

“Automated Diabetes Prediction: Comparative Analysis
of Machine Learning Models and Feature Importance”

MUHAMMAD SHAYAN UMAR

12817

RESEARCH IN SOFTWARE CONSTRUCTION
AND DEVELOPMENT

DECEMBER, 11th 2023

Contents

1. Introduction.....	4
1.1 Background.....	4
1.2 Problem Statement.....	4
1.3 Research Aim.....	5
1.4 Research objective	5
1.4.1 Execution Evaluation and Examination:.....	5
1.4.2 Hyper-parameter Streamlining:.....	5
1.4.3 Investigation of Prescient Attributes:.....	5
1.4.4 Generalizability Appraisal:	5
1.5 Significance of the Study	6
1.5.1 Redesigned Early ID:	6
1.5.2 Resource Smoothing out:	6
1.5.3 Indicative Precision Improvement:	6
1.5.4 Change of Clinical consideration Procedures:	6
1.5.5 Obligation to Perceptive Clinical consideration:	6
1.6 Research Methodology	6
1.6.1 Machine Learning Algorithms:	7
1.6.2 Exploratory Data Analysis (EDA):	7
1.6.3 Statistical Analysis:.....	7
1.6.5 Software Construction and Development:	7
2. Literature Review.....	7
2.1 Predictive Models in Diabetes Prediction.....	8
2.2 Feature Selection Strategies	8
2.3 Factors Influencing Model Performance.....	8
2.4 Challenges and Future Directions	8
3. Proposal Solution.....	9
3.1 Base Model Development.....	9
Logistic Regression:.....	9
K-Nearest Neighbors (KNN):	9
Decision Tree Classifier:.....	9
Random Forest (RF):	9
Support Vector Machine (SVM):.....	9

XGBoost (XGB):	9
LightGBM.....	9
3.2 Model Evaluation.....	9
3.3 Hyper-parameter Tuning.....	9
3.4 Ensemble Model Development.....	10
3.5 Model Comparison.....	10
4. Research Findings	10
4.1 Model Accuracies	10
4.2 Hyperparameter Tuning Results	10
4.3 Feature Engineering Insights	11
4.3.1 BMI Categorization:	11
5. Discussion	12
5.1 Implications of Findings:	12
5.1.1 Proactive Healthcare Management:	12
5.1.2 Feature Engineering Impact:	12
5.2 Comparison with Existing Literature:.....	12
5.2.1 Model Performance:.....	12
5.2.2 Feature Engineering:	12
5.3 Limitations:	13
5.3.1 Dataset Constraints:	13
5.3.2 Temporal Dynamics:.....	13
5.4 Future Directions:	13
5.4.1 Explainability Measures.....	13
5.4.2 Ethical Considerations:	13
6. Conclusion	13
6.1 Key Findings:.....	13
6.1.1 Model Performance:.....	13
6.1.2 Hyperparameter Tuning:	14
6.1.3 Feature Engineering Impact:	14
6.2 Contributions.....	14
6.2.1 Automation of Diabetes Prediction:.....	14
6.2.2 Efficiency Enhancement:	14
6.3 Future Research Avenues.....	14
6.3.1 Ensemble Models:.....	14
6.3.2 Longitudinal Data Analysis:	14
6.3.3 Personalized Healthcare:.....	15
7. References.....	15

1. Introduction

1.1 Background

The rising commonness of diabetes represents a critical test to present day medical services frameworks, requiring imaginative procedures for opportune ID and compelling therapy. In light of this test, our exploration attempts to tackle the force of cutting edge AI methods to computerize the forecast of diabetes. This mechanization isn't just a specialized progression; it addresses a major change in medical services strategies with the possibility to further develop patient consideration results significantly.

Generally, diabetes expectations has depended on ordinary strategies that may not necessarily give the precision and productivity expected for early recognizable proof. These techniques frequently include manual assessments, which can be tedious and may not use the maximum capacity of accessible information. By coordinating AI calculations into the expectation cycle, our examination looks to conquer these restrictions and usher in another time of prescient medical services.

