
A Case Study of US Mortgage Approvals

Max Melchior Langg



University of Oxford

Oxford, 2023

Contents

1	Introduction	2
2	Exploratory Data Analysis	2
3	Modelling	3
3.1	Model selection	3
3.2	Model diagnostics	4
3.3	Results & Interpretation	4
3.3.1	Marginal Effects	7
3.4	Dispersion	7
4	Limitations & Outlook	8
5	Conclusion	8
A	Supplementary EDA Tables	9
B	Modelling	10
B.1	Model summary of initial model	10
C	Outlier Table	11
C.1	Soft Threshold	11
D	Further Modelling Approaches	11
D.1	Oversampling	11
D.2	Weighted GLM	11
D.3	LASSO Regression and n-fold Cross-Validation	12
E	R Code	13

1 Introduction

This dataset assesses mortgage applications from a (unknown) 1990 U.S. city, focusing on the binary outcome of approval of the mortgage application. It includes financial ratios such as housing expenses to income (**hir**), other debts to income (**odir**), loan-to-value ratio (**lvr**), and a mortgage credit score (**lvr**) ranging from 1 (best) to 4. Demographic details cover self-employment (**self**), marital status (**single**), and **ethnicity** (white or black), alongside the 1989 state unemployment rate in the applicant's industry (**uria**). The analysis aims to deeply understand factors affecting mortgage approvals.

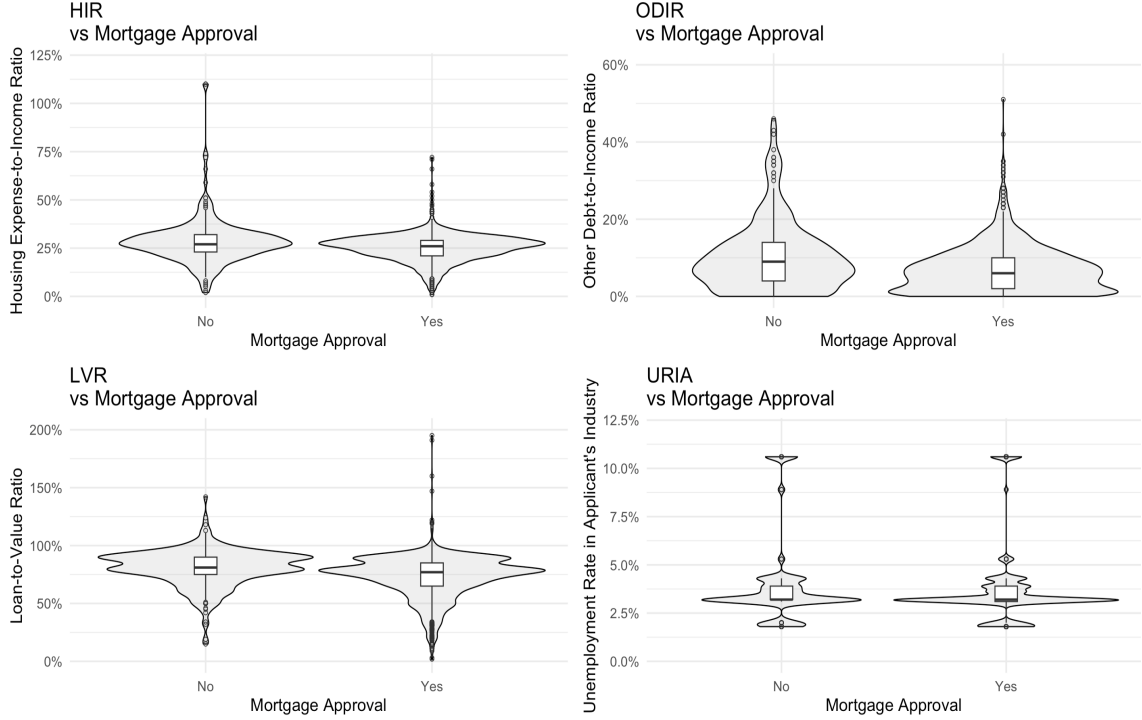


Figure 1: Violin and Box Plot Distribution of Financial Ratios and Unemployment Rate by Mortgage Approval Status

2 Exploratory Data Analysis

The dataset analyzed in this report is a collection of 1902 mortgage applications, with no missing entries, skewed towards approved outcomes (87.5%). For better interpretation, the ratios within the dataset have been transformed to a percentage scale. While the range of values for variables like the Housing Expense to Income Ratio (HIR) and Loan-to-Value ratio (LVR) is broad, they remain within plausible limits (see figure 1). The data reflects a predominance of white, non-single applicants, a lesser proportion of self-employed individuals, and a concentration of applicants with good credit scores. The HIR across applicants averages at 25.5%, with a wider spread of values beyond the interquartile range, suggesting varied housing expenses relative to income among applicants. The Other Debt-to-Income Ratio (ODIR) is skewed, with a median higher for declined (9%) than approved applications (6%). LVR, indicating borrowing risk, also displays asymmetry, with a lower median for approved loans (77%). Outliers on both ends of the spectrum do not imply data issues but rather high variability in loan amounts relative to property values. Similarly, the Unemployment Rate in Industry of Applicant (URIA) is comparable for both approved and denied applications (median for both 3.2 %), with only ten unique values within the dataset. Creditworthiness, as indicated by the Mortgage Credit Score (MCS), shows most applicants in the low-risk categories (1 and 2). Specifically, applicants with the most favorable credit score of 1 have the highest approval rate (92.6%), which marginally decreases with each subsequent credit score category, except from 3 to 4 (2: 85.5%, 3: 75.8%, 4: 78.9%) (see figure 2). Self-employed applicants have a lower approval rate (81.5%) than those not self-employed (88.3%), and single

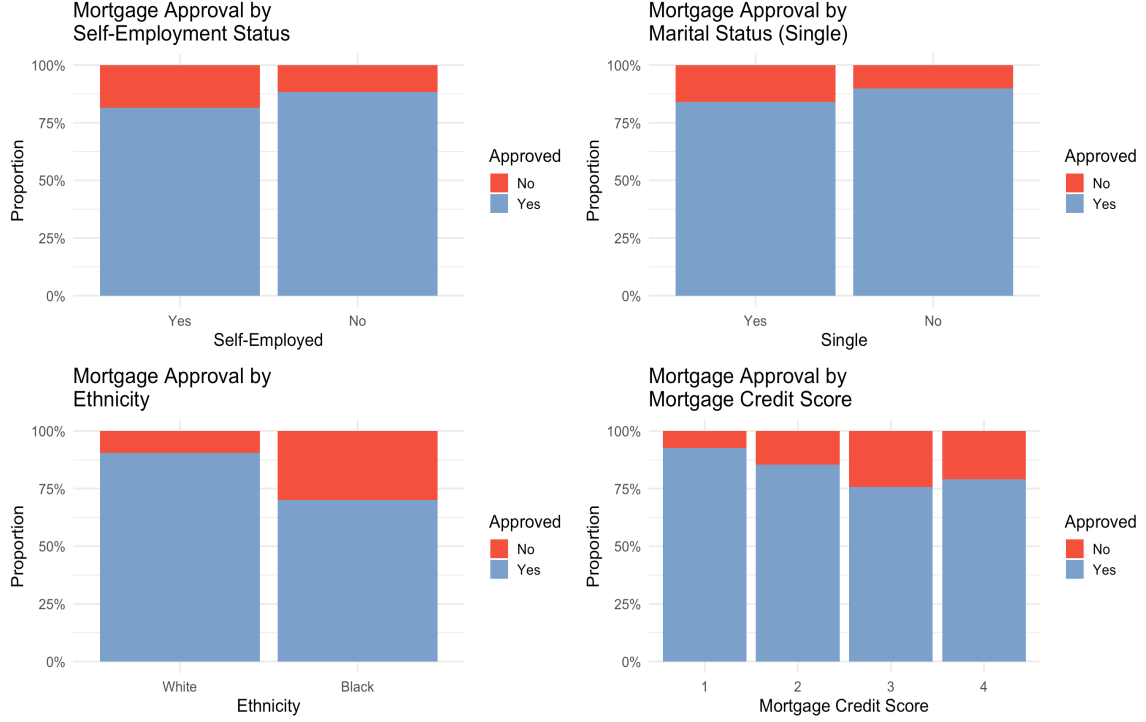


Figure 2: Stacked Bar Charts of Mortgage Approval Rates by Demographic and Financial Criteria.

applicants have a lower approval rate (84.1%) compared to those who are not single (89.8%). Notably, there’s a stark contrast in approval rates between white applicants (90.5%) and non-white applicants (70.0%), highlighting potential disparities. The distribution of mortgage approval rates across different financial indicators shows that variables like HIR, LVR, and ODIR display variable patterns, with approval rates dropping at higher deciles (see figure 3). The ODIR plot demonstrates consistent approval rates across most deciles, with a modest downturn in the highest decile. HIR is characterized by a dynamic range, with approval rates peaking in the first decile and decreasing for higher ratio values. Detailed summary tables, in addition to the graphics in figure 1 and 2, can be found in the appendix (see A).

