Emmanuel Pedernal

## 1. Create dummy variables for categorical data before including in the model.

JASP automatically does this e.g. Feature_name (Feature_category) Sex(Male)

*Coefficients* ▼

| Model | | Estimate | Standard Error | Standardized" | Odds Ratio | z | Wald Test Wald Statistic | df | p | 95% Confidence interval Lower bound | Upper bound |
|---|---|---|---|---|---|---|---|---|---|---|---|
| $M_0$ | (Intercept) | 0.053 | 0.062 | 0.053 | 1.054 | 0.843 | 0.711 | 1 | 0.399 | −0.070 | 0.175 |
| $M_1$ | (Intercept) | 6.678 | 2.428 | 6.310 | 794.340 | 2.750 | 7.561 | 1 | 0.006 | 1.918 | 11.437 |
| | age | 0.027 | 0.014 | 0.244 | 1.027 | 1.924 | 3.704 | 1 | 0.054 | −0.000 | 0.054 |
| | resting_blood_pressure | −0.025 | 0.007 | −0.438 | 0.975 | −3.821 | 14.600 | 1 | < .001 | −0.038 | −0.012 |
| | cholestoral | −0.005 | 0.002 | −0.282 | 0.995 | −2.367 | 5.604 | 1 | 0.018 | −0.010 | −0.001 |
| | Max_heart_rate | 0.022 | 0.007 | 0.499 | 1.022 | 3.324 | 11.051 | 1 | < .001 | 0.009 | 0.034 |
| | oldpeak | −0.403 | 0.132 | −0.474 | 0.668 | −3.053 | 9.318 | 1 | 0.002 | −0.662 | −0.144 |
| | sex (Male) | −1.992 | 0.314 | −1.992 | 0.136 | −6.341 | 40.208 | 1 | < .001 | −2.608 | −1.377 |

## 2. Generate and discuss the following:

-Table of coefficients with odds ratios and associated p-values.

Backward pass and Forward pass give high p value (fail to reject H0) while Enter method reject H0

H0: Fits the data just as well as the more complex model

H1: M1 is better model

## Backward

*Model Summary - target* ▼

| Model | Deviance | AIC | BIC | df | ΔX² | p | McFadden R² | Nagelkerke R² | Tjur R² | Cox & Snell R² |
|---|---|---|---|---|---|---|---|---|---|---|
| $M_0$ | 606.815 | 652.815 | 766.261 | 1002 | | | 0.000 | 0.000 | 0.641 | 0.000 |
| $M_1$ | 608.244 | 652.244 | 760.758 | 1003 | 1.429 | 0.232 | −0.002 | −0.003 | 0.640 | −0.001 |
| $M_2$ | 612.009 | 652.009 | 750.658 | 1005 | 3.765 | 0.152 | −0.008 | −0.011 | 0.638 | −0.005 |

## Forward

| Model | Deviance | AIC | BIC | df | ΔX² | p | McFadden R² | Nagelkerke R² | Tjur R² | Cox & Snell R² |
|---|---|---|---|---|---|---|---|---|---|---|
| M₀ | 1420.240 | 1422.240 | 1427.173 | 1024 | | | 0.000 | | 0.000 | |
| M₁ | 1123.995 | 1131.995 | 1151.724 | 1021 | 296.246 | < .001 | 0.209 | 0.335 | 0.274 | 0.251 |
| M₂ | 933.739 | 949.739 | 989.199 | 1017 | 190.255 | < .001 | 0.343 | 0.504 | 0.410 | 0.378 |
| M₃ | 777.967 | 799.967 | 854.224 | 1014 | 155.773 | < .001 | 0.452 | 0.621 | 0.530 | 0.466 |
| M₄ | 714.041 | 740.041 | 804.162 | 1012 | 63.926 | < .001 | 0.497 | 0.664 | 0.575 | 0.498 |
| M₅ | 678.727 | 706.727 | 775.781 | 1011 | 35.314 | < .001 | 0.522 | 0.687 | 0.595 | 0.515 |
| M₆ | 655.001 | 685.001 | 758.988 | 1010 | 23.725 | < .001 | 0.539 | 0.702 | 0.612 | 0.526 |
| M₇ | 640.753 | 672.753 | 751.672 | 1009 | 14.248 | < .001 | 0.549 | 0.710 | 0.621 | 0.533 |
| M₈ | 628.941 | 662.941 | 746.792 | 1008 | 11.813 | < .001 | 0.557 | 0.717 | 0.629 | 0.538 |
| M₉ | 621.017 | 657.017 | 745.801 | 1007 | 7.923 | 0.005 | 0.563 | 0.722 | 0.633 | 0.541 |
| M₋ | 615.217 | 653.217 | 746.934 | 1006 | 5.800 | 0.016 | 0.567 | 0.726 | 0.635 | 0.544 |
| M₋ | 612.009 | 652.009 | 750.658 | 1005 | 3.208 | 0.073 | 0.569 | 0.727 | 0.638 | 0.545 |

