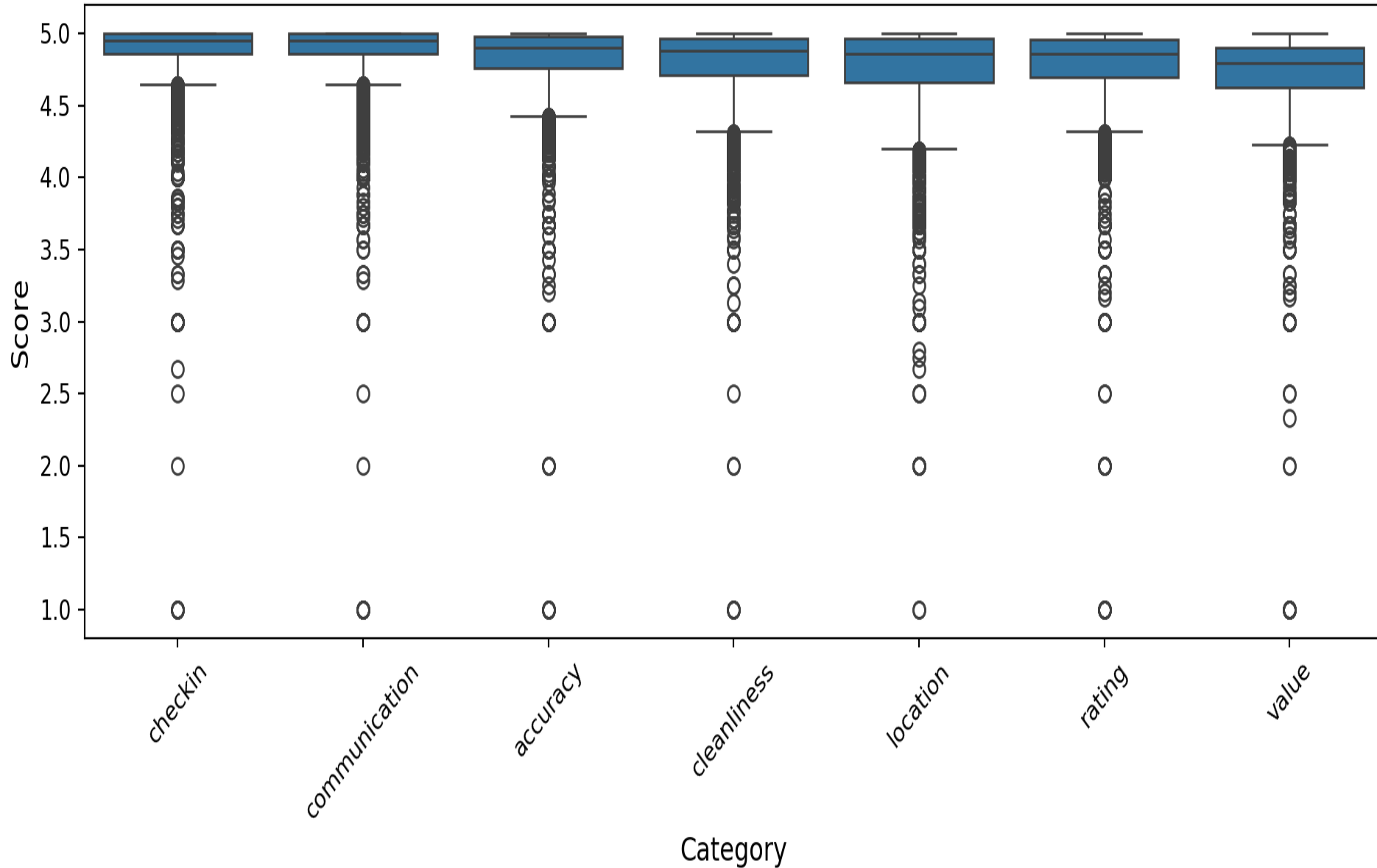


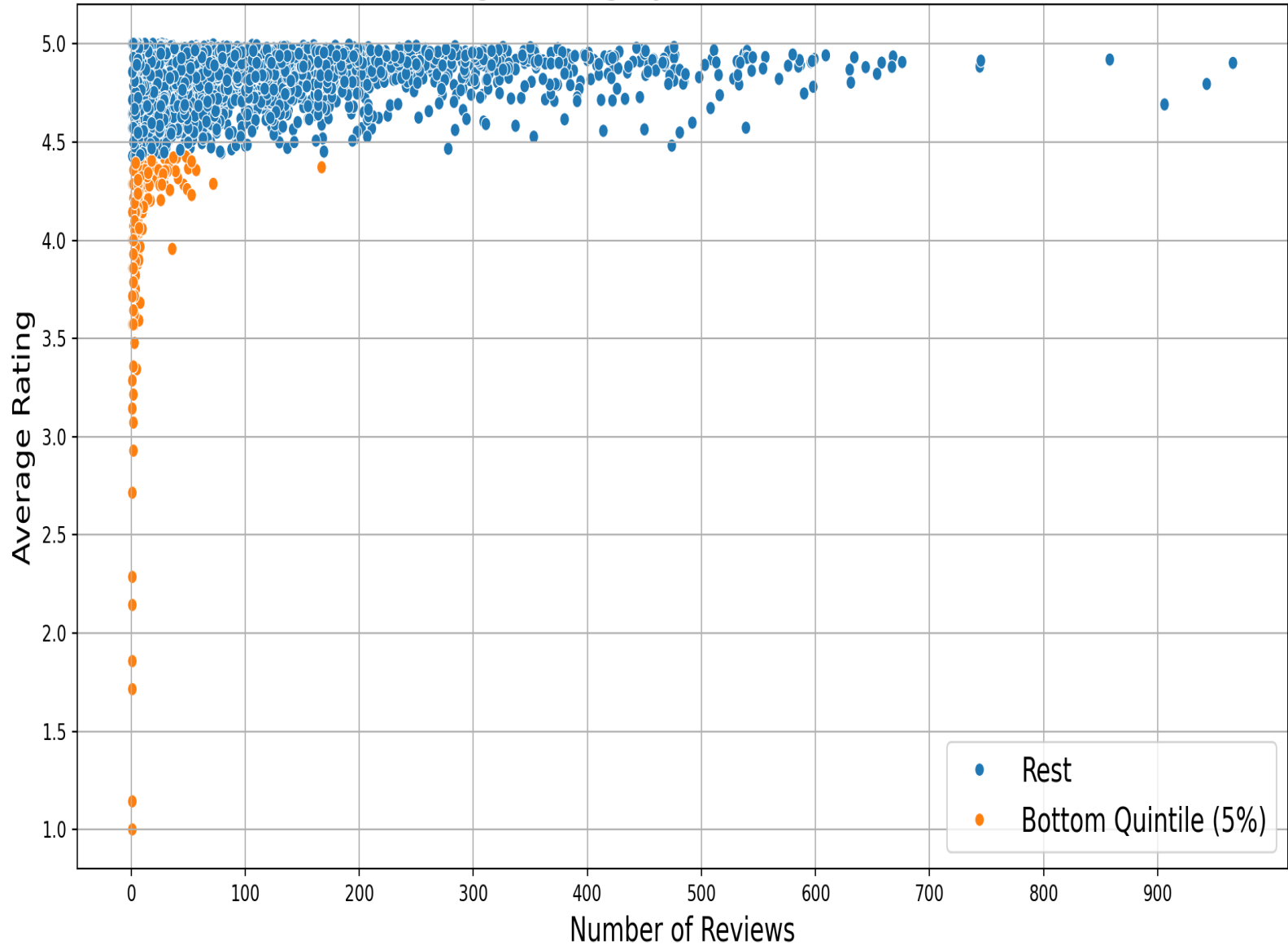
Predicting Poorly Rated AirBnB Listings in the DC Area Using Natural Language Processing