3 Modelling

Given that our main objective is to better understand how the probability of mortgage approvals depends on the explanatory variables, we will leverage a logistic regression model with the logit link to model the binary outcome of mortgage approvals. As given by the task at hand, we constrained our scope to interactions with the *self* variable — denoting self-employment status.

3.1 Model selection

In refining the model selection, I applied the stepwise AIC algorithm that utilizes both forward and backward selection. This approach is applied to the model with all initial variables along with their interaction terms with *self* (see B.1). We are aiming to find a reduced model that balances complexity against its goodness of fit. We retrieve the following linear predictor:

$$\eta = \beta_0 + \beta_1 x_{\text{hir}} + \beta_2 x_{\text{odir}} + \beta_3 x_{\text{lvr}} + \beta_4 x_{\text{mcs}} + \beta_5 x_{\text{self}} + \beta_6 x_{\text{single}} + \beta_7 x_{\text{white}} + \beta_8 x_{\text{uria}} + \beta_9 x_{\text{odir}} \times x_{\text{self}} + \beta_{10} x_{\text{self}} \times x_{\text{white}} + \beta_{11} x_{\text{self}} \times x_{\text{uria}}. \quad (1)$$

Further, I test whether or not the reduced model misfits the data compared to the initial model via the likelihood ratio test with test statistic

$$D_{\text{reduced}} - D_{\text{initial}} \sim \chi^2(p_{\text{initial}} - p_{\text{reduced}}) = \chi^2(20 - 14) = \chi^2(6). \quad (2)$$

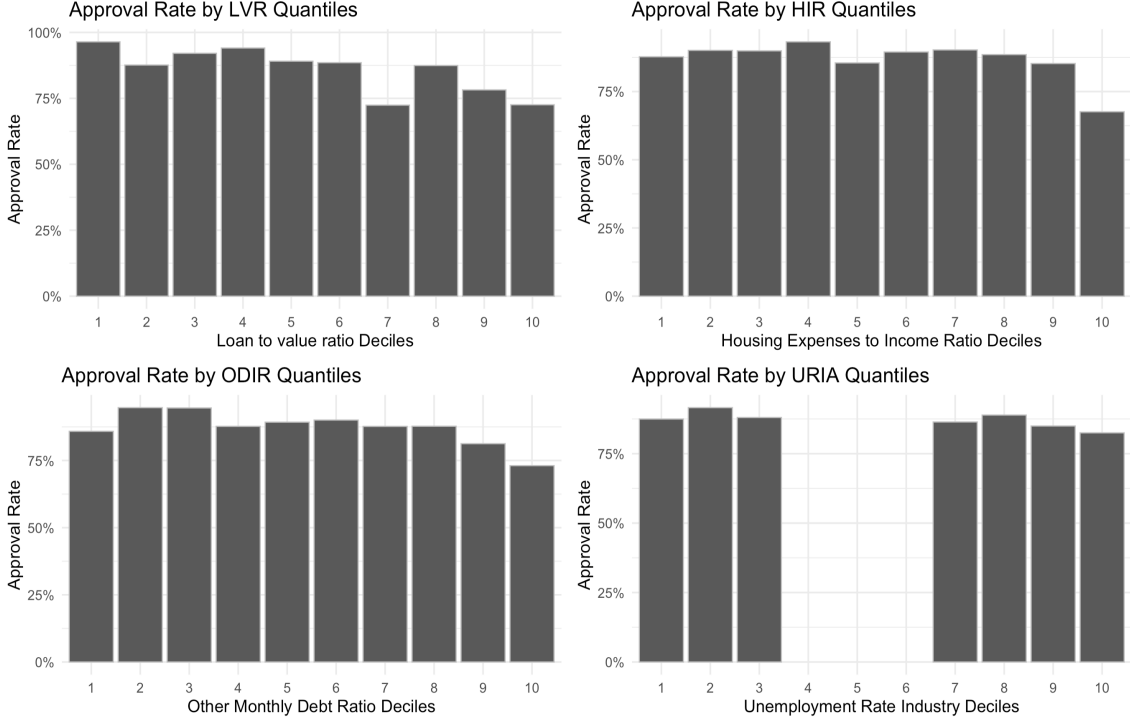


Figure 3: Mortgage Approval Rates across Financial Ratio and Unemployment Rate Deciles.

D_{reduced} is the residual deviance of our reduced model and D_{initial} is the residual deviance for the initial model with all variables and interactions of the self-employment status. We retrieve a p-value of 0.54 showing no significant change in deviance, such that we can simplify the model to our reduced model obtained through the stepwise AIC algorithm. In our models, we encoded the variable mortgage credit score (MCS) as a factor. This is due to the fact that our findings from our exploratory data analysis in 2 suggest that increasing the credit score from category 1 to category 2 is not as severe as increasing the credit score from 2 to 3 or 4. Therefore, we decided to include the variable as a categorical variable (factor) with dummy encoding (reference category MCS = 2) instead of a linear (numeric) effect, which would assume the same effect of changing credit score category between all four categories. Moreover, one can argue that because of the ordinal scale of the credit score, we can only make statements of higher, equal, or lower, but we cannot interpret the differences between categories in a meaningful way and, if included as numeric the relabelling of for instance category 4 to category 100 would influence the estimates and outcome, justifying the factor encoding in our model. The reference categories for the remaining binary/categorical variables are `single` = "Yes", `self` = "No" and `ethnicity` = "White".

3.2 Model diagnostics

In assessing model influence using Cook's Distance, as shown in Figure 4, we identified several influential observations using a soft threshold of 0.02, and the common hard threshold of $8/(n-2p)$. These points, notable enough to warrant scrutiny, were deemed legitimate data rather than errors and are detailed in Table C.1 in the appendix. The model fit visualized in Figure 5 and the associated generalized R-squared ($R_{KL}^2 = 0.1542$) indicate a modest explanatory power of the model, typical for logistic regression analyses.

3.3 Results & Interpretation

The logistic regression model output in table 1 gives an overview of the estimated coefficients of our selected model in 3.1. As already pointed out in section 2, we scaled our ratio variables by a factor of 100 such that they are on a percentage point scale for better interpretation.

We observe that the housing expense to income ratio (HIR) shows a negative association with the log-odds of approval, with an estimate of -0.053, with all other variables held fixed (*ceteris paribus*).

