

1. Data Normalization

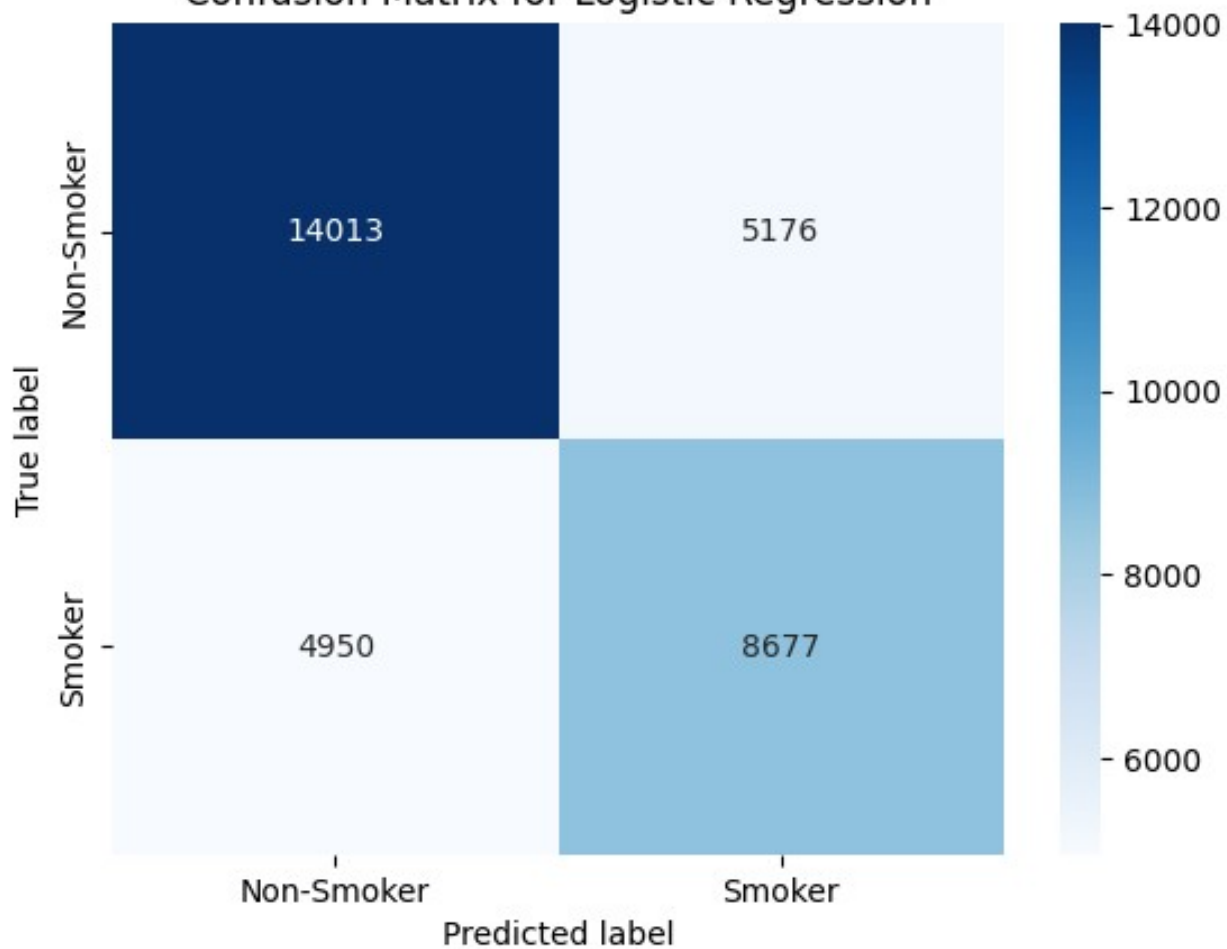
- transforms data into a standard scale without distorting differences in ranges.
- ensures all features contribute equally to the model.
- Common techniques include:
 - Min-Max Scaling: Scales features to a $[0,1]$ range.
 - Standardization: Centers data by subtracting the mean and dividing by the standard deviation.

StandardScaler is used here to standardize data, making it suitable for machine learning algorithms sensitive to feature magnitudes.

