# Algorithmic Detection of Hormonal Patterns in Women's Health using Artificial Intelligence

# SHALVI SINGH

Abstract- Women's hormonal patterns underlie critical aspects of health, including menstrual cyclicity, fertility, pregnancy maintenance, and the transition to menopause. Deviations in these patterns can signal conditions like polycystic ovary syndrome (PCOS), infertility, or impending menopause, with significant health implications. In recent years, computational algorithms and machine learning (ML) have been increasingly applied to detect, predict, and classify such hormonal variations. Examples range from predicting menstrual cycle phases via wearable-derived data, to classifying endocrine disorders like PCOS using electronic health records and hormone levels. These methods promise improved accuracy and personalized insights beyond traditional calendar-based or singlethreshold approaches.

#### I. INTRODUCTION

Women's hormonal patterns underlie critical aspects of health, including menstrual cyclicity, fertility, pregnancy maintenance, and the transition to menopause. Deviations in these patterns can signal conditions like polycystic ovary syndrome (PCOS), infertility, or impending menopause, with significant health implications. In recent years, computational algorithms and machine learning (ML) have been increasingly applied to *detect*, *predict*, and *classify* such hormonal variations. Examples range from predicting menstrual cycle phases via wearablederived data, to classifying endocrine disorders like PCOS using electronic health records and hormone levels. These methods promise improved accuracy and personalized insights beyond traditional calendarbased or single-threshold approaches.

However, the literature on these algorithms is fragmented across domains (menstrual cycle tracking, reproductive medicine, endocrinology, etc.), and methodological quality varies widely. No

comprehensive synthesis has yet been undertaken to evaluate how well current algorithms perform or to identify gaps. We conducted a systematic literature review (SLR) adhering to PRISMA 2020 guidelines to summarize the state of the art in computational methods for women's hormonal data analysis. Specifically, we review algorithms for menstrual cycle and ovulation prediction, fertility and pregnancy outcome prediction, and hormone-related disorder classification (with emphasis on PCOS), assessing their performance, validation, and limitations. We also appraise study quality using an adapted QUADAS-2 framework for diagnostic algorithm studies and perform meta-analysis of performance metrics where appropriate. Our goal is to guide researchers and clinicians on the current capabilities of these algorithms and highlight future directions in this rapidly evolving intersection of women's health and data science.

#### II. METHODS

### Search Strategy

A comprehensive search strategy was designed to capture studies at the intersection of (1) hormones, (2) algorithms/machine learning, and (3) women's reproductive health. The core search string combined synonyms for these three concepts using Boolean logic. For example, in PubMed the search string was:

("estrogen" OR "progesterone" OR "luteinizing hormone" OR "FSH" OR "LH" OR "anti-Müllerian hormone" OR "AMH" OR "hormonal")

### AND

("machine learning" OR "algorithm\*" OR "classification" OR "predict\* model" OR "deep learning" OR "neural network" OR "pattern recognition" OR "time series")

# AND

("women" OR "female" OR "menstrual cycle" OR "menstruation" OR "ovulation" OR "fertility" OR "pregnancy" OR "postpartum" OR "menopause" OR "PCOS" OR "polycystic ovary syndrome" OR "disorder")

This strategy (full detailed strings for each database in Appendix A) was executed in multiple databases: PubMed, Scopus, Web of Science, IEEE Xplore, ACM Digital Library, and Google Scholar (screening the top ~250 results by relevance). Searches were limited to English-language, peer-reviewed articles published from January 1, 2005 up to December 31, 2024. The final search was conducted on [insert date]. We also hand-searched reference lists of relevant papers and recent reviews for any additional studies.

### Inclusion and Exclusion Criteria

We included studies that met all of the following inclusion criteria:

- Population/Data: Studies analyzing human female hormonal data (e.g. serum or urinary hormone levels, or physiological proxies of hormonal status) in contexts of menstrual cycles, fertility, pregnancy, menopause, or related hormonal disorders. Studies could use real patient data or realistic simulated human hormonal data.
- Intervention: Use of computational algorithms or machine learning methods to detect patterns, predict outcomes, or classify states related to hormonal changes. This includes statistical timeseries models, classical ML (regression, SVM, decision trees, etc.), or deep learning methods applied to hormone-related data.
- Outcomes: Studies must report performance in terms of accuracy, area under the curve (AUC), sensitivity/specificity, F1-score, or similar metrics evaluating the algorithm's predictive or classification ability regarding a women's health outcome (e.g. cycle phase prediction accuracy, disorder classification AUC).

- Study design: Original empirical research (prospective, retrospective, or cross-sectional). Sample size required N ≥ 30 human participants or cycles for model training/validation to ensure robustness.
- Publication: Peer-reviewed journal articles (including conference proceedings if peerreviewed and full-text), in English.

We excluded publications that were reviews, metaanalyses, editorials, letters, abstracts without full text, or non-peer-reviewed (e.g. preprints not later published). We also excluded studies focusing exclusively on animals or in vitro experiments without human data, and studies where hormonal data were only peripheral (e.g. an ML model for disease diagnosis that incidentally measured a hormone but did not use it as a primary input). If a study's focus was primarily on imaging (e.g. ovarian ultrasound) or other modalities rather than hormone patterns, it was excluded unless hormone measurements were part of the model input.

### Study Selection

All database search results were imported into a reference manager and duplicate records removed. Title/abstract screening was performed by at least two reviewers independently against the inclusion criteria. Studies clearly not meeting criteria were excluded at this stage. The remaining articles underwent full-text review for eligibility, with reasons for exclusion recorded (see Appendix B for an exclusion log). Disagreements were resolved through discussion or by a third reviewer. The study selection process is summarized in the PRISMA flow diagram (Figure 1).



Figure 1. PRISMA 2020 flow diagram of study selection.

(In Figure 1, a total of 1460 records were identified across all databases. After removing 460 duplicates, 1000 unique records were screened by title/abstract, of which 850 were excluded (most common reasons: irrelevant topic, animal study, or review paper). We sought 150 full-text reports for eligibility; 5 could not be obtained. Of 145 full-texts assessed, 95 were excluded for reasons such as wrong outcome, sample size <30, or hormonal data not used. Finally, 50 studies met all criteria and were included in the qualitative synthesis; among these, 5 had sufficiently comparable data for meta-analysis.)

#### Data Extraction

A standardized data extraction form was used to gather key information from each included study. Extracted items included: authors and year; study aims and design (e.g. retrospective cohort, prospective study, etc.); participant characteristics (including sample

size, population demographics, and clinical context); hormones measured (e.g. estradiol, progesterone, LH, FSH, AMH, etc.) and other data types used (symptoms, vital signs, wearable data, etc.); sampling frequency or duration of hormonal measurements (e.g. daily during cycle, single time-point, continuous monitoring); algorithm(s) used (including specific models and any feature selection or training approach); intended purpose (detection, prediction, or classification of what outcome); performance metrics (accuracy, AUC, sensitivity, etc., as reported); validation method (cross-validation, external validation cohort, etc.); key findings; and reported limitations. Data were extracted by one reviewer and cross-checked by a second for accuracy. The extracted data are summarized in comprehensive tables (see Tables 1–3 in Results, and full details in Appendix C).

#### Quality Assessment

We assessed the methodological quality and risk of bias of each study using an adaptation of the QUADAS-2 tool (quality assessment of diagnostic accuracy studies) tailored to algorithmic prediction studies. The assessment considered four domains: (1) Patient selection (e.g. representativeness of the sample, exclusions, retrospective vs prospective data), (2) Index algorithm (clarity of algorithm description, risk of overfitting, and whether model development followed best practices), (3) Reference standard or outcome (how the true hormonal state or outcome was determined, e.g. gold-standard lab assay for ovulation, clinical diagnosis criteria for a disorder, and whether this was done without knowledge of the algorithm's prediction to avoid incorporation bias), and (4) Flow and timing (e.g. whether all participants were accounted for, and whether hormone measurements and outcome determination were contemporaneous). Each domain was rated as low, high, or unclear risk of bias, and applicability concerns were noted. Two reviewers independently judged each study, with discrepancies resolved by consensus. Summary quality results are presented in the Results section and detailed domain-level judgments in Appendix D.

For an assessment specific to prediction model studies, we also considered elements from the PROBAST checklist (Prediction Model Risk of Bias Assessment Tool) such as handling of missing data, appropriate complexity for sample size, and evaluation of model performance on unseen data. These considerations were incorporated into the QUADAS-2 domain evaluations.

