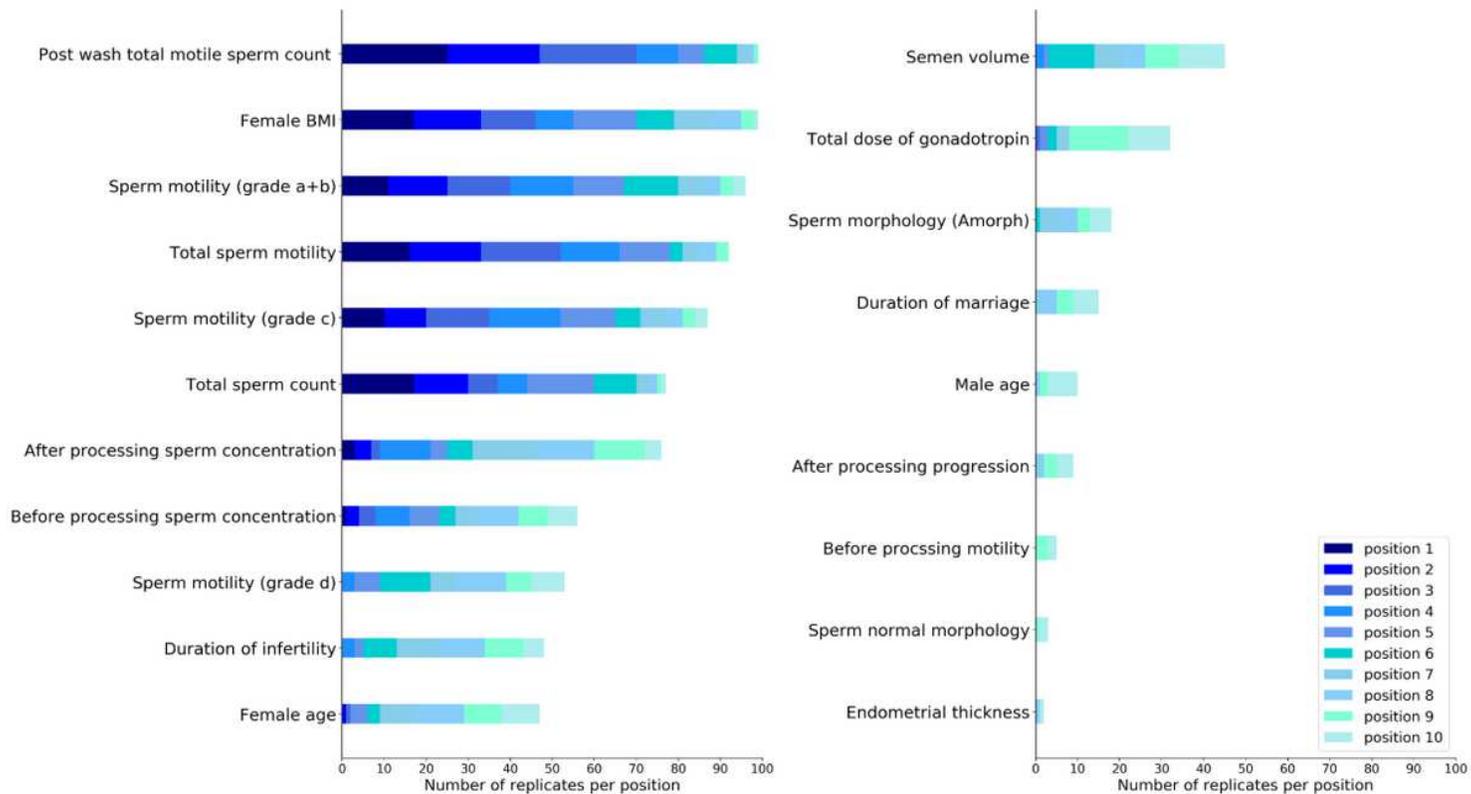


Figure 7

Overview of top features ranked based on CNFE-SE

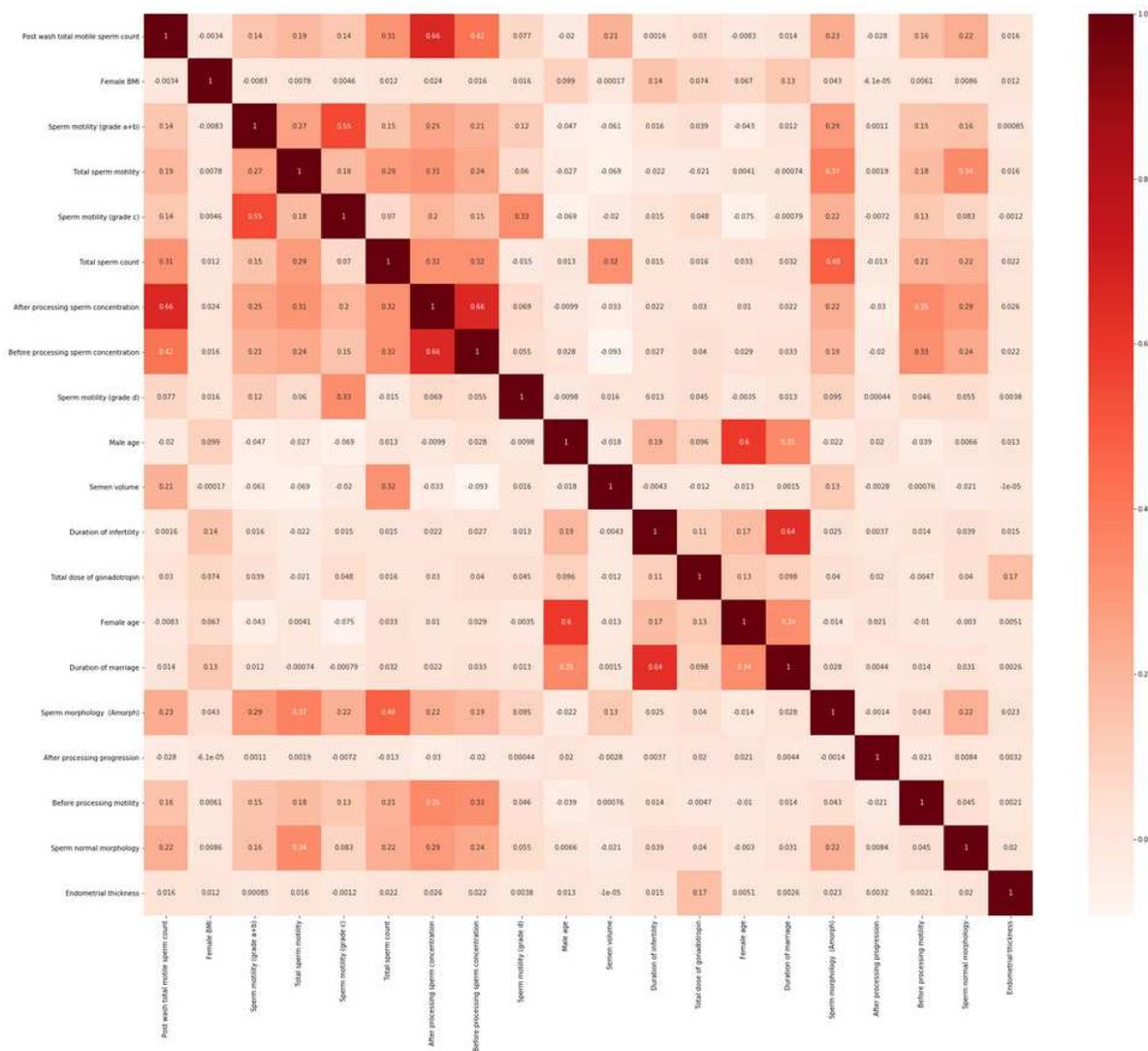


Figure 8

the pairwise correlation analysis of 20 most important features

