

1 **Predicting polycystic ovary syndrome (PCOS) with machine learning algorithms from electronic**
2 **health records**

3
4 Zahra Zad, BS^{1*}; Victoria S. Jiang, MD^{2*}; Amber T. Wolf, BA³; Taiyao Wang, PhD¹; J. Jojo Cheng, BA⁴;
5 Ioannis Ch. Paschalidis, PhD^{1,5}; Shruthi Mahalingaiah, MD, MS^{2,6~}

6
7 ¹Division of Systems Engineering, Center for Information and Systems Engineering (CISE), Boston
8 University, 15 St. Mary's Street, Brookline, MA 02446, USA

9 ²Division of Reproductive Endocrinology and Infertility, Department of Obstetrics and Gynecology,
10 Massachusetts General Hospital, 55 Fruit Street, Yawkey 10, Boston, MA 02114, USA

11 ³Icahn School of Medicine at Mount Sinai, 1 Gustave L. Levy Place, New York, NY 10029, USA

12 ⁴Department of Biostatistics and Medical Informatics, University of Wisconsin, West Johnson Street,
13 Madison, WI 53792, USA

14 ⁵Department of Electrical & Computer Engineering, Department of Biomedical Engineering, and Faculty
15 for Computing & Data Sciences, Boston University, 8 St. Mary's Street, Boston, MA 02215, USA

16 ⁶Department of Environmental Health, Harvard T.H. Chan School of Public Health, 665 Huntington
17 Avenue, Boston, MA 02115, USA

18
19 * Zahra Zad and Victoria S. Jiang contributed equally to this manuscript and are co-first authors of this
20 manuscript.

21
22 ~Corresponding Author & reprint requests:

23 Shruthi Mahalingaiah, MD, MS. **ORCID: 0000-0002-5527-5787**

24 Email: shruthi@hsph.harvard.edu

25
26 **Study funding/competing interest(s):** This study was partially supported by National Science
27 Foundation grants CCF-2200052, IIS-1914792, and DMS-1664644, by the NIH under grants R01
28 GM135930 and UL54 TR004130, and by the Boston University Kilachand Fund for Integrated Life
29 Science and Engineering

30
31 **Disclosure Summary:** The authors declare no conflict of interest and nothing to disclose.

32 **Abstract:**

33 **Introduction:** Predictive models have been used to aid early diagnosis of PCOS, though existing models
34 are based on small sample sizes and limited to fertility clinic populations. We built a predictive model
35 using machine learning algorithms based on an outpatient population at risk for PCOS to predict risk and
36 facilitate earlier diagnosis, particularly among those who meet diagnostic criteria but have not received a
37 diagnosis.

38
39 **Methods:** This is a retrospective cohort study from a SafetyNet hospital's electronic health records
40 (EHR) from 2003-2016. The study population included 30,601 women aged 18-45 years without
41 concurrent endocrinopathy who had any visit to Boston Medical Center for primary care, obstetrics and
42 gynecology, endocrinology, family medicine, or general internal medicine. Four prediction outcomes
43 were assessed for PCOS. The first outcome was PCOS ICD-9 diagnosis with additional model outcomes
44 of algorithm-defined PCOS. The latter was based on Rotterdam criteria and merging laboratory values,
45 radiographic imaging, and ICD data from the EHR to define irregular menstruation, hyperandrogenism,
46 and polycystic ovarian morphology on ultrasound.

47
48 **Results:** We developed predictive models using four machine learning methods: logistic regression,
49 supported vector machine, gradient boosted trees, and random forests. Hormone values (follicle-
50 stimulating hormone, luteinizing hormone, estradiol, and sex hormone binding globulin) were combined
51 to create a multilayer perceptron score using a neural network classifier. Prediction of PCOS prior to
52 clinical diagnosis in an out-of-sample test set of patients achieved AUC of 85%, 81%, 80%, and 82%,
53 respectively in Models I, II, III and IV. Significant positive predictors of PCOS diagnosis across models
54 included hormone levels and obesity; negative predictors included gravidity and positive bHCG.

55
56 **Conclusions:** Machine learning algorithms were used to predict PCOS based on a large at-risk
57 population. This approach may guide early detection of PCOS within EHR-interfaced populations to
58 facilitate counseling and interventions that may reduce long-term health consequences. Our model
59 illustrates the potential benefits of an artificial intelligence-enabled provider assistance tool that can be
60 integrated into the EHR to reduce delays in diagnosis. However, model validation in other hospital-based
61 populations is necessary.

62
63 **Keywords:** Polycystic ovary syndrome (PCOS), disease prediction, predictive model, machine learning,
64 artificial intelligence.

65
66 **Word count:** 5,712

67
68 **Number of figures and tables:** 7 (& 7 supplementary tables)

69 **Introduction**

70 Polycystic ovary syndrome (PCOS) is the most common type of ovulation disorder and
71 endocrinopathy among reproductive age women. PCOS is a diagnosis of exclusion after other
72 endocrinopathies known to affect ovulation have been evaluated including thyroid, adrenal, and pituitary
73 related disease. Based on the Rotterdam criteria, PCOS is diagnosed when two of the three following
74 criteria are exhibited: clinical or biochemical hyperandrogenism, oligo-anovulation, and polycystic ovary
75 morphology (PCOM) on transvaginal or transabdominal ultrasound. PCOS has a population prevalence of
76 5-15%, depending on the diagnostic criteria used (1).

77 PCOS is associated with multiple health issues and increased morbidity and mortality, including a
78 high chronic disease burden that is also very costly for individuals with PCOS and insurers (2). PCOS is
79 the leading cause of anovulatory infertility in reproductive-aged women. In fact, over 90% of anovulatory
80 women who present to infertility clinics have PCOS (3). PCOS patients have an increased risk of
81 endometrial hyperplasia and endometrial cancer (4) due to anovulatory cycles leading to long periods of
82 exposure to the effects of unopposed estrogen. PCOS has been associated with the development of
83 metabolic syndrome (5), diabetes (6), cerebrovascular disease and hypertension (7), compared to women
84 without PCOS. Despite these serious health consequences, PCOS frequently goes undiagnosed due to the
85 wide range of symptom severity on presentation, leading to delayed treatment and potentially more severe
86 clinical sequelae due to lack of preventive care, health management, and counseling (4). Even when
87 PCOS is diagnosed, it is often very delayed. One study found that over one-third of women with PCOS
88 waited over two years and were seen by three or more providers before finally receiving the diagnosis (8).

89 Predictive models can play a significant role in aiding earlier diagnosis of PCOS, though several
90 include only those women presenting for fertility care. One model used serum anti-Müllerian hormone
91 (AMH) and androstenedione levels, menstrual cycle length, and BMI to predict the development of PCOS
92 in Chinese women (9). Another model used only AMH and BMI to predict a diagnosis of PCOS or other
93 ovulatory dysfunction disorders (10). Other studies have created predictive models for certain outcomes
94 among women with PCOS such as pregnancy outcomes (11,12) and insulin resistance (13). In this study,
95 we use clinical and socioeconomic variables among 30,601 women aged 18 to 45 years within the
96 electronic health records (EHR) to develop predictive model utilizing machine learning algorithms with
97 the goal of earlier detection and treatment of PCOS.

99 **Materials and Methods**

100 **Data acquisition**

101 The dataset was created by querying de-identified patient data from female patients aged 18 to 45
102 years who had or were considered at risk for PCOS diagnosis by having had any one of the three testing
103 procedures for PCOS in their EHR. Included within the initial sample were those patients who had any
104 visit to Boston Medical Center (BMC) for primary care, obstetrics and gynecology, endocrinology, family
105 medicine, or general internal medicine and received: 1) a pelvic/transvaginal ultrasound for any reason, 2)
106 androgen lab assessment, or had clinical symptoms of androgen excess, 3) an ICD-9 label for irregular
107 periods, or 4) a PCOS diagnosis, between October 2003 to December 2016 within the BMC Clinical Data
108 Warehouse (CDW). The start-date was selected to reflect the first day that ICD-9 codes were used and
109 recorded at BMC. The end date reflected cessation of use of the ICD-9 codes and transition to ICD-10
110 codes within BMC. To avoid misidentifying an ovulation disorder caused by another endocrinopathy,
111 exclusion criteria included diagnosis of concurrent endocrinopathy, such as thyroid disorders,
112 hyperaldosteronism, Cushing's syndrome, other adrenal gland disorders, or malignancy based on ICD-9
113 codes as listed in Supplementary Table 1.