**Enter**

*Model Summary - target ▼*

| Model | Deviance | AIC | BIC | df | ΔX² | p | McFadden R² | Nagelkerke R² | Tjur R² | Cox & Snell R² |
|---|---|---|---|---|---|---|---|---|---|---|
| M₀ | 1420.240 | 1422.240 | 1427.173 | 1024 | | | 0.000 | | 0.000 | |
| M₁ | 606.815 | 652.815 | 766.261 | 1002 | 813.425 | < .001 | 0.573 | 0.731 | 0.641 | 0.548 |

*Note.* M₁ includes age, resting_blood_pressure, cholestoral, Max_heart_rate, oldpeak, sex, chest_pain_type, fasting_blood_sugar, rest_ecg, exercise_induced_angina, slope, vessels_colored_by_flourosopy, thalassemia

For this model we'll use the following as reference for each feature (will be dropped to avoid collinearity)

| Feature | Reference (to be drop) |
|---|---|
| Sex | Female |
| Chest_pain | Asymptomatic |
| Fasting | Greater than 120mg |
| Rest_ecg | Normal |
| Exercise_induced | No |
| Slope | Flat |
| Flouroscopy | Four |
| Thalassemia | Normal |

*Coefficients*

| Model | | Estimate | Standard Error | Standardized* | Odds Ratio | z | Wald Test | | |
|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | Wald Statistic | df | p |
| M₀ | (Intercept) | 0.053 | 0.062 | 0.053 | 1.054 | 0.843 | 0.711 | 1 | 0.399 |
| M₁ | (Intercept) | 6.475 | 1.891 | 6.107 | 648.582 | 3.424 | 11.724 | 1 | < .001 |
| | age | 0.027 | 0.014 | 0.244 | 1.027 | 1.924 | 3.704 | 1 | 0.054 |
| | resting_blood_pressure | −0.025 | 0.007 | −0.438 | 0.975 | −3.821 | 14.600 | 1 | < .001 |
| | cholestoral | −0.005 | 0.002 | −0.282 | 0.995 | −2.367 | 5.604 | 1 | 0.018 |
| | Max_heart_rate | 0.022 | 0.007 | 0.499 | 1.022 | 3.324 | 11.051 | 1 | < .001 |
| | oldpeak | −0.403 | 0.132 | −0.474 | 0.668 | −3.053 | 9.318 | 1 | 0.002 |
| | sex (Male) | −1.992 | 0.314 | −1.992 | 0.136 | −6.341 | 40.208 | 1 | < .001 |
| | chest_pain_type (Atypical angina) | −1.523 | 0.442 | −1.523 | 0.218 | −3.450 | 11.903 | 1 | < .001 |
| | chest_pain_type (Non-anginal pain) | −0.403 | 0.383 | −0.403 | 0.668 | −1.052 | 1.106 | 1 | 0.293 |
| | chest_pain_type (Typical angina) | −2.410 | 0.392 | −2.410 | 0.090 | −6.148 | 37.795 | 1 | < .001 |
| | fasting_blood_sugar (Lower than 120 mg/ml) | −0.380 | 0.320 | −0.380 | 0.684 | −1.189 | 1.414 | 1 | 0.234 |
| | rest_ecg (Left ventricular hypertrophy) | −0.800 | 1.537 | −0.800 | 0.449 | −0.521 | 0.271 | 1 | 0.603 |
| | rest_ecg (ST-T wave abnormality) | 0.397 | 0.218 | 0.397 | 1.488 | 1.823 | 3.322 | 1 | 0.068 |
| | exercise_induced_angina (Yes) | −0.750 | 0.249 | −0.750 | 0.472 | −3.016 | 9.099 | 1 | 0.003 |
| | slope (Downsloping) | 1.395 | 0.272 | 1.395 | 4.036 | 5.133 | 26.345 | 1 | < .001 |
| | slope (Upsloping) | 0.596 | 0.472 | 0.596 | 1.814 | 1.262 | 1.592 | 1 | 0.207 |
| | vessels_colored_by_flourosopy (One) | −3.900 | 0.966 | −3.900 | 0.020 | −4.036 | 16.288 | 1 | < .001 |
| | vessels_colored_by_flourosopy (Three) | −3.854 | 1.055 | −3.854 | 0.021 | −3.651 | 13.333 | 1 | < .001 |
| | vessels_colored_by_flourosopy (Two) | −5.163 | 1.052 | −5.163 | 0.006 | −4.908 | 24.092 | 1 | < .001 |
| | vessels_colored_by_flourosopy (Zero) | −1.566 | 0.930 | −1.566 | 0.209 | −1.683 | 2.833 | 1 | 0.092 |
| | thalassemia (Fixed Defect) | −0.392 | 0.442 | −0.392 | 0.676 | −0.888 | 0.788 | 1 | 0.375 |
| | thalassemia (No) | −2.797 | 1.466 | −2.797 | 0.061 | −1.908 | 3.639 | 1 | 0.056 |
| | thalassemia (Reversable Defect) | −1.806 | 0.436 | −1.806 | 0.164 | −4.145 | 17.182 | 1 | < .001 |