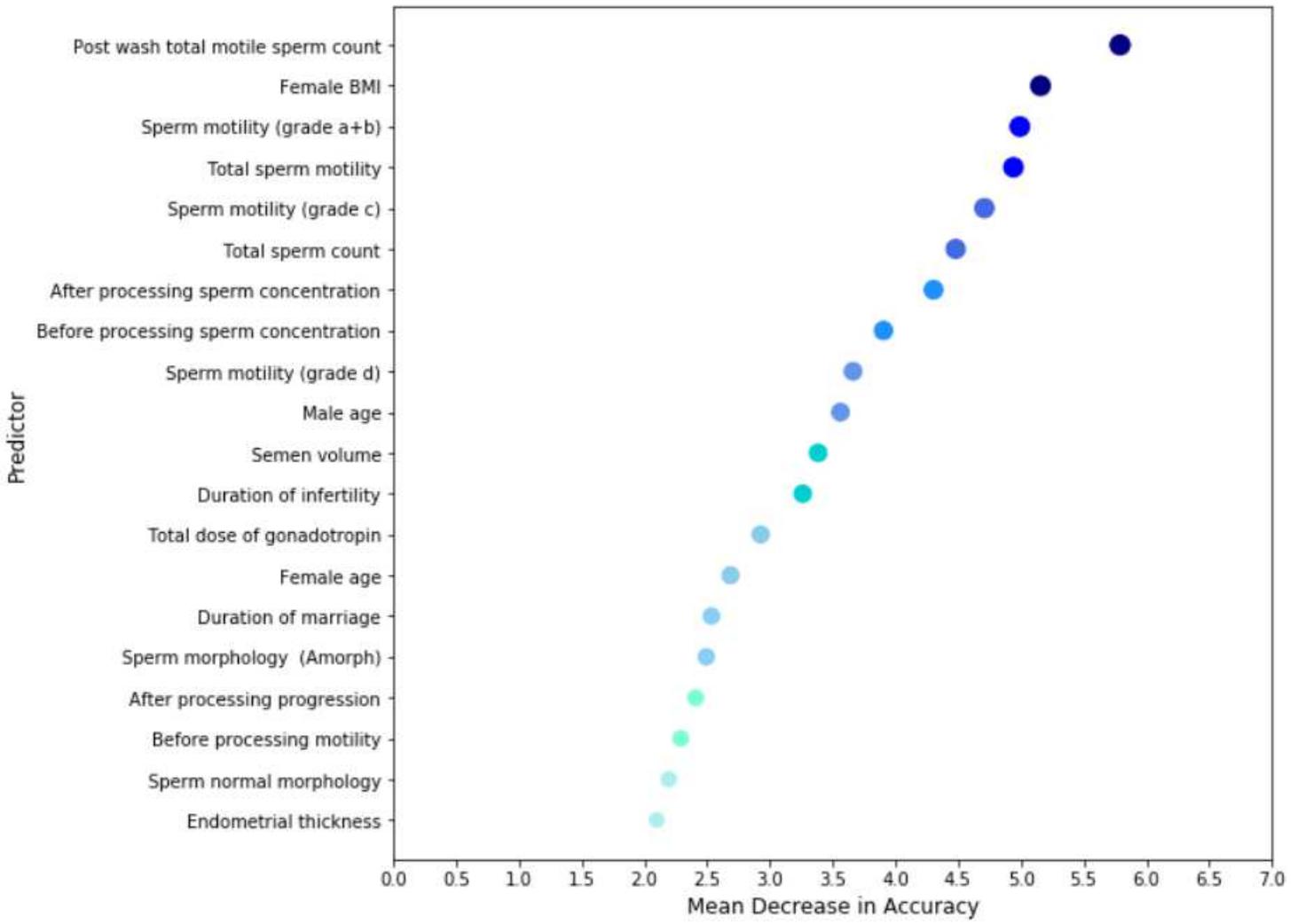


Figure 9

MDA values of top-20 features in this study

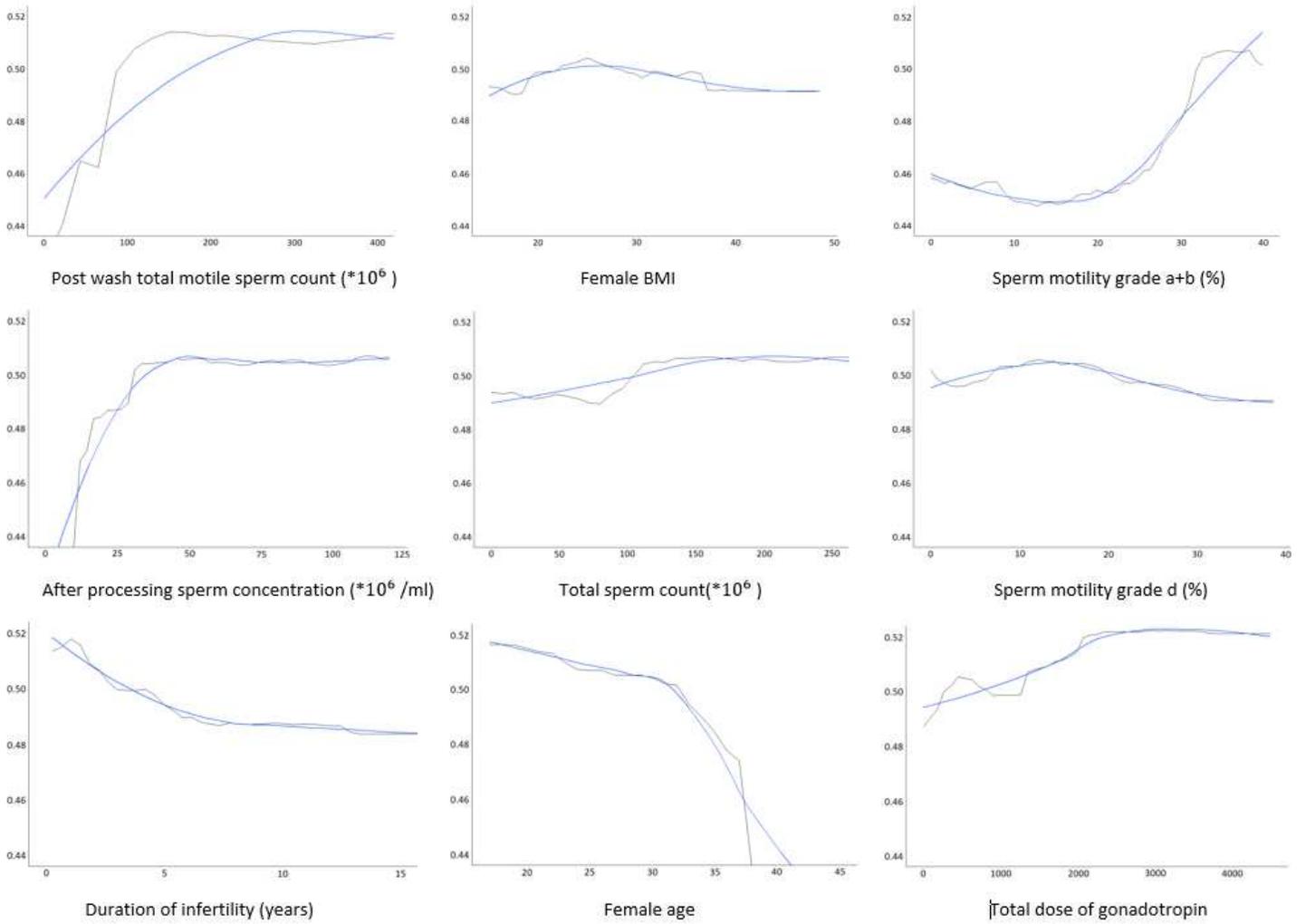


Figure 10

Partial dependency plots of nine features among the important features which the blue curves indicate locally weighted smoothing. It shows pregnancy variation obtained by CNFE-SE (y-axis) as a function of a feature (x-axis) in IUI.