115 **Ethical approval**

116 The study was approved by the Institutional Review Board of Boston University School of
117 Medicine and the Harvard T.H. Chan School of Public Health (Protocol # H35708) and is considered
118 non-human subjects research.

119

120 **Reference label definitions**

121 **Individual predictors**

122 Time-varying predictor variables with a date stamp before that of the outcome of interest were
123 included in our models. We considered the following predictor variables:

124 *Socioeconomic and lifestyle demographic variables:* age, race (White/Caucasian, Black/African
125 American, Hispanic/Latina, Asian, Native Hawaiian/Pacific Islander, Middle Eastern, Other/Unknown),
126 smoking status (yes/no), marital status (single, married, separated, divorced, widowed, other),
127 homelessness (yes/no), and highest level of education (8th grade or less, some high school, high school
128 graduate, some college/technical/vocational training, graduated college/technical school/vocational
129 training, declined to answer, other).

130 *Anthropometrics:* Body mass index (BMI, kg/m²) was either calculated from height and weight or
131 abstracted as the listed BMI variable associated with each visit. BMI was then categorized into three
132 categories: normal (BMI < 25 kg/m²); overweight (BMI between 25-30 kg/m²); and obese (BMI > 30
133 kg/m²). To further capture the obesity population in the absence of height/weight/BMI data, the obese
134 category also included any patient with an ICD-9 code for unspecified obesity (278.00), morbid obesity
135 (278.01), localized adiposity (278.1), and/or a history of gastric bypass.

136 *Cardiovascular health:* To include blood pressure as a predictor variable, we defined a
137 categorical hypertension variable by using systolic (SBP) and diastolic (DBP) blood pressure readings
138 and ICD-9 diagnostic codes for unspecified essential hypertension (401.9), benign essential hypertension
139 (401.1), and essential primary hypertension (401.0). Blood pressure was categorized into three groups:
140 normal, defined by no ICD-9 codes for hypertension recorded and SBP < 120 mmHg, and DBP < 80
141 mmHg; elevated, defined by no ICD-9 codes for hypertension recorded and SBP was 120-129 mmHg or
142 DBP < 80 mmHg; hypertension, defined by any ICD-9 code for hypertension recorded or SBP ≥ 140
143 mmHg or DBP ≥ 90 mmHg.

144 *Reproductive endocrine predictive variables:* beta human chorionic gonadotropin (bHCG) level
145 (negative bHCG < 5 mIU/mL, positive bHCG ≥ 5 mIU/mL), HIV status (negative/positive), age at
146 menarche, pelvic inflammatory disease diagnosis (614.9), history of hysterosalpingogram, and gravidity
147 (history of present or prior pregnancy within obstetric history). Endocrine and metabolic lab values
148 included: TSH, glycosylated hemoglobin (A1c) as a marker for diabetes, low-density lipoprotein (LDL),
149 high density lipoprotein (HDL), and diagnosis of hypercholesterolemia (272.0). Of note, our model did
150 not include androgen precursors such as DHEA or androstenedione as, according to Monash guidelines,
151 these values provide limited additional information in the diagnosis of PCOS (14,15).

152 **Combined predictors**

153 Expecting a nonlinear relationship between many reproductive hormones and a PCOS diagnosis,
154 we used a multilayer perceptron (MLP) neural network to map follicle-stimulating hormone (FSH),
155 luteinizing hormone (LH), sex hormone binding globulin (SHBG), and estradiol (E2) values to a
156 composite metric we call MLP score. The MLP score was repetitively trained and the hyperparameters
157 were tuned to generate a predictive probability associated with PCOS diagnosis for each predictive
158 model, as described with further detail below.

159 **Outcomes**

160 *Defining PCOS:* PCOS diagnosis was assigned for any patient who had an ICD-9 code for PCOS
161 (256.4) or met the Rotterdam criteria (16), according to which a positive diagnosis is made in the
162 presence of two out of the following three features: (i) irregular menses (IM) as defined by rare menses,
163 oligo-ovulation, or anovulation; (ii) hyperandrogenism (HA) as defined by clinical or biochemical
164 androgen excess; and (iii) polycystic ovarian morphology (PCOM) noted on transabdominal or
165 transvaginal ultrasound. Based on these three criteria, we defined three auxiliary variables IM, HA, and
166 PCOM to use in the definition of our labels. PCOM was captured through diagnostic radiology text
167 reports from ovarian ultrasound imaging for the subset that had ultrasound imaging (17).

170 *Defining Irregular Menstruation (IM)*: IM was defined with the following ICD-9 codes: absence
171 of menstruation (626.0), scanty or infrequent menstruation (626.1), irregular menstrual cycle (626.4),
172 unspecified disorders of menstruation and abnormal bleeding from female genital tract (626.9), and
173 infertility, female associated with anovulation (628.0) (3).

174 *Defining Hyperandrogenism (HA)*: HA was assigned to a patient if any of the androgen lab
175 testing for bioavailable testosterone, free testosterone, or total testosterone was greater than clinical
176 thresholds of 11 ng/dL, 5 pg/mL, 45 ng/dL, respectively. In addition, HA was assigned if ICD-9 codes
177 for hirsutism (704.1) or acne (706.1 or 706.0) were recorded for a patient.

178 *Defining Ultrasound characteristics for polycystic ovarian morphology (PCOM)*: Among those
179 with an ultrasound in this dataset, PCOM was identified on ultrasound reports using natural language
180 processing (NLP) with complete methods detailed by Cheng and Mahalingaiah (17), to report PCOM as
181 identified (PCOM present), unidentified (PCOM absent), or indeterminate (PCOM unidentifiable based
182 on source report data).

183 We considered four models to predict the following: **Model I**: patients with ICD-9 diagnosis of
184 PCOS (256.4) within the EHR; **Model II**: patients diagnosed with PCOS by Rotterdam criteria having IM
185 and HA without a specific ICD-9 PCOS code; **Model III**: patients diagnosed with PCOS by Rotterdam
186 criteria having two out of the three conditions IM/HA/PCOM and without a specific ICD-9 PCOS code;
187 **Model IV**: all patients with PCOS using either Model I or Model III criteria. ICD-9 codes were abstracted
188 from the billing code and diagnosis code associated with each encounter within the EHR. Model I
189 included all patients who were diagnosed with PCOS. Model II and its superset Model III was composed
190 of patients who did not have a PCOS diagnosis code but met diagnostic criteria of PCOS based on
191 Rotterdam criteria, representing the patient population with undiagnosed PCOS. Model IV essentially
192 captures all women who were diagnosed or met criteria for PCOS within our population. Supplementary
193 Table 2 details model definitions and includes the count and percent of patients in each category. The date
194 of diagnosis was assigned by the date of PCOS ICD-9 code (256.4) for Model I, the date of the latest
195 diagnostic criteria met for Model II and III, and the earlier date associated with Model I and Model III, for
196 Model IV.

197

198 **Predictive models**

199 **Classification methods**

200 We explored a variety of supervised classification methods, both linear and nonlinear. Linear
201 methods included logistic regression (LR) and support vector machines (SVM) (18) and were fitted with
202 an additional regularization term: an L1-norm of the coefficient vector to inject robustness (19) and
203 induce sparsity. Regularization added a penalty to the objective function, thereby minimizing the sum of a
204 metric capturing fitness to the data and a penalty term that is equal to some multiple of a norm of the
205 model parameters. Sparsity was motivated by the earlier works (20–23), where it was shown that sparse
206 classifiers can perform almost as well as very sophisticated classification methods. Nonlinear methods,
207 including gradient boosted trees (GBT/XGBoost) (24) and random forests (RF) (25) which produce large
208 ensembles of decision trees, may yield better classification performance, but are not interpretable or
209 explainable to enable a safety check by a clinician. Specifically, the RF is a large collection of decision
210 trees and it classifies by averaging the decisions of these trees. The GBT/XGBoost, also called gradient
211 boosting machine (GBM), similarly combines decisions by many decision trees. We used LightGBM
212 which is a fast, high-performance GBM framework (26). We tuned GBM's hyperparameters through
213 cross-validation.

