

Maternal Health Risk Prediction

A Comparative Study of Machine Learning Model Performance

Mengen Liu, Roy Oh, Kim Tan Palanca, & Nicolas Zhu

Summary

This project aims to use machine learning to predict maternal health risk based on key physiological factors by comparing the performances between two models. Maternal health are a major concern in healthcare. Hence, early risk assessment is crucial for the welfare of both the mother and child. By leveraging structured data from the Maternal Health Risk Dataset, we trained and evaluated machine learning models to classify maternal health risk levels into three different classes: low, mid, and high risk.

The dataset contains 1,014 records with features such as age, blood pressures, blood sugar levels, body temperature and heart rate. Our goal will be to determine whether machine learning algorithms could accurately predict risk levels based on these physiological markers. We will first conduct an exploratory data analysis before building any of the models. Once we have trained the models, their accuracies will be evaluated as measures of performance.

Introduction

Maternal health can be defined as “the health condition of women during pregnancy, childbirth, and the postnatal period (WHO 2025). This is a critical area of healthcare, as complications during pregnancy and childbirth can lead to severe consequences for both mothers and newborns. According to the (WHO 2024), around 800 women died each day in 2020 due to preventable causes related to maternal health, further emphasizing the need for risk assessment measures.

Historically, risk assessment have been carried out by medical professionals that relied heavily on clinical expertise and constant monitoring. However, traditional approaches to monitoring basic physiological indicators often lacked efficiency in identifying potential complications (Mu, Yan, and Zhu 2023). Since the boom of machine learning (ML), many members of the academe have explored the use of ML in maternal health risk prediction, offering data-driven approaches

to enhance early detection and intervention to offload the burden on overworked medical professionals (Bajaj, Kumari, and Bansal 2023; Mu, Yan, and Zhu 2023; Ukrit et al. 2024).

The analysis will utilize the [Maternal Health Risk](#) dataset sourced from the UC Irvine Machine Learning Repository (Ahmed 2020). Consisting of 1014 observations, this dataset includes the following 7 features:

- **Age:** Age of the patient (in years).
- **SystolicBP:** Systolic Blood Pressure (mmHg).
- **DiastolicBP:** Diastolic Blood Pressure (mmHg).
- **BS (Blood Sugar Level):** Blood sugar concentration (mmol/L).
- **BodyTemp:** Body temperature (°F).
- **HeartRate:** Heart rate (beats per minute).
- **RiskLevel:** The target variable, categorized into low risk, mid risk, and high risk.

Project Question

To contribute to this discourse, this research aims to conduct a comparative study on the performance of two ML techniques in predicting maternal health risk, assessing each model's reliability in identifying risk levels. The following research question guides this analysis:

- *Which machine learning model most accurately predicts the maternal health risk level (low, medium, or high risk) based on physiological indicators such as blood sugar levels, body temperature, and other relevant health factors?*

Methods

For this analysis, the data will first be loaded into the notebook then cleaned to handle any possible missing values and ensure its usability for the various models. Following the data cleaning stage will be an exploratory data analysis (EDA) to gain a comprehensive view of the data. This step will include visualizing the summary statistics, distributions, and correlations between variables to determine any patterns in the data prior to the model development.

This study will implement 3 ML classification models:

1. Baseline (Majority Class)
2. Logistic Regression
3. Random Forest

Each model will be evaluated based on the appropriate classification metric to compare their relative performance in maternal health risk prediction.

Table 1: Check for Missing Values

```
# A tibble: 7 x 2
  feature      na
  <chr>      <dbl>
1 Age        0
2 SystolicBP 0
3 DiastolicBP 0
4 BS         0
5 BodyTemp   0
6 HeartRate   0
7 RiskLevel   0
```

Table 2: Distinct Risk Levels

```
# A tibble: 3 x 1
  RiskLevel
  <chr>
1 high risk
2 low risk
3 mid risk
```

Wrangling and Cleaning the Data

From Table 1, we find that there are no missing values in the dataset.

Furthermore, we find that the features `Age`, `SystolicBP`, `DiastolicBP`, `BS`, `BodyTemp`, and `HeartRate` are numeric variables, while `RiskLevel` is currently a character variable. Moreover, Table 2 shows that there are three categories under `RiskLevel`: high risk, mid risk, and low risk. Given the three distinct categories under the target feature, we will modify `RiskLevel` to a factor variable to appropriately reflect its categorical nature in further analysis.