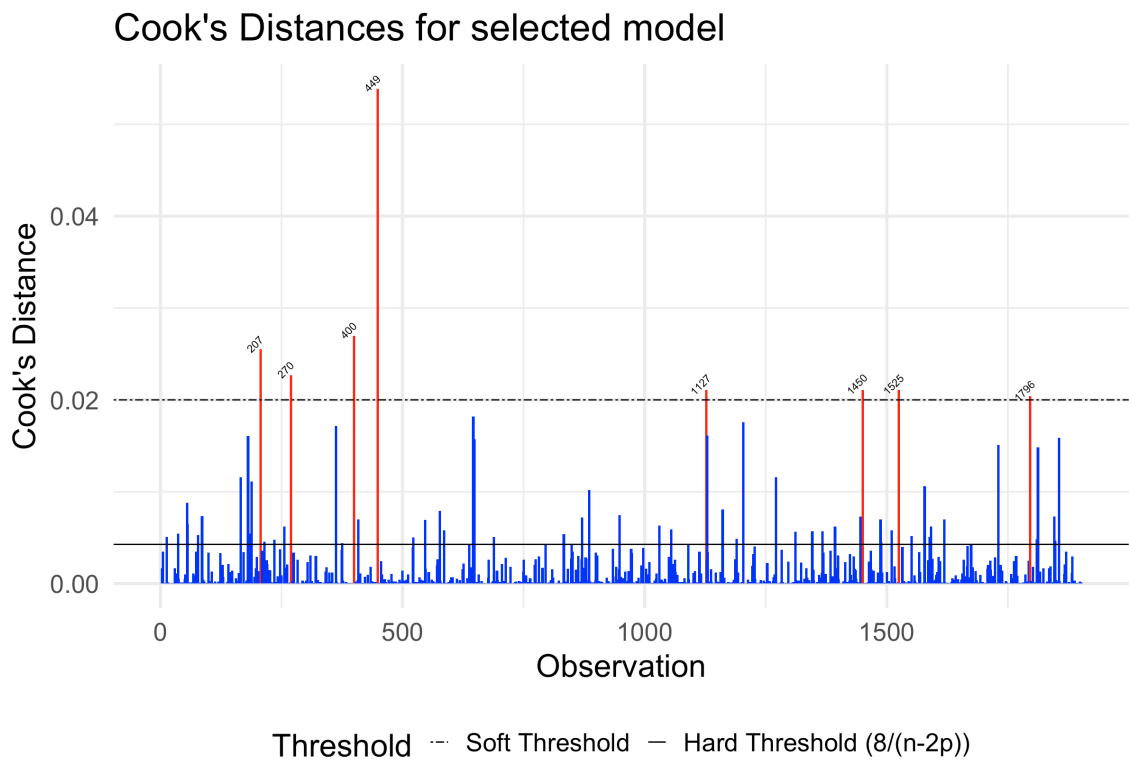


Figure 4: Cook's Distance for Influence Diagnostics in selected model. Vertical bars representing each observation's Cook's distance. The dashed lines indicate the soft and hard thresholds for identifying potentially influential points, with those exceeding the soft threshold highlighted in red.

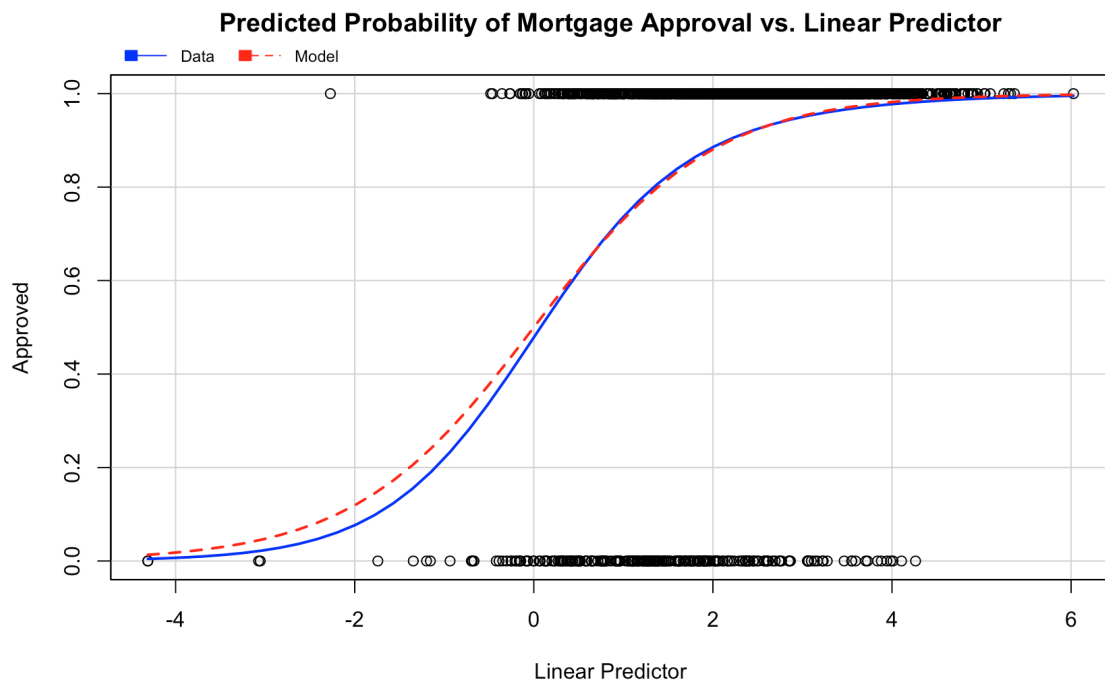


Figure 5: Calibration Curve for Mortgage Approval Predictions. Predicted probability of mortgage approval against linear predictor with the actual data points are marked in blue, while the red dashed line represents the model's predictions.

(c.p)). If the HIR increases by one percentage point, the odds of approval decrease on average by an estimated multiplicative factor of 0.95 c.p.. By considering the confidence interval, we can also make statements about the significance, as if the 95% confidence interval covers the value zero (or one if exponentiated), our effect is not significant at α -level 0.05.

For the mortgage credit score (mcs) variables, we exhibit varying influences. While category 1 shows a positive relationship with the odds for approval, showing that if the credit score category is 1 the estimated odds of approval increase by multiplicative factor of 1.72 (on average) compared to our reference category 2, while the odds for MCS3 and MCS4 decrease with factors 0.63 and 0.59.

For demographic variables, not being single slightly increases the odds of approval, with an OR of 1.39. It is worth mentioning that for the variables included in the interaction terms seen in table 1 we have to distinguish between interpreting effects for self-employed applicants and non self-employed applicants. We included enhancing effects on, for example, the ODIR through our interaction terms. The effect of the ODIR for self-employed applicants is therefore $\beta_{\text{odir}} + \beta_{\text{selfYes:odir}}$ and confidence interval

$$\left(\hat{\beta}_{\text{odir}} + \hat{\beta}_{\text{selfYes:odir}} \right) \pm z_{\alpha/2} \sqrt{\widehat{\text{var}}_{\text{odir}} + \widehat{\text{var}}_{\text{selfYes:odir}} + 2 \cdot \widehat{\text{Cov}} \left(\hat{\beta}_{\text{odir}}, \hat{\beta}_{\text{selfYes:odir}} \right)} \quad (3)$$

Note that we need to consider the estimated covariance of the coefficients when calculating this interval. The results for the coefficients of self-employed applicants can be seen in table 2, including the odds ratios, which are simply $\exp(\beta_{\text{odir}} + \beta_{\text{selfYes:odir}})$. We notice that if the applicant is self-employed, the increase of one percentage point in the ODIR does have a multiplicative effect closer to an OR of one (OR = 0.969) than for applicants that are self-employed (OR = 0.917). For ethnicity, we observe that self-employed black individuals have a point estimate of 0.05 with an odds ratio of 1.051, indicating a minimal effect on approval. For non-self-employed black individuals, the estimate is notably negative (-1.340), with a very low odds ratio of 0.262. This indicates a substantial decrease in the odds of approval by an estimated multiplicative factor of 0.262 if the individual is not self-employed and black compared to not self-employed whites.

For self-employed individuals, the estimate of the effect of URIA corresponds to an odds ratio of 1.113, therefore suggesting that an increase of one percentage point in URIA results in an increase in the odds for approval by an estimated average multiplicative of 1.113. The estimated coefficient associated with self-employment alone represents the unique contribution of being self-employed to the log-odds of the outcome when all other variables in the model are held constant at their reference levels. However, interpreting the 'selfYes' effect is challenging due to its involvement in several interaction terms with continuous variables like URIA and ODIR, hence implying that the impact of being self-employed changes with varying unemployment rates and debt ratios. In this context, marginal effects become particularly insightful.