214

215 **Performance metrics**

216 To assess model performance, we obtained the Receiver Operating Characteristic (ROC) curve.
217 The ROC is created by plotting the true positive rate, which is indicative of sensitivity or recall, against
218 the false positive rate (equal to one minus specificity) at various thresholds. The c-statistic or the area
219 under the ROC curve (AUC), is used to evaluate the prediction performance. A perfect predictor is

220 defined by generating an AUC score of 1, and a predictor which makes random guesses has an AUC
221 score of 0.5. We also used the weighted-F1 score to evaluate the models. The weighted-F1 score is the
222 average of the F1 scores of each class weighted by the number of participants in each class. The class-
223 specific F1 scores are computed as the harmonic mean of precision and recall of a classifier which
224 predicts the label of the given class. The weighted-F1 score is between 0 to 1, and a higher value
225 represents a better model. The AUC is more easily interpretable, and the weighted F1-score is more
226 robust to class imbalance (27).

227

228 **Statistical feature selection (SFS)**

229 Categorical variables were converted into dummy/indicator variables. To avoid collinearity, we
230 dropped the missing or unclassified data (NaN) category. For continuous variables, missing values were
231 imputed by the median value for that variable. A summary of the missing variables for each model is
232 provided in Supplementary Table 3. Variables with very low variability ($SD < 0.0001$) were assessed for
233 removal from the models, however none were noted in any model. We applied statistical feature selection
234 (SFS) to reduce the less informative features and simplify the models. For each of the four models’
235 outcomes, the chi-squared test was applied for binary variables and the Kolmogorov-Smirnov statistic for
236 continuous variables; the variables for which we could not reject the null hypothesis of the same
237 distribution for each class ($p\text{-value} > 0.01$) were removed. Representative aggregated patient-level
238 statistics for each model are shown in Supplementary Table 4. We also removed one from each pair of
239 highly correlated variables (with absolute value of the correlation coefficient > 0.8) to avoid redundant
240 variables. Highly correlated variables and the retained variable are provided in Supplementary Table 5.
241 For all models we standardized the corresponding features by subtracting the mean and scaling to unit
242 variance.

243

244 **Training-test splitting**

245 We split the dataset into five random parts, where four parts were used as the training set, and the
246 remaining part was used for testing. We used the training set to tune the model hyperparameters via 5-fold
247 cross-validation, and we evaluated the performance metrics on the testing set. We repeated training and
248 testing five times, each time with a different random split into training/test sets. The mean and standard
249 deviation of the metrics on the test sets over the five repetitions are reported.

250

251 **Development of the MLP score**

252 For every model, there was a considerable difference between the AUC of linear models and non-
253 linear models. To improve the performance of our linear models, we utilized nonlinear models to capture
254 intricate relationships between features. We utilized Gradient Boosted Trees (GBT) to find which features
255 most commonly appeared together among decision trees. We found FSH, LH, SHBG, and estradiol levels
256 to be a meaningful group of features which are all reproductive hormones and continuous variables that
257 appeared together among trees for all our models. We subsequently used these four features as input
258 features into a multilayer perceptron (MLP) neural network model with three hidden layers, each
259 employing the rectified linear unit (ReLU) activation function. The neural network was trained using the
260 training set to classify PCOS. We used the output probability of the MLP model, which we called “MLP
261 score,” as a new feature into our original predictive models.

262

263 **Recursive feature elimination (RFE)**

264 We also used a recursive feature elimination approach with L1-penalized logistic regression (L1-
265 regularized RFE) to extract the most informative features and develop parsimonious models. Specifically,
266 after running the L1-penalized logistic regression (L1-LR), we obtained weights associated with the
267 variables (i.e., the coefficients of the model, denoted by β), and we eliminated the variable with the
268 smallest absolute weight in each turn. We iterated in this fashion, eliminating one variable at a time, to
269 select a model that maximizes a metric equal to the mean AUC minus the standard deviation (SD) of the

270 AUC in a validation dataset (using 5-fold cross-validation on the training set to obtain an average of this
271 metric over five repetitions).

272

273 **Final predictive models**

274 We computed the performance of the following models: L1-penalized logistic regression (LR-
275 L1), support vector machine (SVM-L1), random forests (RF), and gradient boosted trees
276 (GBT/XGBoost). We calculated each variable's LR coefficient with a 95% confidence interval (β
277 [95% CI]), the correlation of the variable with the outcome (Y-correlation), the p-value of each variable
278 (p-value), the mean of the variable (Y1-mean) in the PCOS labeled patients, the mean of the variable (Y0-
279 mean) in the patients without the PCOS label, and the mean and standard deviation of the variable over all
280 patients (All-mean and All-SD). Ranking predictor variables by the absolute value of their coefficients in
281 the logistic regression model amounts to ranking these variables by how much they affect the predicted
282 probability of the outcome. A positive coefficient implies that the larger the value of the variable within
283 the range specified by the data, the higher the chance of having a PCOS diagnosis as defined by the model
284 outcome.

285

286 **Results**

287 **Results of data acquisition and data pre-processing**

288 After inclusion and exclusion criteria were applied to all 65,431 women within the initial data
289 pool, 30,601 patient records were available for this analysis and defined populations are included in
290 Figure 1. There were 1,329 patients (4.5%) with a PCOS ICD-9 diagnosis code (Model I). 1,465 patients
291 had records with PCOM results as present, absent, or unidentifiable. There were 1,056 patients (3.6%)
292 with undiagnosed PCOS (Model II), and a total of 1,116 (3.8%) of patients with no ICD 256.4 indication
293 and two out of IM/HA/PCOM positive criteria (Model III). Finally, there were 2,445 PCOS patients
294 (8.0%) in the combined analysis (Model IV). The total number of records in each model are included in
295 Supplementary Table 2. In the total cohort, the patients were predominantly Black/African American
296 (40.3%) and White (26.5%), with an average age of 33.6 years (SD = 6.6). Complete demographic
297 characteristics are described in Table 1.

298 There were 43 categorical variables and 12 continuous variables retained as predictors after the
299 data pre-processing procedures. There were four pairs of highly correlated variables and one variable
300 from each correlated pair included in the final model as noted in Supplemental Table 5. Supplementary
301 Table 4 describes all 51 variables used by the predictive models.

302

303 **Model Performance**

304 Tables 2, 3, 4 and 5 display the parsimonious models that use the MLP score (LR-L2-MLP score)
305 and show the most significant variables in the prediction of the outcome for Models I, II, III, and IV,
306 respectively. All p-values were less than 0.05, which was set as the significance level.

307 For Model I, the parsimonious predictive model achieved an AUC (SD) of 82.3% (1.7). The MLP
308 score ($\beta = 0.71$) and obesity ($\beta = 0.45$) were positively correlated with PCOS diagnosis. Pregnancy
309 (gravidity $\beta = -0.53$; positive pregnancy test $\beta = -0.50$), normal BMI ($\beta = -0.24$), smoking ($\beta = -0.18$), age
310 ($\beta = -0.16$), and Hispanic race ($\beta = -0.10$) were inversely correlated with PCOS diagnosis as shown in
311 Table 2.

312 For Model II, the parsimonious predictive model achieved an AUC (SD) of 77.6% (1.3). The
313 MLP score ($\beta = 0.61$), obesity ($\beta = 0.21$), normal BMI ($\beta = 0.15$), normal blood pressure ($\beta = 0.16$),
314 negative pregnancy test ($\beta = 0.12$), and normal HDL ($\beta = 0.08$) were positively correlated with
315 undiagnosed PCOS. Age ($\beta = -0.27$), pregnancy (gravidity $\beta = -0.26$; positive pregnancy test $\beta = -0.19$),
316 and Hispanic race ($\beta = -0.18$) were inversely correlated with undiagnosed PCOS as show in Table 3.