Table 3 shows a snapshot of the cleaned data.

Table 3: Sample of the Cleaned Data

```
# A tibble: 6 x 7
  Age SystolicBP DiastolicBP BS BodyTemp HeartRate RiskLevel
  <dbl>      <dbl>      <dbl> <dbl>      <dbl>      <dbl> <chr>
1    25        130        80 15          98        86 high risk
2    35        140        90 13          98        70 high risk
3    29         90        70 8           100       80 high risk
4    30        140        85 7           98        70 high risk
5    35        120        60 6.1         98        76 low risk
6    23        140        80 7.01        98        70 high risk
```

Table 4: Summary Statistics

```
# A tibble: 6 x 6
  Feature      min    max  mean median    sd
  <chr>      <dbl> <dbl> <dbl> <dbl> <dbl>
1 Age          10     70  29.2    25  13.8
2 SystolicBP   70    160 111.    120  17.9
3 DiastolicBP  49    100  75.4    80  13.8
4 BS           6     19   8.35    7.5  2.83
5 BodyTemp     98    103  98.7    98   1.41
6 HeartRate    7     90  73.9    76   8.16
```

Exploratory Data Analysis

Summary Statistics

Age Distributions

Since age is an important factor in maternal health, we visualize the age distribution by risk level. From the visualization, high risk individuals have a higher median age around 35 years old. Additionally, the interquartile range indicates that the high risk group has more variation in age. We observe some outliers in the low and mid risk groups. Based on the visualization, older aged individuals seem more associated with maternal health risks.