Variable	Estimate	Std. Error	95% CI	OR ($\exp(\hat{\beta})$)	95% CI (OR)
(Intercept)	6.186***	0.551	(5.106, 7.266)	485.86	(164.61, 1437.82)
HIR	-0.053***	0.010	(-0.073, -0.033)	0.95	(0.93, 0.97)
ODIR	-0.087***	0.013	(-0.111, -0.062)	0.92	(0.89, 0.94)
LVR	-0.021***	0.005	(-0.030, -0.011)	0.98	(0.97, 0.99)
MCS1	0.540**	0.195	(0.157, 0.923)	1.72	(1.17, 2.52)
MCS3	-0.454	0.474	(-1.384, 0.475)	0.63	(0.25, 1.61)
MCS4	-0.529	0.604	(-1.715, 0.656)	0.59	(0.18, 1.93)
Single No	0.331*	0.153	(0.031, 0.631)	1.39	(1.03, 1.88)
Ethnic Black	-1.340***	0.179	(-1.690, -0.989)	0.26	(0.18, 0.37)
URIA	-0.105**	0.038	(-0.180, -0.030)	0.90	(0.83, 0.97)
Self Yes	-2.352***	0.497	(-3.326, -1.378)	0.095	(0.036, 0.252)
Self Yes:odir	0.054*	0.025	(0.006, 0.103)	1.06	(1.01, 1.11)
Self Yes:Black	1.389*	0.656	(0.103, 2.675)	4.01	(1.11, 14.48)
self Yes:URIA	0.212*	0.087	(0.042, 0.383)	1.24	(1.04, 1.47)

*** $p < 0.001$; ** $p < 0.01$; * $p < 0.05$

Table 1: Model coefficients, standard errors, 95% confidence intervals, exponentiated coefficients (OR), and their 95% confidence intervals with significance codes

Variable	Employment	Estimate	Std. Error	95% CI	OR ($\exp(\hat{\beta})$)
ODIR	<i>Self</i> = "Yes"	-0.032	0.021	(-0.074, 0.010)	0.969
	<i>Self</i> = "No"	-0.087***	0.013	(-0.111 -0.062)	0.917
Ethnicity Black	<i>Self</i> = "Yes"	0.05	0.633	(-1.192, 1.291)	1.051
	<i>Self</i> = "No"	-1.340***	0.179	(-1.690, -0.989)	0.262
URIA	<i>Self</i> = "Yes"	0.107	0.078	(-0.046, 0.261)	1.113
	<i>Self</i> = "No"	-0.105**	0.038	(-0.180, -0.030)	0.900

*** $p < 0.001$; ** $p < 0.01$; * $p < 0.05$

Table 2: Model coefficients, standard errors, 95% confidence intervals, exponentiated coefficients (OR) for variables included in interaction terms.

3.3.1 Marginal Effects

Marginal effects are particularly valuable when interpreting complex models with interaction terms between categorical and continuous variables, hence for the interpretation of the role of self-employment, they seem worth exploring. Table 3 presents the average marginal effects (AME) for each variable in the logistic regression model estimated via the delta method, which quantifies the expected average change in the log odds of mortgage approval associated with a one-unit change in the predictor variable, while other variables are held at their fixed average values. We can now obtain the average marginal effect (AME) of being self-employed of -0.9425, which can be interpreted as follows for self-employed applicants, the odds of an approved mortgage application decrease multiplicatively by factor $\exp(-0.9425) = 0.3897$ c.p. compared to a person who is not self-employed. The set of plots in figure 6 illustrates the marginal effects of the interactions in

Factor	AME	$\exp(\text{AME})$	SE	p	CI
HIR	-0.0532	0.9483	0.0101	<0.0001	(-0.073, -0.033)
LVR	-0.0206	0.9795	0.0048	<0.0001	(-0.030, -0.011)
MCS1	0.5399	1.7161	0.1954	0.0057	(0.157, 0.923)
MCS3	-0.4544	0.6348	0.4742	0.3380	(-1.384, 0.475)
MCS4	-0.5294	0.5891	0.6048	0.3815	(-1.715, 0.656)
ODIR	-0.0803	0.9230	0.0114	<0.0001	(-0.103, -0.058)
Self Yes	-0.9425	0.3897	0.2195	<0.0001	(-1.373, -0.512)
Single No	0.3307	1.3922	0.1530	0.0307	(0.031, 0.631)
URIA	-0.0806	0.9226	0.0352	0.0220	(-0.150, -0.012)
Ethnicity Black	-1.1818	0.3072	0.1748	<0.0001	(-1.524, -0.839)

Table 3: Average Marginal Effects (AME) and Exponentiated AME for Mortgage Approval Factors. AME of various factors on mortgage approval, their exponentiated forms representing odds ratios, standard errors (SE), p-values, and 95% confidence intervals (CI). Negative AME values suggest a decrease, and positive values suggest an increase in the probability of mortgage approval.

our reduced model on the probability scale, which I want to touch on shortly, in addition to our odds-scale interpretations. In the first plot, we see the predicted probability of mortgage approval as a function of ODIR, stratified by self-employment status. For non-self-employed individuals ('No'), there is a steeper decline in approval probability as ODIR increases.

The second plot depicts the effect of the URIA on approval probability, again differentiated by self-employment status. For those not self-employed, a higher unemployment rate in the applicant's industry corresponds to a lower probability of approval.

Finally, the third plot shows the approval probability by ethnicity. Non-self-employed black individuals have a lower predicted probability of approval than their white counterparts, which is less the case for self-employed individuals.

3.4 Dispersion

We obtain a dispersion parameter estimate of 1.102 by using the formula $\hat{\phi} = \frac{1}{n-p} \sum_{i=1}^n \frac{(y_i - \hat{\mu}_i)^2}{V(\hat{\mu}_i)}$. This difference from the assumed $\phi = 1$ suggests minor overdispersion, which means that the observed variation in the response variable is greater than the variation predicted and assumed

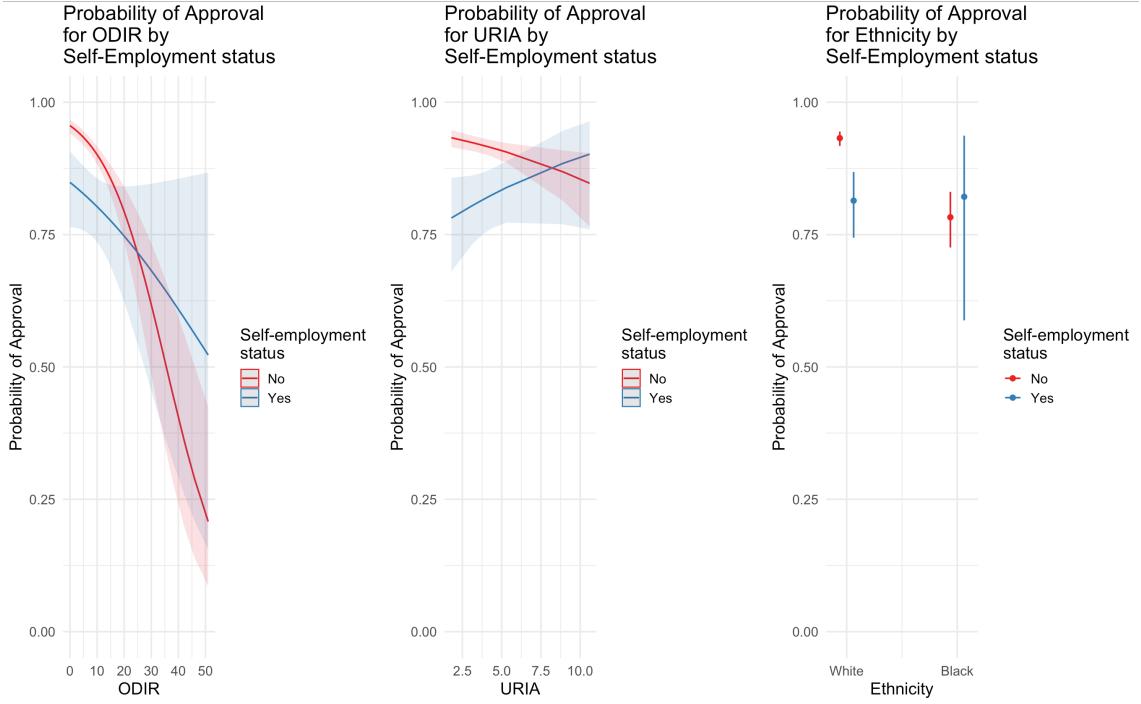


Figure 6: Conditional Effects of ODIR, URIA, and Ethnicity on Mortgage Approval Probability by Self-employment status. Lines show the predicted probability of mortgage approval at values of ODIR, URIA, and Ethnicity—focal terms in the statistical model—while holding non-focal variables constant. Ribbons show the 95% confidence intervals.

in fitting the model, leading to potentially attenuated variance estimates. It would be possible to tackle overdispersion by using a quasibinomial model, which adjusts the variance to account for the extra variation observed in the response variable.

4 Limitations & Outlook

This analysis, while detailed, encounters limitations such as slight overdispersion and the absence of variables like wealth, which may influence the precision of effect estimation. Addressing non-linear predictor-outcome relationships through generalized additive models could also refine the understanding of the data. For those seeking predictive approaches, the appendix presents methodologies like oversampling and LASSO regression as alternatives.