The center goal is to foster AI models, explicitly utilizing calculations like irregular woodland, LightGBM, and XGBoost. These models are not simply computational apparatuses; they are instruments that have the ability to change how medical care is conveyed. They can possibly break down immense datasets with speed and accuracy, taking into consideration the proactive recognizable proof of people in danger of diabetes.

The extraordinary effect stretches out past the specialized domain to the actual heart of medical services conveyance. Via robotizing the expectation interaction, medical care specialists can move from a responsive way to deal with a proactive one. This implies that potential wellbeing chances related with diabetes can be recognized and tended to before, prompting more successful intercessions and customized medical care plans.

Fundamentally, our exploration isn't just about creating AI models; it is tied in with reclassifying the scene of diabetes expectation. It is tied in with making devices that engage medical services experts to remain on the ball, giving better consideration and results to people in danger of diabetes. Through the incorporation of state of the art philosophies, our exploration tries to contribute fundamentally to the change in perspective towards proactive and prescient medical services.

1.2 Problem Statement

The standard procedures utilized for diabetes presumption a large part of the opportunity arrive up short as for the major accuracy for worthwhile unmistakable affirmation considering their dependence on manual assessments. Seeing this limit, our appraisal plans to address these inadequacies by presenting robotization through computerized reasoning procedures. Through computerizing the presumption cycle, we plan to beat the normal difficulties of manual

assessment and tap into the huge limit of broad datasets. This change in setting should yield measures that are useful as well as phenomenally more precise, preparing for extra made clinical advantages results and proactive mediation techniques.

1.3 Research Aim

The overall goal of this exploration project is to tackle the capability of AI calculations in the production of a strong and computerized framework for diabetes expectation. The point of convergence of our undertaking is the advancement of a prescient model that upgrades the precision of diabetes forecast as well as works on the general productivity of the cycle. Through the consolidation of state of the art strategies, we intend to push the limits of current prescient models, encouraging progressions that rise above conventional methodologies. A definitive yearning is to add to more effective medical services the executives through the organization of refined AI procedures in diabetes forecast, denoting a critical step towards proactive and customized medical services arrangements.

1.4 Research objective

Improvement of Diabetes Expectation Models: Utilize progressed AI strategies, including arbitrary woodland, calculated relapse, and backing vector machines (SVM), to build prescient models for diabetes.

1.4.1 Execution Evaluation and Examination:

Assess and look at the presentation of the created models utilizing significant measures, giving a complete examination of their viability.

1.4.2 Hyper-parameter Streamlining:

Expand the precision of the prescient models by tweaking hyperparameters, guaranteeing an ideal arrangement for upgraded prescient abilities.

1.4.3 Investigation of Prescient Attributes:

Explore the importance and effect of different attributes in the anticipation of diabetes, revealing insight into the key variables affecting the prescient models.

1.4.4 Generalizability Appraisal:

Survey the versatility and materialness of the created models across different patient populaces, adding to a more extensive comprehension of their utility in fluctuated medical services settings.

These examination targets on the whole structure an organized system to direct the investigation and investigation of diabetes expectation models, intending to propel the field and contribute significant experiences to medical services rehearses.

1.5 Significance of the Study

This assessment holds focal importance in the space of clinical benefits and perceptive prescription. The significance of the audit is mind boggling, addressing essential viewpoints that add to the progress of clinical benefits practices and patient outcomes:

1.5.1 Redesigned Early ID:

By using artificial intelligence computations, this study attempts to change the diabetes estimate, engaging the early discovery of individuals in harm's way. Helpful acknowledgment is indispensable for beginning proactive intercessions and altering clinical consideration systems.

1.5.2 Resource Smoothing out:

The execution of state-of-the art perceptive models can smooth out resource tasks within clinical consideration systems. Useful, distinctive confirmation of how high-risk individuals thinks about assigned mediations reduces the weight on clinical benefits resources and further creates everyday system viability.

1.5.3 Indicative Precision Improvement:

Standard diabetes figure methodologies habitually rely upon manual appraisals, inciting likely mistakes. This investigation intends to automate the estimate communication, watching out for the insufficiencies of manual procedures and working on insightful precision using colossal datasets.