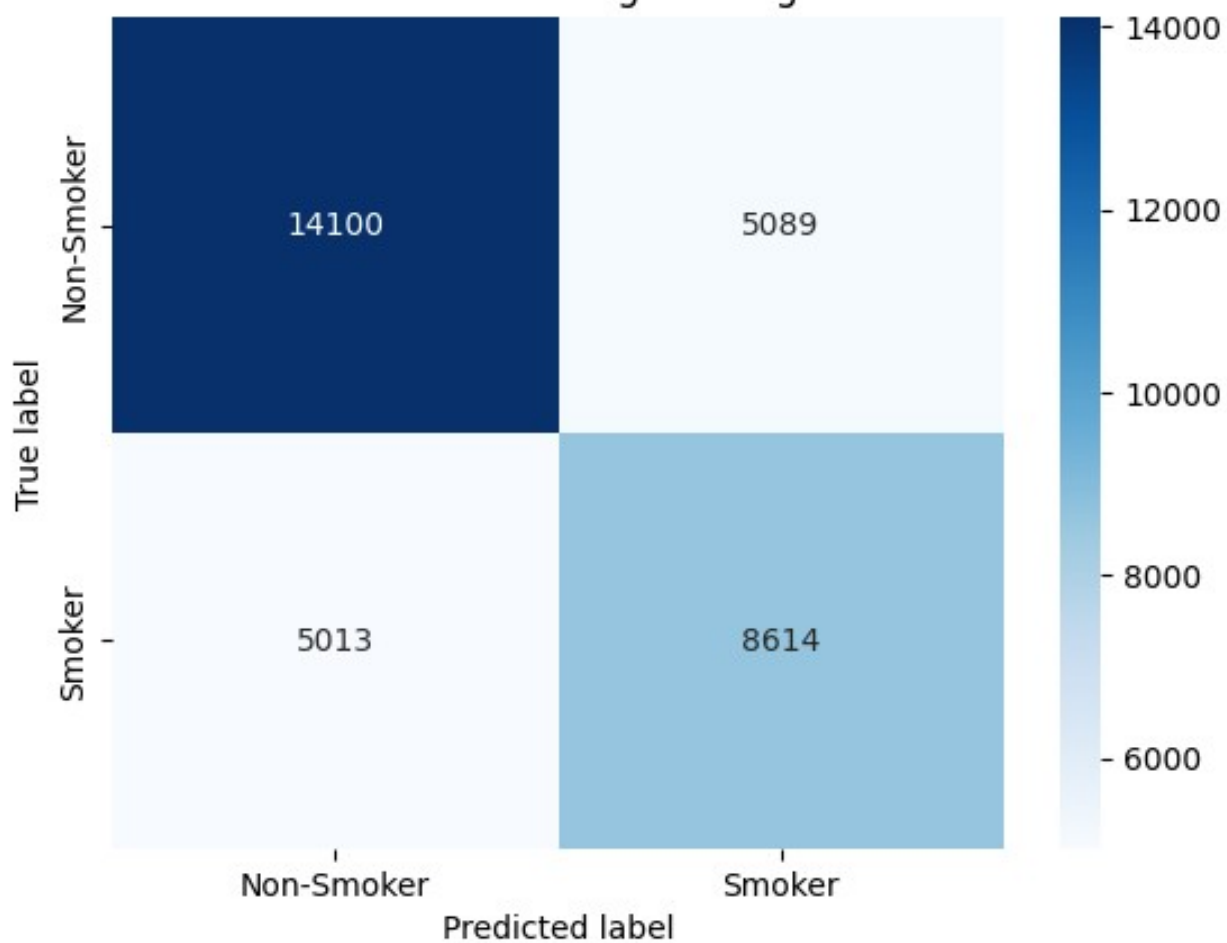
2. Visualization and Data Analysis

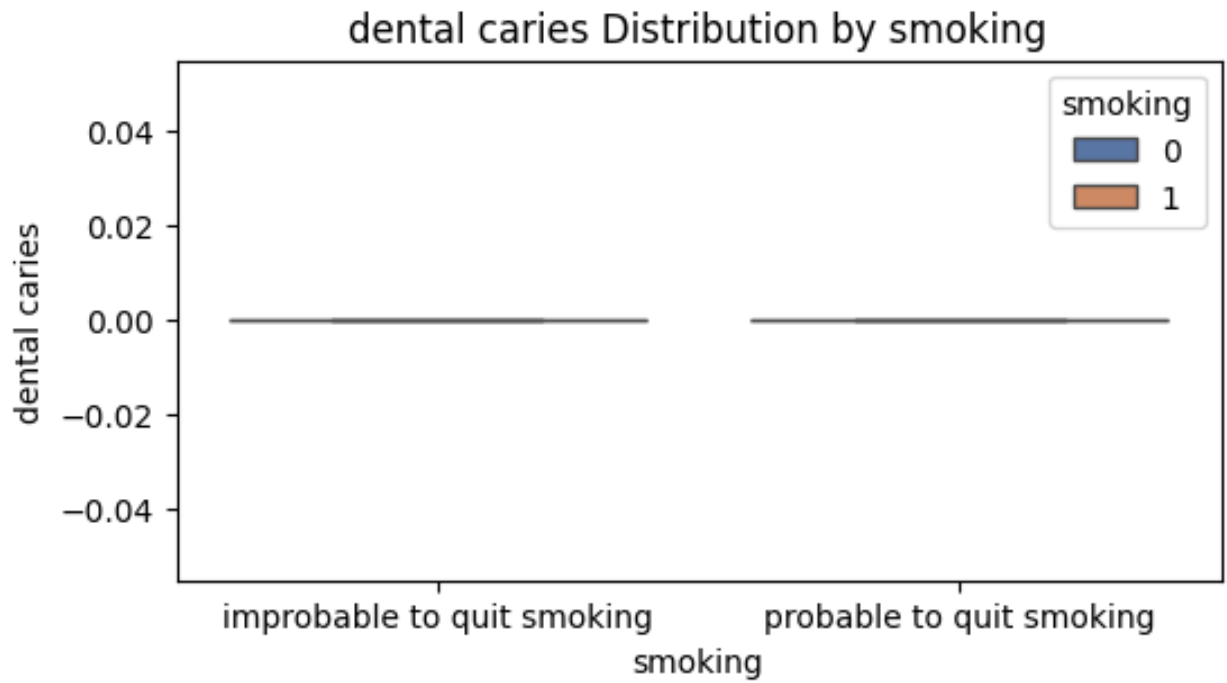
Univariate:

Confusion Matrix for Logistic Regression

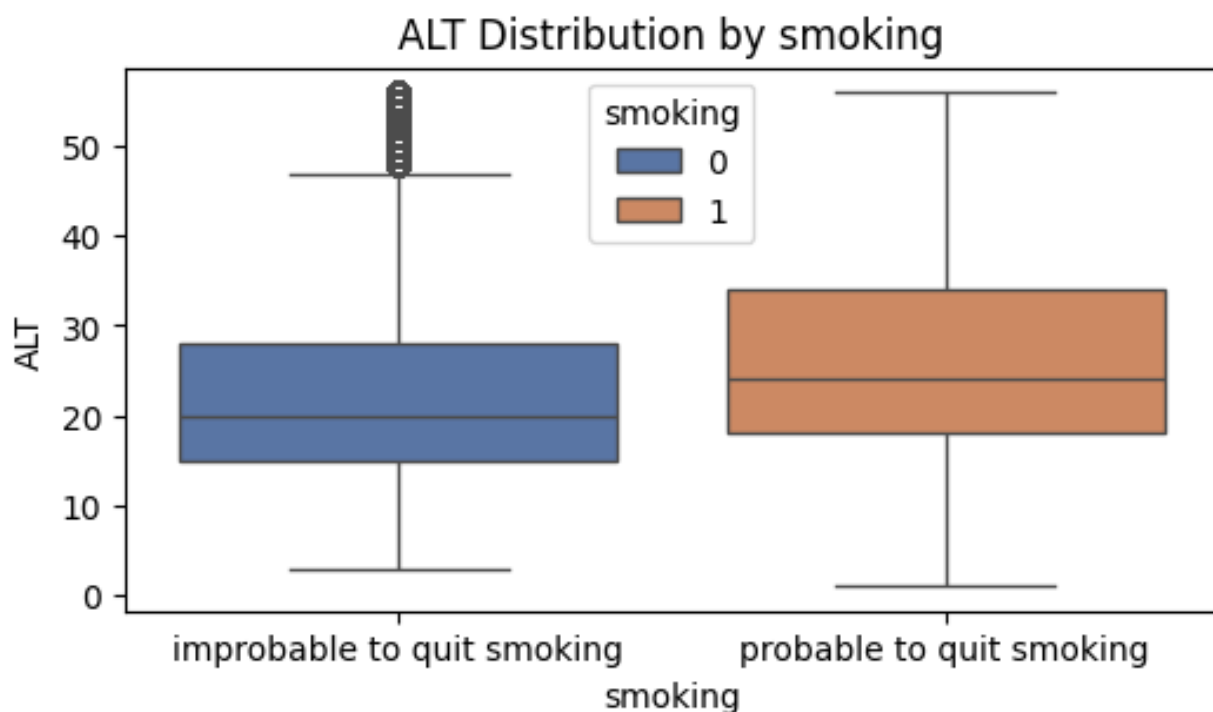
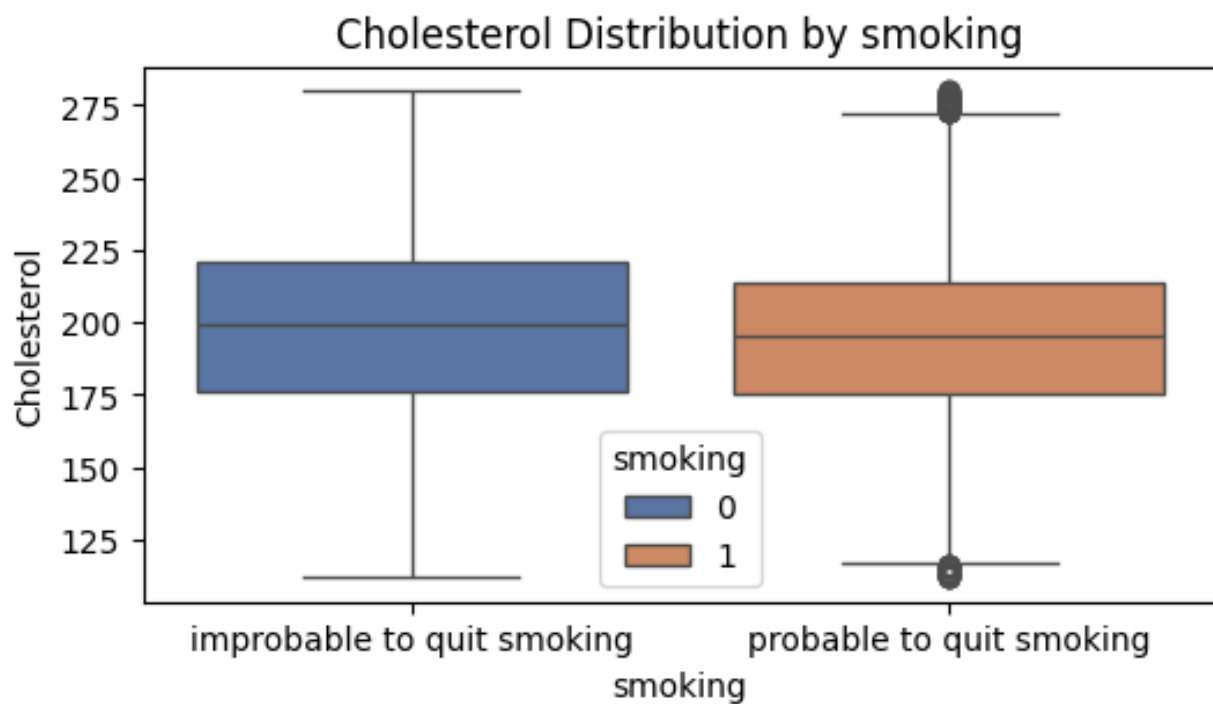