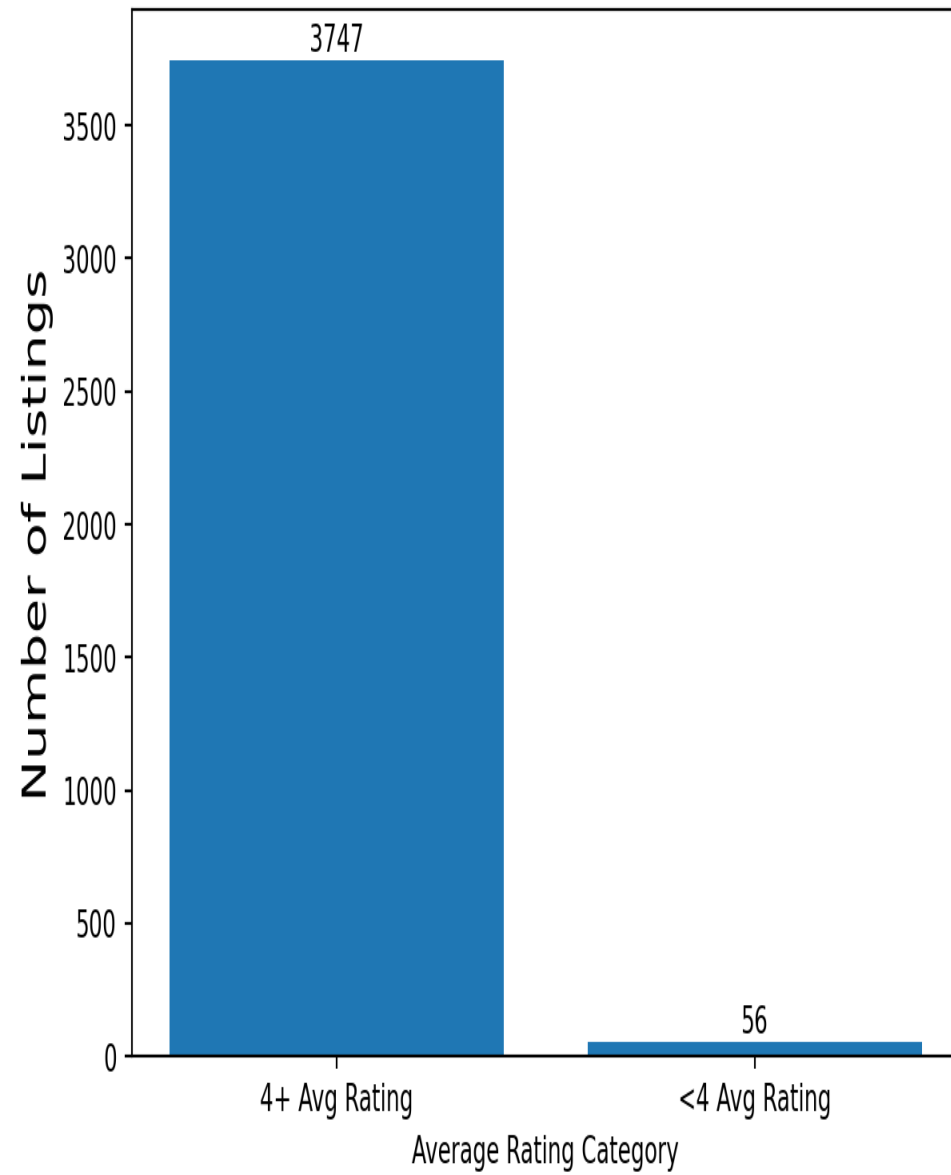
Distribution of Ratings by Category



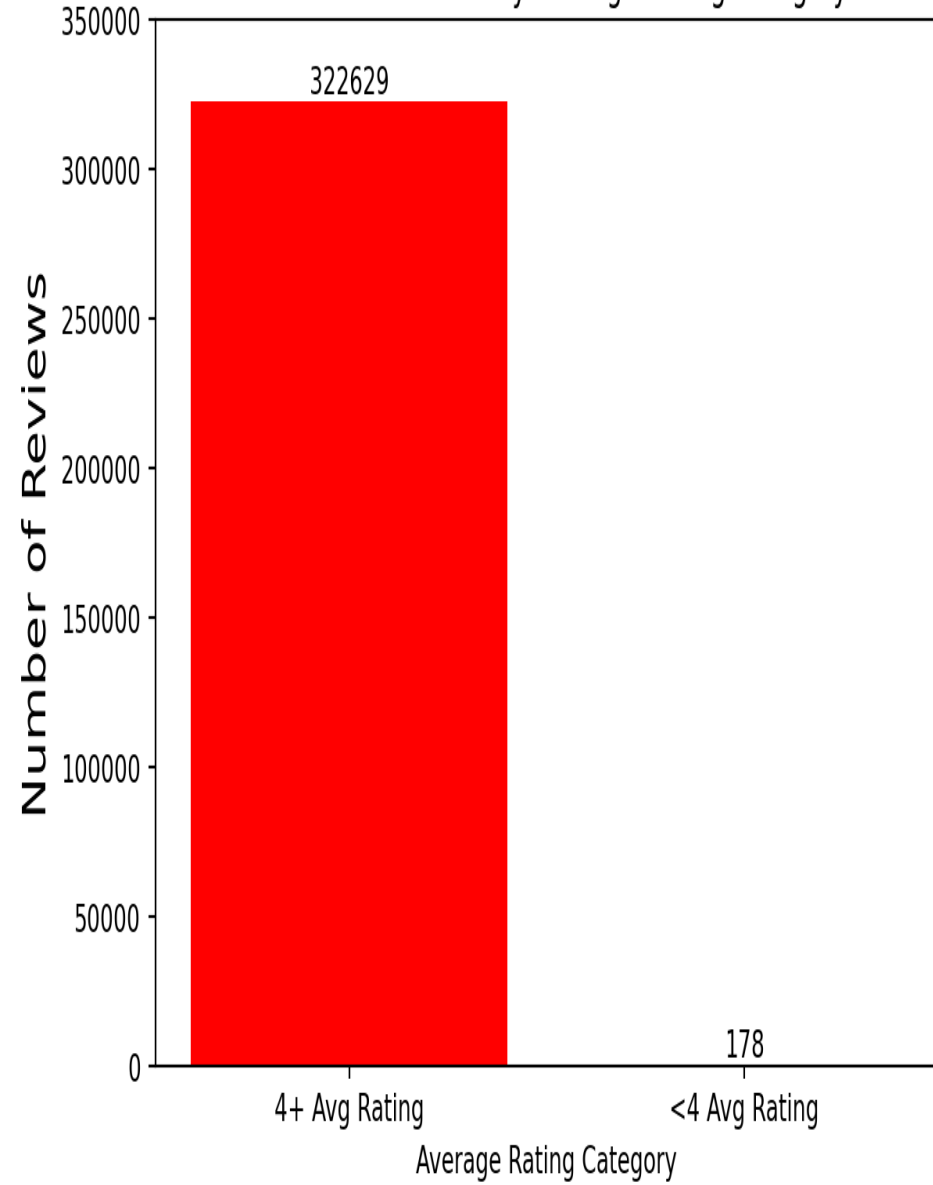