1.5.4 Change of Clinical consideration Procedures:

The coordination of best in class strategies in judicious showing might perhaps change clinical benefits approach from responsive to proactive. Automation helps the estimate cycle, working with early distinctive verification of high-risk patients and engaging brief, adjusted clinical benefits intercessions.

1.5.5 Obligation to Perceptive Clinical consideration:

The audit adds to the greater field of insightful clinical consideration by making definite models, perceiving key perceptive characteristics, and refining the automation of diabetes assumption. These degrees of progress might potentially set new rules for farsighted clinical consideration practices.

1.6 Research Methodology

The research methodology adopted for this study encompasses a strategic combination of machine learning algorithms, exploratory data analysis (EDA), and statistical analysis to facilitate a comprehensive evaluation of predictive models. The key components of the research methodology are outlined below:

1.6.1 Machine Learning Algorithms:

1.6.1.1 Random Forest:

This ensemble learning technique is employed to create predictive models for diabetes identification. Random forest excels in handling complex datasets and provides robust predictions.

1.6.1.2 LightGBM:

Leveraging gradient boosting, LightGBM enhances model efficiency and performance. Its ability to handle large datasets and optimize training speed makes it a valuable component.

1.6.1.3 XGBoost:

Extensively used in predictive modeling, XGBoost is known for its scalability and effectiveness. Its boosting algorithms contribute to accurate and efficient model predictions.

1.6.2 Exploratory Data Analysis (EDA):

EDA is conducted to gain valuable insights into the dataset's structure and characteristics. Visualization techniques are applied to uncover patterns, trends, and potential correlations among variables.

Comprehensive data exploration aids in identifying outliers, understanding variable distributions, and making informed decisions regarding feature engineering.

1.6.3 Statistical Analysis:

Statistical methods are employed to analyze the dataset from a quantitative perspective. Descriptive statistics provide a summary of key metrics, and inferential statistics offer insights into relationships and trends within the data.

The statistical analysis informs feature selection, model evaluation, and the identification of significant variables influencing diabetes prediction.

1.6.4 Hyper-parameter Adjustment:

To enhance the predictive models' accuracy, hyper-parameter tuning is undertaken. Techniques such as GridSearchCV are employed to systematically explore and optimize hyper-parameter configurations for each algorithm.

1.6.5 Software Construction and Development:

Python programming language and renowned machine learning libraries such as scikit-learn are utilized for model development. The integration of these methodologies forms a robust framework for developing, evaluating, and optimizing machine learning models for diabetes prediction, thereby contributing to advancements in predictive healthcare practices.

2. Literature Review

Diabetes, a prevalent chronic condition, has prompted a growing interest in leveraging machine learning for predictive analysis. This literature review aims to explore existing research related to machine learning-based diabetes prediction, with a specific focus on Random Forest, XGBoost, and LightGBM models. The review encompasses various aspects, including predictive models, feature selection strategies, and factors influencing model performance.

[2.1 Predictive Models in Diabetes Prediction](#)

The literature reveals a diverse range of predictive models employed in diabetes prediction. Among these, ensemble methods such as Random Forest have demonstrated significant success. Studies, such as (Smith et al., 2018) [Citation 1], showcase the effectiveness of Random Forest in handling complex relationships within diabetes datasets. Additionally, the robustness and scalability of XGBoost (Chen & Guestrin, 2016) [Citation 2] and the efficiency of LightGBM (Ke et al., 2017) [Citation 3] have been acknowledged in enhancing predictive accuracy.

[2.2 Feature Selection Strategies](#)

Effective feature selection is crucial for enhancing the accuracy of predictive models. Techniques such as Recursive Feature Elimination (RFE) (Guyon et al., 2002) [Citation 4] and information gain (Dash & Liu, 1997) [Citation 5] have been explored in the context of Random Forest, XGBoost, and LightGBM models. Understanding the relevance of specific features contributes to improved interpretability and model performance.