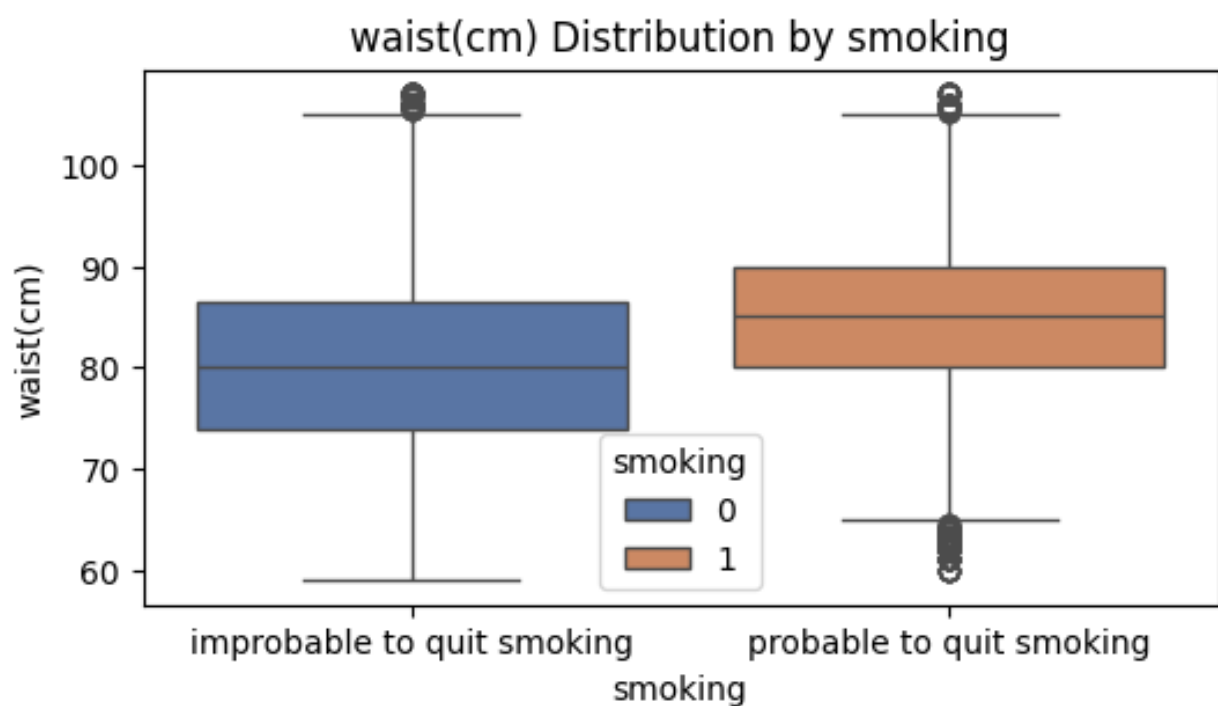
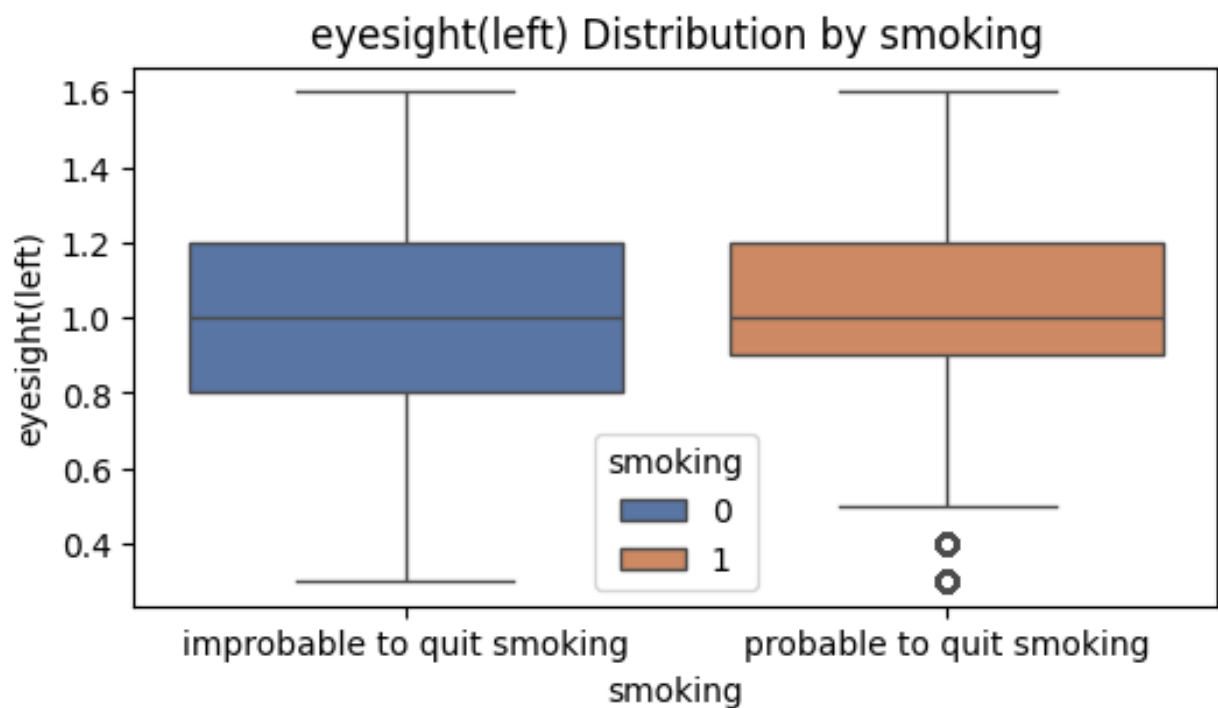
Confusion Matrix for Logistic Regression