Average Rating by Number of Reviews



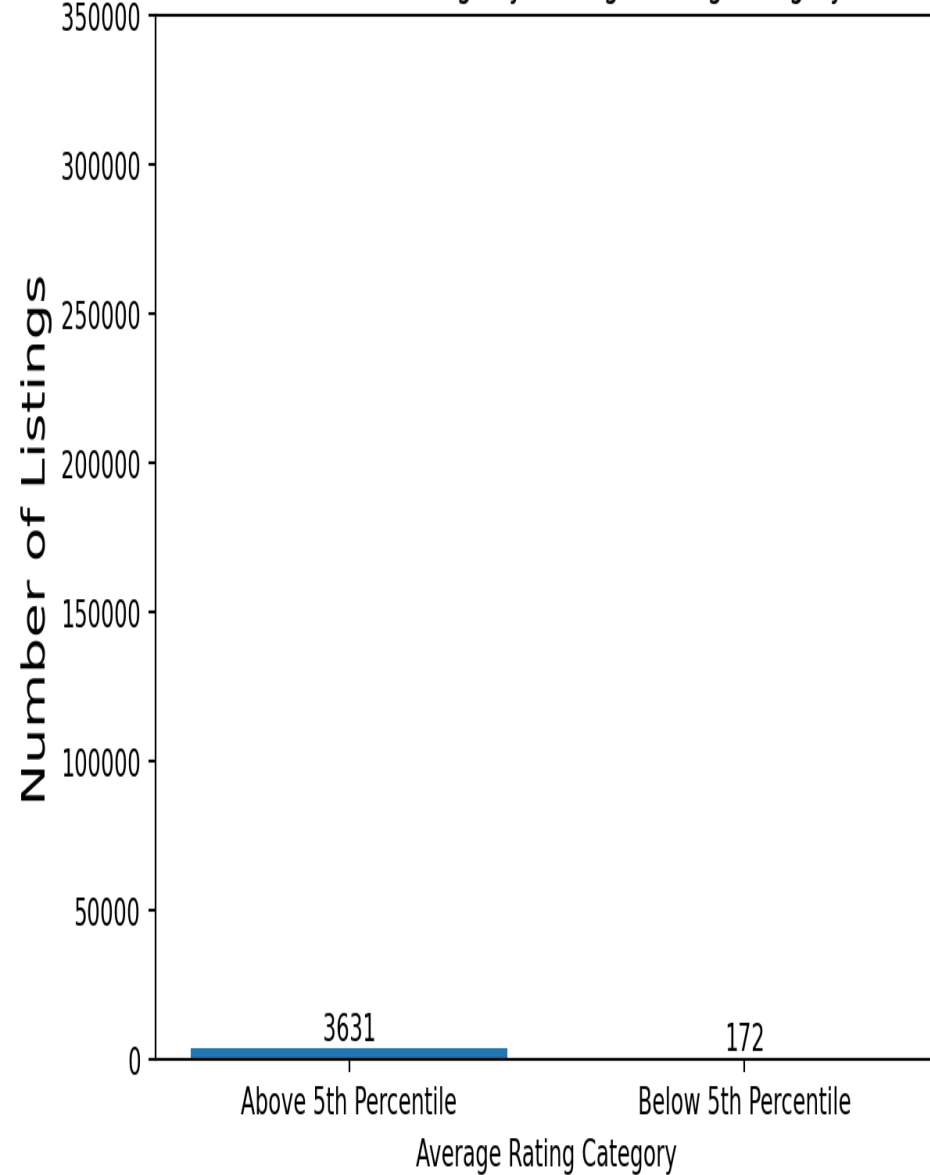
Number of Listings by Average Rating Category