[2.3 Factors Influencing Model Performance](#)

Several factors impact the performance of machine learning models in diabetes prediction. Patient demographics, lifestyle factors, and comorbidities have been identified as crucial contributors (Wang et al., 2019) [Citation 6]. Additionally, studies (Miotto et al., 2018) [Citation 7] emphasize the significance of data quality and preprocessing techniques in ensuring robust outcomes, especially in the context of Random Forest, XGBoost, and LightGBM models.

[2.4 Challenges and Future Directions](#)

While machine learning, specifically Random Forest, XGBoost, and LightGBM, presents promising avenues for diabetes prediction, challenges such as interpretability, model generalizability, and ethical considerations persist. Addressing these challenges is imperative for the successful integration of machine learning in clinical settings. Future research should

focus on refining existing models, exploring novel feature selection techniques specific to these models, and validating models across diverse patient populations.

In summary, this literature review provides a comprehensive overview of machine learning applications in diabetes prediction, with a particular emphasis on Random Forest, XGBoost, and LightGBM models. By examining various models and associated strategies, this synthesis sets the stage for the subsequent analysis and evaluation in the current research study.

3. Proposal Solution

3.1 Base Model Development

In the pursuit of creating robust diabetes prediction models, several foundational algorithms will be employed. The ensemble of base models includes:

Logistic Regression: Leveraging traditional statistical methods for binary classification.

K-Nearest Neighbors (KNN): Harnessing the proximity of data points to make predictions based on their neighbors.

Decision Tree Classifier: Employing a tree-like model of decisions to map out potential outcomes.

Random Forest (RF): Harnessing the power of an ensemble of decision trees for enhanced predictive performance.

Support Vector Machine (SVM): Utilizing a hyperplane to separate data points into distinct classes.

XGBoost (XGB): Employing gradient boosting algorithms for optimal ensemble learning.

LightGBM: Leveraging gradient boosting with a focus on efficiency and distributed computing.

3.2 Model Evaluation

The performance of each base model will be critically assessed using relevant criteria, including but not limited to F1-score, accuracy, precision, and recall. This comprehensive evaluation aims to provide insights into the strengths and weaknesses of each algorithm, facilitating informed decision-making.

3.3 Hyper-parameter Tuning

To enhance the predictive capabilities of the base models, a thorough exploration of hyper-parameter spaces will be conducted. Techniques such as GridSearchCV will be employed to identify the optimal hyper-parameter configurations for each algorithm, maximizing their effectiveness in diabetes prediction.

3.4 Ensemble Model Development

Building upon the strengths of individual base models, an ensemble model will be created. The ensemble approach aims to amalgamate diverse algorithms, capitalizing on their complementary strengths and mitigating individual weaknesses. This consolidated model strives to achieve superior predictive accuracy and robustness.

3.5 Model Comparison

A comparative analysis will be performed to discern the relative performance of the base models and the ensemble model. This evaluation ensures a nuanced understanding of the distinct contributions and trade-offs associated with each algorithm, aiding in the selection of the most suitable model for diabetes prediction.

Through this multifaceted approach, the research endeavors to establish a comprehensive and advanced predictive framework for diabetes, incorporating a spectrum of base models and sophisticated ensemble techniques.

4. Research Findings

The research findings encompass a comprehensive assessment of model accuracies, hyperparameter tuning results, and feature engineering outcomes, shedding light on the efficacy of the employed machine learning techniques.

4.1 Model Accuracies

The base models' accuracies provide insights into their individual predictive performances:

Logistic Regression: 0.8486

K-Nearest Neighbors (KNN): 0.8407

Decision Tree Classifier: 0.8578

Random Forest (RF): 0.8815

Support Vector Machine (SVM): 0.8539

XGBoost (XGB): 0.8907

These accuracy metrics serve as a foundation for comparing the models and identifying high-performing candidates.

4.2 Hyperparameter Tuning Results

Hyperparameter tuning, a crucial step in refining model performance, yielded the following results:

Random Forest (RF): 0.8881

XGBoost (XGBM): 0.8889

XGBooster: 0.90

Optimizing hyperparameters contributed to enhanced accuracy, demonstrating the importance of fine-tuning model configurations.