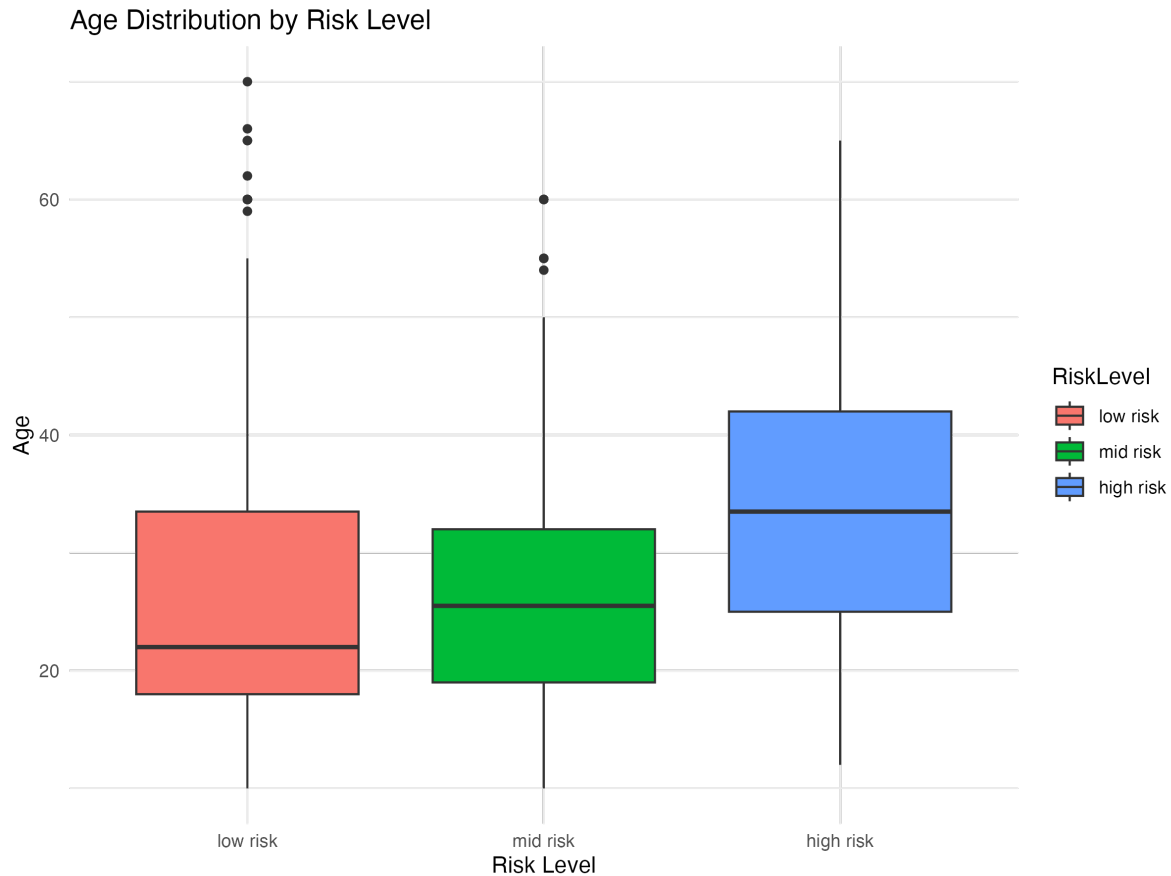


Figure 1: Age Distribution by Risk Level1

Since the target classes seem relatively balanced, it would be appropriate to use **accuracy** as the main scoring metric. Accuracy is given by the number correct prediction out of all predictions made

Correlation Matrix

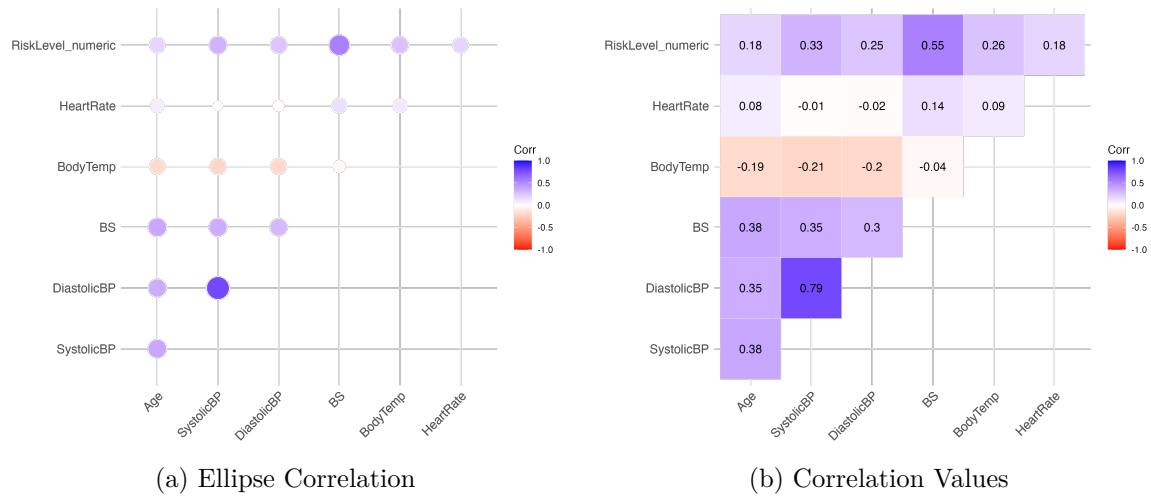


Figure 2: Correlation Matrix2

All of the variables have a positive correlation with RiskLevel, indicating that increases in these variables generally correspond to a higher maternal health risk. BS (Blood Sugar level) has the strongest correlation of 0.7900021, suggesting it is likely to be the most influential factor. We thought age would have a stronger correlation with RiskLevel, however, systolic blood pressure and diastolic blood pressure seems to have a stronger correlation with RiskLevel than age. Additionally, there are no signs of concern for multicollinearity issues in this dataset.

These findings may be a possible reason for the outliers observed above. Younger individuals with high blood pressures or sugar levels may be classified into higher risk levels. This indicates the importance of other factors.

Classification Model Building

Train/Test Splitting

The cleaned data will be split into training and testing sets, with 80% allocated for training and 20% for testing.

Baseline Model (Majority Class)

The baseline model has shown an accuracy of 0.5168539.