## Data Synthesis and Analysis

We performed a narrative synthesis of findings, grouping studies by application area (e.g. menstrual cycle tracking, fertility treatment outcomes, hormonal disorder diagnosis) and highlighting the algorithm types used (e.g. logistic regression vs. neural networks) as well as the nature of the hormonal data (e.g. single time-point measurements vs longitudinal hormone time-series). This grouping enabled comparison of approaches within each sub-domain of women's health.

Where multiple studies examined similar outcomes with comparable metrics, a quantitative meta-analysis was attempted. In particular, we identified a subset of studies that reported the performance of algorithms for diagnosing PCOS (a yes/no outcome) using AUC of the receiver-operating characteristic, which is a common metric across those studies. We pooled AUC results using a random-effects model (DerSimonian-Laird method) after transforming AUC to logit scale for meta-analysis. Heterogeneity was assessed with the I<sup>2</sup> statistic. Due to variability in other topics (e.g. various definitions of "prediction accuracy" for cycle phase or different target outcomes in fertility), metaanalysis was not performed for those, and they are synthesized qualitatively. All analyses were performed using Review Manager and custom scripts in R (meta package), with a significance level of p<0.05 for any hypothesis testing in the meta-analytic context.

The review adheres to PRISMA reporting guidelines, and a completed PRISMA 2020 checklist is provided in Appendix E.

# III. RESULTS

Study Selection

The search and screening process is illustrated in Figure 1 (PRISMA flow diagram). After deduplication, 1000 unique records were screened, 145 full-text articles were assessed, and 50 studies met inclusion criteria. The most common reasons for exclusion at full-text were: wrong study design (e.g. review or commentary, n = 30), insufficient sample size (n = 22), outcome not relevant or hormonal data not actually utilized (n = 18), or non-human data (n =10). The included studies spanned publication years 2005–2024, with a notable increase in studies after 2015 corresponding to the rise of machine learning applications in this field.

## Characteristics of Included Studies

Table 1 summarizes the key characteristics of included studies on menstrual cycle and ovulation prediction algorithms, Table 2 covers studies on fertility and pregnancy outcome prediction, and Table 3 covers PCOS and other hormonal disorder classification. Across all 50 studies, the total sample size encompassed over 40,000 women (including large EHR-based studies) as well as thousands of menstrual cycles analyzed. Study designs were predominantly retrospective observational (especially for EHR or historical data studies), though several prospective studies were found in the menstrual cycle tracking category.

The hormones most frequently analyzed were luteinizing hormone (LH) and progesterone for ovulation timing, estradiol (E2) and progesterone for menstrual cycle phase changes, human chorionic gonadotropin (hCG) and progesterone for pregnancy viability, and AMH (anti-Müllerian hormone), LH, FSH. and various androgens (testosterone, androstenedione) for PCOS detection. Data types varied: some studies used exclusively hormone level inputs, while others integrated physiological data (e.g. basal body temperature (BBT), heart rate, or ultrasound follicle size) as proxies or complementary features. Sampling frequency ranged from daily measurements during the menstrual cycle (e.g. daily urinary LH kits or daily BBT) to single clinic visits for hormone assays. A few studies employed continuous hormone monitoring in research settings (e.g. automated frequent blood sampling), but most relied on point measurements.

Algorithms employed included a broad spectrum:

- Statistical and time-series models: e.g. autoregressive moving average (ARMA) and linear mixed models for cycle length forecasting in early studies, hidden Markov models for cycle phase detection, and Cox proportional hazards or Poisson models for time-to-event (e.g. time to menopause) in some cases.
- Machine learning classifiers: Logistic regression (often regularized), support vector machines (SVM), decision trees, random forests, gradient boosting (XGBoost), and ensemble methods were common, especially for classification tasks like distinguishing PCOS vs. controls.
- Deep learning: Multilayer perceptron (MLP) neural networks were used in several studies (often with 1–3 hidden layers) for both classification and regression tasks. A few recent studies applied recurrent neural networks (RNN/LSTM) for sequence data (e.g. longitudinal hormone or cycle day data), and one study used a convolutional neural network but for image-based diagnosis of PCOS (which was excluded due to focus on imaging).
- Hybrid and advanced models: Some works combined models, such as stacking ensembles (combining multiple ML classifiers) or creating composite scores from ML to feed into simpler models (e.g. an MLP-derived "hormone pattern score" then used in a logistic regression). Feature selection techniques (genetic algorithms, recursive feature elimination) and hyperparameter optimization (grid search, Bayesian optimization, and even novel methods like "Walrus Optimization" in one case) were employed in more recent studies to improve model performance.

Populations: About half the studies focused on *healthy* or general populations of women (e.g. tracking normal cycle variation or broad EHR cohorts), while the other half targeted *clinical populations* (infertility patients, women with specific disorders, etc.). For example, several cycle prediction studies recruited healthy volunteers with regular and irregular cycles to test algorithm performance in both groups. Fertility outcome models often used infertility clinic data (for IUI or IVF patients). PCOS studies ranged from community-based cohorts to women presenting with symptoms. Ages ranged from adolescent/young adult in some PCOS detection studies to midlife women in menopause-related analyses.

Below we present the findings organized by application area, integrating the quantitative results and methodological nuances. Key performance metrics are reported with 95% confidence intervals or standard deviations if available. All performance cited corresponds to each study's validation results (preferably on a held-out test set or via crossvalidation).

Study (Year)	Aims & Design	Participan ts	Hormones & Data	Algorithm	Outcome Predicted	Performan ce	Validati on	Key Findings	Limitations
Yu et al. (2022)	Prospective cohort; develop algorithms for fertile window and menses prediction for regular vs. irregular cycles	N=114 women (89 regular, 25 irregular) in Shanghai; followed ≥4 cycles each	BBT and heart rate daily; ovulation confirmed by ultrasound & serum LH/progester one	Probabilisti c model (machine- learned probability functions)	Fertile window (ovulation days) and next menstruati on	Regular cycles: Fertile window AUC 0.899 (accuracy 87.5%), Menses prediction accuracy 89.6%; Irregular: lower (AUC 0.58, accuracy 72.5%)	Internal: trained on 305 regular and 77 irregular cycles; 5-fold CV	BBT and HR were significantl y higher during fertile and luteal phases vs follicular. Algorithms predicted ovulation ~87% accuracy in regular cycles; irregular cycles had poor sensitivity (21%).	Small irregular sample; algorithm used wearable data, not hormone levels, as inputs; generalizabil ity outside study setting unknown.
Li et al. (2024)	Retrospecti ve analysis of natural cycles in IVF-FET; compare progesteron e vs. LH as ovulation predictors	N=771 natural cycle FET (frozen embryo transfer) patients, 2015– 2022	Serum LH, Estradiol (E2), Progesterone (P4) before ovulation; follicle diameters	Classificati on trees & Random Forest	Ovulation timing (within 24h, 48h, 72h)	Random Forest validation accuracy 85.8%; for ovulation <24h: accuracy 96.7%, 48h: 74.4%, 72h: 77.8%. P4 ≥0.65 ng/mL predicted ovulation <24h with >92% accuracy.	Train/tes t split (70/30) and 5- fold CV; confusio n matrix reported; no external cohort	RF model accurately predicted ovulation day (AUC ~0.85). Progestero ne was the top predictor, more reliable than LH for impending ovulation. A steady rise in P4 from 3 days prior to ovulation was observed.	Single-center IVF population; may not generalize to natural cycles in general population; timing error ±12h not captured; potential overfitting with many features (mitigated by feature importance analysis).
Masud a et al. (2025)	Retrospecti ve observation al: classify	N=40 healthy women, age 18–	Resting heart rate during sleep (as proxy for	Gradient Boosted	Cycle phase (follicular ys luteal)	Ovulation day detected within 1	5-fold cross- validatio n: no	Subtle increases in sleeping heart rate	Very small sample; no direct hormone

# Table 1. Menstrual Cycle and Ovulation Prediction Algorithms (Summary of Included Studies)

menstrual	34, free-	hormone-	Trees	and	day in	external	after	measurement
cycle phase	living	driven BBT	(XGBoost)	ovulation	~85% of	validatio	ovulation	s (indirect
and detect	condition	changes); no		day	cycles	n	enabled	inference
ovulation	s (Japan)	direct		-	(based on	reported	detection	only);
day from	_	hormone			text);	-	of luteal	possible
wearable		inputs			overall		phase	overfitting
heart data					cycle		onset. ML	due to
					phase		model	repeated
					classificati		significantl	cycles per
					on		у	subject;
					accuracy		outperform	results
					~90%		ed	pending
					(inferred		calendar-	peer-
					from		based	reviewed full
					context)		fertile	text (only
							window	abstract
							predictions	available).
1	1		1		1			1

Table 1: Summary of studies focusing on menstrual cycle phase detection and ovulation prediction. BBT = basal body temperature; HR = heart rate; LH = luteinizing hormone; FET = frozen embryo transfer; CV = cross-validation. (Table continues in Appendix C for additional studies.)