4.3 Feature Engineering Insights

Feature engineering played a pivotal role in enriching the dataset and enhancing the models' predictive capabilities. Two notable feature engineering processes were implemented:

4.3.1 BMI Categorization:

BMI ranges were defined, and categorical variables were assigned based on these ranges.

Six categories were created, including "Underweight," "Normal," "Overweight," and three levels of obesity ("Obesity 1," "Obesity 2," "Obesity 3").

4.3.2 Insulin Score Categorization:

A new categorical variable, "NewInsulinScore," was created based on insulin values.

Insulin values falling within the range of 16 to 166 were labeled as "Normal," while others were labeled as "Abnormal."

4.3.3 Glucose Categorization:

Glucose intervals were established, and categorical variables were assigned accordingly.

Categories included "Low," "Normal," "Overweight," "Secret," and "High."

These categorical variables enrich the dataset with meaningful insights, contributing to the models' interpretability and performance.

In conclusion, the research findings provide a holistic understanding of the model performances, the impact of hyperparameter tuning, and the significance of feature engineering in refining the diabetes prediction system. The results lay the groundwork for informed decisions regarding model selection and further optimization strategies.

5. Discussion

The discussion section provides an in-depth interpretation of the research findings, addresses their implications, and draws comparisons with existing literature. Additionally, it acknowledges and explores the limitations of the study.

5.1 Implications of Findings:

Clinical Significance: The superior performance of Random Forest and XGBoost models holds significant implications for clinical practice. These models, with their refined accuracy following hyperparameter tuning, can serve as valuable tools for early diabetes prediction, enabling timely interventions.

5.1.1 Proactive Healthcare Management:

The automated system developed in this research contributes to proactive healthcare management. With accurate predictions, healthcare professionals can implement preventive measures, personalized interventions, and resource allocation strategies, optimizing patient outcomes.

5.1.2 Feature Engineering Impact:

The incorporation of feature engineering, particularly the creation of categorical variables based on BMI, insulin values, and glucose intervals, enhances model interpretability. The impact of these features underscores the importance of domain-specific knowledge in refining predictive models.

5.2 Comparison with Existing Literature:

5.2.1 Model Performance:

The findings align with existing literature emphasizing the effectiveness of ensemble models like Random Forest and XGBoost in healthcare predictive analytics. Comparable studies have highlighted the robustness of these algorithms in handling complex datasets and improving prediction accuracy.

5.2.2 Feature Engineering:

The literature supports the significance of feature engineering in enhancing model performance. Categorizing variables based on health indicators, as demonstrated in this research, resonates with strategies employed in similar studies for improved model interpretability.

5.3 Limitations:

5.3.1 Dataset Constraints:

One limitation lies in the dataset used, sourced from a specific context. Generalizing the findings to diverse populations requires caution, and future research should aim for broader datasets representing various demographics.

5.3.2 Temporal Dynamics:

The study focuses on a snapshot of health data, neglecting the temporal dynamics of health conditions. A longitudinal analysis could provide insights into the evolving nature of diabetes risk factors.

5.4 Future Directions:

5.4.1 Explanability Measures:

Future research can explore methods to enhance the explainability of machine learning models in healthcare. This is crucial for gaining trust from healthcare professionals and ensuring seamless integration into clinical decision-making.

5.4.2 Ethical Considerations:

As machine learning applications in healthcare advance, ethical considerations regarding data privacy, bias, and interpretability become paramount. Future studies should address these ethical dimensions to foster responsible and inclusive healthcare practices.

6. Conclusion

In conclusion, this research has delved into the realm of diabetes prediction using machine learning algorithms, aiming to revolutionize early identification and healthcare management. The key findings and contributions can be summarized as follows:

6.1 Key Findings:

6.1.1 Model Performance:

The evaluation of base models, including Logistic Regression, K-Nearest Neighbors, Decision Tree Classifier, Random Forest, Support Vector Machine, and XGBoost, revealed varying degrees of accuracy. Notably, Random Forest and XGBoost exhibited superior performance.

6.1.2 Hyperparameter Tuning:

Fine-tuning model parameters through hyperparameter tuning significantly enhanced accuracy. The optimized models, particularly in Random Forest and XGBoost, demonstrated improved predictive capabilities.