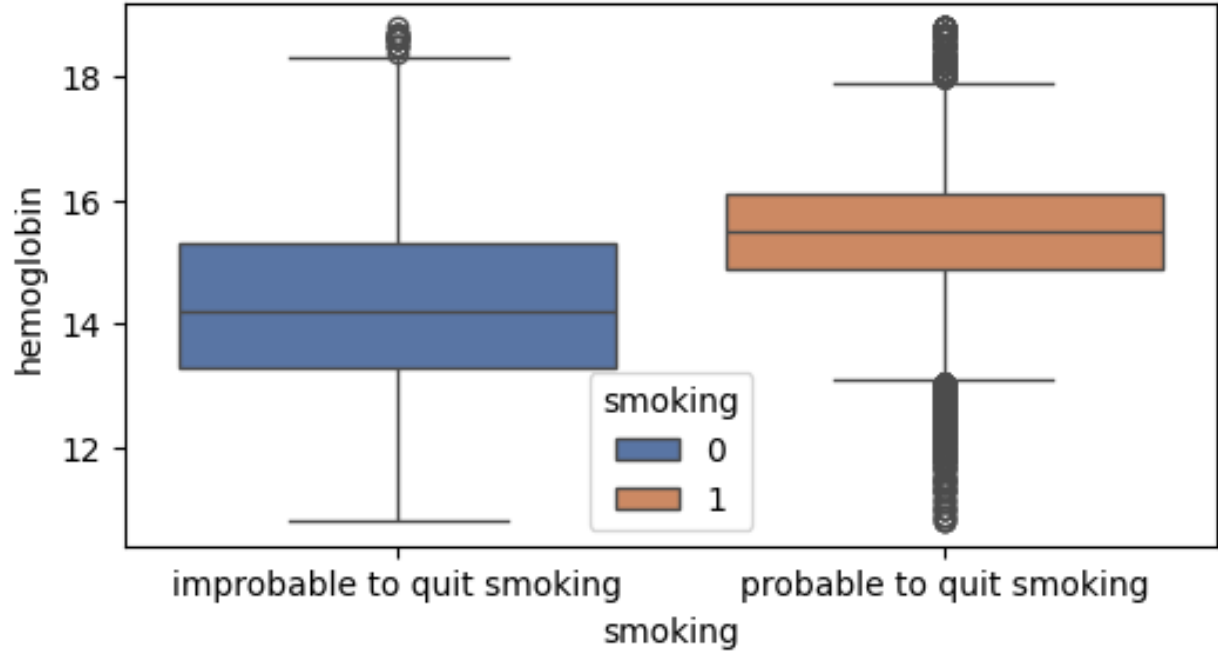


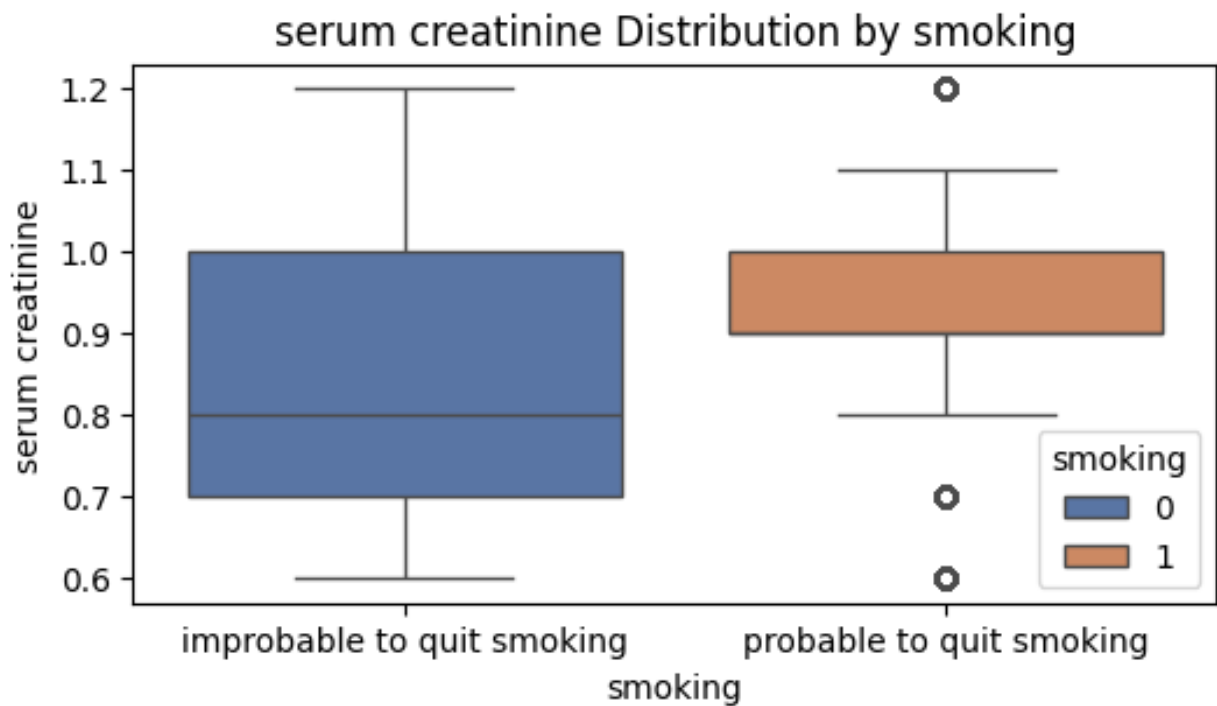
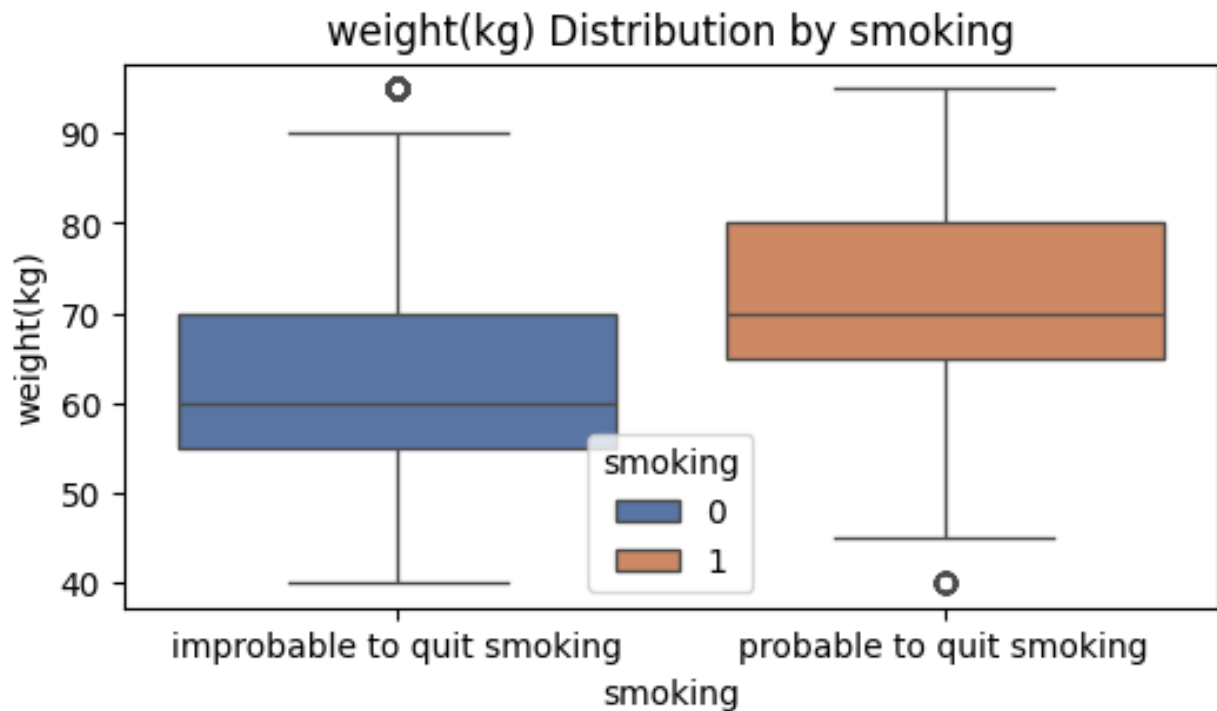
- The distribution for these 3 graphs shows a uniform spread across the range, indicating consistent representation without notable skewness or clustering.