5 Conclusion

Our logistic regression model discerns key determinants of mortgage approval. Financial metrics like the housing expense-to-income ratio (HIR) and loan-to-value ratio (LVR) negatively impact approval odds, whereas a superior mortgage credit score markedly boosts them. Demographics influence outcomes too; singles and people of black ethnicity face disadvantages, hinting at broader issues not encapsulated by the model. The inclusion of additional data and the application of alternative modeling techniques could enhance the robustness of future models, offering a more comprehensive view of the factors influencing mortgage approvals.

A Supplementary EDA Tables

Continuous Variables				
Variable	Mean	St. Dev.	Min	Max
hir	0.255	0.079	0.010	1.100
odir	0.075	0.066	0.000	0.510
lvr	0.737	0.180	0.020	1.950
uria	3.763	2.015	1.800	10.600

Table 4: Summary of Continuous Variables

Categorical/Binary Variables			
Variable	Levels	<i>n</i>	%
approved	Yes	1665	87.5
	No	237	12.5
mcs	1	609	32.0
	2	1241	65.2
	3	33	1.7
	4	19	1.0
self	Yes	216	11.4
	No	1686	88.6
single	Yes	753	39.6
	No	1149	60.4
white	White	1625	85.4
	Black	277	14.6

Table 5: Summary of Categorical/Binary Variables

B Modelling

B.1 Model summary of initial model

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	6.0374	0.6087	9.92	<0.0001
HIR	-0.0530	0.0117	-4.53	<0.0001
ODIR	-0.0851	0.0127	-6.71	<0.0001
LVR	-0.0186	0.0054	-3.45	0.0006
MCS1	0.5379	0.2206	2.44	0.0148
MCS3	-0.6383	0.5399	-1.18	0.2371
MCS4	-0.9631	0.6376	-1.51	0.1309
Self Yes	-1.6640	1.2988	-1.28	0.2001
Single No	0.2708	0.1676	1.62	0.1061
Ethnic Black	-1.3599	0.1794	-7.58	<0.0001
URIA	-0.1015	0.0384	-2.64	0.0082
hir:selfYes	-0.0069	0.0242	-0.29	0.7742
odir:selfYes	0.0517	0.0258	2.00	0.0451
lvr:selfYes	-0.0098	0.0119	-0.82	0.4105
MCS1:selfYes	0.0613	0.4812	0.13	0.8987
MCS3:selfYes	0.7142	1.1123	0.64	0.5208
MCS4:selfYes	14.0039	428.5445	0.03	0.9739
singleNo:selfYes	0.3862	0.4269	0.90	0.3656
ethnBlack:selfYes	1.5179	0.6829	2.22	0.0262
uria:selfYes	0.2022	0.0898	2.25	0.0243

Table 6: Model summary of initial model before applying stepwise AIC selection

C Outlier Table

C.1 Soft Threshold

	approved	hir	odir	lvr	mcs	self	single	white	uria
207	No	10.00	0.00	67.00	2	Yes	Yes	White	10.60
270	No	18.00	5.00	61.00	1	Yes	No	White	10.60
400	No	31.00	25.00	57.00	3	Yes	No	White	10.60
449	Yes	44.00	51.00	55.00	3	Yes	No	White	10.60
1127	No	13.00	8.00	80.00	2	Yes	Yes	Black	3.20
1450	No	11.00	36.00	90.00	1	Yes	No	White	3.20
1525	No	6.00	12.00	80.00	4	No	No	Black	3.20
1796	No	30.00	11.00	72.00	4	No	No	White	3.20

Table 7: Observation with Cooks distance of over 0.02 in the selected model

D Further Modelling Approaches

D.1 Oversampling

Oversampled the minority class of unapproved applications, however, results should be interpreted with caution as this assumes a representative representation of the minority class in the data.

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	3.7498	0.2993	12.53	<0.001
hir	-0.0411	0.0060	-6.80	<0.001
odir	-0.0794	0.0068	-11.67	<0.001
lvr	-0.0184	0.0028	-6.60	<0.001
mcs1	0.4437	0.1054	4.21	<0.001
mcs3	-0.4961	0.3134	-1.58	0.1134
mcs4	-0.3584	0.4228	-0.85	0.3967
selfYes	-0.8089	0.7745	-1.04	0.2963
singleNo	0.4007	0.0858	4.67	<0.001
BlackYes	-1.3740	0.1041	-13.20	<0.001
uria	-0.1266	0.0207	-6.13	<0.001
hir:selfYes	-0.0333	0.0143	-2.34	0.0195
odir:selfYes	0.0395	0.0155	2.54	0.0110
lvr:selfYes	-0.0104	0.0078	-1.33	0.1835
mcs1:selfYes	0.2287	0.2699	0.85	0.3969
mcs3:selfYes	0.3197	0.7075	0.45	0.6513
mcs4:selfYes	14.4137	250.4227	0.06	0.9541
selfYes:singleNo	0.7702	0.2524	3.05	0.0023
selfYes:BlackYes	1.2852	0.4043	3.18	0.0015
selfYes:uria	0.1727	0.0506	3.41	0.0006

Table 8: Results for model fitted with oversampled / balanced data set

D.2 Weighted GLM

$$w(x) = \begin{cases} \frac{\sum_{i=1}^N 1_{\{approved_i = \text{Yes}\}}}{N} & \text{if } approved = \text{Yes} \\ \frac{\sum_{i=1}^N 1_{\{approved_i = \text{No}\}}}{N} & \text{if } approved = \text{No} \end{cases} \quad (4)$$

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	6.1859	0.1945	31.80	<0.001
hir	-0.0532	0.0036	-14.87	<0.001
odir	-0.0865	0.0044	-19.54	<0.001
lvr	-0.0206	0.0017	-12.19	<0.001
mcs1	0.5399	0.0690	7.83	<0.001
mcs3	-0.4544	0.1674	-2.71	0.0066
mcs4	-0.5294	0.2135	-2.48	0.0132
selfYes	-2.3519	0.1754	-13.41	<0.001
singleNo	0.3307	0.0540	6.12	<0.001
whiteBlack	-1.3395	0.0631	-21.22	<0.001
uria	-0.1047	0.0135	-7.74	<0.001
odir:selfYes	0.0545	0.0087	6.27	<0.001
selfYes:whiteBlack	1.3892	0.2316	6.00	<0.001
selfYes:uria	0.2122	0.0307	6.91	<0.001

Table 9: Results for weighted GLM with weights calculated as described in 4.

D.3 LASSO Regression and n-fold Cross-Validation

20-fold cross-validation to select the optimal regularization parameter, lambda. Figure D.3 visualizes the cross-validation process with the model’s binomial deviance (a goodness-of-fit measure) against different values of $\log(\text{lambda})$.

	Feature	Coefficient
1	hir	-0.02
2	odir	-0.03
3	lvr	-0.01
4	whiteBlack	-0.84
5	hir.selfYes	-0.01

Table 10: Non-zero coefficients for lambda which is one standard error above the minimum cross-validated error

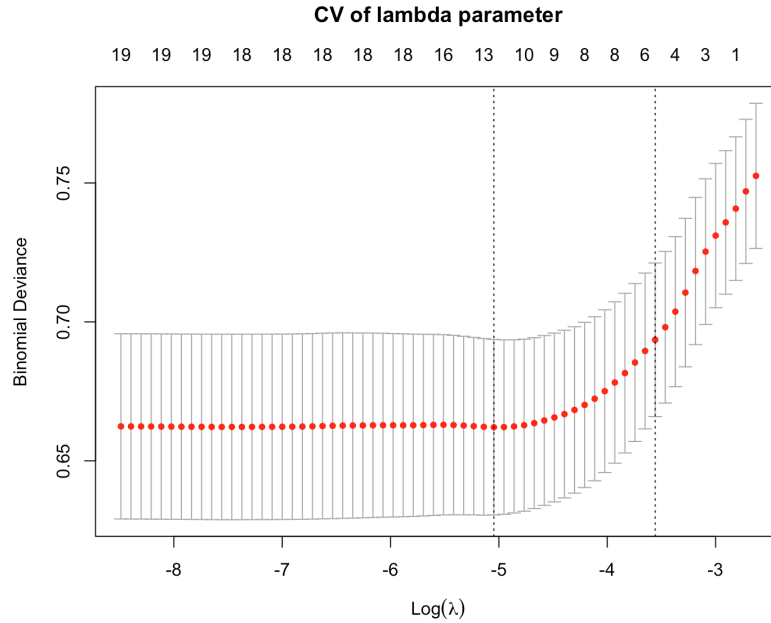


Figure 7: Cross-Validation of Lambda in LASSO Regression. Dotted lines represent the lambda values yielding minimum deviance and the most regularized model within one standard error of the minimum.