317 For Model III, the parsimonious predictive model achieved an AUC (SD) of 77.4% (1.6). The
318 MLP score ($\beta = 0.60$), obesity ($\beta = 0.19$), normal blood pressure ($\beta = 0.17$), normal BMI ($\beta = 0.14$), Black
319 race (0.13), negative pregnancy test ($\beta = 0.12$), and normal HDL ($\beta = 0.09$) were positively correlated
320 with undiagnosed PCOS. Age ($\beta = -0.25$), pregnancy (gravidity $\beta = -0.24$; positive pregnancy test $\beta = -$

321 0.20), and Hispanic race ($\beta = -0.15$) were inversely correlated with undiagnosed PCOS as show in Table
322 4.

323 For Model IV, the parsimonious predictive model achieved an AUC (SD) of 79.1% (1.1). The
324 MLP score ($\beta = 0.7$), obesity ($\beta = 0.31$), normal BMI ($\beta = 0.15$), hypertension ($\beta = 0.07$) and some higher
325 degree of education, such as college or vocational/technical school ($\beta = 0.06$) were positively correlated
326 with PCOS diagnosis. Age ($\beta = -0.21$), pregnancy (gravidity $\beta = -0.37$; positive pregnancy test $\beta = -0.34$;
327 negative pregnancy test $\beta = -0.05$), Hispanic race ($\beta = -0.12$), and smoking ($\beta = -0.08$) were inversely
328 correlated with PCOS diagnosis as shown in Table 5.

329 GBT models had the highest performance. Predictions of PCOS in a test set of patients not used
330 during algorithm training achieved 85%, 81%, 80%, and 82% AUC for Models I, II, III, and IV,
331 respectively. We also report the performance with the logistic regression model (LR-L1) after SFS and
332 the performance when using our developed MLP score alongside variables selected via recursive feature
333 elimination (LR-L2-MLP score). Supplementary Table 6 displays features for each model, associated
334 with LR-L1 algorithm after SFS. As we hypothesized, developing models using the MLP score (LR-L2-
335 MLP score) leads to improvement of the performance of linear models (LR-L1) for Models I, II, III, and
336 IV, respectively from 79%, 72%, 73%, and 75% AUC to 82%, 78%, 77%, and 79% AUC. Table 6 details
337 the models with the best performance (highest AUC) using all 51 features before and after statistical
338 feature selection (SFS). In Table 6, the means and standard deviations of AUC and weighted-F1 scores on
339 the test set over the five repetitions are listed. Supplementary Table 7 displays the performance of all
340 models and all algorithms, before and after statistical feature selection (SFS).

341

342 **Discussion**

343 Evaluating an at-risk population for PCOS is essential for early diagnosis and initiating multi-
344 disciplinary care with the goal of reducing health risks (endometrial hyperplasia/cancer), infertility and
345 pregnancy complications, and chronic disease burden including cardiometabolic disorders associated with
346 PCOS. Retrospective analysis of the at-risk population within an urban health center allows for
347 assessment of factors predictive of diagnosis. Of note, the study sample represents a population of
348 patients who had any visit to BMC for primary care, obstetrics and gynecology, endocrinology, family
349 medicine, or general internal medicine and does not represent a random sample. While this is not a
350 population level assessment, our model is applicable to patients with high suspicion for PCOS who
351 interact with the healthcare system.

352 The ranked list of variables, from the most predictive to the least predictive of the PCOS
353 outcome, informed the main drivers of the predictive models. For example, non-gravidity, high levels of
354 LH, low levels of FSH, obesity, and higher BMI increase the likelihood of PCOS. These variables are
355 consistent with key variables from other models and in the pathophysiology of PCOS. The overall
356 predictive accuracy was high for all models, suggesting that a predictive model may assist in early
357 detection of PCOS within those at risk in an electronically interfaced medical record. Furthermore, we
358 found that non-linear models had superior predictive capacity compared to linear models for all four
359 model outcomes, potentially allowing for inclusion of non-linear reproductive hormone relationships.

360 When assessing patients who received a diagnosis of PCOS (Model I), the most predictive factors
361 related to diagnosis were hormone levels (as captured by the MLP score) and obesity, a clinical factor in
362 supporting a PCOS diagnosis. Specifically, there is a non-linear relationship between reproductive
363 hormones such as FSH, LH, and estradiol. Often these hormonal lab tests are obtained randomly in those
364 with oligomenorrhea, and it is also common to find an elevated FSH to LH ratio. A concern may also be
365 the misclassification of hypothalamic amenorrhea into the group classified as PCOS where the FSH and
366 LH levels would be low or suppressed, or in the setting of premature ovarian insufficiency, notable by an
367 elevated FSH and low estradiol. The MLP score allows for the diversity of relationships of these hormone
368 levels and was trained using a neural network to appropriately classify PCOS. Additionally, prior
369 pregnancy (gravidity) and a positive pregnancy test were negatively associated with a diagnosis of PCOS,
370 consistent with the underlying increased risk of infertility due to oligo-ovulation. Normal BMI and
371 smoking, a known ovarian toxicant, were negatively associated with the presence of a PCOS diagnosis,

372 which may indicate patient characteristics that increase risk of a delayed PCOS diagnosis. These
373 identified variables demonstrate the robustness of the model towards predicting phenotypic traits of
374 patients with PCOS, which is aligned with the performance accuracy. While the significant factors such
375 as hormone levels, gravidity, bHCG, and obesity identified in the model are already known to be
376 associated with PCOS, the true impact of our model lies within the implementation of such a tool within
377 the EHR. For example, a real-world application of this model in the clinical setting would entail
378 integration of our model into the electronic health record system that would provide the probability of
379 PCOS diagnosis or set a threshold for suspicion for each patient to aid a provider's evaluation. This
380 would lead to more timely diagnosis and optimize referrals for downstream follow-up for known clinical
381 sequelae associated with PCOS.

382 When assessing patients who met diagnostic criteria without the ICD-9 label of PCOS (Models II
383 and III), predictive factors both supported the underlying PCOS diagnosis and alluded towards factors
384 that may contribute to missing the diagnosis despite meeting Rotterdam criteria. Similar to Model I,
385 gravidity and a positive pregnancy test were negatively associated with Models II and III diagnosis, while
386 obesity was positively associated with Models II and III diagnosis, consistent with Model I. Interestingly,
387 distinct positive predictors among Models II and III were normal BMI, normal blood pressure, and
388 normal HDL. These patients may present as the "lean" phenotype of PCOS or those with mild features,
389 leading to underdiagnosis of PCOS. Diagnosing "lean" PCOS can be more nuanced, potentially delaying
390 diagnosis or requiring more specialized consultation (28). Within our cohort, 1,116 individuals were
391 identified by the model without the ICD-9 code that met Rotterdam PCOS diagnostic criteria (Model III),
392 suggesting the predictive value of our models to identify at risk groups within a large health system and
393 reduce delays in diagnosis. Given that women often wait over two years and see numerous health
394 professionals before receiving a diagnosis of PCOS, the integration of high-quality AI-based diagnostic
395 tools with the EHR could significantly contribute to more timely diagnosis (8).

396 Consistent with Models I, II, and III, positive pregnancy test and gravidity were both negatively
397 associated with PCOS diagnosis in Model IV while obesity and presence of hypertension were both
398 positively associated with the Model IV combined PCOS outcome. Some higher degree of education,
399 such as college or vocational/technical school, was also positively associated with the outcomes of
400 undiagnosed PCOS and combined PCOS (Models II, III, and IV), which may suggest that education
401 status and patient's self-advocacy for seeking care within a medical system may be implicated specifically
402 in under-diagnosed individuals. Of note, we dropped insurance status after finding that the null was a
403 strong predictor of PCOS, though it is interesting to note that 83% of 331 patients in this dataset with
404 missing insurance have PCOS. Insurance status alludes to socioeconomic barriers such as access to care,
405 which can result in a delay in timely diagnosis through either inability to seek evaluation or follow
406 through with testing. While the implications of insurance status and social determinants of health are
407 beyond the scope of this paper, it is important to note that persistence in seeking treatment within a
408 fractionated health care system can be challenging financially and psychologically, as patients may need
409 multiple evaluation or specialist's consultation to reach the right diagnosis.