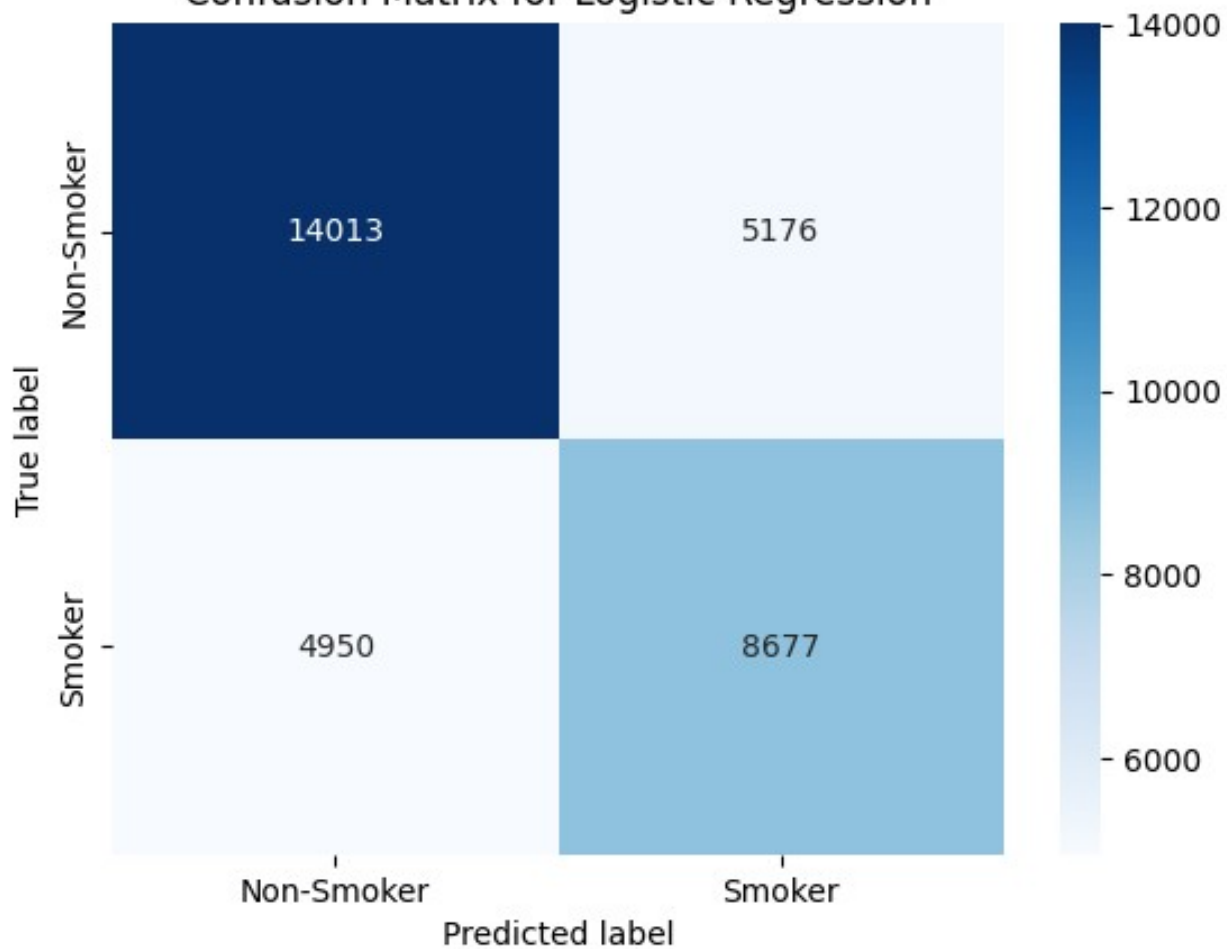
hemoglobin Distribution by smoking



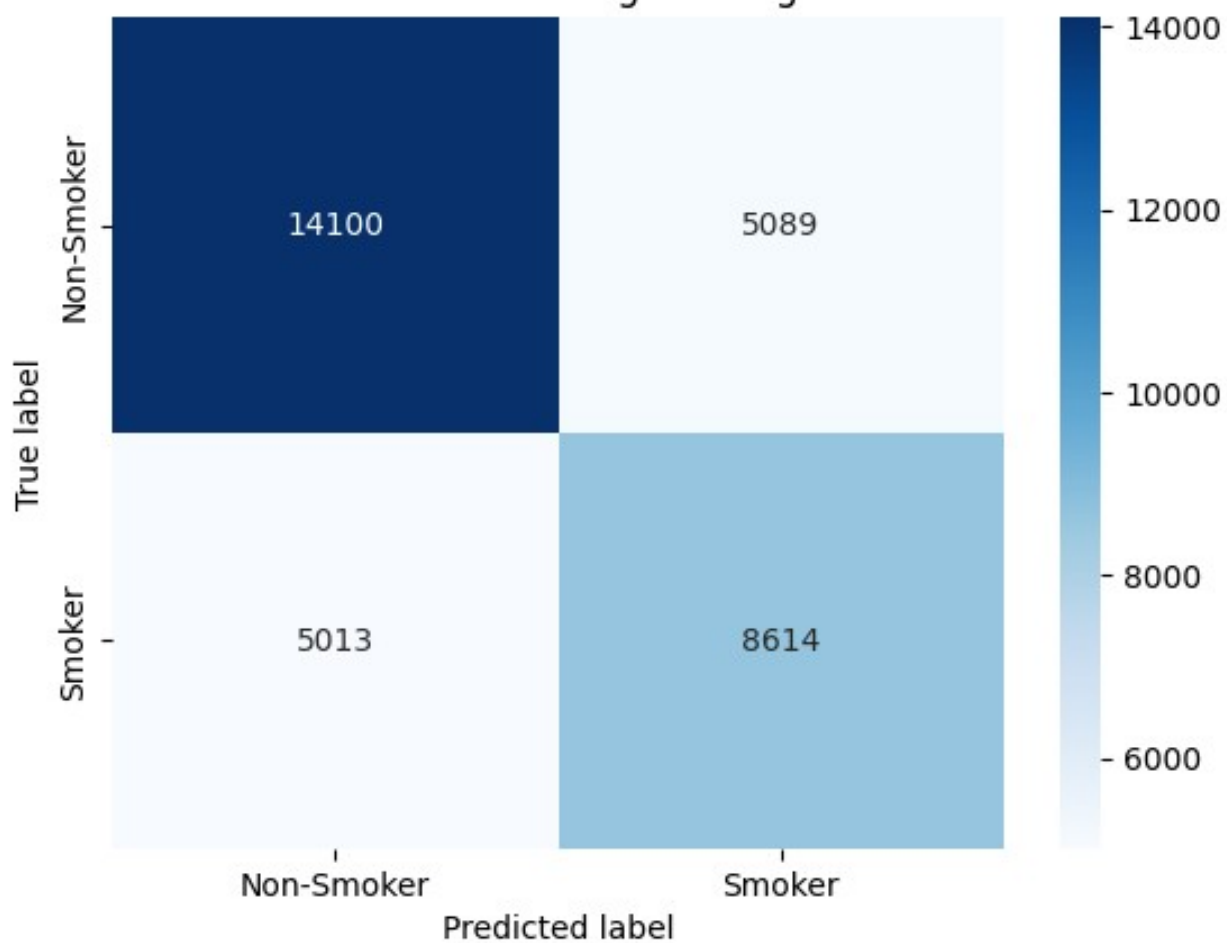


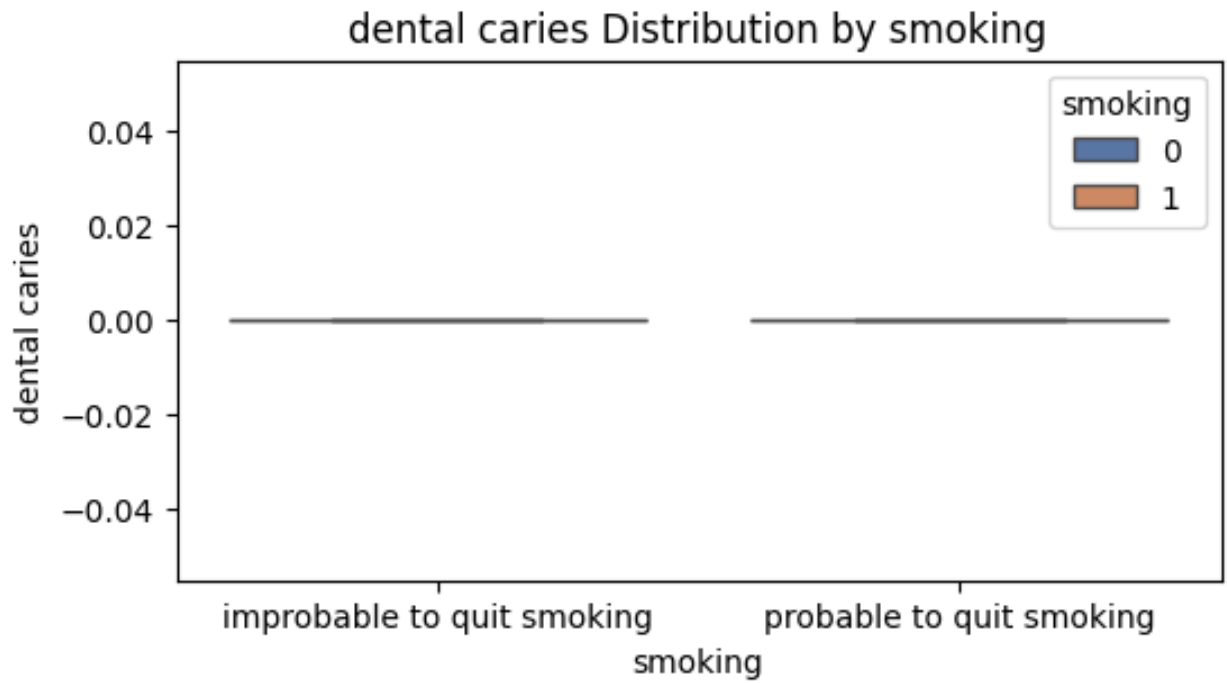
- These graphs are close to a gaussian distribution and most of them have almost zero mean.