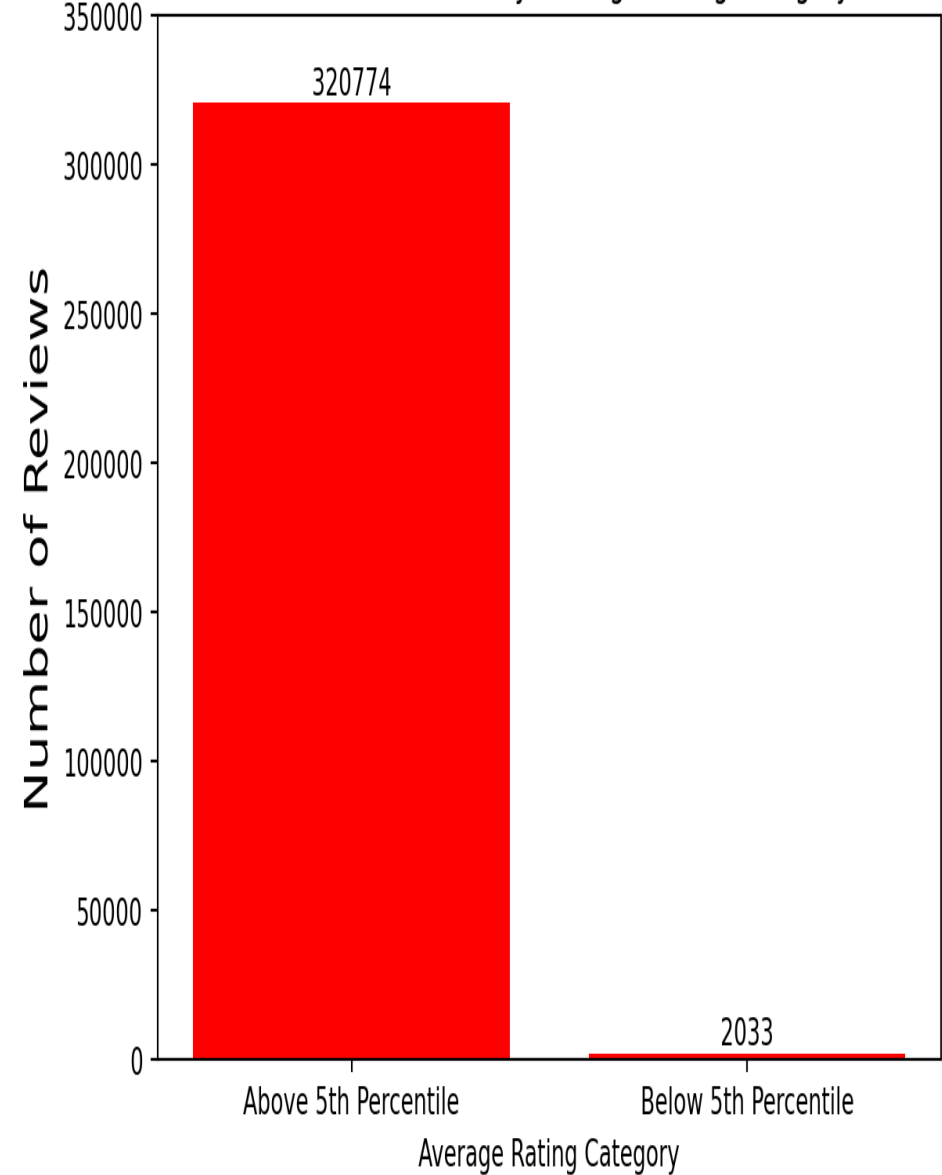
Number of Reviews by Average Rating Category



Number of Listings by Average Rating Category



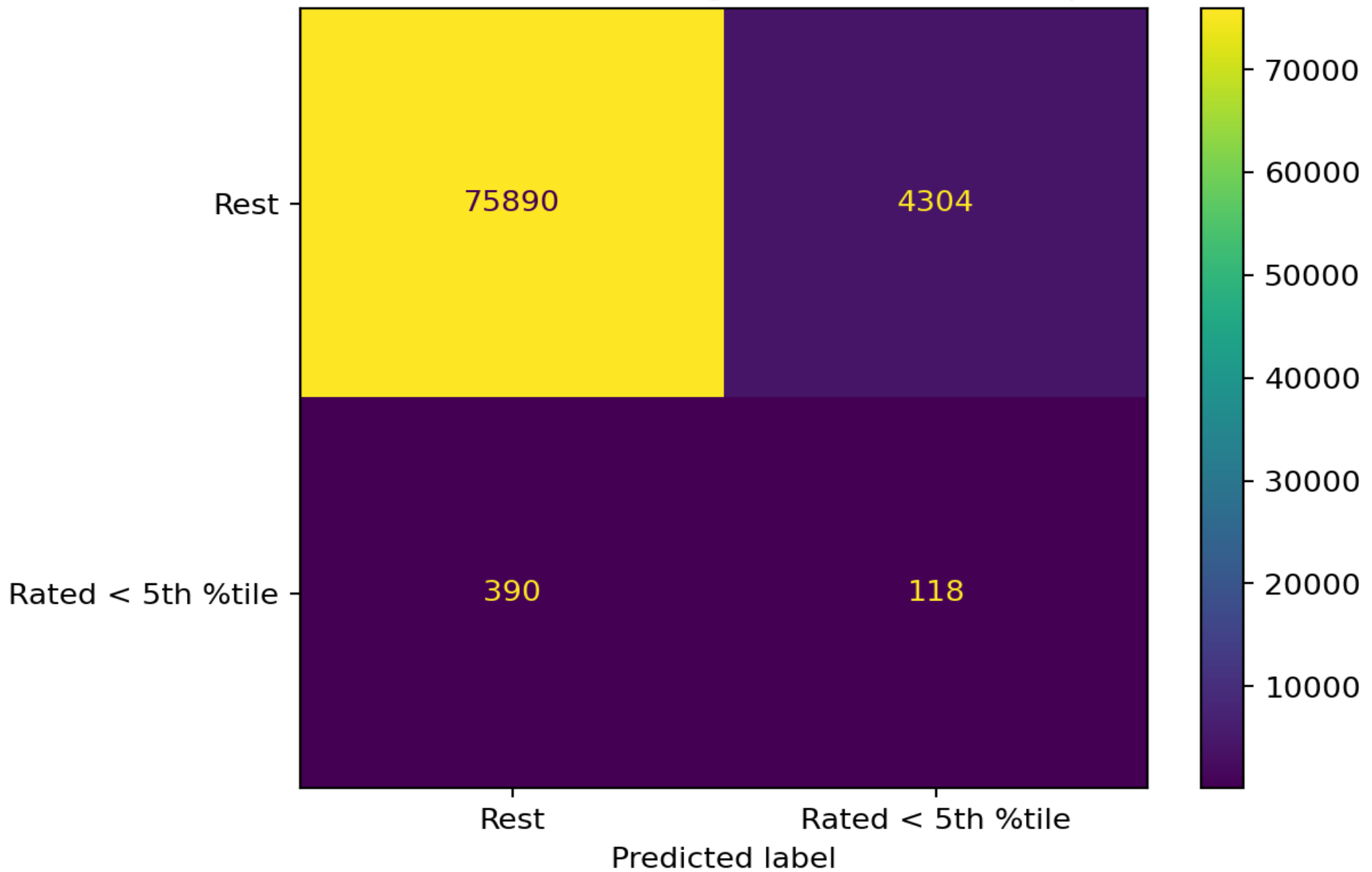
Number of Reviews by Average Rating Category