410 A recent systematic review investigated the utility of artificial intelligence and machine learning
411 in the diagnosis or classification of PCOS (29). Their search ultimately included 31 studies with sample
412 sizes ranging from 9 to 2,000 patients with PCOS. Methods employed by these models included support
413 vector machine, K-nearest neighbor, regression models, random forest, and neural networks. Only 19% of
414 included studies performed all major steps of training, testing, and validating their model. Furthermore,
415 only 32% of included studies used standardized diagnostic criteria such as the Rotterdam criteria or NIH
416 criteria. The authors found that the ROC of included studies ranged from 73-100%. Only one study
417 sourced their data from electronic health records to build their model (30). Despite the lack of
418 standardized model training and diagnostic criteria used in these studies, the review concluded that
419 artificial intelligence and machine learning provide promise in detecting PCOS, allowing for an avenue
420 for early diagnosis.

421 Outside of the machine learning models included in the systematic review, other predictive
422 models have been created for earlier detection of PCOS as well as for predicting long-term health

423 outcomes among women with a diagnosis of PCOS. One such model was created from 11,720 ovarian
424 stimulation cycles at Peking University Third Hospital. The model used serum antimullerian hormone
425 (AMH) and androstenedione levels, BMI, and menstrual cycle length to predict a diagnosis of PCOS. The
426 algorithm was then developed into an online platform that is able to calculate one's risk of PCOS given
427 certain indicators that are inputted into the model, allowing for better screening abilities in the clinic (31).
428 Another study created a similar model, taking into account AMH and BMI to predict a diagnosis of PCOS
429 or other ovulatory dysfunction disorders among 2,322 women (10). They found that in women with
430 higher BMIs and lower AMH levels could be used to predict PCOS compared to normal-weight or
431 underweight women. Deshmukh et al. created a simple four-variable model which included free androgen
432 index (FAI), 17-hydroxyprogesterone, AMH, and waist circumference for predicting risk of PCOS in a
433 cross-sectional study involving 111 women with PCOS and 67 women without PCOS (32). Lastly, Joo et
434 al. used polygenic and phenotypic risk scores to develop a PCOS risk prediction algorithm (33). They
435 found high degrees of association between PCOS and various metabolic and endocrine disorders
436 including obesity, type 2 diabetes, hypercholesterolemia, disorders of lipid metabolism, hypertension, and
437 sleep apnea (33).

438 In addition to the goal of improved screening for PCOS, models have been created to predict
439 long-term clinical outcomes in women with PCOS, such as ovulation, conception, and live birth (11,12).
440 Given the increased risk of insulin resistance in women with PCOS, Gennarelli et al. created a
441 mathematical model to predict insulin sensitivity based on variables such as BMI, waist and hip
442 circumferences, truncal-abdominal skin folds, and serum concentrations of androgens, SHBG,
443 triglycerides, and cholesterol (13). Models to predict non-alcoholic fatty liver disease risk among young
444 adults with PCOS have also been generated (34). Combining earlier detection with more accurate risk
445 stratification of clinical sequelae through predictive modeling can significantly improve the long-term
446 health outcomes of women with PCOS. Application of our models to predict other downstream health
447 risks after the diagnosis of PCOS is a future area of research.

448 Beyond the long-term health impacts of PCOS, the condition also carries a significant economic
449 cost for our healthcare system. A study by Riestenberg et al (2022) recently estimated the total economic
450 burden of PCOS, as well as the cost specifically for pregnancy-related complications and long-term health
451 morbidities (2). The authors estimated the annual economic burden of PCOS to be \$8 billion as of 2020 in
452 the United States. Furthermore, the excess cost of pregnancy-related comorbidities such as gestational
453 hypertension, gestational diabetes, and preeclampsia attributable to PCOS totals \$375 million USD
454 annually. Outside of pregnancy, the cost of long-term comorbidities associated with PCOS including
455 stroke and type 2 diabetes mellitus was estimated at \$3.9 billion USD. Meanwhile, the cost for diagnostic
456 evaluation of PCOS was less than 2% of the total economic burden. This estimated financial burden
457 suggests that predictive models aiding earlier diagnosis could not only reduce long-term health
458 consequences of PCOS but also alleviate significant healthcare costs associated with the condition.

459 Given the high prevalence, significant healthcare burden, and heterogeneity in clinical
460 presentation of PCOS, AI-based tools are well suited for earlier diagnosis of PCOS. Our study had many
461 strengths. First, our machine learning models, which were highly accurate and robust in PCOS diagnosis
462 prediction, were created using the largest sample size to date (29). Second, our model was tested and
463 trained on a diverse Safety-Net hospital-sourced population not restricted to the context of fertility care.
464 Third, it is the only model that incorporated three data streams (ICD-9 codes, clinical laboratory findings,
465 and radiologic findings) and an MLP score. Fourth, the parsimonious and interpretable models were very
466 close in achieving full model predictive accuracy, performing relatively closely to the best-performing
467 non-linear models. Essentially, our parsimonious models "isolate" nonlinearities in hormone levels
468 (captured by the MLP score) and linearly combine that score with other variables. Most models evaluate
469 reproductive hormones (FSH, estradiol, LH, and SHBG) as individual variables within linear models,
470 which does not account for the high inter- and intra-patient variability. By using non-linear mapping of
471 the hormone values, we were able to generate a composite variable allowing for a linear function that
472 correlates with the likelihood of an accurate prediction. Last, our variables are easily accessible in an
473 electronic health dataset, rendering the models helpful for clinical prediction. Our study did not evaluate

474 AMH as a predictive variable because it was not widely utilized during the time window of this data
475 extraction corresponding with ICD-9 codes.

476 Despite these strengths, our model is not without limitations. First, it is only directly applicable to
477 those who interact with the medical system and those deemed “at-risk” for a PCOS diagnosis, which
478 would not facilitate population-based prediction. Additional studies need to be conducted in other patient
479 populations or unselected community-based populations to validate the use of these models, especially
480 expanding to the entire population within a health system to evaluate the accuracy of our models (35).
481 Second, we must interpret our data within the limitations of informative presence in EHR data.
482 Informative presence is defined as data that is present and informed with respect to the health outcome, in
483 this case PCOS, as well as behavioral patterns of interaction with healthcare institutions which may be
484 additionally impacted by marginalization (36). This is an important consideration for interpreting
485 predictive models using EHR data (36,37). Nevertheless, we were able to extract over 1000 patients who
486 were undiagnosed with PCOS among the population, suggesting the predictive value of the modelling in
487 identifying diagnosis gaps among specific populations within a large health system. Third, it is possible
488 that additional examination of the medical record beyond ICD-9 diagnosis may allow for more
489 clarification of risk in the presumed PCOS group. Last, our exclusion of concurrent endocrinopathies was
490 chosen to avoid incorrectly including ovulation disorders caused by other endocrinopathies, but it is
491 possible that this was an overly strict exclusion criterion.

492 In conclusion, this novel machine learning algorithm incorporates three data streams from a large
493 EHR dataset to assess PCOS risk. This model can be integrated into the EHR to aid clinicians in earlier
494 diagnosis of PCOS and connect patients to interventions and healthcare providers across their
495 reproductive lifespan with the goal of health optimization and risk reduction.

496 **Acknowledgements**

497 We would like to acknowledge Linda Rosen, the research manager of the Clinical Data Warehouse at the
498 Boston Medical Center for procuring the dataset, and Alexis Veiga, the research assistant who provided
499 administrative support for this project.

500

501 **Author contributions**

502 ZZ performed the analysis and co-wrote the manuscript with a focus on the methods. VJ interpreted the
503 findings and drafted the initial manuscript. AW conducted a literature review, contributed to
504 interpretation of data, writing, and editing the manuscript. TW initiated the analytical approach to the
505 research question. JC curated the initial dataset and reviewed the analysis and manuscript drafts. SM and
506 ICP designed the study, oversaw analysis, interpretation of the findings, and manuscript drafting and
507 revision process. All authors met ICJME criteria for authorship.