Confusion Matrix for Logistic Regression

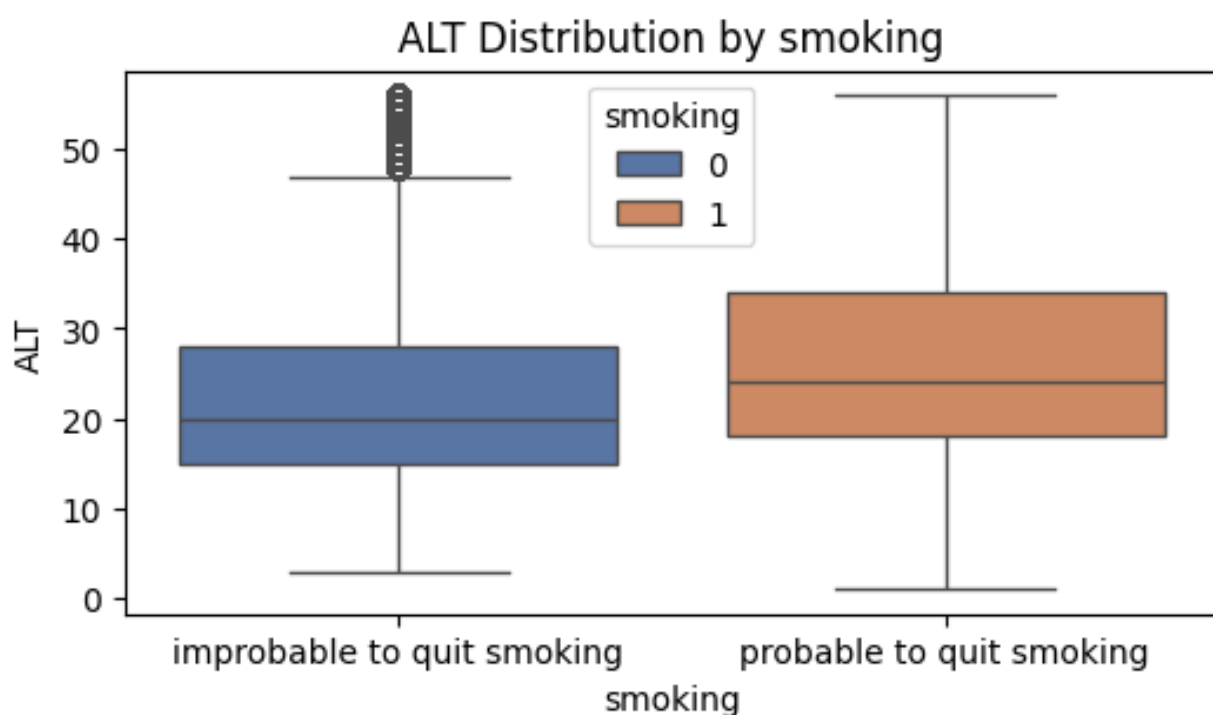
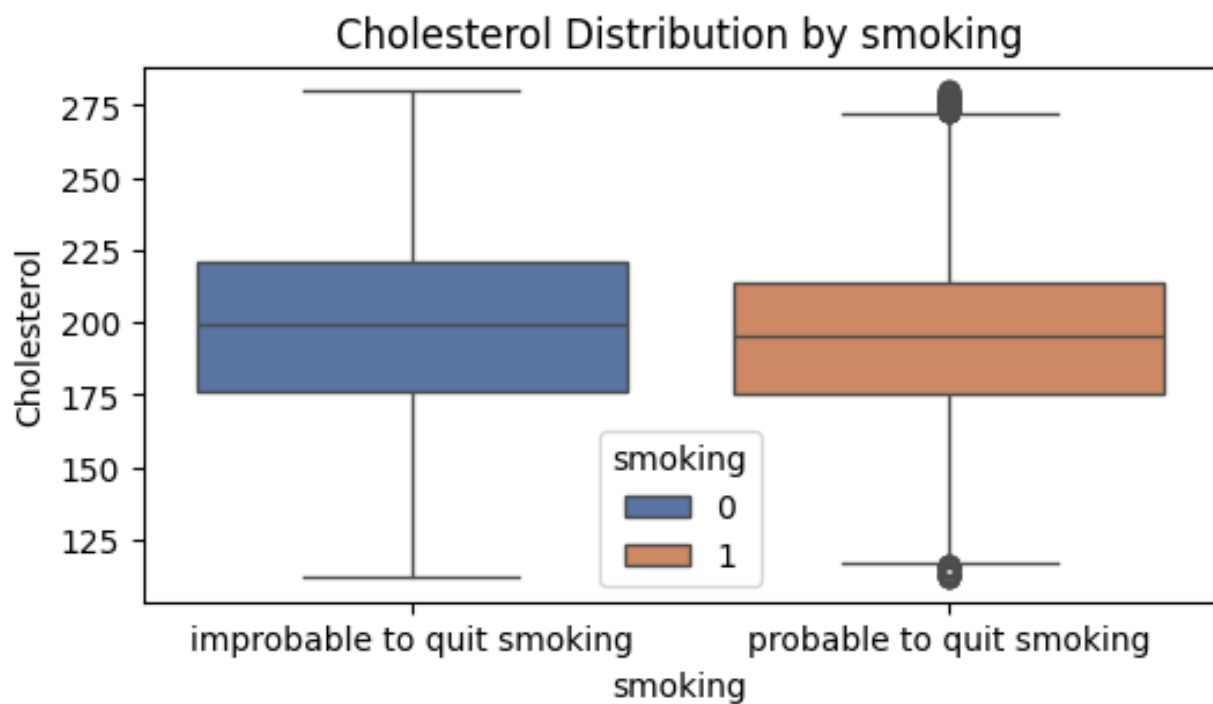


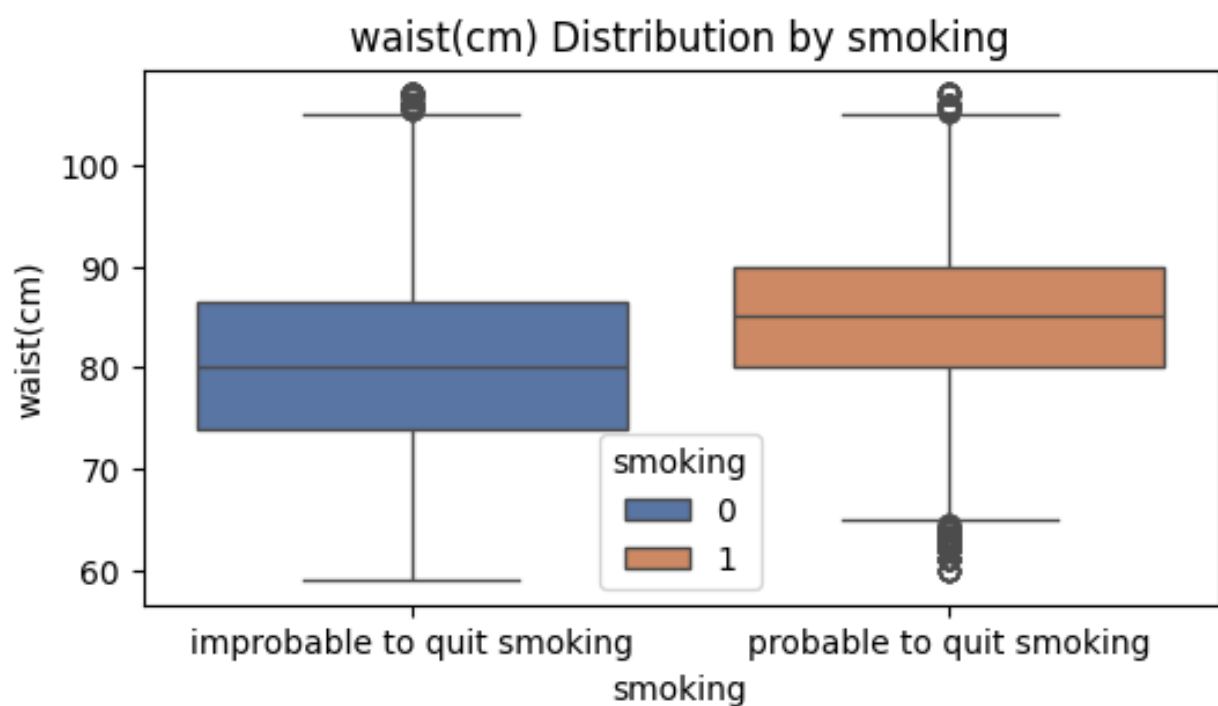
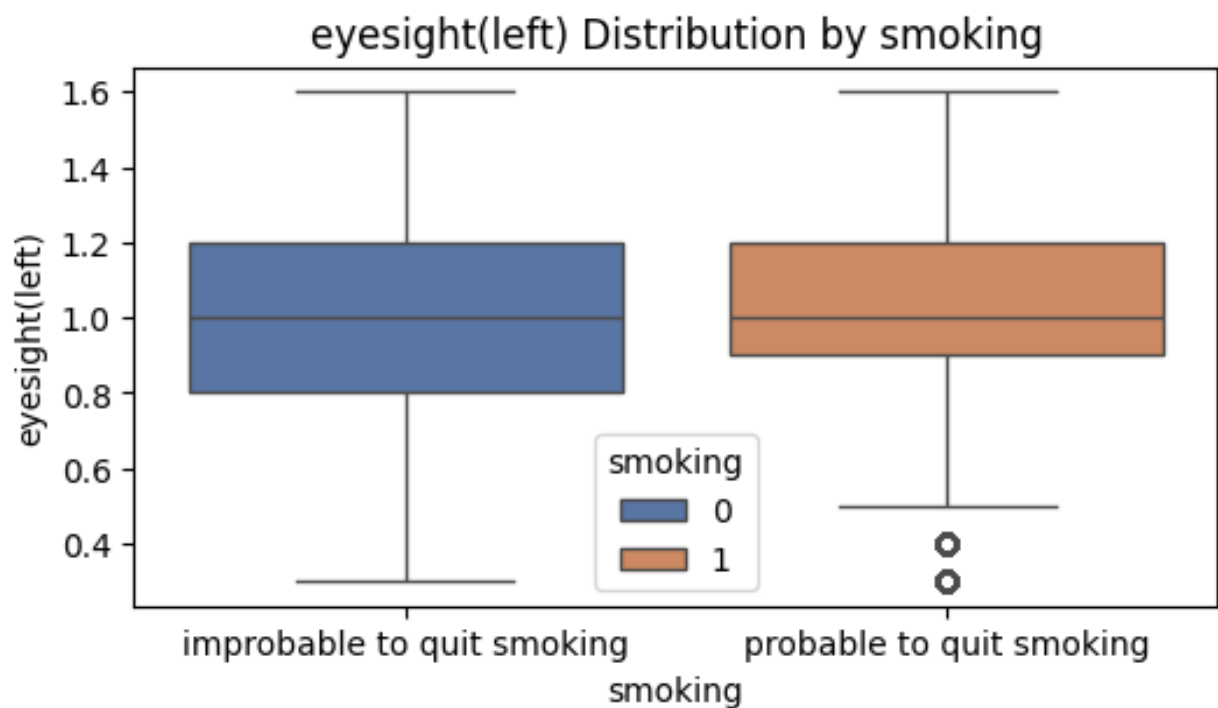
Confusion Matrix for Logistic Regression