E R Code

```

1  # Libraries ----
2  set.seed(42)
3  library(ggplot2)
4  library(gridExtra)
5  library(dplyr)
6  library(car)
7  library(rsq)
8  library(gbm)
9  library(boot)
10 library(xgboost)
11 library(margins)
12 library(ggeffects)
13 library(stargazer)
14 # Data read-in and data definition ----
15 ## Read and definition ----
16 mortg_raw <- read.csv("data/mortg.csv")
17 mortg <- read.csv("data/mortg.csv")
18 mortg$approved <- factor(mortg$approved,
19                           levels = c("0", "1"),
20                           labels = c("No", "Yes"))
21 mortg$mcs <- factor(mortg$mcs)
22 mortg$self <- factor(
23   as.character(mortg$self),
24   levels = c("1", "0"),
25   labels = c("Yes", "No")
26 )
27 mortg$single <- factor(
28   as.character(mortg$single),
29   levels = c("1", "0"),
30   labels = c("Yes", "No")
31 )

```

```

32 mortg$white <- factor(
33   as.character(mortg$white),
34   levels = c("1", "0"),
35   labels = c("White", "Black")
36 )
37
38 ##Scale ratios to percentage point scale for interpretability ----
39 mortg$hir <- mortg$hir * 100
40 mortg$odir <- mortg$odir * 100
41 mortg$lvr <- mortg$lvr * 100
42
43 # Exploratory Data Analysis ----
44 ## Overview ----
45 summary(mortg)
46 table(mortg$approved, mortg$self)
47 table(mortg$approved, mortg$single)
48 table(mortg$approved, mortg$white)
49 prop.table(table(mortg$approved, mortg$self), 1)
50 prop.table(table(mortg$approved, mortg$single))
51 prop.table(table(mortg$approved, mortg$white), 2)
52 prop.table(table(mortg$white, mortg$self), 2)
53
54 ## Bivarite Plots ----
55 ### Boxplots (continous data) ----
56 create_vio_box_plot <-
57   function(data,
58     variable,
59     title,
60     xlab = "",
61     limits = c(0, 100)) {
62     ggplot(data, aes_string(x = variable, y = "approved")) +
63       geom_violin(
64         width = 1,
65         color = "black",
66         fill = "darkgrey",
67         alpha = 0.2
68       ) +
69       geom_boxplot(width = 0.1,
70         outlier.size = 1.4,
71         outlier.shape = 1) +
72       labs(title = title,
73         x = xlab,
74         y = "Mortgage Approval") +
75       coord_flip() +
76       theme_minimal(base_size = 15) +
77       scale_x_continuous(labels = scales::percent_format(scale = 1),
78         limits = limits)
79   }
80 create_scaled_stacked_bar_plot <-
81   function(data,
82     category,
83     xlab = category,
84     ylab = "Percentage",
85     fill = c("#F05039", "#7CA1CC"),
86     title = paste("Mortgage Approval by", category, "Status")) {
87     data %>%
88       group_by(!!sym(category), approved) %>%
89       summarise(Count = n()) %>%

```

```

90     mutate(Total = sum(Count)) %>%
91     mutate(Percent = (Count / Total) * 100) %>%
92     ggplot(aes_string(x = category,
93                       y = "Percent", fill = "factor(approved)")) +
94     geom_bar(stat = "identity", position = "fill") +
95     labs(
96       title = title,
97       x = xlab,
98       y = ylab,
99       fill = "Approved"
100    ) +
101    scale_fill_manual(values = fill) +
102    scale_y_continuous(labels = scales::percent_format()) +
103    theme_minimal(base_size = 15)
104  }
105
106
107 plot_hir <-
108   create_vio_box_plot(mortg,
109                       "hir",
110                       "HIR\nvs Mortgage Approval",
111                       xlab = "Housing Expense-to-Income Ratio",
112                       limits = c(0, 120))
113
114 plot_odir <-
115   create_vio_box_plot(
116     mortg,
117     "odir",
118     "ODIR\nvs Mortgage Approval",
119     xlab = "Other Debt-to-Income Ratio",
120     limits = c(0, 60)
121   )
122 plot_lvr <-
123   create_vio_box_plot(mortg,
124                       "lvr",
125                       "LVR\nvs Mortgage Approval",
126                       xlab = "Loan-to-Value Ratio",
127                       limits = c(0, 200))
128 plot_uria <-
129   create_vio_box_plot(
130     mortg,
131     "uria",
132     "URIA\nvs Mortgage Approval",
133     xlab = "Unemployment Rate in Applicant's Industry",
134     limits = c(0, 12)
135   )
136 plot_self <-
137   create_scaled_stacked_bar_plot(mortg,
138                                   "self",
139                                   xlab = "Self-Employed",
140                                   ylab = "Proportion",
141                                   title = "Mortgage Approval by\nSelf-Employment Status")
142 plot_single <-
143   create_scaled_stacked_bar_plot(mortg,
144                                   "single",
145                                   xlab = "Single",
146                                   ylab = "Proportion",
147                                   title = "Mortgage Approval by\nMarital Status (Single)")

```



```

148 plot_white <-
149   create_scaled_stacked_bar_plot(mortg,
150                                   "white",
151                                   xlab = "Ethnicity",
152                                   ylab = "Proportion",
153                                   title = "Mortgage Approval by\nEthnicity")
154 plot_mcs <-
155   create_scaled_stacked_bar_plot(mortg,
156                                   "mcs",
157                                   xlab = "Mortgage Credit Score",
158                                   ylab = "Proportion",
159                                   title = "Mortgage Approval by\nMortgage Credit Score")
160
161 grid.arrange(plot_hir, plot_odir, plot_lvr, plot_uria, ncol = 2)
162 grid.arrange(plot_self, plot_single, plot_white, plot_mcs, ncol = 2)
163
164 ## Decile Plots ----
165 mortg_unscaled <- mortg_raw
166 plot_approval_rate_by_quantiles <-
167   function(data,
168            variable_name,
169            xlab = paste(variable_name, "Deciles"),
170            title = paste("Approval Rate by", toupper(variable_name), "Quantiles")) {
171     # Compute the quantile bins for the variable
172     data <- data %>%
173       mutate(quantile_bin = findInterval(.data[[variable_name]],
174                                           quantile(.data[[variable_name]],
175                                                       probs = seq(0.1, 1, 0.1)),
176                                           left.open = TRUE) + 1)
177
178     # Calculate the mean approval rate by quantile bin
179     approval_rate <- data %>%
180       group_by(quantile_bin) %>%
181       summarise(approval_rate = mean(approved, na.rm = TRUE)) %>%
182       ungroup()
183
184     # Create the bar plot
185     ggplot(approval_rate, aes(x = as.factor(quantile_bin), y = approval_rate)) +
186       geom_bar(stat = "identity",
187               colour = "grey",
188               position = position_dodge()) +
189       scale_x_discrete(limits = as.character(1:10), drop = FALSE) +
190       labs(title = title,
191            x = xlab,
192            y = "Approval Rate") +
193       scale_y_continuous(labels = scales::percent_format()) +
194       theme_minimal(base_size = 15)
195   }
196
197
198
199 decile_approval_lvr <-
200   plot_approval_rate_by_quantiles(mortg_raw, "lvr",
201                                   xlab = "Loan to value ratio Deciles")
202 decile_approval_hir <-
203   plot_approval_rate_by_quantiles(mortg_raw, "hir",
204                                   xlab = "Housing Expenses to Income Ratio Deciles")
205 decile_approval_odir <-