508

509 **Data availability**

510 All datasets generated during and/or analyzed during the current study are not publicly available but are
511 available from the corresponding author on reasonable request.

512

513 **Competing interests**

514 The authors declare no competing interests.

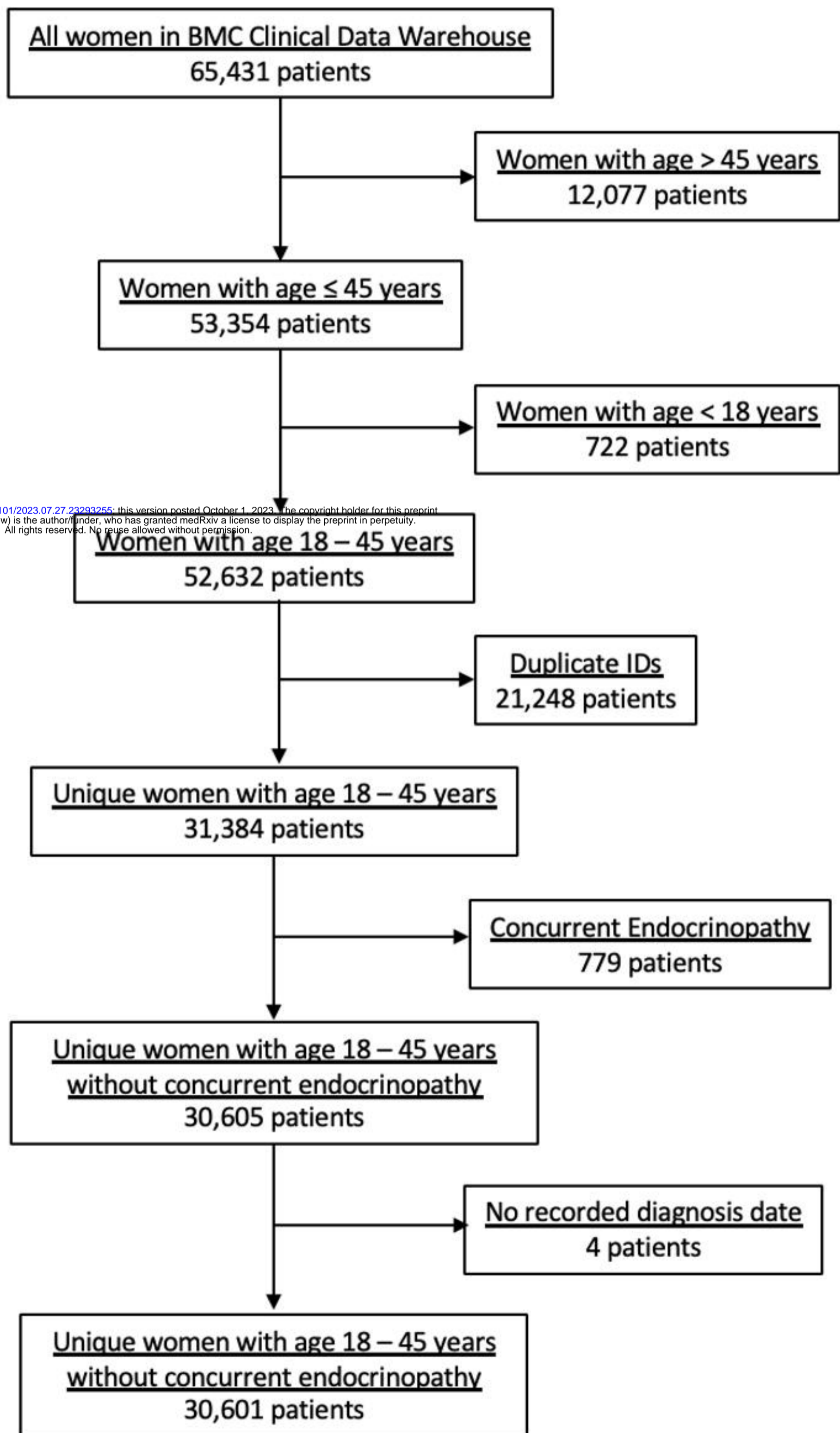
515 **References**

- 516 1. Azziz R, Carmina E, Dewailly D, Diamanti-Kandarakis E, Escobar-Morreale HF, Futterweit W, et al.
517 The Androgen Excess and PCOS Society criteria for the polycystic ovary syndrome: the complete task
518 force report. *Fertil Steril*. 2009 Feb;91(2):456–88.
- 519 2. Riestenberg C, Jagasia A, Markovic D, Buyalos RP, Azziz R. Health Care-Related Economic Burden of
520 Polycystic Ovary Syndrome in the United States: Pregnancy-Related and Long-Term Health
521 Consequences. *J Clin Endocrinol Metab*. 2022 Jan 18;107(2):575–85.
- 522 3. Sirmans SM, Pate KA. Epidemiology, diagnosis, and management of polycystic ovary syndrome. *Clin*
523 *Epidemiol*. 2013 Dec 18;6:1–13.
- 524 4. Barry JA, Azizia MM, Hardiman PJ. Risk of endometrial, ovarian and breast cancer in women with
525 polycystic ovary syndrome: a systematic review and meta-analysis. *Hum Reprod Update*.
526 2014;20(5):748–58.
- 527 5. Lim SS, Kakoly NS, Tan JWJ, Fitzgerald G, Bahri Khomami M, Joham AE, et al. Metabolic syndrome
528 in polycystic ovary syndrome: a systematic review, meta-analysis and meta-regression. *Obes Rev*.
529 2019;20(2):339–52.
- 530 6. Anagnostis P, Tarlatzis BC, Kauffman RP. Polycystic ovarian syndrome (PCOS): Long-term metabolic
531 consequences. *Metabolism*. 2018;86:33–43.
- 532 7. Wekker V, Van Dammen L, Koning A, Heida KY, Painter RC, Limpens J, et al. Long-term
533 cardiometabolic disease risk in women with PCOS: a systematic review and meta-analysis. *Hum*
534 *Reprod Update*. 2020;26(6):942–60.
- 535 8. Gibson-Helm M, Teede H, Dunaif A, Dokras A. Delayed Diagnosis and a Lack of Information
536 Associated With Dissatisfaction in Women With Polycystic Ovary Syndrome. *J Clin Endocrinol*
537 *Metab*. 2017 Feb 1;102(2):604–12.
- 538 9. Xu H, Feng G, Alpadi K, Han Y, Yang R, Chen L, et al. A Model for Predicting Polycystic Ovary
539 Syndrome Using Serum AMH, Menstrual Cycle Length, Body Mass Index and Serum
540 Androstenedione in Chinese Reproductive Aged Population: A Retrospective Cohort Study. *Front*
541 *Endocrinol*. 2022;13:821368.
- 542 10. Vagios S, James KE, Sacha CR, Hsu JY, Dimitriadis I, Bormann CL, et al. A patient-specific model
543 combining antimüllerian hormone and body mass index as a predictor of polycystic ovary syndrome
544 and other oligo-anovulation disorders. *Fertil Steril*. 2021;115(1):229–37.
- 545 11. Kuang H, Jin S, Hansen KR, Diamond MP, Coutifaris C, Casson P, et al. Identification and replication
546 of prediction models for ovulation, pregnancy and live birth in infertile women with polycystic ovary
547 syndrome. *Hum Reprod*. 2015;30(9):2222–33.
- 548 12. Jiang X, Liu R, Liao T, He Y, Li C, Guo P, et al. A Predictive Model of Live Birth Based on Obesity
549 and Metabolic Parameters in Patients With PCOS Undergoing Frozen-Thawed Embryo Transfer. *Front*
550 *Endocrinol*. 2021;12.
- 551 13. Gennarelli G, Holte J, Berglund L, Berne C, Massobrio M, Lithell H. Prediction models for insulin
552 resistance in the polycystic ovary syndrome. *Hum Reprod*. 2000;15(10):2098–102.