*Note.* target level '1' coded as class 1.

* Standardized estimates represent estimates where the continuous predictors are standardized (X-standardization).

Based from the Coefficients

**H0: Feature has no effect on the outcome**

**H1: Feature does have effect**

| Feature | Insights |
|---|---|
| age | Age has 2.7% increase in odds of having disease although Age has 0.054 p value, we'll consider it due medical field where age has positive correlation with heart diseases |
| resting_blood_pressure | is associated with a 2.5% decrease in odds of disease |
| cholestoral | Higher cholesterol slightly decreases odds by 0.5% per unit this should be reviewed since it does not make sense that a high cholesterol is better hence greater than 5% p value |
| Max_heart_rate | Each additional unit increases the odds by 2.2%. With a p value less than 5% means this feature have an effect in determining if a patient have disease |
| oldpeak | Higher oldpeak slightly decreases odds by 33% per unit this should be reviewed since it |

| | does not make sense that a high oldpleak is better hence greater than 5% p value. |
|---|---|
| sex (Male) | The odds of having heart disease for Male is lower 13.6% compared to women 86.4% |
| chest_pain_type | Atypical and Typical angina significantly lower odds compared to Asymptomatic (Reference) |
| fasting_blood_sugar (Lower than 120 mg/ml) | Has lower odds of having (32%) disease vs 68% of greater than 120mg/ml fasting blood sugar. But has greater than 5% p val meaning we have weak evidence to tell if this is by chance or not (small dataset) |
| rest_ecg (LVH) | About 55% lower odds of heart disease compared to those with a normal ECG, but this is not statistically significant high p val. So, no strong evidence LVH is associated with heart disease risk in this model. |
| rest_ecg (ST-T wave abnormality) | About 49% higher odds of heart disease compared to those with a normal ECG. High p value also so no strong evidence with heart disease for this dataset |
| exercise_induced_angina (Yes) | About 53% lower odds of having heart disease compared to those without exercise-induced angina (No). With p val lower than 5%, exercise-induced angina is associated with lower odds of heart disease in this model. |
| slope (Downsloping) | odds ratio of 4.036 indicates that patients with a downsloping slope have about 4 times higher odds of having heart disease. P val less than 5% means this is significant. |
| slope (Upsloping) | odds ratio of 1.814 suggests these patients have about 1.8 times higher odds of heart disease relative to the reference. With high p value the evidence is not strong enough |
| vessels_colored_by_flourosopy vs Four | having 1, 2, or 3 vessels colored is associated with much lower odds (0.006) of heart disease compared to 4 vessels.<br><br>Zero vessels also seem to have lower odds, but this result is not quite statistically significant but is insignificant since it has high p val. |