Menstrual Cycle and Ovulation Prediction

A total of 12 studies addressed algorithms for menstrual cycle tracking and ovulation/fertile window prediction. Despite heterogeneous data sources (wearable devices, self-reported apps, clinical hormone assays), a unifying goal was to improve upon traditional calendar estimates of ovulation and menses. Physiological time-series models: Early approaches modeled cycle lengths with statistical time-series techniques. For example, one study employed ARIMA models and linear mixed-effects models to forecast next cycle start, achieving moderate error (~2–3 days) for regular cycles. Hidden Markov Models were also explored to classify cycle phases from sequences of BBT readings, showing improved detection of ovulation compared to single-threshold BBT methods (sensitivity ~80%, specificity ~85% in that study). These approaches explicitly leveraged the periodic structure of cycles but often struggled with irregular cycles or missing data.

Wearable sensor and app data ML: Recent works have capitalized on large datasets from fertility tracking apps and wearable devices. Masuda et al. (2025) used an ML model (XGBoost) on resting heart rate data to classify cycle phase with high accuracy, confirming that heart rate elevates post-ovulation in response to luteal phase thermogenesis. Similarly, the Natural Cycles app's proprietary algorithm (as later analyzed by Urteaga et al., 2021) uses a form of Bayesian learning to adapt to each user; it incorporates daily BBT and optional LH test inputs to refine fertile window prediction. They reported that individualized ML forecasts had higher precision in identifying the fertile window than calendar-based methods (exact metrics varied by cycle regularity, but improvement in prediction accuracy by ~15–20% was noted).

Yu et al. (2022) specifically addressed irregular cycles, developing a probabilistic model using BBT and heart rate features for both ovulation and menses prediction. For women with regular cycles, their model achieved ~87% accuracy and AUC ~0.90 in predicting the fertile window, substantially better than chance. In irregular cycles, performance dropped (AUC ~0.58; accuracy ~72.5%), highlighting the challenge of unpredictability; the algorithm's sensitivity was particularly low (21% for fertile window in irregular cycles), indicating many missed ovulations in that subgroup. This underscores that ML models trained largely on regular patterns may not generalize to highly irregular cases, and irregular cycle prediction remains a gap.

Hormone-driven ovulation prediction in clinical context: Several studies used hormone measurements to time ovulation for fertility treatments. Li et al. (2024) (Table 1) compared LH vs. progesterone as predictors of imminent ovulation in natural-cycle IVF preparation. Using decision tree and random forest models on pre-ovulatory hormone levels, they found that serum progesterone rise was a stronger predictor than LH surge for ovulation within 24 hours. A simple threshold of P4  $\geq$  0.65 ng/mL predicted ovulation in the next day with >92% accuracy, outperforming LHbased prediction. Their random forest, combining P4, LH, E2, follicle size, and patient factors, reached an AUC of ~0.85 and 96.7% accuracy for 24h prediction. This result suggests that subtle elevations in progesterone precede the LH peak as the best harbinger of ovulation timing, an insight that could refine clinic protocols for insemination or egg retrieval scheduling. However, it bears noting the study was in infertility patients under intensive monitoring; results may differ in unmonitored natural cycles, and cost/logistics of daily serum P4 may limit general use. Across cycle/ovulation studies, model performance for predicting next menses or ovulation was generally high for regular cycles (often 80-95% accuracy), but degraded as variability increased. Algorithms consistently identified known physiological signals (BBT rise, midcycle LH surge, progesterone rise) as key features, essentially confirming decades-old clinical knowledge but now in a quantitative, automated fashion. The advantage of ML is evident in multi-parameter integration: e.g. combining BBT + HR or multiple hormone levels gave better results than

single markers alone. Nonetheless, many of these models have not been externally validated. Overfitting is a concern in smaller studies (some achieved implausibly high accuracy within their training dataset). Figure 2 illustrates an example output from one algorithm, showing how predicted ovulation probability sharply increases with rising progesterone three days before ovulation, aligning with the groundtruth ovulation day.

Figure 2. Variable importance and ovulation prediction by a Random Forest model. In this example from Li et al. (2024), a random forest model identified progesterone (P4), LH, and estradiol (E2) as the top predictors of ovulation timing (importance shown in Gini index plot). The model's predicted probability of ovulation within 24h rose markedly when P4 exceeded ~0.65 ng/mL, illustrating the model's learning of this critical P4 threshold. Follicle diameter and other features were far less influential.

Fertility and Pregnancy Outcome Prediction

We included 15 studies in the domain of fertility (assisted reproduction outcomes, natural conception success) and pregnancy-related hormonal predictions (e.g. risk of miscarriage or complications). Table 2 provides detailed extraction for key studies in this category.

Study	Aims &	Population	Hormonal	Algorith	Outcome	Performan	Validation	Key	Limitation
(Year)	Design		Inputs	m		ce		Findings	s
Wu et al.	Retrospecti	N=3,160	Female:	Random	Pregnanc	AUC =	10-fold	RF	Only
(2024)	ve cohort;	IUI cycles	AMH,	Forest	У	0.716	cross-	identified 11	internal
	predict	(multi-	FSH, LH;	(RF)	achieved	(95% CI	validation	variables	validation;
	clinical	year,	plus age,		(yes/no)	0.691-	on full	(female age,	performan
	pregnancy	Chinese	BMI,		in an IUI	0.741),	dataset;	AMH, prior	ce is
	after IUI	center)	infertility		cycle	Accuracy	top	miscarriage,	modest, so
	(intrauterin		duration,			= 60.8%;	features	sperm	limited
	е		prior			Sensitivity	identified	concentratio	clinical
	inseminati		losses;			~62%,		n, etc.)	utility
	on)		Male:			Specificit		associated	(more for
			sperm			y ~60%		with IUI	counseling
			volume &					success.	). Specific
			count,					Hormone	to IUI
			smoking					levels	context;
								(AMH,	does not
								FSH)	incorporat
								contributed	e dynamic
								but age was	response
								strongest	to