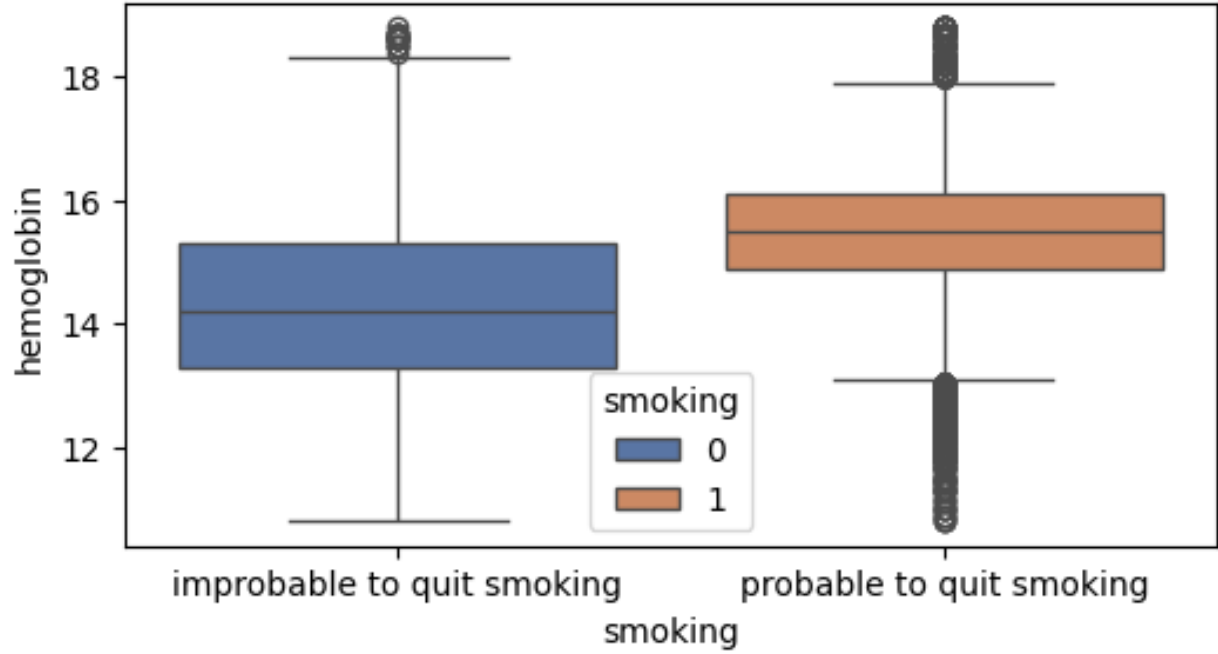


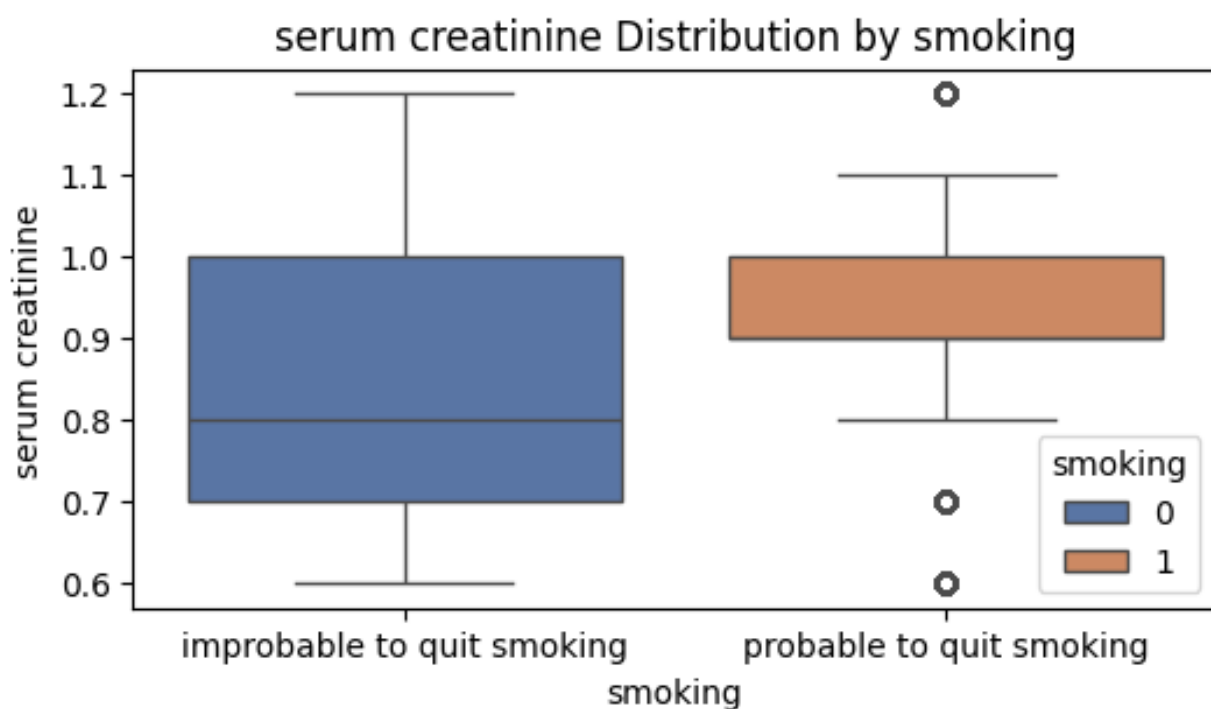
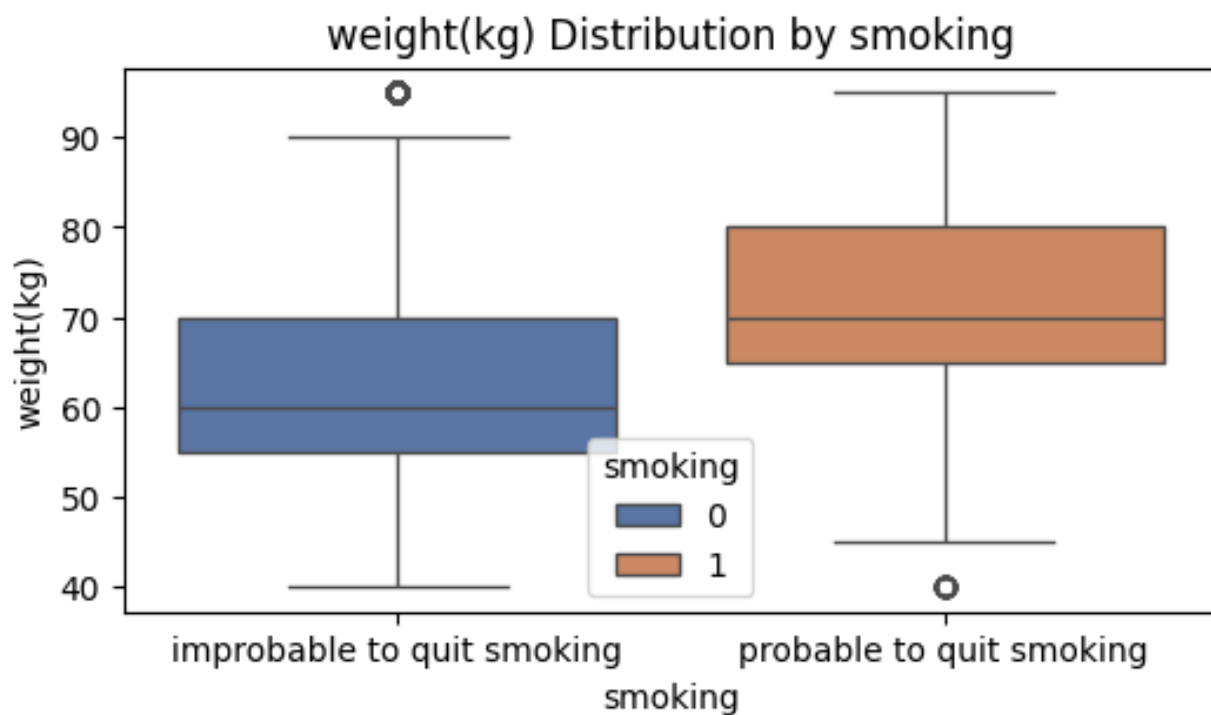
- Those graphs show uniform distribution indicating no relation to smoking