Table 5: Multinomial Logistic Regression Model

Call:

```
nnet::multinom(formula = RiskLevel ~ ., data = train_data, trace = FALSE)
```

Coefficients:

	(Intercept)	Age	SystolicBP	DiastolicBP	BS	BodyTemp
mid risk	-52.00101	-0.003909453	0.06037105	-0.03544363	0.4016684	0.4357209
high risk	-88.06475	-0.018588524	0.04783398	0.03051463	0.8807465	0.7050298

	HeartRate
mid risk	0.01943005
high risk	0.04539915

Residual Deviance: 535.76

AIC: 563.76

Multinomial Logistic Regression

The MLR model is as follows:

Model Testing

The multinomial logistic regression has only given us an slightly better accuracy than the baseline with a score of 0.7078652.

Lastly, we plot multinomial logistic regression graphs to visualize how predicted probabilities changes among different variable levels.

We observe from Figure 3 that as blood sugar level rises the probability of high risk increases, while mid and low risk decreases.

Random Forest

The random forest model is as follows:

Here are the parameters that have been passed through the randomForest object:

- `RiskLevel ~ .`: Predicts RiskLevel based on all other features
- `ntree = 500`: Uses 500 trees in the forest
- `importance = TRUE`: Computes feature importance

Table 6: Exponentiated coefficients to transform log-odds to odds ratio

```
# A tibble: 14 x 6
```

	y.level <chr>	term <chr>	estimate <dbl>	std.error <dbl>	statistic <dbl>	p.value <dbl>
1	mid risk	(Intercept)	2.61e-23	0.000346	-150195.	0
2	mid risk	Age	9.96e- 1	0.0116	-0.338	7.36e- 1
3	mid risk	SystolicBP	1.06e+ 0	0.0136	4.43	9.42e- 6
4	mid risk	DiastolicBP	9.65e- 1	0.0171	-2.07	3.85e- 2
5	mid risk	BS	1.49e+ 0	0.153	2.63	8.45e- 3
6	mid risk	BodyTemp	1.55e+ 0	0.0208	21.0	1.83e- 97
7	mid risk	HeartRate	1.02e+ 0	0.0182	1.07	2.85e- 1
8	high risk	(Intercept)	5.67e-39	0.000264	-333932.	0
9	high risk	Age	9.82e- 1	0.0165	-1.13	2.59e- 1
10	high risk	SystolicBP	1.05e+ 0	0.0182	2.63	8.54e- 3
11	high risk	DiastolicBP	1.03e+ 0	0.0230	1.33	1.84e- 1
12	high risk	BS	2.41e+ 0	0.156	5.65	1.63e- 8
13	high risk	BodyTemp	2.02e+ 0	0.0254	27.8	5.00e-170
14	high risk	HeartRate	1.05e+ 0	0.0235	1.93	5.38e- 2

Table 7: Multinomial Logistic Regression Predicted Probabilities

```
# A tibble: 89 x 5
```

	ID	Predicted_Class <chr>	`low risk` <dbl>	`mid risk` <dbl>	`high risk` <dbl>
1	1	high risk	0.00168	0.0461	0.952
2	2	high risk	0.0102	0.114	0.876
3	3	low risk	0.730	0.147	0.122
4	4	high risk	0.0000465	0.00390	0.996
5	5	mid risk	0.421	0.500	0.0782
6	6	low risk	0.855	0.0918	0.0534
7	7	low risk	0.546	0.373	0.0813
8	8	low risk	0.762	0.166	0.0718
9	9	low risk	0.653	0.267	0.0804
10	10	low risk	0.914	0.0796	0.00624

```
# i 79 more rows
```