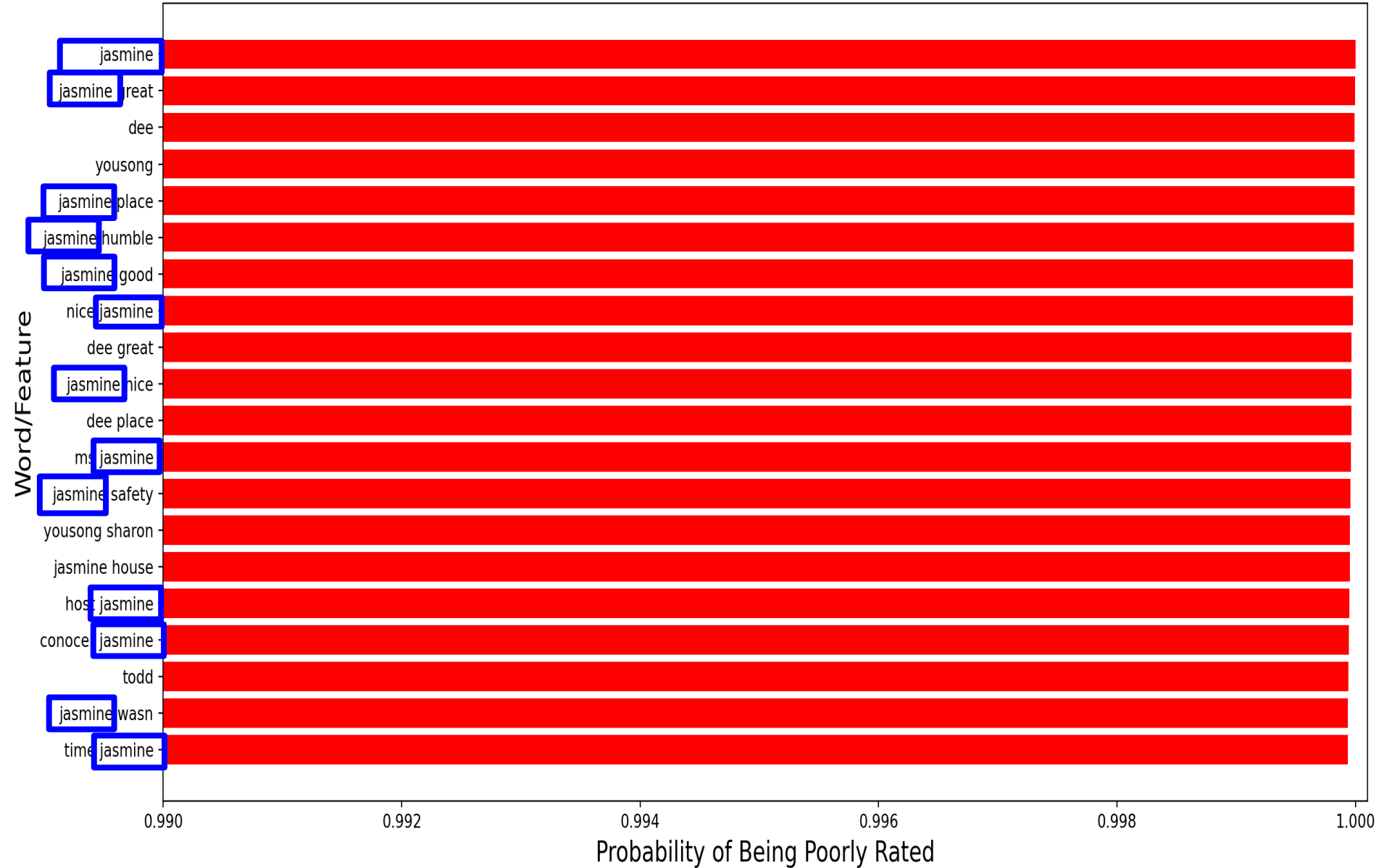
Model Iteration

Rating	Model	Recall Score	F1 Score
+/- 4.0 Avg Rating	Initial Model: TFIDF Vectorizer and Logistic Regression	0.000	0.000
+/- 4.0 Avg Rating	TFIDF & Logistic Regression <i>with Balanced Classes</i>	0.111	0.009
<i>5th Percentile</i>	TFIDF & Logistic Regression with Balanced Classes	0.343	0.059
5th Percentile	TFIDF & Logistic Regression <i>with Multinomial Naïve Bayes</i>	0.000	0.000
5th Percentile	<i>Count Vectorizer</i> and <i>Logistic Regression</i> with Balanced Classes	0.242	0.047
5th Percentile	<i>TFIDF Vectorizer</i> and Logistic Regression with <i>Random Over Sampler</i>	0.467	0.048
5th Percentile	TFIDF Vectorizer and Logistic Regression with Balanced Classes and <i>Custom Stop Words</i>	0.252	0.052
5th Percentile	TFIDF Vectorizer, Logistic Regression Scoring and Random Over Sampler, but <i>Changing Sampling Strategy and Scoring to Balanced Accuracy</i>	0.254	0.068