- 553 14. Villarroel C, López P, Merino PM, Iñiguez G, Sir-Petermann T, Codner E. Hirsutism and
554 oligomenorrhea are appropriate screening criteria for polycystic ovary syndrome in adolescents.
555 *Gynecol Endocrinol*. 2015 Aug 3;31(8):625–9.
- 556 15. Monash University. International evidencebased guideline for the assessment and management of
557 polycystic ovary syndrome. 2018.
- 558 16. ESHRE TR, Group ASPCW. Revised 2003 consensus on diagnostic criteria and long-term health risks
559 related to polycystic ovary syndrome. *Fertil Steril*. 2004;81(1):19–25.
- 560 17. Cheng JJ, Mahalingaiah S. Data mining polycystic ovary morphology in electronic medical record
561 ultrasound reports. *Fertil Res Pract*. 2019;5(1):1–7.
- 562 18. Hastie T, Tibshirani R, Friedman J. The elements of statistical learning: data mining, inference and
563 prediction. Vol. 1. Springer series in statistics Springer, Berlin; 2001.
- 564 19. Chen R, Paschalidis IC. Distributionally Robust Learning. *Found Trends® Optim*. 2020 Dec 22;4(1–
565 2):1–243.
- 566 20. Brisimi TS, Xu T, Wang T, Dai W, Adams WG, Paschalidis IC. Predicting Chronic Disease
567 Hospitalizations from Electronic Health Records: An Interpretable Classification Approach. *Proc*
568 *IEEE*. 2018 Apr;106(4):690–707.
- 569 21. Brisimi TS, Xu T, Wang T, Dai W, Paschalidis IC. Predicting diabetes-related hospitalizations based
570 on electronic health records. *Stat Methods Med Res*. 2019 Dec 1;28(12):3667–82.
- 571 22. Chen R, Paschalidis IC. Robust Grouped Variable Selection Using Distributionally Robust
572 Optimization. *J Optim Theory Appl*. 2022 Sep 1;194(3):1042–71.
- 573 23. Chen R, Paschalidis IC, Hatabu H, Valtchinov VI, Siegelman J. Detection of unwarranted CT radiation
574 exposure from patient and imaging protocol meta-data using regularized regression. *Eur J Radiol Open*.
575 2019;6:206–11.
- 576 24. Chen T, Guestrin C. XGBoost: A Scalable Tree Boosting System. In: Proceedings of the 22nd ACM
577 SIGKDD International Conference on Knowledge Discovery and Data Mining [Internet]. San
578 Francisco, California, USA: Association for Computing Machinery; 2016 [cited 2020 Jun 19]. p. 785–
579 94. (KDD '16). Available from: <https://doi.org/10.1145/2939672.2939785>
- 580 25. Breiman L. Random forests. *Mach Learn*. 2001;45(1):5–32.
- 581 26. Ke G, Meng Q, Finley T, Wang T, Chen W, Ma W, et al. Lightgbm: A highly efficient gradient boosting
582 decision tree. *Adv Neural Inf Process Syst*. 2017;30.
- 583 27. Saito T, Rehmsmeier M. The Precision-Recall Plot Is More Informative than the ROC Plot When
584 Evaluating Binary Classifiers on Imbalanced Datasets. Brock G, editor. *PLOS ONE*. 2015 Mar
585 4;10(3):e0118432.
- 586 28. Toosy S, Sodi R, Pappachan JM. Lean polycystic ovary syndrome (PCOS): an evidence-based practical
587 approach. *J Diabetes Metab Disord*. 2018 Nov 13;17(2):277–85.

- 588 29. Barrera FJ, Brown EDL, Rojo A, Obeso J, Plata H, Lincango EP, et al. Application of machine learning
589 and artificial intelligence in the diagnosis and classification of polycystic ovarian syndrome: a
590 systematic review. *Front Endocrinol* [Internet]. 2023 [cited 2023 Sep 21];14. Available from:
591 <https://www.frontiersin.org/articles/10.3389/fendo.2023.1106625>
- 592 30. Castro V, Shen Y, Yu S, Finan S, Pau CT, Gainer V, et al. Identification of subjects with polycystic
593 ovary syndrome using electronic health records. *Reprod Biol Endocrinol RBE*. 2015 Oct 29;13:116.
- 594 31. Xu H, Feng G, Alpadi K, Han Y, Yang R, Chen L, et al. A Model for Predicting Polycystic Ovary
595 Syndrome Using Serum AMH, Menstrual Cycle Length, Body Mass Index and Serum
596 Androstenedione in Chinese Reproductive Aged Population: A Retrospective Cohort Study. *Front*
597 *Endocrinol*. 2022;13.
- 598 32. Deshmukh H, Papageorgiou M, Kilpatrick ES, Atkin SL, Sathyapalan T. Development of a novel risk
599 prediction and risk stratification score for polycystic ovary syndrome. *Clin Endocrinol (Oxf)*.
600 2019;90(1):162–9.
- 601 33. Joo YY, Actkins K, Pacheco JA, Basile AO, Carroll R, Crosslin DR, et al. A polygenic and phenotypic
602 risk prediction for polycystic ovary syndrome evaluated by phenome-wide association studies. *J Clin*
603 *Endocrinol Metab*. 2020;105(6):1918–36.
- 604 34. Carreau AM, Pyle L, Garcia-Reyes Y, Rahat H, Vigers T, Jensen T, et al. Clinical prediction score of
605 nonalcoholic fatty liver disease in adolescent girls with polycystic ovary syndrome (PCOS-HS index).
606 *Clin Endocrinol (Oxf)*. 2019;91(4):544–52.
- 607 35. Azziz R, Woods KS, Reyna R, Key TJ, Knochenhauer ES, Yildiz BO. The Prevalence and Features of
608 the Polycystic Ovary Syndrome in an Unselected Population. *J Clin Endocrinol Metab*. 2004 Jun
609 1;89(6):2745–9.
- 610 36. Harton J, Mitra N, Hubbard RA. Informative presence bias in analyses of electronic health records-
611 derived data: a cautionary note. *J Am Med Inform Assoc JAMIA*. 2022 Jun 14;29(7):1191–9.
- 612 37. Sisk R, Lin L, Sperrin M, Barrett JK, Tom B, Diaz-Ordaz K, et al. Informative presence and observation
613 in routine health data: A review of methodology for clinical risk prediction. *J Am Med Inform Assoc*
614 *JAMIA*. 2020 Nov 9;28(1):155–66.
- 615

Figure 1. Flow of patients from the BMC CDW into the dataset used by the study.



medRxiv preprint doi: <https://doi.org/10.1101/2023.07.27.23293255>; this version posted October 1, 2023. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted medRxiv a license to display the preprint in perpetuity. All rights reserved. No reuse allowed without permission.

Table 1. Demographic characteristics of the study population and by model.

Variable	Model I	Model II	Model III	Model IV
Age, Mean years (SD)	33.6 (6.6)	33.7 (6.6)	33.7 (6.6)	33.6 (6.6)
Race, n (%)				
Black/African American	11881 (40.3)	11824 (40.5)	11861 (40.5)	12395 (40.5)
White/Caucasian	7812 (26.5)	7733 (26.5)	7741 (26.4)	8086 (26.4)
Hispanic/Latina	2858 (9.7)	2837 (9.7)	2841 (9.7)	2929 (9.6)
Asian	1350 (4.6)	1354 (4.6)	1354 (4.6)	1406 (4.6)
Middle Eastern	175 (0.6)	176 (0.6)	176 (0.6)	184 (0.6)
American Indian/Native American	163 (0.6)	162 (0.6)	162 (0.6)	168 (0.5)
Native Hawaiian/Pacific Islander	17 (0.1)	18 (0.1)	18 (0.1)	18 (0.1)
Other	979 (3.3)	966 (3.3)	966 (3.3)	1023 (3.3)
Unknown	4250 (14.41)	4146 (14.19)	4153 (14.19)	4392 (14.4)
Marital Status				
Single	22325 (75.7)	22155 (75.8)	22199 (75.8)	23224 (75.9)
Married	5833 (19.8)	5753 (19.7)	5767 (19.7)	6018 (19.7)
Separated	392 (1.3)	391 (1.3)	392 (1.3)	401 (1.3)
Divorced	388 (1.3)	379 (1.3)	380 (1.3)	397 (1.3)
Widowed	35 (0.1)	35 (0.1)	35 (0.1)	35 (0.1)
Other	502 (1.7)	489 (1.7)	489 (1.7)	516 (1.7)
Unknown	10 (0.03)	10 (0.03)	10 (0.03)	10 (0.03)
Body Mass Index				
Normal (BMI < 25)	7534 (25.6)	7685 (26.3)	7697 (26.3)	7902 (25.8)
Overweight (BMI between 25-30)	5694 (19.3)	5689 (19.5)	5707 (19.5)	5941 (19.4)
Obese (BMI ≥ 30)	7645 (25.9)	7369 (25.2)	7387 (25.2)	7985 (26.1)
Unknown	8612 (29.2)	8469 (29.0)	8481 (29.0)	8,773 (28.7)

Table 2. Most significant variables for PCOS diagnosis prediction in Model I.