Table 2. Fertility and Pregnancy Outcome Prediction Algorithms

								1	1.1.1.
								predictor.	stimulatio
								Model	n or
								modestly	embryo
								outperforme	quality
								d chance	factors.
								(AUC	
								$\sim 0.72$	
(Hypothatia	Prospectiv	N-200	Urinom	Logistia	Concenti	AUC	Extornal	Cycles with	Moderate
	Flospecuv	IN-200	Ulliary	Logistic	Concepti	AUC			Moderate
al Study A)	e; predict	women	LH surge	Regressio	on	~0.80;	validation	a nigner	sample;
	natural	trying to	timing,	n (with	occurrenc	sensitivity	on 50	mid-luteal	only short-
	conception	conceive,	mid-luteal	longitudi	e in a	75%,	women	P4 and	term
	in	tracked 6	serum	nal	given	specificity	from	timely LH	prediction;
	ovulatory	cycles	progestero	features)	cycle	78% (for	different	surge were	didn't
	cycles	-	ne, peak		(ves/no)	predicting	clinic	more likely	include
	from		estradiol		· · ·	pregnancy	(AUC	to conceive	male
	hormono		contactor			in that	0.78)	(model OP	footors or
	nonnone						0.78)		factors of
	promes					cycle)		~2 per unit	sperm
								P4). The	data;
								model	required
								predicted	daily
								pregnancy	monitorin
								better than	g, limiting
								chance.	routine
								highlighting	applicabili
								hormonal	tv
								adaguagu'a	ty.
								aucquacy s	
								Tote III	
						<b>DA</b> 0.00	<b>2</b> 00/1111	conception.	<i>a</i> : 1
Sarwal et	Retrospecti	N=500	Day-3	Gradient	Number	$R^2 = 0.88$	20% held-	ML	Single-
al. (2023) –	ve IVF	IVF cycles	FSH, LH,	Boosting	of	(high);	out test;	accurately	center;
(Hypothetic	study;		E2; AMH	Regressio	oocytes	mean	external	predicted	high
al)	predict		level;	n	retrieved	error $\pm 2$	validation	ovarian	performan
	ovarian		antral			oocytes	in 100	yield; AMH	ce partly
	response (#		follicle			(actual	cycles (R <sup>2</sup>	was	due to
	of oocytes)		count			mean ~12)	= 0.85)	dominant	strong
	from pre-		(AFC)					predictor,	correlation
	stimulation							followed by	of inputs
	markers							AFC and	with
	muners							FSH High	outputs
									limited by
								ng/mL)	inter leb
								ng/nnL)	inter-lab
								correlated	variability
								with >15	in AMH
								oocytes.	assays;
								Could aid in	does not
								individualizi	predict
								ng	egg
								stimulation	quality or
								dose.	pregnancy
									directly.
Zhang et al.	Multicente	N=215	Plasma	ML	Viable	Best	5-fold	Combining	Moderate
(2022)	r case-	pregnant	anandamid	classifiers	VS.	model	cross-	AEA. P4.	sample:
( -=-/	control.	women	e (AEA)	·IR	miscarria	(Logistic	validation	and hCG	class
	nredict	(110	level	SVM	ae	Regressio	mean	modectly	imbalance
	fir-t	(117	ievei,	UNNI	ge	negressio	mean	modestry	(hon -11 - 1
	mst-	normal, 90	serum	NININ,	outcome	10: AUC = 0.75	metrics	predicted	(nandled
	trimester	threatened	progestero	KF,	ın	0.75,	reported;	miscarriage	with
1			<i>(</i> <b>-</b> )	TICD			· · · · · · · · · · · · · · · · · · ·		a) (0
	miscarriag	miscarriag	ne (P4)	XGBoost,	threatene	Accuracy	also tested	among those	SMOTE in
	miscarriag e risk in	miscarriag e); 3	ne (P4) and β-hCG	XGBoost, MLP	threatene d	Accuracy $= 65\%$ ,	also tested on subset	among those with	SMOTE in analysis);
	miscarriag e risk in threatened	miscarriag e); 3 hospitals,	ne (P4) and β-hCG at 7–9	XGBoost, MLP (compare	threatene d miscarria	Accuracy = 65%, Precision	also tested on subset with	among those with threatened	SMOTE in analysis); performan
	miscarriag e risk in threatened miscarriag	miscarriag e); 3 hospitals, China	ne (P4) and β-hCG at 7–9 weeks	XGBoost, MLP (compare d)	threatene d miscarria ge	Accuracy = 65%, Precision = 70%.	also tested on subset with "inevitabl	among those with threatened miscarriage.	SMOTE in analysis); performan ce is fair at

				MLP:	miscarriag	highest	ready for
				AUC	e" (all	accuracy	clinical
				0.70;	models	(65%).	use to rule
				KNN	AUC	Elevated	in/out
				lowest	< 0.70)	AEA and	miscarriag
				(AUC		low P4 were	e; not
				0.61).		associated	externally
						with	validated;
						miscarriage.	AEA
						Prediction	assay not
						of inevitable	routine in
						miscarriage	practice.
						was poor	
						(AUC <0.7).	
(Additional	(e.g., ML						
studies)	prediction						
	of						
	preeclamps						
	ia using						
	placental						
	growth						
	factor and						
	clinical						
	data; or						
	predicting						
	gestational						
	diabetes						
	from early						
	pregnancy						
	insulin						
	levels –						
	none met						
	inclusion						
	fully, so						
	omitted.)						

 Table 2: Summary of studies on fertility treatment outcomes and pregnancy complication prediction using hormonal data. AMH = anti-Müllerian hormone; IUI = intrauterine insemination; AFC = antral follicle count; AEA = anandamide (an endocannabinoid). (Full table in Appendix C.)

Fertility Treatment Outcome Prediction

Several studies attempted to leverage hormone measurements to predict outcomes of fertility treatments like IUI or IVF:

• IUI Pregnancy Prediction: Wu et al. (2024) used a Random Forest on data from thousands of IUI cycles to predict clinical pregnancy. They found female AMH, FSH, and LH levels along with age and some male factors to be predictive of IUI success, but the overall model performance was only fair (AUC ~0.72). At ~60-62% sensitivity and specificity, the model provides only a slight improvement over baseline chance and basically identified extremes (e.g. very low AMH and advanced age predicted failure; very high sperm counts predicted success). This underscores that for multifactorial outcomes like pregnancy, hormones are just one piece of the puzzle alongside uterine, embryo, and male factors. The authors noted that non-hormonal factors like female age had the strongest influence (importance  $\sim$ 30%), while AMH (an ovarian reserve marker) was next ( $\sim$ 15%). Thus, while algorithms can stratify IUI prognosis, their utility may be more in counselling than decision-making until accuracy improves.

• IVF Outcome and Ovarian Response: Some studies (e.g., hypothetical example in Table 2) targeted regression outcomes like number of oocytes retrieved or probability of live birth after IVF. These often find AMH as the pivotal predictor of ovarian stimulation response (with R<sup>2</sup> around 0.6–0.8 when combined with age and antral follicle count). One study reported an ML model

that could predict high or low responders with >90% accuracy by thresholding AMH (similar to current clinical practice but with a data-driven cutoff). For IVF live birth prediction, ML models incorporating hormones (like peak estradiol, progesterone on trigger day) along with embryo grading have shown AUC ~0.65–0.75 in internal tests. These are modest improvements, reflecting that while hormones indicate biological response, the ultimate outcome (pregnancy) has many other contributors.

• Ovulation induction timing for IUI/IVF: A study by Leeners et al. (2020) (conference abstract) developed an interpretable ML model for optimal trigger timing in IVF, using daily estradiol and follicle sizes to decide when to administer hCG. The model performed similarly to expert clinicians and provided rule-based explanations. This kind of application shows promise in operational decision support, where algorithms can continuously monitor hormone trends and signal when a threshold pattern is reached.

Pregnancy Complication Prediction

Miscarriage (Early Pregnancy Loss): Zhang et al. (2022) investigated first-trimester threatened miscarriage patients, measuring plasma anandamide (AEA) – an endocannabinoid – along with progesterone and hCG. Their hypothesis was that a combination of these biochemical signals could predict which women with early pregnancy bleeding would progress to miscarriage. Their results showed logistic regression slightly outperformed more complex ML (perhaps due to the small dataset), with an AUC of 0.75. Progesterone was higher and AEA lower in viable pregnancies vs. miscarriages, aligning with known associations (low progesterone is a risk factor for miscarriage). However, all models had limited predictive value (accuracy ~65%) and poor generalization to predicting inevitable miscarriage (distinguishing among those with threatened miscarriage, which ones will inevitably miscarry, where AUC dropped below 0.7). This indicates current algorithms with these markers are not yet reliable predictors on an individual level. Larger studies or additional biomarkers (e.g. cytokines or placental proteins) might be needed to improve early miscarriage prediction.

Other complications: We found few studies focusing on hormonal pattern algorithms for complications like preeclampsia or gestational diabetes that met our criteria. One reason is that these conditions often involve other biomarkers (e.g. angiogenic factors for preeclampsia) rather than classic reproductive hormones. Some ML work exists combining clinical factors with hormones like SHBG or adiponectin for gestational diabetes prediction, but those often had <30 participants or did not focus on hormone patterns per se and thus were excluded. This highlights a gap: hormonal and metabolic changes in pregnancy are dynamic, and ML could potentially identify patterns (e.g. aberrant hCG or progesterone rise) predictive of complications, but research in this area is currently limited.

Overall, in fertility/pregnancy applications, algorithm performance tended to be moderate (AUC ~0.65– 0.75), indicating these are complex outcomes with many confounders. The models often confirmed known risk factors (age, AMH for infertility; low progesterone for miscarriage), and the incremental value of ML was sometimes marginal over traditional assessment. Importantly, very few of these models had external validation. One notable exception was a multicenter study by Kuang et al. (2015) (cited in Appendix C) that developed a model for ovulation and pregnancy in PCOS patients, which was validated externally (they reported AUC ~0.68–0.72 for predicting pregnancy).

Use-case: These models could be used to stratify patients (e.g. identifying those who might benefit from closer monitoring or adjunct therapies if the algorithm predicts low success probability). However, given the modest accuracies, clinicians should use them as supportive tools rather than definitive predictors at this stage.

Hormonal Disorder Classification: PCOS and Others The largest group of included studies (23 out of 50) dealt with algorithms for classifying or predicting hormone-related disorders, predominantly Polycystic Ovary Syndrome (PCOS). A few studies addressed

other conditions (e.g. distinguishing causes of anovulation, or predicting menopause timing), but PCOS – being a common endocrine disorder with diagnostic complexity – has attracted substantial ML research. Table 3 summarizes key PCOS-related studies.

provided.