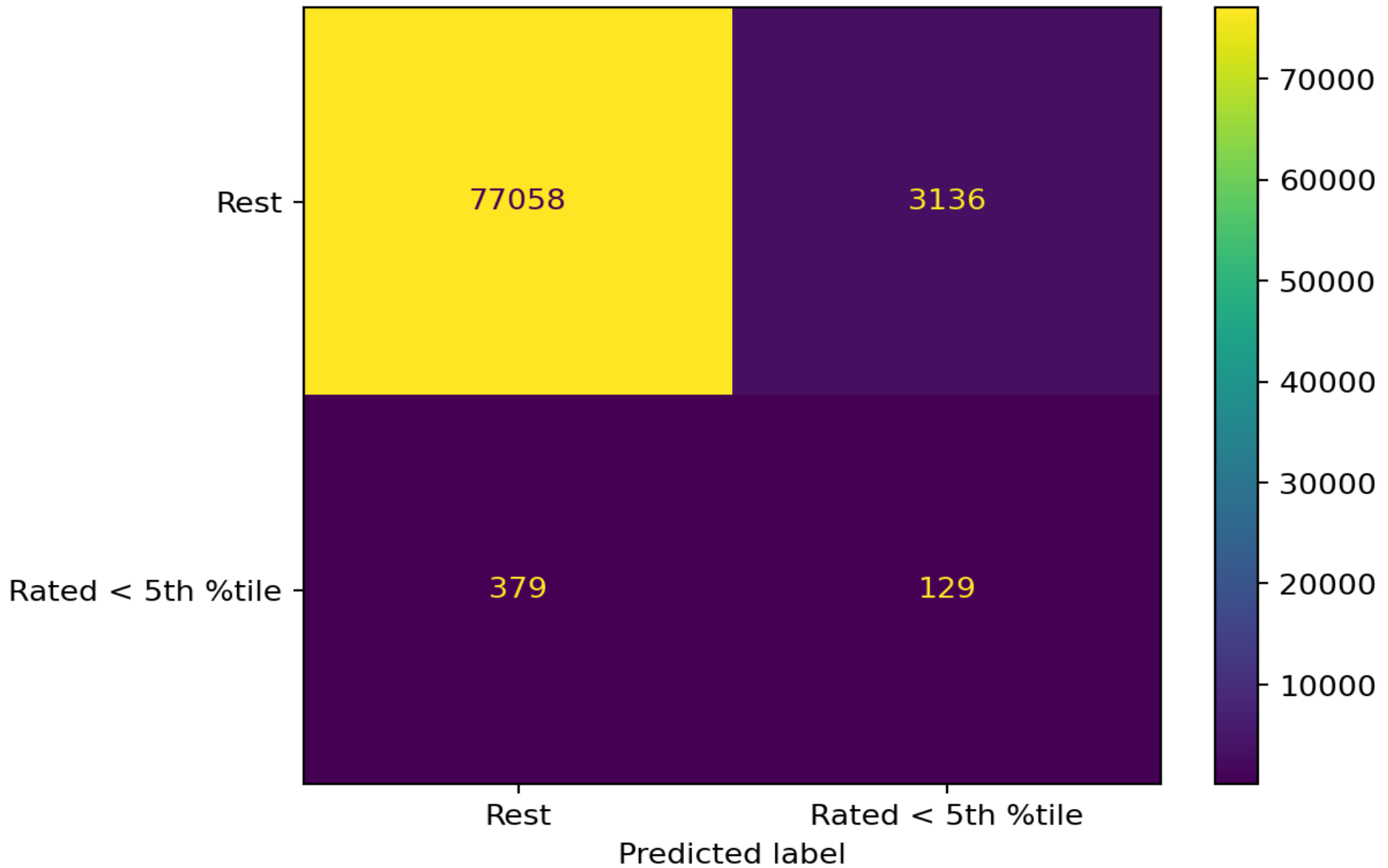
Confusion Matrix using Random Oversampler