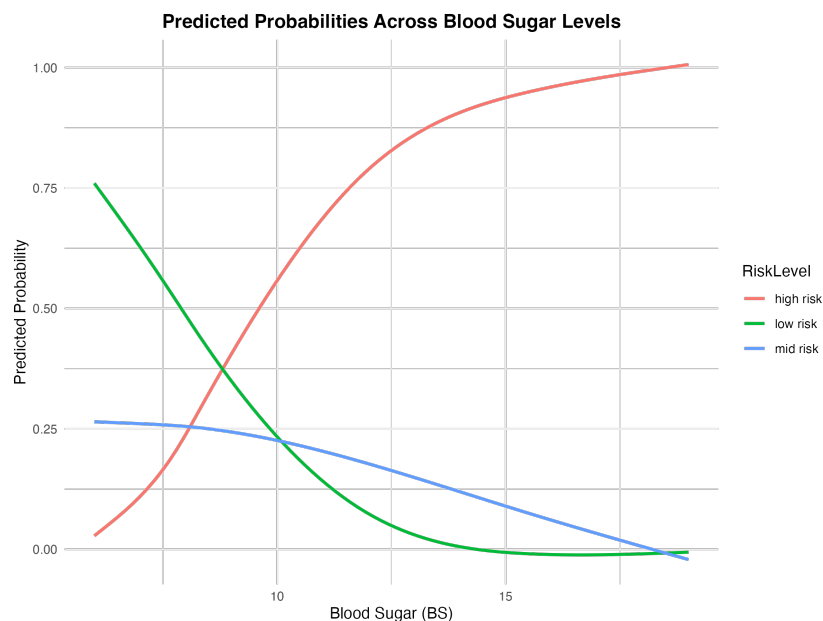



Figure 3: Predicted Probabilities Across Blood Sugar Levels3

Table 8: Random Forest Model

Call:

```
randomForest(formula = RiskLevel ~ ., data = train_data, ntree = 500, importance = TRUE)
```

Type of random forest: classification

Number of trees: 500

No. of variables tried at each split: 2

OOB estimate of error rate: 31.13%

Confusion matrix:

	high risk	low risk	mid risk	class.error
high risk	65	12	13	0.2777778
low risk	10	160	18	0.1489362
mid risk	13	47	25	0.7058824

Model Testing

Now that the model is trained using our train set, we can make predictions on the test set. We find that the random forest model provides a better accuracy score than the multinomial regression model with an accuracy of 0.6516854. Moreover, we can assess feature importances with a random forest model.

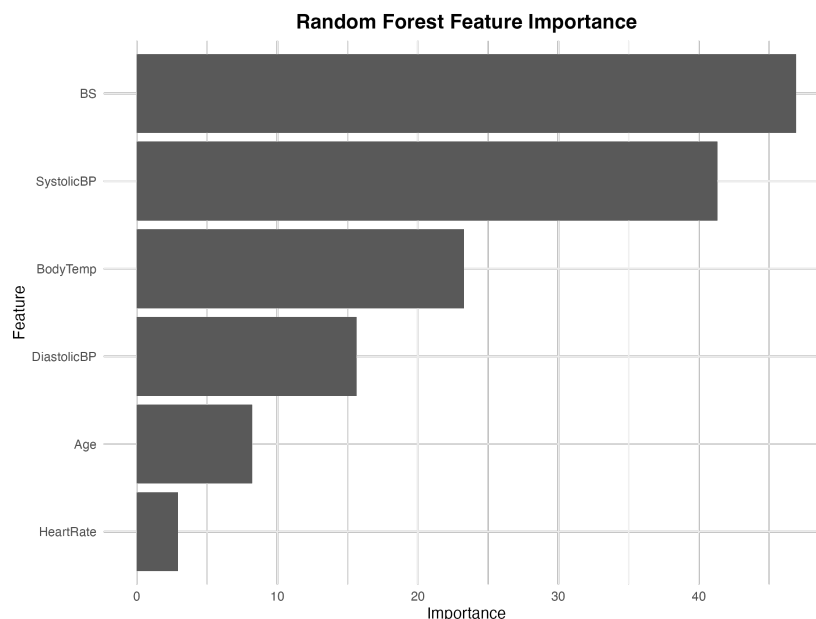


Figure 4: Random Forest Feature Importances4

Based on Figure 4 above, we see that the model identifies blood sugar, systolic blood pressure, and body temperature as the top predictors of maternal health risk (i.e., have the most predictive power). With blood sugar specifically, we see that its feature importance reaches over 100 indicating that it is highly influential in predicting maternal health risk compared to the other features.

Results

Comparison of Results

From Figure 6, random forest has yielded a 0.6516854 accuracy score, which is -0.0561798 higher than the multinomial logistic regression accuracy score and 0.1348315 higher than of the baseline model.

