



# Endometriosis early detection

Daniel Moshe, Adi Haber





## Contents

Background .....	3
Problem Statement .....	3
Proposed Solution .....	4
Introduction .....	4
Machine Learning .....	4
Endometriosis .....	4
UK BioBank [6] .....	5
Related Works .....	5
Solution Description .....	6
Stages .....	6
Extracting Data .....	6
Tools and Technologies .....	7
Anticipated Challenges .....	7
Acknowledgments .....	8
Bibliography .....	9



## Background

Endometriosis, a chronic inflammatory condition, primarily manifests through symptoms such as pain and infertility [1]. It occurs when tissue resembling the uterine lining grows outside the uterus, adhering to pelvic organs, and occasionally other areas of the body. This abnormal adhesion triggers inflammation and the formation of scar tissue, resulting in debilitating pain and, in some instances, infertility. Endometriosis predominantly affects women of reproductive age, with research suggesting that approximately 5-10% of this demographic, totalling around 180 million individuals globally, are affected [2].

Main known indicators of the endometriosis include:

- Pelvic pain and/or lower abdominal pain
- Painful menstrual cramps
- Abnormal menstrual bleeding pattern (either by amount or irregularity)
- Family history [3]
- Infertility [4]

Diagnosing endometriosis presents a challenge since adhesions are not always detectable through imaging techniques like ultrasound or MRI. Typically, a definitive diagnosis necessitates undergoing laparoscopic surgery [3] and a subsequent biopsy.

## Problem Statement

Our primary objective revolves around the prompt identification of endometriosis. Remarkably, 60% of women dealing with endometriosis navigate consultations with three or more clinicians before receiving a diagnosis, leading to an average delay of seven years before definitive diagnosis [5]. This prolonged delay intensifies symptoms, lowers overall quality of life, and contributes to enduring reproductive health challenges. Conventional diagnostic methods, predominantly reliant on invasive procedures and subjective assessments, further complicate the diagnostic process.



## Proposed Solution

This project endeavors to aid in diagnosing endometriosis by analyzing patient data. We will collect data on endometriosis and healthy patients from the UK Biobank and select a group of features (symptoms and risk factors) from which we will try to detect the existence of endometriosis. Our primary goal is to build the optimal machine-learning model for accurate endometriosis detection based on the features we found.

## Introduction

### Machine Learning

Machine learning, a subset of artificial intelligence, revolutionizes medical research by extracting insights from vast datasets to enhance diagnostic accuracy, treatment efficacy, patient outcomes, and identifying risk factors. There are two primary subcategories of machine learning - supervised and unsupervised learning.

Supervised learning algorithms use labeled data to train models to predict outcomes or classify instances, offering valuable insights into disease detection and prognosis.

Unsupervised learning techniques uncover hidden patterns within unlabeled data, enabling researchers to identify unexplored disease subtypes or biomarkers.

Deep learning, a subset of machine learning, utilizes neural networks with multiple layers to automatically extract complex features from raw data, paving the way for advanced image analysis, genomic sequencing, and complex medical issues.

With the integration of these machine learning paradigms, medical researchers unlock unprecedented opportunities to unravel the complexities of diseases, revolutionizing healthcare delivery.

### Endometriosis

Endometriosis, a prevalent chronic gynecological condition reliant on estrogen, concerns the presence of uterine endometrial tissue outside its normal cavity. This



disorder is characterized by the presence of endometrial tissue outside the uterus, leading to pelvic pain and fertility issues.

## UK BioBank [6]

UK Biobank is a large-scale biomedical database and research resource, containing in-depth, de-identified genetic and health information from half a million UK participants. The database, which is regularly augmented with additional data, is globally accessible to approved researchers and scientists undertaking vital research into the most common and life-threatening diseases. UK Biobank's research resource is a major contributor to the advancement of modern medicine and treatment and has enabled several scientific discoveries that improve human health.

## Related Works

Several research papers have applied machine learning techniques to predict endometriosis. These studies serve as a starting point for our research, allowing us to refine our unique research question and build upon existing knowledge in the field.

Study	Link	Date of Publication
Revisiting the Risk Factors for Endometriosis: A Machine Learning Approach	<a href="https://www.mdpi.com/1716284">https://www.mdpi.com/1716284</a>	Jul-2022
Diagnosis of Endometriosis Based on Comorbidities: A Machine Learning Approach	<a href="https://www.mdpi.com/2554716">https://www.mdpi.com/2554716</a>	Nov-2023
Machine learning algorithms as new screening approach for patients with endometriosis	<a href="https://doi.org/10.1038/s41598-021-04637-2">https://doi.org/10.1038/s41598-021-04637-2</a>	Jan-2022



## Solution Description

### Stages

In the initial phase of our research, we comprehensively explore the realm of endometriosis research, delving into the medical domain to better understand the features embedded within our dataset. We extract relevant data from the extensive UK Biobank dataset, excluding irrelevant features. We meticulously curate our cohort, ensuring the inclusion of only relevant data points and features. To comprehensively explore endometriosis research, it is essential to delve into the medical domain and gain a precise understanding of the parameters within our dataset. By doing so, we can discern the significance of each parameter and optimize the dataset accordingly. This optimization process enables us to harness familiar tools for analyzing our predictive models and developing a specialized model tailored to address the specific challenges in endometriosis detection. Ultimately, we aim to assess and compare the results of various machine learning models to recommend the most comprehensive and precise approach for predicting endometriosis based on symptoms.