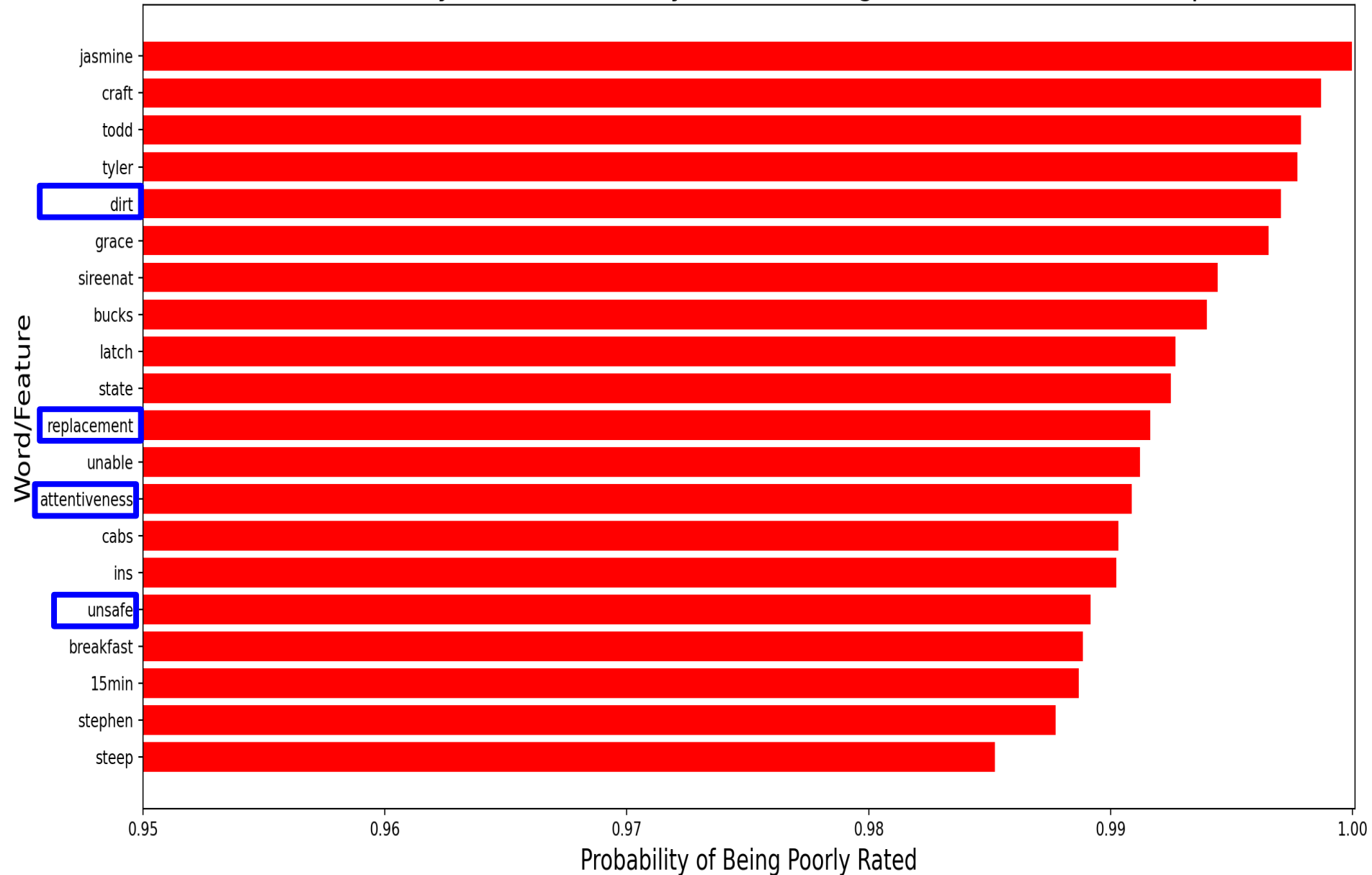
Features Most Likely to Predict Poorly Rated Listings in Random Oversample Model



Confusion Matrix using Balanced OverSampler



Features Most Likely to Predict Poorly Rated Listings in Balanced Oversample Model



Conclusions and Recommendations

- The issue with imbalanced classes made it difficult to produce a model that had a good recall score and balanced accuracy. Models had an issue with false positives due to over-learning the names of hosts/listings, particularly when over-sampling the minority class of poorly reviewed listings.

Caveats

- The rating is for the listing's average rating across ALL reviews, not the rating each reviewer gave to their stay. Someone who stayed at a listing that is overall poorly reviewed but could give a positive review themselves (though it should be cancelled out by other negative reviews). There's no data on how each individual reviewer rated the listing.

Future Considerations

- Started looking at SMOTE to synthesize minority class of poor reviews.
- Additionally, might look at lemmatizing/stemming words (to make “safety” and “unsafe” both be “safe”) and/or spaCy to develop context around words.
- Get data from multiple cities, especially on poorly reviewed listings. This could help the model avoid learning about specific poorly reviewed listings/hosts in one city and instead look at what the commonalities are about poorly reviewed listings across multiple cities.
- Was reluctant to do random under sampling and get rid of so much data from positively rated listings. But the solution may be to look at just listings with under 100 reviews. Since all poorly rated listings have under 100 stays, looking at poorly rated listings with under 100 stays and well rated listings with under 100 stays may be a better comparison.
- Created binary classification by drawing a semi-arbitrary line. But since most listings are positively rated, difference is less between above/below 4.0 or above/below bottom quartile, and more sliding scale on difference between a 4.9 and 4.8. Tobit regression may work better to determine what features/words improve/reduce a listing’s rating rather than classification model.