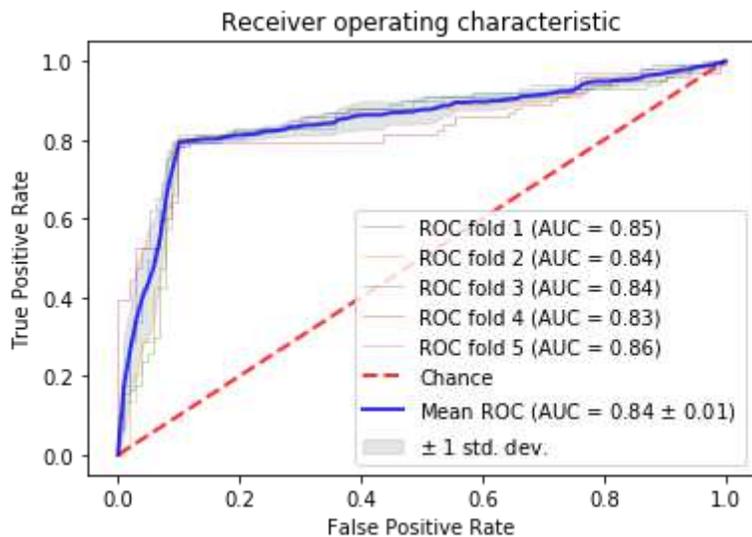


Figure 11

ROC curve for CNFE-SE trained with all features

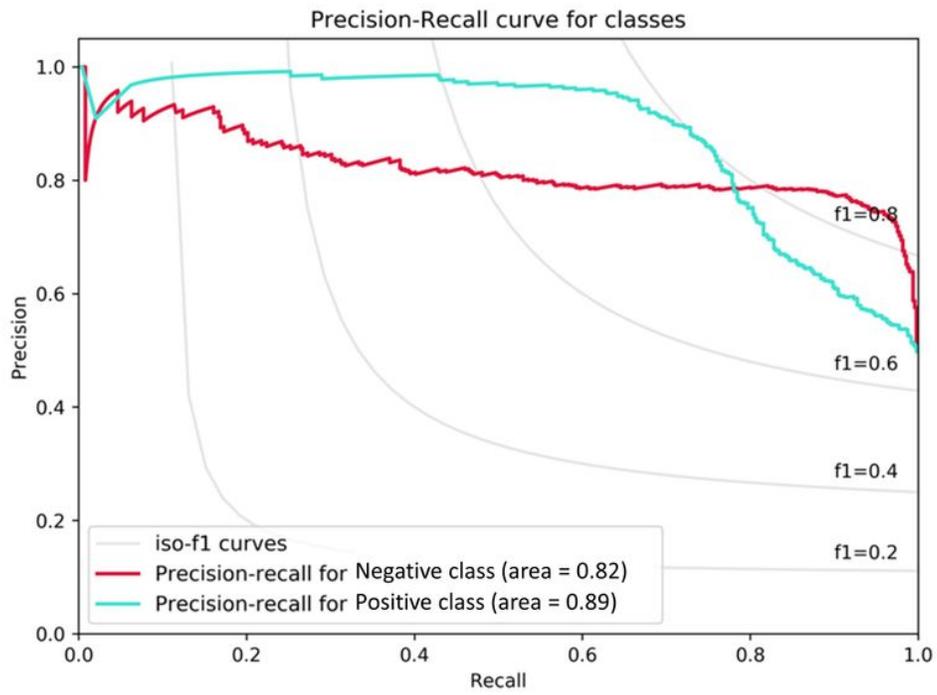


Figure 12

Precision-recall curves for CNFE-SE trained using top-20 features