Table 9: Baseline Confusion Matrix

```
# A tibble: 9 x 4
  True Predicted Frequency Percentage
  <chr> <chr>      <dbl>      <dbl>
1 high risk high risk      0         NA
2 low risk  high risk     22        24.7
3 mid risk  high risk      0         NA
4 high risk low risk      0         NA
5 low risk  low risk     46        51.7
6 mid risk  low risk      0         NA
7 high risk mid risk      0         NA
8 low risk  mid risk     21        23.6
9 mid risk  mid risk      0         NA
```

Table 10: Multinomial Logistic Regression Confusion Matrix

```
# A tibble: 9 x 4
  True Predicted Frequency Percentage
  <chr> <chr>      <dbl>      <dbl>
1 high risk high risk     16        72.7
2 low risk  high risk      4         6.7
3 mid risk  high risk      2        28.6
4 high risk low risk      0         0
5 low risk  low risk     44        73.3
6 mid risk  low risk      2        28.6
7 high risk mid risk      6        27.3
8 low risk  mid risk     12         20
9 mid risk  mid risk      3        42.9
```

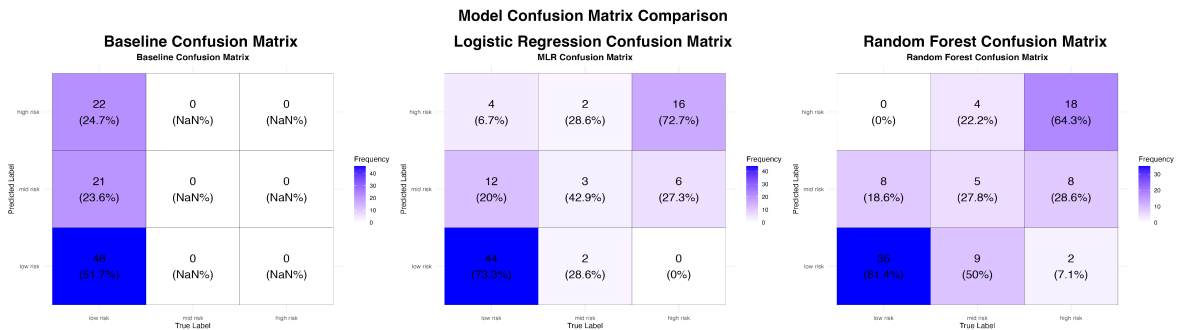


Figure 5: ML Model Confusion Matrices

Table 11: Random Forest Confusion Matrix

A tibble: 9 x 4

	True <chr>	Predicted <chr>	Frequency <dbl>	Percentage <dbl>
1	high risk	high risk	18	64.3
2	low risk	high risk	0	0
3	mid risk	high risk	4	22.2
4	high risk	low risk	2	7.1
5	low risk	low risk	35	81.4
6	mid risk	low risk	9	50
7	high risk	mid risk	8	28.6
8	low risk	mid risk	8	18.6
9	mid risk	mid risk	5	27.8