Rank	Variables	β	$\beta - \%95$ CI	Y- correlation	p-value	Y1- mean	Y0- mean	All- mean	All- std
0	MLP Score	0.71	0.028	0.33	6.80E-197	0.17	0.04	0.05	0.08
1	Intercept	-0.68	-	-	-	-	-	-	-
2	Gravidity	-0.53	0.018	-0.12	4.55E-78	1.28	2.08	2.04	1.39

3	Positive bHCG	-0.5	0.019	-0.09	1.50E-48	0.05	0.23	0.22	0.42
4	Obesity	0.45	0.017	0.11	1.38E-81	0.51	0.27	0.28	0.45
5	Normal BMI	-0.24	0.017	-0.05	3.57E-16	0.15	0.26	0.26	0.44
6	Smoker	-0.18	0.017	-0.03	6.62E-05	0.09	0.14	0.14	0.34
7	Age	-0.16	0.016	-0.08	1.70E-25	31.34	33.79	33.68	6.61
8	Hispanic/Latina Race	-0.1	0.016	-0.02	1.82E-03	0.07	0.10	0.10	0.30

Table 3. Most significant variables for PCOS diagnosis prediction in Model II.

Rank	Variables	β	β - %95 CI	Y-correlation	p-value	Y1-mean	Y0-mean	All-mean	All-std
0	MLP Score	0.61	0.023	0.26	2.13E-142	0.12	0.04	0.04	0.06
1	Intercept	-0.44	-	-	-	-	-	-	-
2	Age	-0.27	0.015	-0.08	2.26E-31	31.01	33.79	33.69	6.61
3	Gravidity	-0.26	0.016	-0.09	2.35E-63	1.42	2.08	2.06	1.39
4	Obesity	0.21	0.016	0.03	9.60E-06	0.34	0.27	0.27	0.44
5	Positive bHCG	-0.19	0.017	-0.06	4.14E-21	0.10	0.23	0.23	0.42
6	Hispanic/Latina Race	-0.18	0.016	-0.02	2.69E-03	0.06	0.10	0.10	0.30
7	Normal BP	0.16	0.015	0.03	1.37E-07	0.60	0.51	0.51	0.50
8	Normal BMI	0.15	0.016	0.03	8.57E-07	0.34	0.26	0.26	0.44
9	Negative bHCG	0.12	0.015	0.06	1.44E-22	0.37	0.23	0.23	0.42
10	HDL	0.08	0.015	0.01	1.03E-10	52.13	51.59	51.61	7.86

Table 4. Most significant variables for PCOS diagnosis prediction in Model III.

Rank	Variables	β	β - %95 CI	Y-correlation	p-value	Y1-mean	Y0-mean	All-mean	All-std
0	MLP Score	0.6	0.023	0.26	7.41E-142	0.10	0.04	0.04	0.05
1	Intercept	-	-	-	-	-	-	-	-
2	Age	-0.25	0.015	-0.08	5.91E-30	31.16	33.79	33.69	6.61
3	Gravidity	-0.24	0.016	-0.09	2.47E-63	1.46	2.08	2.06	1.39
4	Positive bHCG	-0.20	0.017	-0.06	3.59E-20	0.11	0.23	0.23	0.42
5	Obesity	0.19	0.016	0.03	2.73E-06	0.34	0.27	0.27	0.44
6	Normal BP	0.17	0.015	0.04	3.94E-08	0.60	0.51	0.51	0.50

7	Hispanic/Latina Race	-0.15	0.016	-0.02	2.00E-03	0.06	0.10	0.10	0.30
8	Normal BMI	0.14	0.016	0.03	6.76E-06	0.33	0.26	0.26	0.44
9	Black/African American Race	0.13	0.015	0.02	2.03E-03	0.46	0.40	0.41	0.49
10	Negative bHCG	0.12	0.015	0.06	2.20E-25	0.37	0.23	0.23	0.42
11	HDL	0.09	0.015	0.01	4.06E-12	52.04	51.59	51.61	7.86

Table 5. Most significant variables for PCOS diagnosis prediction in Model IV.

Rank	Variables	β	β - %95 CI	Y-correlation	p-value	Y1-mean	Y0-mean	All-mean	All-std
0	MLP Score	0.7	0.024	0.36	0.00E-01	0.20	0.07	0.08	0.10
1	Intercept	-0.44	-	-	-	-	-	-	-
2	Gravidity	-0.37	0.017	-0.14	2.17E-135	1.36	2.08	2.02	1.39
3	Positive bHCG	-0.34	0.017	-0.10	2.23E-65	0.08	0.23	0.22	0.41
4	Obesity	0.31	0.015	0.10	2.86E-66	0.43	0.27	0.28	0.45
5	Age	-0.21	0.015	-0.10	1.91E-52	31.26	33.79	33.59	6.62
6	Hispanic/Latina Race	-0.12	0.015	-0.03	2.34E-06	0.07	0.10	0.10	0.29
7	Smoker	-0.08	0.015	-0.02	3.00E-04	0.11	0.14	0.14	0.34
8	Hypertension	0.07	0.015	0.04	3.63E-12	0.28	0.21	0.22	0.41
9	Education – Some College/Technical/ Vocational School	0.06	0.014	0.03	1.55E-04	0.18	0.15	0.15	0.36
10	Negative bHCG	-0.05	0.015	0.05	2.29E-16	0.31	0.23	0.24	0.42

Table 6. Model performance over the test set, in the format of mean percentage (SD percentage) over 5 repetitions.

	Model I		Model II		Model III		Model IV	
	AUC	F1-weighted	AUC	F1-weighted	AUC	F1-weighted	AUC	F1-weighted
Best full models before SFS	XGBoost (51 features)		XGBoost (51 features)		XGBoost (51 features)		XGBoost (51 features)	
	85.2 (1.8)	94.5 (0.2)	80.6 (0.5)	95.1 (0.2)	80.4 (0.7)	94.8 (0.1)	81.8 (1.4)	91.1 (0.4)
Best full models after SFS	XGBoost (14 features)		XGBoost (16 features)		XGBoost (17 features)		XGBoost (17 features)	
	83.6 (1.7)	94.5 (0.2)	80.5 (0.7)	95.1 (0.2)	79.8 (1.1)	94.8 (0.1)	81.1 (1.3)	90.9 (0.3)

	LR-L1 (14 features)		LR-L1 (16 features)		LR-L1 (17 features)		LR-L1 (17 features)	
	79.2 (1.9)	93.9 (0.2)	71.7 (0.9)	94.7 (0.1)	72.9 (2.1)	94.4 (0.1)	74.8 (1.1)	89.7 (0.3)
	Parsimonious models LR-L2-MLP score (8 features)		Parsimonious models LR-L2-MLP score (10 features)		Parsimonious models LR-L2-MLP score (11 features)		Parsimonious models LR-L2-MLP score (10 features)	
	82.3 (1.7)	94.5 (0.1)	77.6 (1.3)	95.1 (0.1)	77.4 (1.6)	94.9 (0.1)	79.1 (1.1)	90.8 (0.3)