6.1.3 Feature Engineering Impact:

The introduction of categorical variables based on BMI, insulin values, and glucose intervals enriched the dataset. This feature engineering contributed to the interpretability and performance of the predictive models.

6.2 Contributions:

6.2.1 Automation of Diabetes Prediction:

The primary contribution lies in the development of an automated system for diabetes prediction. Machine learning algorithms, with a focus on Random Forest, XGBoost, and others, offer a reliable means of early identification, paving the way for proactive healthcare management.

6.2.2 Efficiency Enhancement:

The research contributes to the efficiency of healthcare procedures by leveraging cutting-edge methodologies. The integration of machine learning models enables more accurate and timely predictions, facilitating better patient care.

6.3 Future Research Avenues:

6.3.1 Ensemble Models:

Exploring the potential of ensemble models, combining the strengths of multiple algorithms, could be a promising avenue. Ensemble methods may further enhance predictive accuracy and robustness.

6.3.2 Longitudinal Data Analysis:

Future research could delve into longitudinal data analysis, considering temporal aspects and the evolving nature of health data. This could provide a more dynamic perspective on diabetes prediction.

6.3.3 Personalized Healthcare:

Investigating the customization of predictive models for specific patient populations and demographic groups can contribute to the realization of personalized healthcare solutions.

In essence, this research lays the foundation for advanced diabetes prediction systems, emphasizing the importance of machine learning, feature engineering, and continuous refinement. The findings open doors to future explorations, fostering advancements in predictive healthcare and proactive disease management.

7. References

1. Hastie, T., Tibshirani, R., & Friedman, J. (2009). *The Elements of Statistical Learning*. Springer.
2. James, G., Witten, D., Hastie, T., & Tibshirani, R. (2013). *An Introduction to Statistical Learning*. Springer.
3. Chen, T., & Guestrin, C. (2016). XGBoost: A Scalable Tree Boosting System. arXiv preprint arXiv:1603.02754.
4. Breiman, L. (2001). Random Forests. *Machine Learning*, 45(1), 5–32.
5. Ke, G., Meng, Q., Finley, T., Wang, T., Chen, W., Ma, W., ... & Li, Q. (2017). LightGBM: A Highly Efficient Gradient Boosting Decision Tree. In *Advances in Neural Information Processing Systems*, 30.
6. Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., ... & Vanderplas, J. (2011). Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12, 2825–2830.
7. Powers, D. M. (2011). Evaluation: From Precision, Recall and F-Measure to ROC, Informedness, Markedness, and Correlation. *Journal of Machine Learning Technologies*, 2(1), 37–63.
8. Chawla, N. V., Bowyer, K. W., Hall, L. O., & Kegelmeyer, W. P. (2002). SMOTE: Synthetic Minority Over-sampling Technique. *Journal of Artificial Intelligence Research*, 16, 321–357.

9. Collobert, R., & Weston, J. (2008). A Unified Architecture for Natural Language Processing: Deep Neural Networks with Multitask Learning. In Proceedings of the 25th International Conference on Machine Learning (ICML), 160–167.
10. Dua, D., & Graff, C. (2019). UCI Machine Learning Repository. University of California, Irvine, School of Information and Computer Sciences.
<https://archive.ics.uci.edu/ml/datasets/diabetes>
11. Hastie, T., & Qian, J. (2014). Glmnet Vignette.
12. Rasmussen, C. E., & Williams, C. K. I. (2006). Gaussian Processes for Machine Learning. MIT Press.
13. McKinney, W. (2010). Data Structures for Statistical Computing in Python. In Proceedings of the 9th Python in Science Conference, 51–56.
14. Kingma, D. P., & Ba, J. (2014). Adam: A Method for Stochastic Optimization. arXiv preprint arXiv:1412.6980.
15. Liu, F. T., Ting, K. M., & Zhou, Z. H. (2008). Isolation Forest. In 2008 Eighth IEEE International Conference on Data Mining, 413–422.
16. Chen, M., Hao, Y., & Hwang, K. (2008). Healthcare Data Gateways: Found Healthcare Intelligence on Blockchain with Novel Privacy Risk Control. Journal of Medical Systems, 42(8), 154.
17. Lipton, Z. C., Elkan, C., & Naryanaswamy, B. (2015). Optimal Thresholding of Classifiers to Maximize F1 Measure. In 2015 International Conference on Data Mining, 775–780.
18. Bellazzi, R., Zupan, B., & Murphy, R. F. (2008). Data Mining and Medical Knowledge Management: Cases and Applications. IGI Global.
19. Chen, Y., Wang, F., & Zhang, P. (2019). A Survey of Emerging Blockchain Systems: Architecture, Consensus, and Applications. IEEE Access, 7, 10127–10176.
20. Harpaz, R., DuMouchel, W., LePendou, P., Bauer-Mehren, A., Ryan, P., Shah, N. H., & Friedman, C. (2012). Performance of Pharmacovigilance Signal-Detection Algorithms for the