| | |
|---|---|
| thalassemia (Fixed Defect) | about 32.4% lower odds of having heart disease compared to Normal thalassemia. difference is not statistically significant high p val |
| thalassemia (No) | about 93.9% lower odds of heart disease compared to normal thalassemia. With p val close to 5% we can say that this is significant |
| thalassemia (Reversable Defect) | about 83.6% lower odds of having heart disease compared to normal thalassemia patients low p val meaning there is strong evidence in this dataset that this feature has lower odds compared to Normal thalassemia |

Based from the dataset a Higher max heart rate, slop(downsloping) and rest_ecg (ST-T wave abnormality) increases the odds of heart disease, while higher resting blood pressure, cholesterol, and oldpeak decreased the odds (might need further investigation because some metrics are inverse of real-life application). Males, those with typical or atypical angina, and patients with multiple-colored vessels and reversible thalassemia defects also had lower odds compared to female counter parts.

**-Model evaluation metrics such as accuracy, precision, recall, and other metrics.**

### Performance Diagnostics

*Confusion matrix*

| Observed | Predicted 0 | 1 | % Correct |
|---|---|---|---|
| 0 | 420 | 79 | 84.168 |
| 1 | 42 | 484 | 92.015 |
| Overall % Correct | | | 88.195 |

*Note. The cut-off value is set to 0.5*

*Performance metrics*

| | Value |
|---|---|
| Accuracy | 0.882 |
| AUC | 0.946 |
| Sensitivity | 0.920 |
| Specificity | 0.842 |
| Precision | 0.860 |
| F-measure | 0.889 |
| Brier score | 0.087 |
| H-measure | 0.707 |

Since this is for medical application, the best metric here would be **sensitivity** or true positive rate it's best to capture actual positive patients.

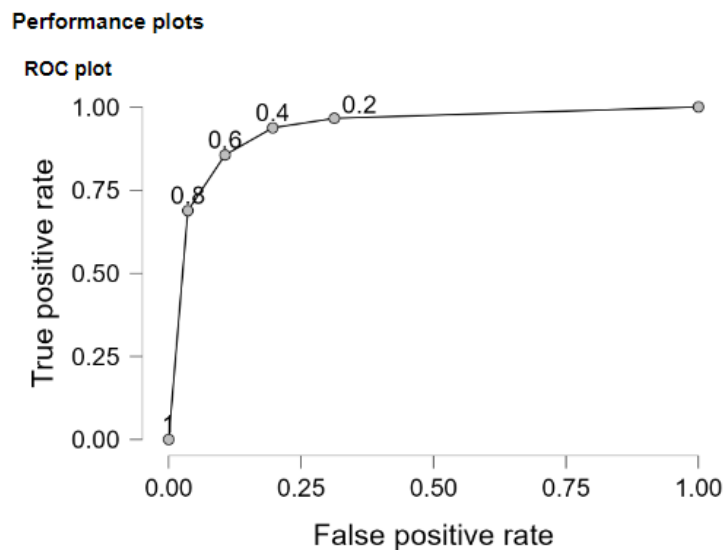| |
|---|
| **True Positives (TP)** = 484 |
| **True Negatives (TN)** = 420 |
| **False Positives (FP)** = 79 |
| **False Negatives (FN)** = 42 |

**Odds ratio**

**TP * TN / FP * FN =**

**484 * 420 / 79 * 42 = 61.27**

The odds of the model correctly predict (true positive or true negative) are 61.27 times higher than the odds of an incorrect prediction (false positive or false negative).

**ROC curve.**

Performance plots

ROC plot



**Thresholds**

0.8 – Low false positives but will miss many positives

0.6 – trade-off between sensitivity/specificity

**0.4** – Has more positives but an increase in false positives

0.2 – highest positives detected but has more false positives

**For medical application we choose 0.4 Threshold for screening cases while 0.2 for high-risk population**

**With a AUC score of 94.6 the model distinguishes the classes well**

- High TPR (Sensitivity = 0.920)
- Low FPR (1 - Specificity = 0.158)