```

```

206     plot_approval_rate_by_quantiles(mortg_raw, "odir",
207                                     xlab = "Other Monthly Debt Ratio Deciles")
208 decile_approval_uria <-
209     plot_approval_rate_by_quantiles(mortg_raw, "uria",
210                                     xlab = "Unemployment Rate Industry Deciles")
211
212 grid.arrange(
213     decile_approval_lvr,
214     decile_approval_hir,
215     decile_approval_odir,
216     decile_approval_uria,
217     ncol = 2
218 )
219
220 # Modelling ----
221 ## Set reference categories ----
222 mortg$mcs <-
223     relevel(mortg$mcs, 2) # MCS 2 as reference category (rc)
224 mortg$self <- relevel(mortg$self, 2) # Not self-employed as rc
225 mortg$single <- relevel(mortg$single, 1) # Single Yes as rc
226 mortg$white <- relevel(mortg$white, 1) # white as rc
227
228 ## Variable selection stepwise AIC ----
229 full_model <-
230     glm(approved ~ . * self,
231         family = binomial(link = "logit"),
232         data = mortg)
233 summary(full_model)
234 full_model_summary <- summary(full_model)
235 reduced_model <- step(full_model, direction = "both")
236 summary(reduced_model)
237 reduced_model_summary <- summary(reduced_model)
238 ## Model selection Analysis of Deviance ----
239 ### LRT ----
240 D_full <- full_model$deviance
241 D_reduced <- reduced_model$deviance
242 p_full <- length(full_model$coefficients)
243 p_reduced <-
244     length(reduced_model$coefficients) # number of parameters in reduced model
245 deviance_change <- D_reduced - D_full
246 pchisq(deviance_change,
247         df = p_full - p_reduced,
248         lower.tail = FALSE)
249 ### via ANOVA function ----
250 anova(reduced_model, full_model, test = "Chisq")
251
252 # Model Diagnosis ----
253 # Cooks Distance ----
254 plot_cooks_distance <-
255     function(model,
256             model_name,
257             alpha = 0.1,
258             threshold = NULL) {
259         # Calculate Cook's distance
260         cooks_d <- cooks.distance(model)
261
262         # If threshold is not set, use the default value
263         if (is.null(threshold)) {

```

```

264     threshold <- 4 / length(cooks_d)
265   }
266
267   # Create a data frame for plotting
268   plot_data <- data.frame(
269     Observation = 1:length(cooks_d),
270     CooksDistance = cooks_d,
271     AboveThreshold = cooks_d > threshold
272   )
273
274   # Create the plot using ggplot2
275   ggplot(plot_data, aes(x = Observation, y = CooksDistance)) +
276     geom_point(aes(color = AboveThreshold), alpha = alpha) +
277     geom_point(
278       data = subset(plot_data, AboveThreshold),
279       aes(color = AboveThreshold),
280       alpha = 1
281     ) + # Points above threshold with full opacity
282     scale_color_manual(values = c("FALSE" = "blue", "TRUE" = "red"),
283       guide = FALSE) +
284     ggtitle(paste("Cook's Distance for", model_name)) + # Title
285     xlab('Observation') + # X-axis label
286     ylab("Cook's Distance") + # Y-axis label
287     geom_text(
288       data = subset(plot_data, AboveThreshold),
289       aes(label = Observation),
290       vjust = -0.5,
291       hjust = "inward",
292       size = 3,
293       angle = 45
294     ) + # Label high points
295     geom_hline(yintercept = threshold,
296       linetype = 'dotted',
297       col = 'black') +
298     geom_hline(
299       yintercept = 4 / length(cooks_d),
300       linetype = 'dotted',
301       col = 'pink'
302     )
303   }
304
305   plot_cooks_distance_barplot <-
306     function(model,
307       model_name,
308       threshold = 4 / length(cooks.distance(model))) {
309     # Calculate Cook's distance
310     cooks_d <- cooks.distance(model)
311
312     # Create a data frame for plotting
313     plot_data <- data.frame(Observation = 1:length(cooks_d),
314       CooksDistance = cooks_d)
315
316     # Create the bar plot using ggplot2
317     ggplot(plot_data, aes(x = Observation, y = CooksDistance)) +
318       geom_bar(stat = "identity",
319         aes(fill = CooksDistance > threshold),
320         width = 5) +
321       scale_fill_manual(values = c("blue", "red"), guide = FALSE) +

```

```

322 ggtitle(paste("Cook's Distances for selected model")) +
323 xlab("Observation") +
324 ylab("Cook's Distance") +
325 geom_text(
326   data = subset(plot_data, CooksDistance > threshold),
327   aes(label = Observation),
328   vjust = -0.5,
329   size = 3,
330   angle = 45
331 ) +
332 geom_hline(aes(yintercept = threshold, linetype = "Soft threshold"),
333            color = 'black') +
334 geom_hline(aes(
335   yintercept = 8 / (length(cooks_d) - 2 * length(coef(reduced_model))),
336   linetype = "Hard threshold"
337 ),
338 color = 'black') +
339 scale_linetype_manual(
340   name = "Threshold",
341   values = c("Soft threshold" = "twodash", "Hard threshold" = "solid"),
342   breaks = c("Soft threshold", "Hard threshold"),
343   labels = c("Soft Threshold", "Hard Threshold (8/(n-2p))")
344 ) +
345 theme_minimal(base_size = 25) +
346 theme(legend.position = "bottom")
347 }
348
349 plot_cooks_distance(reduced_model, "Reduced", 0.1, 0.02)
350 plot_cooks_distance_barplot(reduced_model, "Reduced", 0.02)
351 cooks_d <- cooks.distance(reduced_model)
352
353
354 # Calibration Curve ----
355 marginalModelPlot(reduced_model,
356                   ylab = "Approved",
357                   cex.lab = 15)
358 title(main = "Predicted Probability of Mortgage Approval vs. Linear Predictor")
359 # Model results ----
360 # Model table ----
361 summary(reduced_model)
362
363 # Confidence Intervals ----
364 confint.default(reduced_model)
365 exp(confint.default(reduced_model))
366 ## Multiple Confidence Intervals -----
367 calculate_interac_confint <-
368   function(model,
369            var1,
370            var2 = "selfYes",
371            sig.level = 0.05) {
372     # Extract variance-covariance matrix from the model
373     coefs_var <- vcov(model)
374     # Extract the coefficients
375     coefs <- coef(model)
376     # Calculate the z-value based on the significance level
377     z <- qnorm(1 - sig.level / 2)
378
379     # Calculate the standard error for the interaction term

```

```

380     # considering the covariance
381     se <-
382     sqrt(coefs_var[var1, var1] + coefs_var[var2, var2]
383           + 2 * coefs_var[var1, var2])
384     ci_lower <- coefs[var1] + coefs[var2] - z * se
385     ci_upper <- coefs[var1] + coefs[var2] + z * se
386     # Return the confidence interval
387
388     return(list(
389       "confint" = c(lower = ci_lower, upper = ci_upper),
390       "se" = se
391     ))
392   }
393
394   calculate_interac_pv <- function(model, var1, var2) {
395     # Extract variance-covariance matrix from the model
396     coefs_var <- vcov(model)
397     # Extract the coefficients
398     coefs <- coef(model)
399
400     # Calculate the standard error for the interaction term
401     # considering the covariance
402     se <-
403     sqrt(coefs_var[var1, var1] + coefs_var[var2, var2]
404           + 2 * coefs_var[var1, var2])
405
406     p <- 2 * (1 - pnorm(abs((coefs[var1] + coefs[var2]) / se)))
407     names(p) <- paste(var1, var2, sep = "+")
408     # Return the confidence interval
409     return(p)
410   }
411
412   ### odir:selfYes ----
413   round(sum(coef(reduced_model)[c("odir", "odir:selfYes")]), 3)
414   calculate_interac_confint(
415     model = reduced_model,
416     var1 = "odir",
417     var2 = "odir:selfYes",
418     sig.level = 0.05
419   )
420   calculate_interac_pv(model = reduced_model,
421                       var1 = "odir", var2 = "odir:selfYes")
422   ### selfYes:whiteBlack ----
423   round(sum(coef(reduced_model)[c("whiteBlack", "selfYes:whiteBlack")]), 3)
424   calculate_interac_confint(
425     model = reduced_model,
426     var1 = "whiteBlack",
427     var2 = "selfYes:whiteBlack",
428     sig.level = 0.05
429   )
430   calculate_interac_pv(model = reduced_model,
431                       var1 = "whiteBlack", var2 = "selfYes:whiteBlack")
432   ### selfYes:uria ----
433   round(sum(coef(reduced_model)[c("uria", "selfYes:uria")]), 3)
434   calculate_interac_confint(
435     model = reduced_model,
436     var1 = "uria",
437     var2 = "selfYes:uria",