### Extracting Data

In our effort to extract meaningful insights from the UK Biobank dataset for predicting endometriosis presence, we encounter several challenges stemming from the dataset's vastness and heterogeneity. One of the primary issues is the abundance of irrelevant data that do not pertain to our research question, requiring meticulous data selection. Additionally, a considerable portion of the records within the dataset contain missing values, further complicating our analysis. As we choose features for our predictive models, we may exclude records that have missing values for other selected features, which could make our dataset incomplete. To overcome this challenge, we will employ various data imputation techniques to address missing values, such as mean or median imputation or predictive imputation using machine learning algorithms (for example, KNN imputation). As part of our data cleaning process, we will meticulously review our dataset, removing any inaccurate or outlier data points. Despite these challenges, our



approach will aim to extract the most relevant and informative data from the entirety of the UK Biobank dataset, enabling us to develop robust predictive models for endometriosis diagnosis.

## Tools and Technologies

In our research, we plan to use various tools and technologies to address the challenge of predicting endometriosis presence from tabular data, framed as a binary classification problem. We will explore classic machine learning models such as logistic regression, decision trees, random forests, and support vector machines. These models provide a solid foundation for understanding the relationships between features and the target variable. Additionally, we plan to experiment with more complex models, like gradient boosting machines (GBM) and ensemble methods, to help capture intricate patterns and interactions within the data. Furthermore, we will employ feature engineering techniques to extract meaningful insights and enhance model performance. Feature engineering may involve creating new features, transforming existing ones, or selecting the most relevant features using techniques like recursive feature elimination or feature importance analysis. Finally, we intend to explore deep learning models utilizing neural networks to uncover deeper layers of abstraction and potentially capture nonlinear relationships present in the data. By employing this diverse array of tools and techniques, we aim to identify the most effective approach for accurately predicting endometriosis presence and ultimately contributing to advancements in medical research and diagnosis using machine learning.

## Anticipated Challenges

In our research, we confront several notable constraints that shape the scope and reliability of our findings. Firstly, while utilizing the UK Biobank dataset provides valuable insights, its demographic skew towards women averaging 50 years old presents a limitation. Given that our objective is to assist in diagnosing endometriosis in younger



women, the dataset may not fully represent the nuances of the condition in this demographic. Moreover, a substantial proportion of cases labelled as 'has endometriosis' are self-diagnosed, comprising approximately 40% of our dataset. This fact introduces a layer of uncertainty regarding the accuracy of the diagnoses, as self-diagnoses may lack the certainty and precision of clinical assessments. Furthermore, the dynamic nature of endometriosis progression and treatment outcomes necessitates longitudinal data, which may be limited in our dataset. Despite these constraints, our research strives to navigate these complexities and contribute towards advancing the understanding and diagnosis of endometriosis.

## Acknowledgments

We thank Adi Shraibman, Dorit Shweiki, and Yonatan Bilu for their mentorship and support throughout the project.

We thank The Academic College of Tel Aviv-Yaffo for providing us access to the UK BioBank data.

We thank Grammarly and ChatGPT for grammar and tone suggestions.





## Bibliography

- [1] Wang, P.-H., Yang, S.-T., Chang, W.-H., Liu, C.-H., Lee, F.-K., & Lee, W.-L. (2022). Endometriosis: Part I. Basic concept. *Taiwanese Journal of Obstetrics and Gynecology*, 61(6), 927–934. <https://doi.org/10.1016/j.tjog.2022.08.002>
- [2] Zondervan, K. T., Becker, C. M., Koga, K., Missmer, S. A., Taylor, R. N., & Viganò, P. (2018). Endometriosis. *Nature reviews. Disease primers*, 4(1), 9. <https://doi.org/10.1038/s41572-018-0008-5>
- [3] Blass, I., Sahar, T., Shraibman, A., Ofer, D., Rappoport, N., & Linial, M. (2022). Revisiting the Risk Factors for Endometriosis: A Machine Learning Approach. *Journal of Personalized Medicine*, 12(7), 1114–1114. <https://doi.org/10.3390/jpm12071114>
- [4] Tom Gunnar Tanbo, & Péter Fedorcsák. (2017). Endometriosis-associated infertility: aspects of pathophysiological mechanisms and treatment options. *Acta Obstetrica et Gynecologica Scandinavica*, 96(6), 659–667. <https://doi.org/10.1111/aogs.13082>
- [5] Horne, A. W., & Missmer, S. A. (2022). Pathophysiology, diagnosis, and management of endometriosis. *BMJ (Clinical research ed.)*, 379, e070750. <https://doi.org/10.1136/bmj-2022-070750>
- [6] UK BioBank (2015). [Ukbiobank.ac.uk](https://www.ukbiobank.ac.uk/learn-more-about-uk-biobank/about-us). <https://www.ukbiobank.ac.uk/learn-more-about-uk-biobank/about-us>