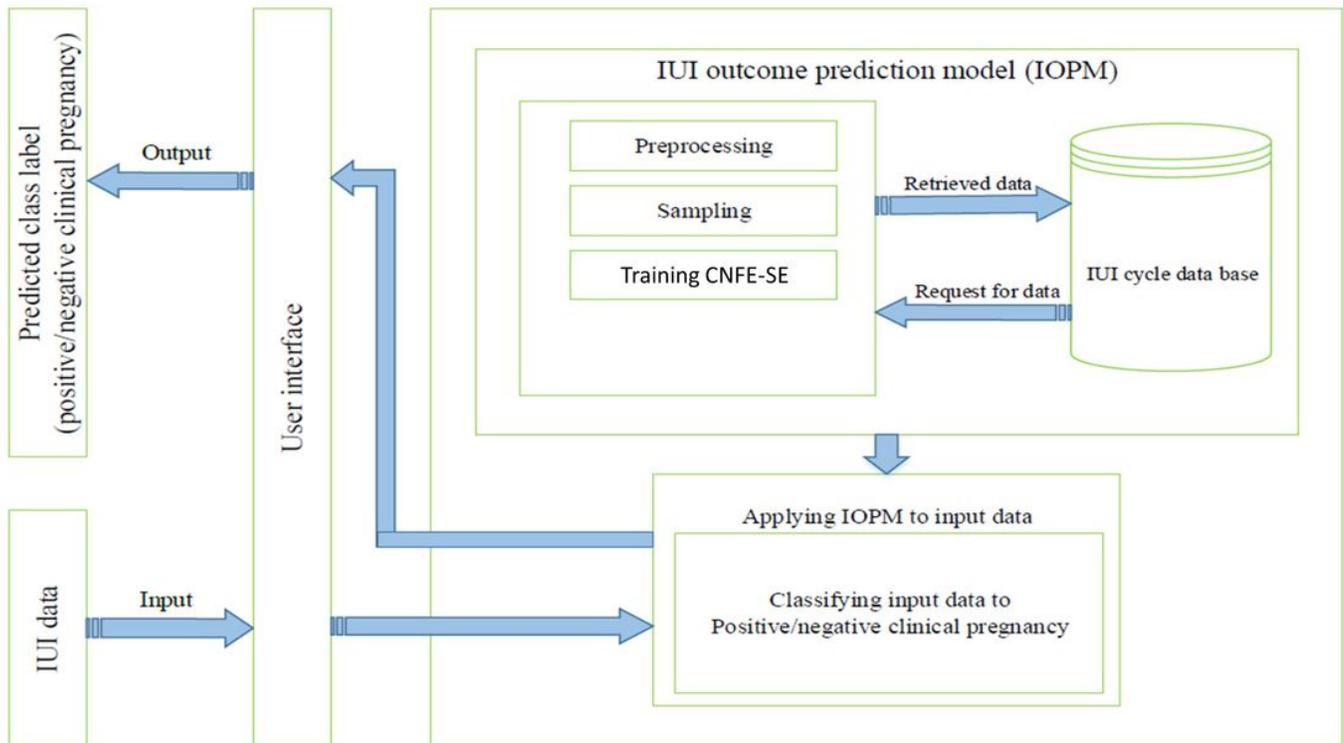


Figure 13

Schematic of our proposed medical decision support system for IUI outcome prediction

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [EditedAppendix.docx](#)