```

```

438   sig.level = 0.05
439 )
440 calculate_interac_pv(model = reduced_model,
441                     var1 = "uria", var2 = "selfYes:uria")
442
443
444
445 # Marginal Effects ----
446 marginal_effects <-
447   margins(
448     reduced_model,
449     vce = "delta",
450     vcov = vcov(reduced_model),
451     type = "link"
452   )
453
454 ## Summary of marginal effects ----
455 summary(marginal_effects) %>% as.data.frame()
456
457
458 marginal_effects_df <- summary(marginal_effects) %>% as.data.frame()
459
460 marginal_effects_df <- marginal_effects_df %>%
461   mutate(AME_exp = exp(AME))
462
463
464
465 marginal_effects_df <- marginal_effects_df %>%
466   mutate(p = 2 * (1 - pnorm(abs(z))),
467          p_rounded = round(2 * (1 - pnorm(abs(
468            z
469          ))), 4))
470
471 ## AME Plot ----
472 ggplot(marginal_effects_df, aes(x = factor, y = AME)) +
473   geom_point() +
474   geom_errorbar(aes(ymin = lower, ymax = upper), width = 0.2) +
475   coord_flip() + # Flip coordinates to have factors on the y-axis
476   labs(title = "Average Marginal Effects (AME) with Confidence Intervals",
477        x = "Average Marginal Effect",
478        y = "Factor") +
479   theme_minimal()
480
481
482 ## Effects on Prob Scale----
483 ### Single effects ----
484 meff_self <- ggpredict(reduced_model, terms = c("self")) |>
485   plot() +
486   ylim(c(0, 1))
487 meff_white <- ggpredict(reduced_model, terms = c("white")) |>
488   plot() +
489   ylim(c(0, 1))
490 meff_single <- ggpredict(reduced_model, terms = c("single")) |>
491   plot() +
492   ylim(c(0, 1))
493 meff_mcs <- ggpredict(reduced_model, terms = c("mcs")) |>
494   plot() +
495   ylim(c(0, 1))

```

```

496 meff_odir <- ggpredict(reduced_model, terms = c("odir[all]")) |>
497   plot() +
498   ylim(c(0, 1))
499 meff_lvr <- ggpredict(reduced_model, terms = c("lvr[all]")) |>
500   plot() +
501   ylim(c(0, 1))
502 meff_hir <- ggpredict(reduced_model, terms = c("hir[all]")) |>
503   plot() +
504   ylim(c(0, 1))
505 meff_uria <- ggpredict(reduced_model, terms = c("uria[all]")) |>
506   plot() +
507   ylim(c(0, 1))
508
509 grid.arrange(
510   meff_odir,
511   meff_lvr,
512   meff_hir,
513   meff_uria,
514   meff_self,
515   meff_white,
516   meff_single,
517   meff_mcs,
518   ncol = 4
519 )
520
521 ## Interaction Plots ----
522 meff_odir_self <-
523   ggeffect(reduced_model, terms = c("odir[all]", "self")) |>
524   plot() +
525   ylim(c(0, 1)) +
526   geom_line() +
527   xlab("ODIR") +
528   ylab("Probability of Approval") +
529   labs(colour = "Self-employment\nstatus") +
530   ggtitle("Probability of Approval\nfor ODIR by\nSelf-Employment status") +
531   theme_minimal(base_size = 15)
532
533
534 meff_uria_self <-
535   ggeffect(reduced_model, terms = c("uria[all]", "self")) |>
536   plot() +
537   ylim(c(0, 1)) +
538   xlab("URIA") +
539   ylab("Probability of Approval") +
540   labs(colour = "Self-employment\nstatus") +
541   ggtitle("Probability of Approval\nfor URIA by\nSelf-Employment status") +
542   theme_minimal(base_size = 15)
543
544
545 meff_white_self <-
546   ggeffect(reduced_model, terms = c("white[all]", "self")) |>
547   plot() +
548   ylim(c(0, 1)) +
549   xlab("Ethnicity") +
550   ylab("Probability of Approval") +
551   labs(colour = "Self-employment\nstatus") +
552   ggtitle("Probability of Approval\nfor Ethnicity by\nSelf-Employment status") +
553   theme_minimal(base_size = 15)

```

```

554
555 grid.arrange(meff_odor_self, meff_uria_self, meff_white_self, ncol = 3)
556
557 # Dispersion Parameter ----
558 E2 <- resid(reduced_model, type = "pearson")
559 N <- nrow(mortg)
560 p <- length(coef(reduced_model))
561 sum(E2 ^ 2) / (N - p)
562
563 check_overdispersion <- function(logit_model) {
564   residual_df <- df.residual(logit_model)
565   pearson_resid <- residuals(logit_model, type = "pearson")
566   chi_squared <- sum(pearson_resid ^ 2)
567   dispersion_ratio <- chi_squared / residual_df
568   p_value <-
569     pchisq(chi_squared, df = residual_df, lower.tail = FALSE)
570   c(
571     chi_sq = chi_squared,
572     disp_ratio = dispersion_ratio,
573     res_df = residual_df,
574     p_val = p_value
575   )
576 }
577
578 round(check_overdispersion(reduced_model), 5)
579
580 #Appendix ---
581 ## Oversampling ----
582 # Splitting the data into majority and minority
583
584 minority_data <- mortg[mortg$approved == "No", ]
585 majority_data <- mortg[mortg$approved == "Yes", ]
586
587 # Oversampling minority class
588 oversampled_minority <-
589   minority_data[sample(nrow(minority_data), nrow(majority_data),
590     replace = TRUE), ]
591
592 # Combine back with majority class
593 balanced_data <- rbind(majority_data, oversampled_minority)
594 balanced_full_model <-
595   glm(approved ~ . * self,
596     family = binomial(link = 'logit'),
597     data = balanced_data)
598 summary(balanced_full_model)
599 balanced_reduced_model <-
600   glm(reduced_model$formula,
601     family = binomial(link = 'logit'),
602     data = balanced_data)
603 summary(balanced_reduced_model)
604
605 ## Lasso Regression ----
606 model_lasso <- glmnet::glmnet(
607   x = model.matrix(~ . * self, data = mortg[, -1]),
608   y = model.frame(mortg) |> model.response(),
609   alpha = 1,
610   family = "binomial"
611 )

```



```

612
613 lasso_cv <- glmnet::cv.glmnet(
614   x = model.matrix(~ . * self, data = mortg[, -1]),
615   y = model.frame(mortg) |> model.response(),
616   alpha = 1,
617   nfolds = 20,
618   family = "binomial"
619 )
620 plot(lasso_cv)
621 lasso_coef <-
622   coef(model_lasso, s = lasso_cv$lambda.1se , digits = 0.3)
623 lasso_coef
624
625 lasso_coef <- coef(model_lasso, s = lasso_cv$lambda.1se)
626
627 coef_df <- as.data.frame(Matrix::as.matrix(lasso_coef))
628 names(coef_df) <- c("Coefficient")
629 coef_df$Feature <- row.names(coef_df)
630 coef_df <- coef_df[-c(1:2),] # remove the intercept row
631
632 # Filter out zero coefficients for a cleaner plot
633 coef_df <- coef_df[coef_df$Coefficient != 0,]
634 row.names(coef_df) <- NULL
635
636 ggplot(coef_df, aes(x = Feature, y = Coefficient)) +
637   geom_hline(yintercept = 0,
638             color = "red",
639             linetype = "dashed") +
640   geom_point() +
641   coord_flip() + # Flip the axes to make it horizontal
642   labs(x = "Features", y = "LASSO Coefficients",
643        title = "Non-zero LASSO Coefficients at lambda.1se") +
644   theme_minimal()
645
646 ## Weighted Logistic Regression ----
647 weights <- ifelse(
648   mortg$approved == 1,
649   nrow(mortg) / sum(mortg$approved == "Yes"),
650   nrow(mortg) / sum(mortg$approved == "No")
651 )
652
653 weighted_full_model <-
654   glm(
655     approved ~ . * self,
656     family = binomial(link = "logit"),
657     data = mortg,
658     weights = weights
659   )
660 summary(weighted_full_model)
661 weighted_reduced_model <-
662   glm(
663     approved ~ hir + odir + lvr + mcs + self + single +
664       odir * self + self * white + self * uria,
665     family = binomial(link = "logit"),
666     data = mortg,
667     weights = weights
668   )
669 summary(weighted_reduced_model)

```