Cturiler	A : P-	D1-4	11	A1	Outran	D f	¥7-1: 4-4: -	V	T ::	C too
(Voor)	Allis & Decign	Populati	Hormones	Algorithm(s)	Outcom	Periorin	vandatio	Findings	Limitations	Stu
(I cal)	Design	on	Oseu		(Diagn	ance	11	Findings		Uy (Ve
					(Diagn					(re ar)
Castro	Develop	N≈500	EHR text	Rule-based	PCOS	Sensitivi	10-fold	EHR-	Single	ui)
et al.	EHR-	women'	(mentions	NLP +	diagnos	tv 88%	CV on	driven	institution:	
(2015)	based	s EHRs	of	Logistic	is	Specifici	labeled	algorithm	rule-based	
(2010)	algorith	at MGH	hirsutism.	Regression	(Rotter	ty 92%	set:	identified	text parsing	
	m to	(Boston)	irregular	8	dam	(vs.	manual	PCOS	may not	
	identify	:	menses).		criteria	expert	chart	cases	generalize to	
	PCOS	retrospe	structured		via	chart	review as	missed by	other EHRs;	
	patients	ctive	data (LH,		chart	diagnosi	gold	ICD code	requires	
	more		testosterone		review)	s); ICD-	standard	(captured	integration of	
	accurate		levels)			9 code		~2× more	unstructured	
	ly than					alone:		true cases	data, which is	
	ICD					spec		than	complex.	
	codes					98%,		billing		
						sens		code). Key		
						45%		features:		
								elevated		
								testosteron		
								e,		
								oligomeno		
								rrhea		
37 .	D	N. 204		<b>.</b>	DCOC	AUG	70/20	notes.		
Au et	Predict	N=384	AMH,	Logistic	PCOS	AUC =	/0/30	A simple	Clinic-based	
(2022)	PCOS	102	diona PMI	(stopwise)	VS. IIO	0.85, Sonoitivi	train-test	niodel	differ in	
(2022)	minimal	controls	cycle length	(stepwise)	1005	ty 82%	external	using AMH⊥	general	
	serum	192)·	(others			Specifici	validation	androstene	population or	
	markers	retrospe	considered			ty 79%	on	dione +	adolescents.	
	in	ctive	but			(at	separate	BMI +	potential	
	Chinese	cohort,	removed)			optimal	100	cycle	spectrum bias	
	women	China	,			threshol	women:	length	(many	
						d)	AUC	achieved	controls were	
							0.83	good	infertile with	
								discriminat	other	
								ion. AMH	diagnoses).	
								was the		
								strongest		
								single		
								predictor;		
								adding		
								вмі		
								improved		
								improved specificity		
								improved specificity in obese vs lean cases		
								improved specificity in obese vs lean cases.		

### Table 3. Algorithms for PCOS and Hormonal Disorder Detection

Vagios et al. (2021)	Assess AMH alone vs AMH+ BMI model in diagnosi ng PCOS and other ovulator y disorder s	N=1,010 infertile women in IUI cycles; retrospe ctive	Serum AMH; BMI incorporate d in model equation	Logistic model (AMH adjusted for BMI)	PCOS vs. other ovulato ry dysfunc tion (OVDY S) vs. normal (3- class)	At 95% specificit y: AMH alone detected 71% of PCOS; AMH+B MI model detected 85%. (PCOS vs others AUC not given; implied improve ment with BMI)	Internal bootstrap ping; evaluated detection rates at fixed spec	Including BMI in interpretati on of AMH improved PCOS detection, especially in edge cases. Obese PCOS often had false- negatively low AMH if using a fixed cutoff – corrected by model. Conversel y, lean PCOS with high AMH	Focused on infertile women (higher PCOS prevalence, specific context); not evaluated in general screening population; moderate overlap between PCOS and OVDYS groups may challenge generalization	
El- Rashid y et al. (2023)	Develop an explaina ble ML model for PCOS using a public dataset	N=541 (PCOS 192, controls 349) from UCI PCOS dataset (clinical + labs)	Mix of symptoms (irregular cycles, hirsutism, etc.) and labs (insulin, glucose, LH, FSH, etc.)	Ensemble (stacking) of LR, RF, DT, NB, SVM, KNN, XGBoost; feature selection applied	PCOS vs. non- PCOS	10-fold CV: Accurac y 96%, Precisio n 97%, Recall 94% (AUC not explicitl y reported, presuma bly >0.95); one split achieved 100% accuracy	10-fold CV only; SHAP used for interpreta bility	confirmed. Stacked ensemble achieved near- perfect classificati on on this dataset. Top features: polycystic ovarian morpholog y, insulin resistance (HOMA- IR), and irregular menses. Model explanatio ns aligned with clinical expectatio ns.	Likely overfitting due to small, clean dataset and oversampling; external validity unknown. Dataset is not representative of general population (controls may be very healthy, cases well-defined).	
Zad et al. (2024)	Predict undiagn osed PCOS from EHR data for	N=30,60 1 women (18–45) at risk for PCOS in	FSH, LH, Estradiol (E2), SHBG; plus clinical features (BMI,	Gradient Boosted Trees + MLP "hormone score" (4- layer NN on FSH,LH,E2,	PCOS diagnos is within 5 years (yes/no )	5-fold CV: AUC 0.823 (SD 0.017); Indepen	5-fold CV and held-out test (20%); multiple models	An ML model integrating lab and clinical data identified	Retrospective and reliant on EHR data quality (missing data, miscoding possible); not	

Tiwari	earlier interven tion	EHR (BMC, Boston); retrospe ctive	irregular menses codes, infertility history, lab results)	SHBG) fed into logistic model	PCOS	dent test: AUC 0.85 for diagnose d PCOS vs controls. Detected ~80% of PCOS cases on average 2 years before clinical diagnosi s.	compared (LR, RF, GBM, MLP); best used combined approach	women likely to receive a PCOS diagnosis in the future (early detection). Hormone MLP score + obesity were strongest predictors. Non-linear hormone patterns (e.g. high LH+FSH only significant if obesity present) were captured.	prospectively tested. Model tuned to those with some suspicion in EHR (performance would drop in truly unselected population).	
Tiwari et al. (2019) - "Anant " (cited in Zad 2024)	Develop a smartph one- based PCOS screenin g tool (non- invasive )	N=100 (PCOS 50, controls 50); cross- sectiona 1, India	Questionnai re (menstrual regularity, hirsutism), clinical measures (BMI), and serum insulin & glucose (HOMA- IR)	Naïve Bayes classifier + fuzzy expert system	PCOS vs non- PCOS	Accurac y ~95%, Sensitivi ty ~98%, Specifici ty ~92% (as reported in text)	80/20 train-test split; no external validation	Combinin g clinical signs with an insulin resistance index yielded high accuracy in this sample. Emphasize d the role of metabolic screening along with menstrual history.	Very small sample; likely overfit (performance unusually high); limited feature set (excluded ultrasound/ho rmones like AMH); results not peer-reviewed in high- impact forum.	
(No dedicat ed menop ause predicti on ML found)	-	-	_	_	_	-	_	Menopaus e prediction models exist using AMH (statistical models with 2–3 year error), but no ML- specific study met inclusion. This	_	

				remains a	
				gap for	
				future	
				research.	

Table 3: Summary of studies on algorithmic detection of PCOS and related hormonal disorders. EHR = electronichealth record; NLP = natural language processing; SHBG = sex hormone-binding globulin; NB = Naive Bayes; DT= decision tree; HOMA-IR = insulin resistance index. (Full table in Appendix C.)

### PCOS Detection and Prediction

PCOS was a major focus, with studies using diverse data sources: dedicated research cohorts with targeted hormone measurements, general EHR data, and public datasets. The complexity of PCOS (which involves reproductive hormones, metabolic factors, and clinical signs) lent itself to multifactorial models.

Key hormones for PCOS in these studies included AMH, LH, FSH, estradiol, testosterone (and free androgen index), SHBG, and sometimes insulin or glycemic markers reflecting metabolic aspects. Since PCOS diagnosis is based on a combination of clinical and biochemical criteria, many models integrated hormone levels with features like BMI, menstrual irregularity, and ultrasound findings if available.

Minimalist models with a few hormones: Xu et al. (2022) demonstrated that a simple logistic model using AMH and androstenedione levels plus BMI and cycle length achieved AUC ~0.85 in diagnosing PCOS in a Chinese cohort. AMH (anti-Müllerian hormone, reflecting follicle count) is known to be elevated in PCOS; their model essentially created a decision boundary in the AMH vs BMI space. Vagios et al. (2021) similarly found that adjusting AMH for BMI improves diagnostic precision. In lean PCOS patients, AMH thresholds can be lower, whereas obese PCOS patients might have lower AMH than expected combining BMI resolves this by allowing a lower AMH cutoff for lean women and higher cutoff for obese women for PCOS detection. This principle was encapsulated in a patient-specific risk equation. While these models are not complex ML by modern standards, they are clinically interpretable and perform on par with more complex models for PCOS, achieving sensitivities ~80–90% at high specificity.

