

Logistic Regression Review

What we are to do when trying to solve a logistic Regression problem:

① Identify Target variable: Binary or ordinal

~~Binary or ordinal~~

② Check for Rare event $\geq 5\%$

If we have a Rare event we can do the following:

- OverSampling: duplicating the rare event observations

- UnderSampling: Randomly sample the non-rare events to make them smaller

Technique to correct the bias from the above methods:

	Model correct	Model misspecified
Small sample $n \leq 1000$	Adjust intercept	weighted observations
Large sample $n > 1000$	Either	weighted observations

③ Split into Train, validation, & Test sets

④ Identify the predictor variables:

- If Continuous:

→ build individual Logistic Regression models for each variable. (MLE)

→ Evaluate if variable is statistically significant (p-values)

- If Ordinal:

→ Check for Separation Concerns

$\begin{matrix} 50 & 0 \\ 0 & 50 \end{matrix}$

Linear separation: When a combo perfectly predicts every outcome

$\begin{matrix} 50 & 23 \\ 0 & 22 \end{matrix}$

Quasi Separation: When a combo perfectly predicts a subset of outcomes

↳ To correct the above cases we can fix it through thresholding

→ Check if variable is significant w/ Mantel-Haenszel χ^2 test.

- If Nominal:

→ Check for Separation Concerns

↳ To correct issues we cluster categories

with issues to their most common category. ~~example~~

$\begin{matrix} A & B & C \\ a & 3 & 11 \\ b & 3 & 11 \\ c & 5 & 0 \end{matrix} \rightarrow \begin{matrix} A/C & B & C \\ a & 13 & 3 \\ b & 3 & 11 \end{matrix}$

→ Pearson χ^2 test to see if variable is significant.

⑤ Variable Selection:

→ Forward, backward, & Stepwise

↳ AIC, BIC, p-value

Logistic Regression Review

2

6. Adding Interactions: (be careful with separation issues)

large β w/ large p-value
↓
Signal something is wrong

7. Check Model Diagnostics:

Deviance = measure of how far our model is from the saturated model.

Through Deviance we can calculate: Cook's D, DFBetas, DfDev

8. Model Evaluation:

- Likelihood based calculations: (AIC, BIC, Generalized R^2)

- Probability Metrics: (Rank order observations)

→ Concordance, discordance, # Ties

→ Discrimination Slope = $(\hat{\beta}_1 - \hat{\beta}_0)$ bigger is better

- Classification Metrics: $\frac{TP}{TP+FN}$ $\frac{TN}{TN+FP}$

→ Sensitivity (TPR) vs. Specificity (TNR)

→ ROC Curve (we want the curve to be as far from the diagonal as possible)

→ KS Statistic (we want the two lines to as far as possible from each other)

→ Precision = $\frac{\# \text{ true positive}}{\# \text{ true positive} + \# \text{ false positive}}$

→ Precision vs. Recall (Sensitivity (TPR))

→ F_1 Score is how we select optimal values for the above.
The highest F_1 score is the optimal F_1 score.

→ Lift curve, Response curve, Gain chart

↓ Lift ↓ CRSP Rate ↓ CAP Rate

→ Lift curve: represents how much better the model is at predicting responses compared to random selection. Higher is better
Ex: "In the top 10% of data our model is X times more likely to identify responses compared to random selection of 10% of the data"

→ Response Curve: the proportion of responders in each bucket
Ex: "In the top 10% of our data our model's precision is X"

→ Gain Chart: Shows how well our model captures the total response across all buckets.

Ex: "By forgetting the top X percent of the data X% of the total responses can be captured using our model."

Optimal Point is
Youden's Index
↓
PR + (TNR - 1)

buckets are sorted by our model's predicted probability

↙ Don't use for model selection, only reporting!

- Accuracy & Error: $1 = \text{Accuracy} + \text{Error}$

↳ Accuracy = % of predictions our model got correct

↳ Error = % of predictions our model got wrong.

⑨ Interpreting Coefficients:

- Odds Ratios: $e^{\hat{\beta}}$ or $100 \times (e^{\hat{\beta}} - 1)\%$

\uparrow
x times more likely...

\uparrow
x% higher odds of...

⑩ Make Final Cut off: Cost/profit should always be how we choose cutoffs!

- If Cost/profit is equal, we choose cutoff from: Youden, KS, F₁

⑪ Create the Report for Stakeholder!