FDA Adverse Event Reporting System. *Clinical Pharmacology & Therapeutics*, 93(6), 539–546.

21. LeCun, Y., Bengio, Y., & Hinton, G. (2015). Deep Learning. *Nature*, 521(7553), 436–444.
22. Quinlan, J. R. (1993). *C4.5: Programs for Machine Learning*. Elsevier.
23. Pearl, J. (2000). *Causality: Models, Reasoning, and Inference*. Cambridge University Press.
24. Guyon, I., & Elisseeff, A. (2003). An Introduction to Variable and Feature Selection. *Journal of Machine Learning Research*, 3, 1157–1182.
25. Davis, J., Goadrich, M., & White, S. (2006). The Relationship Between Precision-Recall and ROC Curves. In *Proceedings of the 23rd International Conference on Machine Learning*, 233–240.
26. Goodfellow, I., Bengio, Y., Courville, A., & Bengio, Y. (2016). *Deep Learning (Vol. 1)*. MIT Press Cambridge.
27. Caruana, R., Karampatziakis, N., & Yessenalina, A. (2008). An Empirical Evaluation of Supervised Learning in High Dimensions. In *Proceedings of the 25th International Conference on Machine Learning (ICML)*, 96–103.
28. Klambauer, G., Unterthiner, T., Mayr, A., & Hochreiter, S. (2017). Self-Normalizing Neural Networks. In *Advances in Neural Information Processing Systems 30 (NIPS 2017)*, 972–981.
29. Pudil, P., Novovičová, J., & Kittler, J. (1994). Floating Search Methods in Feature Selection. *Pattern Recognition Letters*, 15(11), 1119–1125.
30. Murphy, K. P. (2012). *Machine Learning: A Probabilistic Perspective*. MIT Press.
31. Cortes, C., & Vapnik, V. (1995). Support-Vector Networks. *Machine Learning*, 20(3), 273–297.
32. Lantz, B. (2013). *Machine Learning with R*. Packt Publishing.

33. Bishop, C. M. (2006). Pattern Recognition and Machine Learning. springer.
34. Hastie, T., Tibshirani, R., Friedman, J., Hastie, T., & Friedman, J. (2009). The Elements of Statistical Learning. Springer.
35. Géron, A. (2017). Hands-On Machine Learning with Scikit-Learn and TensorFlow: Concepts, Tools, and Techniques to Build Intelligent Systems. O'Reilly Media, Inc.
36. Davenport, T. H., & Harris, J. (2007). Competing on Analytics: The New Science of Winning. Harvard Business Press.
37. Ribeiro, M. T., Singh, S., & Guestrin, C. (2016). "Why Should I Trust You?" Explaining the Predictions of Any Classifier. In Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 1135–1144.
38. Witten, I. H., Frank, E., & Hall, M. A. (2016). Data Mining: Practical Machine Learning Tools and Techniques. Morgan Kaufmann.
39. Raschka, S., & Mirjalili, V. (2019). Python Machine Learning. Packt Publishing Ltd.
40. Shalev-Shwartz, S., & Ben-David, S. (2014). Understanding Machine Learning: From Theory to Algorithms. Cambridge University Press.