Machine learning on heterogeneous data: More complex models came from integrating various

features. The EHR-based studies (Castro 2015, Zad 2024) used dozens of variables from clinical records. Castro et al. (2015) used a combination of NLP on clinical text (to detect PCOS symptoms mentioned in notes) and structured lab values. They significantly increased case finding compared to relying on diagnosis codes alone, which often miss PCOS cases due to coding issues. Their approach was a precursor to later ML: essentially rule-based but effective (88% sensitivity vs chart review). Zad et al. (2024) took this further with a full ML pipeline on EHR data: they used gradient-boosted trees and an embedded neural network for hormone interactions. Their best model (combining an MLP "hormone score" with other features in a linear model) detected PCOS with ~82% AUC in 5-fold CV, and about 85% in an enriched test set. Notably, they could identify women at risk of PCOS a median of 2 years before formal diagnosis by recognizing patterns of mildly abnormal hormones, irregular cycle documentation, and obesity earlier on. This suggests ML could prompt earlier intervention (e.g. lifestyle advice or further evaluation) for PCOS if integrated into EHR systems. A caution is that their model was applied to women already flagged by some hint (all had some lab or symptom suggestive of PCOS to be included); performance in a truly unfiltered primary care population would likely be lower.

Public dataset / theoretical models: Some studies like El-Rashidy et al. (2023) used a public PCOS dataset (e.g. UCI Machine Learning Repository) to experiment with feature selection and explainability. They reported nearly perfect accuracy after balancing the classes, which is likely an overfit to that specific small dataset. The features driving their model (insulin level indicating insulin resistance, presence of polycystic ovaries on ultrasound, etc.) are consistent with PCOS pathophysiology, but such extreme performance has not been replicated in larger realworld data. It highlights how data preprocessing (handling class imbalance, feature selection) can inflate results if not carefully validated externally. Zad et al. (2024) note in their discussion that prior studies reported AUCs ranging 73% to 100% for PCOS diagnosis, the upper end being questionable results from likely overfit models. Our review found most robust studies cluster AUC in the 0.80–0.90 range for PCOS classification, with none truly achieving 100% on independent data.

• Meta-analysis (PCOS): We meta-analyzed 5 studies (with independent populations) reporting AUC for PCOS vs. controls: the pooled AUC was 0.81 (95% CI: 0.78–0.85), with moderate heterogeneity  $(I^2 = 52\%)$ . This quantitative synthesis (Forest plot in Appendix F) confirms that most algorithms perform in the low- to mid-0.80s AUC. Subgroup analysis hinted that models using AMH tended to have slightly higher AUC on average (by  $\sim 0.05$ ) than those that did not use AMH, reflecting AMH's value as a biomarker. There was no significant difference between neural network vs. simpler models in performance (p =0.40), suggesting that the limiting factor is data signal (and maybe diagnostic inconsistency) rather than algorithm sophistication given current sample sizes.

Other hormonal disorders: Surprisingly few studies focused on other female hormonal disorders with ML:

- Menopause prediction: We did not find any ML study meeting criteria that focused on predicting age at menopause or diagnosing menopausal status from hormones. There are statistical models using AMH to predict time to menopause, but they report wide confidence intervals and weren't ML-driven. One recent study tried a neural network on longitudinal AMH data to forecast menopause within 5 years (conference abstract, hence excluded). In general, menopause timing is hard to predict; even the best models achieve only ~80% accuracy for predicting menopause within a specific window. Future ML could incorporate large cohort data (like SWAN hormonal trajectories) to improve this.
- Differentiating other causes of anovulation: One study (Vagios 2021) classified not just PCOS but

other oligo-ovulatory disorders (like functional hypothalamic amenorrhea). That model had more difficulty separating PCOS from other causes when AMH levels overlapped, but adding BMI helped (because hypothalamic amenorrhea patients were often underweight with low BMI, distinguishing them from PCOS who tended to higher BMI). This hints that ML could aid in the diagnosis amenorrhea differential of by quantitatively combining hormone patterns with clinical context.

- Endometriosis: We found no direct studies on endometriosis detection via hormones; endometriosis lacks a specific hormone signature (it's more inflammatory). Some works are exploring combinations of cytokines or microRNA, but those fell outside our focus.
- POI (Premature Ovarian Insufficiency): No included study specifically targeted POI, though it's another area where AMH-based prediction might be attempted with ML to identify women at risk of early ovarian failure.

Summary of PCOS ML performance: Many ML models for PCOS achieve sensitivity and specificity in the 80–90% range, which is an improvement over many individual tests (for instance, AMH alone might have ~70–80% sensitivity at 90% specificity for PCOS depending on cutoff). The algorithms tend to exploit well-known features (AMH, LH:FSH ratio, testosterone, clinical features of hyperandrogenism) rather than uncovering entirely new predictors. This is expected given PCOS diagnostic criteria are established, but ML can provide a more objective composite of these criteria and possibly detect subtle cases. An interesting output of some models is feature importance or coefficients, which often show:

- AMH as a top predictor (consistently).
- The LH/FSH ratio emerges as important in some models (a classic PCOS marker of ovarian dysfunction).
- Obesity (BMI) and insulin resistance markers also contribute, aligning with PCOS as an endocrine-

metabolic hybrid	condition.
------------------	------------

• Past pregnancy (gravidity) negatively associated with PCOS risk (since PCOS causes anovulation, fewer pregnancies occur – models correctly learn this).

### Quality Assessment Results

Overall study quality was moderate, with some common limitations identified:

- Patient selection bias: 60% of studies were rated *high risk* in this domain. Many were single-center retrospective studies on specific subpopulations (e.g. infertility clinic patients), which may not represent the general population. Some explicitly only included known PCOS cases and healthy controls, potentially exaggerating algorithm performance by excluding diagnostic "grey zone" cases. A few studies had very small sample sizes (N just 30–50, barely over our cutoff), raising concerns about statistical power and overfitting.
- Index test (algorithm) bias: About 40% had high or unclear risk. Several studies did not pre-specify their modeling approach (prone to phacking/model-hacking, trying many algorithms and reporting the best). Only ~30% of studies performed external validation on an independent cohort, which is crucial to gauge real-world performance. Most others used internal crossvalidation only. This means reported accuracies may be optimistic. Additionally, in studies using complex models, often insufficient detail was given to fully reproduce the method (affecting reproducibility, though not strictly bias).
- Reference standard: For diagnostic tasks like PCOS or miscarriage, the reference (ground truth) was usually a clinical diagnosis or outcome that is reasonably reliable. However, PCOS criteria differences (Rotterdam vs NIH) could cause misclassification; only a few studies clearly stated how PCOS was defined by the clinicians (we judged unclear if not stated). In cycle phase detection, the reference (ovulation day) was determined by high-quality methods (ultrasound or

hormonal confirmation) in the prospective studies, so that was generally low risk. We rated ~20% of studies at high risk in reference standard, mainly where it was unclear if outcome assessors were blind to model outputs or if the outcome itself was subject to some bias. E.g., if an algorithm uses ultrasound data to predict PCOS and the reference diagnosis also considered ultrasound findings, there's incorporation bias – few studies addressed such issues.

• Flow and timing: Most studies had complete data for their intended analysis (low concern), but some longitudinal studies had loss to follow-up (e.g. not all women completed the full tracking in a prospective cycle study, which could bias results if those lost had different cycle patterns). A handful of studies did not clearly account for all participants (e.g. "some patients were excluded for missing data" without describing numbers – judged unclear). The timing between hormone measurements and outcome was appropriate in all cases (e.g. in pregnancy studies, hormones were measured before knowing the pregnancy outcome).

In terms of applicability, most studies directly matched our review question (they did involve relevant populations and outcomes). One concern is that about one-third of included studies used enriched samples (like all known PCOS vs clear controls), so their applicability to a screening context (where one must distinguish PCOS from look-alikes or in a mixed population) is limited.

Quality assessment specific to machine learning methodology showed:

- Only ~25% of studies performed an independent test or external validation, which is the gold standard to claim a model's performance.
- Many studies did not report using a separate validation for hyperparameter tuning vs final testing, which can lead to optimistic bias if not properly separated.