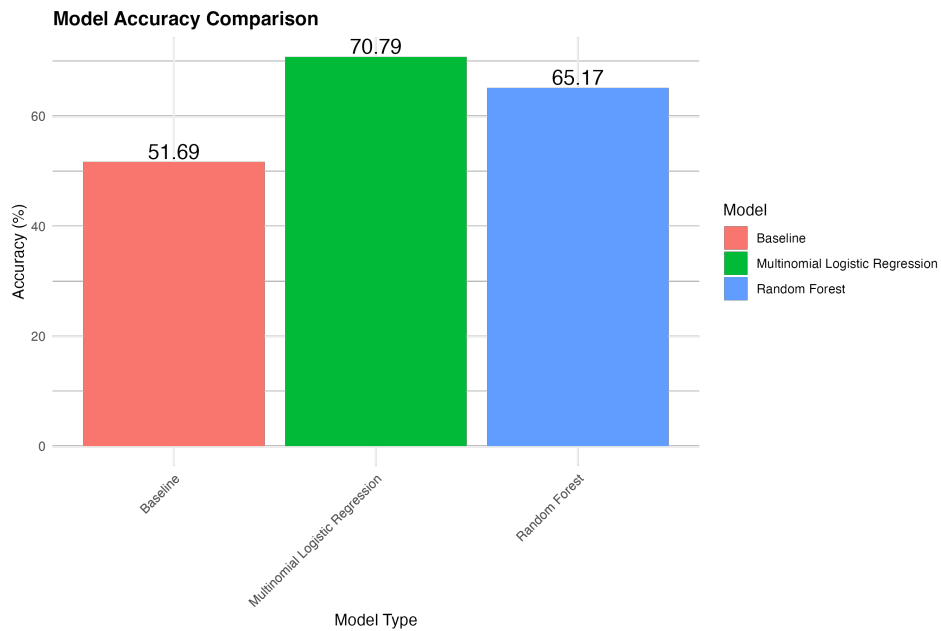


Figure 6: ML Model Accuracies (%)6

Discussion

Best performing model

As shown in the bar graph above, our analyses suggest that maternal health risk during pregnancy can be predicted with up to a 65.17% accuracy using our random forests decision trees model. The second best model was the multinomial logistic regression that had a 70.79% accuracy. Both models performed better than the baseline (51.69%). Both models performing better than the baseline is expected as well as the random forest performing better than the logistic regression as past literature have found similar results (Mu, Yan, and Zhu 2023; Ukrit et al. 2024).

Interpretation

The multinomial logistic regression, while less accurate, is more interpretable and gives us an idea of how a 1-unit increase in a variable is associated with a change in the odds of being in a certain risk category. For example, the multinomial logistic regression suggests that a 1 unit increase in Blood Sugar level is associated with an increase in the odds of being in **high risk** compared to **low risk** by a factor of 2.4127001.

The best predictors for **high risk** compared to **low risk** were body temperature ($OR = 2.0239071$, $p < .001$) and blood sugar ($OR = 2.4127001$, $p < .001$). A one unit increase in both was associated with a more than double increase in the odds of being in the **high risk** category.

The best predictors for **medium risk** compared to **low risk** were also body temperature ($OR = 1.5460773$, $p < .001$) and blood sugar ($OR = 1.4943157$, $p < .001$).

Both the multinomial logistic regression model and the Random Forest model performed best when predicting **high risk**.

Impact

Our analyses show that body temperature and blood sugar are both relatively strongly associated with increasing maternal health risk. We do however acknowledge that our models do not necessarily imply a causal effect such that reducing blood sugar or body temperature will reduce your maternal health risk. Additionally, our models are limited by the number of variables we accounted for. Other factors such as age, parental health conditions, and many more would improve the generalizability of our models.

Our analyses should therefore not be used as guidelines for pregnant mothers. Now that we have further evidence that blood sugar and body temperature are associated with maternal health risk, future research could explore the potential causal mechanism of these

relationships. Future research may also explore whether this effect remains constant across age or whether certain age groups are more susceptible to the effects of blood sugar/body temperature on maternal health risk.

References

- Ahmed, Marzia. 2020. “Maternal Health Risk.” UCI Machine Learning Repository.
- Bajaj, Disha, Ritika Kumari, and Poonam Bansal. 2023. “Risk Level Prediction for Maternal Health Using Machine Learning Algorithms.” *2023 International Conference on Communication, Security and Artificial Intelligence (ICCSAI)*, November, 405–9. <https://doi.org/10.1109/iccsai59793.2023.10421156>.
- Mu, Chenyu, Zexuan Yan, and Yidi Zhu. 2023. “Prediction of Maternal Health Risk Based on Physiological Indicators.” *Proceedings of the 2023 4th International Symposium on Artificial Intelligence for Medicine Science*, October, 578–84. <https://doi.org/10.1145/3644116.3644212>.
- Ukrit, M.Ferni, R. Beaulah Jeyavathana, Aluru Leela Rani, and Vasa Chandana. 2024. “Maternal Health Risk Prediction with Machine Learning Methods.” *2024 Second International Conference on Emerging Trends in Information Technology and Engineering (ICETITE)*, February, 1–9. <https://doi.org/10.1109/ic-etite58242.2024.10493737>.
- WHO. 2024. “Maternal Mortality Fact Sheet.” World Health Organization. <https://www.who.int/news-room/fact-sheets/detail/maternal-mortality>.
- . 2025. “Maternal Health.” World Health Organization. https://www.who.int/health-topics/maternal-health#tab=tab_1.

Figure Scripts

1. Generated by `scripts/04-eda.R`
2. Generated by `scripts/04-eda.R`
3. Generated by `scripts/11-graph-mlr.R`
4. Generated by `scripts/14-rf_feature_importances.R`
5. Generated by `scripts/18-model_comparison.R`
6. Generated by `scripts/18-model_comparison.R`