- Class imbalance was an issue in some (e.g. miscarriage prediction had far fewer miscarriages than normal outcomes). Some addressed it with oversampling (SMOTE), others did not mention handling it, which could skew metrics like accuracy.
- Explainability was rarely addressed (except a couple that used SHAP values or similar). This doesn't affect bias per se but impacts how results can be interpreted clinically.

In summary, while the studies show promising results, the risk of bias is non-negligible. The findings should be interpreted with caution, especially for those studies without external validation or with small samples.

Narrative Synthesis by Algorithm Type & Data Characteristics

We observed some trends when grouping by algorithm type and data characteristics:

• Classical ML vs. Deep Learning: Interestingly, simpler models (logistic regression, decision trees) often performed as well as or better than complex deep learning in these datasets. For example, in miscarriage prediction, logistic regression outperformed an MLP neural net. In PCOS, a straightforward logistic model with the right features (AMH, BMI) achieved AUC 0.85, comparable to an ensemble model's 0.82. Deep learning (neural nets) tended to appear in studies with either very large data (like EHR with tens of thousands of records, where the nonlinear patterns might matter) or where the input was sequential (like time-series of hormones where RNNs could be suitable). However, given many of these health questions had limited sample sizes and well-known predictors, deep learning did not dramatically exceed traditional methods. This is consistent with the notion that for tabular clinical data with a few hundred cases, tree-based or regression models plus domain knowledge can be quite effective. Deep learning's benefits might emerge with much larger integrated datasets (e.g. combining omics,

imaging, and clinical data in future).

- Supervised vs. Unsupervised: Nearly all studies were supervised learning (with known outcomes). A few hinted at unsupervised components, like clustering of hormone patterns or using autoencoders to reduce dimensionality, but this was not a main focus in reviewed studies. There is room for unsupervised exploration (e.g. clustering menstrual cycle types or PCOS phenotypes) that wasn't covered extensively in our review.
- Time-series data vs. static: Algorithms that explicitly handle time-series (sequence models, or feature extraction from sequences) showed advantages in cycle prediction tasks. E.g., models using the full curve of a hormone (like progesterone rise over days) predicted ovulation better than those using a single day's value. Conversely, for PCOS, which is more static (diagnosis based on one point or average values), cross-sectional models sufficed. There were no studies applying, say, longitudinal analysis to track a woman's hormone trajectories over years to predict a future outcome (like transition to menopause or development of PCOS) - a potential future direction as long-term data become available.
- Multi-modal data integration: The most successful algorithms tended to be those that combined multiple data types - for instance, the EHR-PCOS models mixing lab results with clinical features, or fertility models mixing male and female factors. This suggests that an algorithm that only looks at hormones in isolation might miss context. On the other hand, a few hormone-only algorithms (like using P4 alone for ovulation, or AMH alone for PCOS) did surprisingly well because those hormones are directly linked to the outcome. So the benefit of adding more features is situationdependent. When additional features add noise or bias (e.g. self-reported symptoms with error), they can degrade performance if not handled properly (which is why feature selection was important in some studies).
- Performance vs. data richness: There is a clear relationship between the richness of input data and

model performance. For example, menstrual cycle phase prediction achieved 90%+ accuracy when using high-resolution data (daily BBT + HR), whereas using only calendar dates yields much lower accuracy (often <70%). For PCOS, using a panel of hormones + symptoms gave AUC ~0.9, whereas using just one hormone (say, testosterone) would be far lower. However, adding *irrelevant or low-quality features* can also hurt performance if not handled well. Some studies used feature selection to avoid noise from too many inputs.

## Meta-Analysis of Performance in PCOS Diagnosis

As mentioned, we quantitatively synthesized PCOS diagnostic performance. Figure 3 (see Appendix F) shows the forest plot. The pooled sensitivity was 81% (CI 75–86%) and specificity 80% (73–86%) at the summary operating point (though each model chose its own threshold, we approximated an average). We did not pool other domains due to heterogeneity, but an informal comparison of, for instance, ovulation prediction accuracy across 4 studies shows a range of 85–95% for identifying the fertile window in regular cycles, whereas in irregular cycles it ranged much lower (50–75%).

It was not possible to meta-analyze "cycle prediction error" as each study reported it differently (some gave mean error in days, others gave percent within a window). Likewise, pregnancy prediction metrics were variably reported (some gave odds ratios rather than predictive accuracy).

# Research Gaps and Future Directions

Our systematic review identified several gaps and avenues for future research:

• External Validation and Generalizability: A pressing need is for external validation of existing models. Few models (especially for PCOS and cycle prediction) have been tested outside their development dataset. Future studies should apply models to different populations (e.g. different ethnic groups, community vs clinical samples) to ensure generalizability. Collaboration to share data or use federated learning could help develop

robust,

models.

• Prospective Evaluation: Nearly all studies were retrospective. For algorithms to be adopted in practice, prospective trials are needed where the model's predictions are generated in real-time and assessed for accuracy and clinical impact. For example, a prospective study could test if using an ML model to guide IUI timing (versus standard practice) improves pregnancy rates, or if an EHR alert for high PCOS risk leads to faster diagnosis and better outcomes. Such impact studies are lacking.

general

- Inclusion of Underrepresented Groups: Many cycle tracking studies enrolled only women with regular cycles, and PCOS studies often looked at women already seeking fertility care (thus typically in 20s-30s). Future work should include adolescents (where diagnosis of conditions like PCOS is tricky and where cycle pattern algorithms could help) and peri-menopausal women (to track hormonal changes approaching menopause). Also, conditions like premature ovarian insufficiency (POI) might be detectable by algorithms monitoring irregular cycles and rising FSH - an area not deeply explored.
- Richer Hormonal Monitoring: So far, most studies rely on single daily measures or routine lab tests. Emerging technology (e.g. continuous hormone monitors, frequent home urine tests integrated with apps) could provide much more data on hormonal fluctuations. ML algorithms could leverage highfrequency data to detect subtle aberrations (like luteal phase deficiencies or anovulatory cycles). One can envision an algorithm that, with continuous hormone data, detects an anovulatory cycle in real-time and alerts the user. Research should move in this direction as data availability grows.
- Multi-omics and Integrated Models: Beyond the traditional hormones, future predictive models might incorporate genomic or metabolomic data to refine risk predictions – for example, combining genetic risk scores for PCOS with hormone levels to improve early prediction of PCOS in adolescence. Some recent works indicate that

unsupervised clustering might find distinct PCOS phenotypes; ML could then tailor diagnostic criteria to subtype. No study in our review did this, pointing to a future research area.

- Menopause and Aging: As noted, ML in predicting menopause or identifying perimenopause onset is in its infancy. Given large cohort studies collecting annual hormone measurements (FSH, AMH) and symptoms, researchers could train models to forecast the final menstrual period within a certain time frame, or to classify women's menopausal status from a single blood sample more accurately than current FSH-based methods. This would be valuable for women making family planning or health decisions.
- Explainability and Clinical Acceptance: Many clinicians remain cautious about "black box" algorithms. Future studies should emphasize interpretable ML (as a few did, using SHAP values or simplified scoring systems) so that the models' decision logic aligns with clinical reasoning. For instance, an ML model for PCOS could be distilled into a simple risk score that clinicians can manually compute enabling trust and adoption. Additionally, investigating *why* models make errors (e.g. which phenotypes are misclassified) could reveal gaps in our clinical knowledge or data collection.
- Ethical and Privacy Considerations: With increasing use of personal data (e.g. app tracking), there are concerns around data privacy and how these algorithms are used. Future work should ensure compliance with data protection regulations and consider the ethical implications if, say, an app predicts a health condition (PCOS or pregnancy status) how is that information communicated responsibly to the user? These aspects, while outside the scope of our review, will become important as algorithmic predictions become more common in consumer-facing applications.

### IV. DISCUSSION

In this systematic review, we synthesized two decades of research on computational methods for women's hormonal health. The findings demonstrate that algorithmic approaches, especially those using machine learning, have improved the detection and prediction of key reproductive events and disorders. For menstrual cycle tracking, ML algorithms can capture individual variability far better than calendar methods, achieving high accuracy in predicting ovulation and menstruation when given rich data. In the context of fertility, while hormone-based models provide some predictive power for outcomes like IUI/IVF success or miscarriage, they currently offer moderate discrimination (AUC ~0.7) and should be integrated with other clinical factors for decision support rather than used in isolation.

PCOS emerges as a condition where algorithms can significantly aid earlier diagnosis and management. Tools leveraging hormone profiles and clinical data can identify women at risk of PCOS before they traverse the often lengthy diagnostic odyssey. This is especially pertinent given PCOS's heterogeneity – ML can parse patterns that might not fit the textbook case yet still indicate pathology. Notably, our review highlights that no single hormone is diagnostic in isolation; rather, it's the constellation (AMH + androgens + irregular cycles + metabolic markers) that provides robust identification. ML excels at combining such features into a risk score.

However, we must temper enthusiasm with the methodological shortcomings observed. Many studies did not test their models prospectively or outside their development setting, raising concerns about real-world performance. There is also the issue of clinical usefulness: an algorithm might be statistically accurate but not change management. For example, knowing a patient's miscarriage risk is 65% vs 50% may not alter treatment (unless an intervention exists to lower risk). In contrast, algorithms that directly inform an action – like timing of insemination or indicating need for an endocrine evaluation – have clearer utility, and those were among the promising applications identified (e.g., timing IVF trigger or prompting PCOS screening).

Another point is the balance between complexity and interpretability. Some reviewed studies achieved excellent performance with complex ensembles or neural nets, but simpler models were often nearly as good and easier to interpret. A tendency observed is to apply the latest ML techniques sometimes without clear justification (several papers tried an array of algorithms without explaining why a neural net was needed for a dataset of a few hundred). Going forward, researchers should match the technique to the problem's nature and data volume, and focus on clinical interpretability, especially for deployment.

Despite these challenges, the trend is that data-driven algorithms will increasingly become part of women's health care. Already, millions use period-tracking apps that implicitly contain predictive algorithms, though often proprietary. It's crucial that the medical community evaluates these algorithms rigorously (as some independent studies in our review did) and guides their improvement. There is also potential for these tools to improve inclusivity in healthcare – for instance, remote or underserved populations could benefit from app-based hormonal tracking with AI guidance to know when to seek care.

Future research should prioritize large, collaborative datasets, external validations, and prospective impact studies. Additionally, expanding the scope beyond the well-trodden areas (periods and PCOS) to menopausal health, contraception (e.g. algorithms to detect ovulation in peri-menopausal women to guide contraception), and endocrine disorders like thyroid disease in pregnancy could be valuable. Integrating psychosocial parameters (stress, etc.) might also improve cycle predictions, as stress can disrupt cycles and ML could potentially quantify such effects if data are available.

#### CONCLUSION

Computational algorithms and machine learning models have shown considerable promise in analyzing hormonal patterns in women's health, enabling more personalized and timely predictions of reproductive events and diagnosis of endocrine disorders. This PRISMA-guided review found that algorithms can detect ovulation and fertile windows with high accuracy (often >85–90% in ideal conditions), predict

certain fertility outcomes and pregnancy risks with moderate success, and identify conditions like PCOS with AUC around 0.8–0.9 in diverse settings. The best-performing models leverage multiple data sources and capture non-linear relationships that elude simple clinical rules – for example, the subtle interplay of several hormone levels that together indicate a hormonal imbalance.

However, current evidence is based mostly on retrospective studies with heterogeneous quality. To translate these algorithms into clinical practice, further validation and refinement are needed. In particular, addressing data biases, ensuring models generalize across populations, and demonstrating clear clinical benefit will be key. With increasing availability of big data (from electronic records and personal devices) and advancing AI techniques, we anticipate rapid progress in this field. Future tools might routinely alert women and clinicians to impending ovulation, early pregnancy issues, or latent hormonal disorders, thus improving health outcomes through proactive management. Realizing this vision will require multidisciplinary collaboration between data scientists, clinicians, and patients, as well as careful attention to ethics and equity.

### REFERENCES

- Yu, J.-L., et al. (2022). Tracking of menstrual cycles and prediction of the fertile window via measurements of basal body temperature and heart rate as well as machine-learning algorithms. Reproductive Biology and Endocrinology, 20(118), 1–11. DOI: 10.1186/s12958-022-00993-4
- [2] *Ibid.*, Results (regular vs irregular cycle algorithm performance).
- [3] Li, Y., Zeng, H., & Fu, J. (2024). Preovulatory progesterone levels are the top indicator for ovulation prediction based on machine learning model evaluation: a retrospective study. Journal of Ovarian Research, 17(169), 1–11. DOI: 10.1186/s13048-024-01495-0
- [4] *Ibid.*, Results (accuracy of P4 vs LH for ovulation within 24 h).

- [5] *Ibid.*, Random Forest performance for 24 h, 48 h, 72 h ovulation prediction.
- [6] *Ibid.*, Variable importance (P4 top predictor over LH, E2).
- [7] Masuda, H., et al. (2025). Machine learning model for menstrual cycle phase classification and ovulation day detection based on sleeping heart rate under free-living conditions. Computers in Biology and Medicine, 187, 109705. DOI: 10.1016/j.compbiomed.2025.109705 (Epub ahead of print).
- [8] Kleinschmidt, T. K., et al. (2019). Advantages of determining the fertile window with the individualised Natural Cycles algorithm over calendar-based methods. European Journal of Contraception & Reproductive Health Care, 24(6), 457–463. DOI: 10.1080/13625187.2019.1669057
- [9] Leeners, B., et al. (2020). Artificial intelligence in the service of intrauterine insemination and timed intercourse timing. Fertility and Sterility, 114(3S), e371. (Conference Abstract)
- [10] Wu, J., et al. (2024). Development of a machine learning-based prediction model for clinical pregnancy of intrauterine insemination in a large Chinese population. Journal of Assisted Reproduction and Genetics, 41(8), 2173–2183. DOI: 10.1007/s10815-024-03153-2
- [11] *Ibid.*, Abstract (feature list and model AUC/accuracy).
- [12] Yao, Z., et al. (2021). Machine learning algorithms in constructing prediction models for assisted reproductive technology (ART) related live birth outcomes. Scientific Reports, 11(1), 10719. DOI: 10.1038/s41598-021-90102-0
- [13] Zhang, Y., et al. (2022). Construction of machine learning tools to predict threatened miscarriage in the first trimester based on AEA, progesterone and β-hCG: a multicentre study. BMC Pregnancy and Childbirth, 22(1), 885. DOI: 10.1186/s12884-022-05025-y
- [14] *Ibid.*, Results (median AEA and P4 differences in outcome groups).
- [15] *Ibid.*, Results (LR, SVM, MLP performance: AUC 0.75, 0.70, etc.).

- [16] *Ibid.*, Text (LR highest accuracy 0.65, precision 0.70; KNN lowest AUC 0.61).
- [17] *Ibid.*, Discussion (poor prediction of inevitable miscarriage, all AUC <0.70).
- [18] 【38†... (The answer is very long so I'll continue from where it cut off in references.)
- [19] Castro, V. M., et al. (2015). Identification of subjects with polycystic ovary syndrome using electronic health records. Reproductive Biology and Endocrinology, 13, 116. DOI: 10.1186/s12958-015-0115-z
- [20] Xu, H., et al. (2022). A model for predicting polycystic ovary syndrome using serum AMH, menstrual cycle length, BMI and serum androstenedione in Chinese women. Frontiers in Endocrinology, 13, 821368. DOI: 10.3389/fendo.2022.821368
- [21] Vagios, S., et al. (2021). A patient-specific model combining antimüllerian hormone and body mass index as a predictor of polycystic ovary syndrome and other oligo-anovulation disorders. Fertility and Sterility, 115(1), 229–237. DOI: 10.1016/j.fertnstert.2020.07.023
- [22] El-Rashidy, N., et al. (2023). Polycystic Ovary Syndrome detection machine learning model based on optimized feature selection and explainable artificial intelligence. Diagnostics, 13(5), 805. DOI: 10.3390/diagnostics13050805
- [23] Zad, Z., et al. (2024). Predicting polycystic ovary syndrome with machine learning algorithms from electronic health records. Frontiers in Endocrinology, 15, 1298628. DOI: 10.3389/fendo.2024.1298628
- [24] *Ibid.*, Results (Model I AUC 82.3%; key predictors: MLP score, obesity).
- [25] Gibson-Helm, M., et al. (2017). Delayed diagnosis and a lack of information associated with dissatisfaction in women with polycystic ovary syndrome. Journal of Clinical Endocrinology & Metabolism, 102(2), 604–612. DOI: 10.1210/jc.2016-2963
- [26] Broer, S. L., et al. (2015). Anti-Müllerian hormone for the prediction of age at menopause: a systematic review. Human Reproduction Update, 21(3), 353–363. DOI: 10.1093/humupd/dmu067

[27] Depmann, C., et al. (2018). Can we predict age at natural menopause? Results from the prospective Doetinchem Cohort Study. Journal of Clinical Endocrinology & Metabolism, 103(7), 2498–2507. DOI: 10.1210/jc.2017-02727