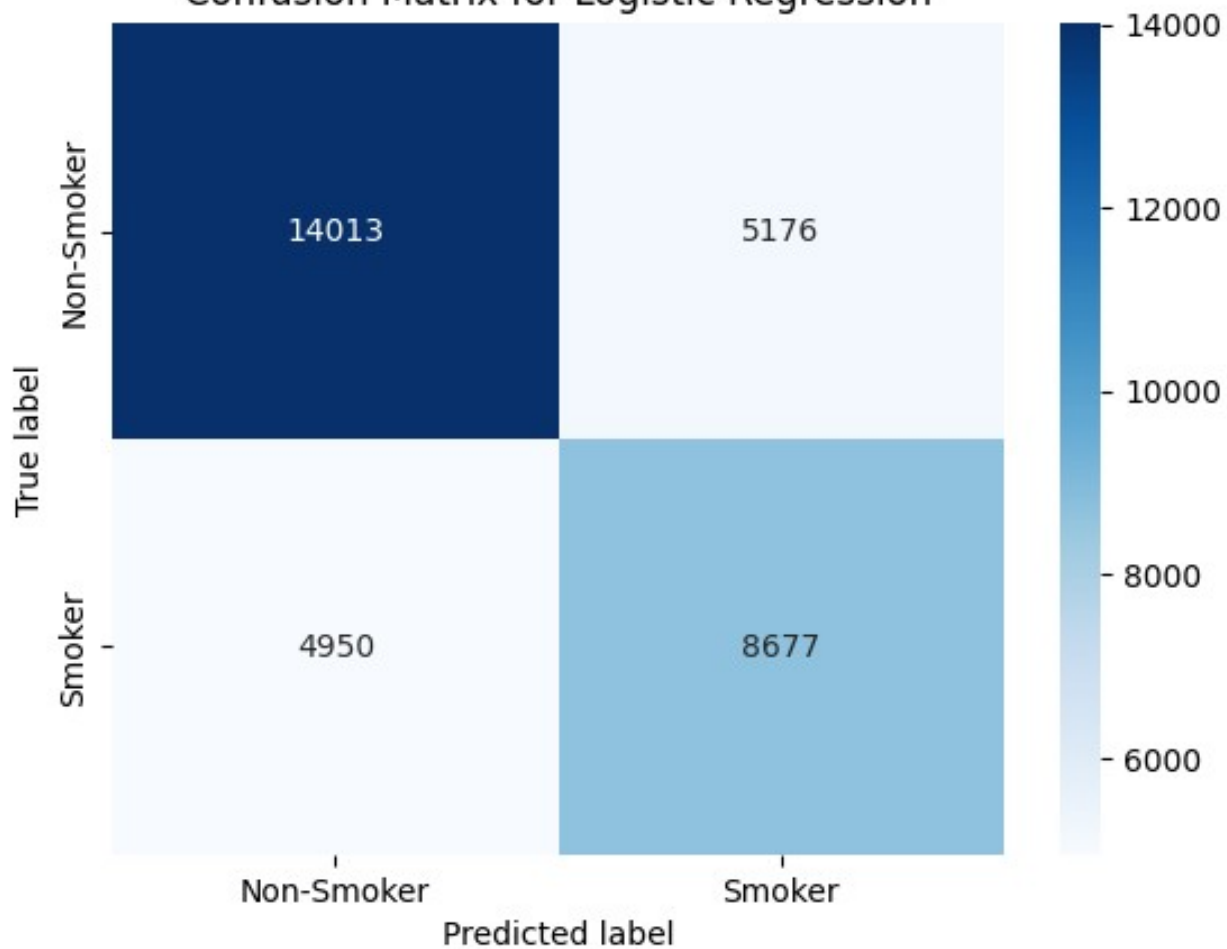


hemoglobin Distribution by smoking

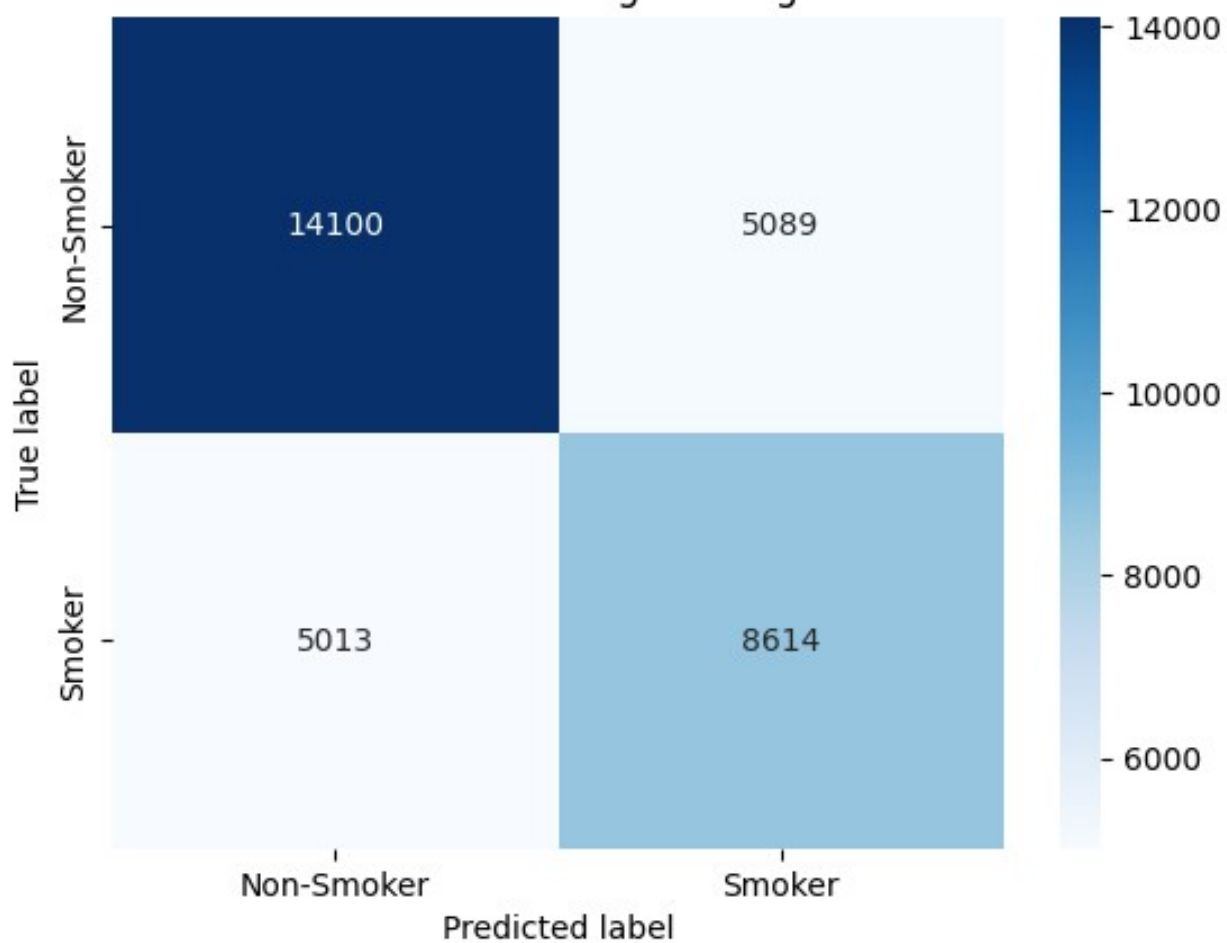


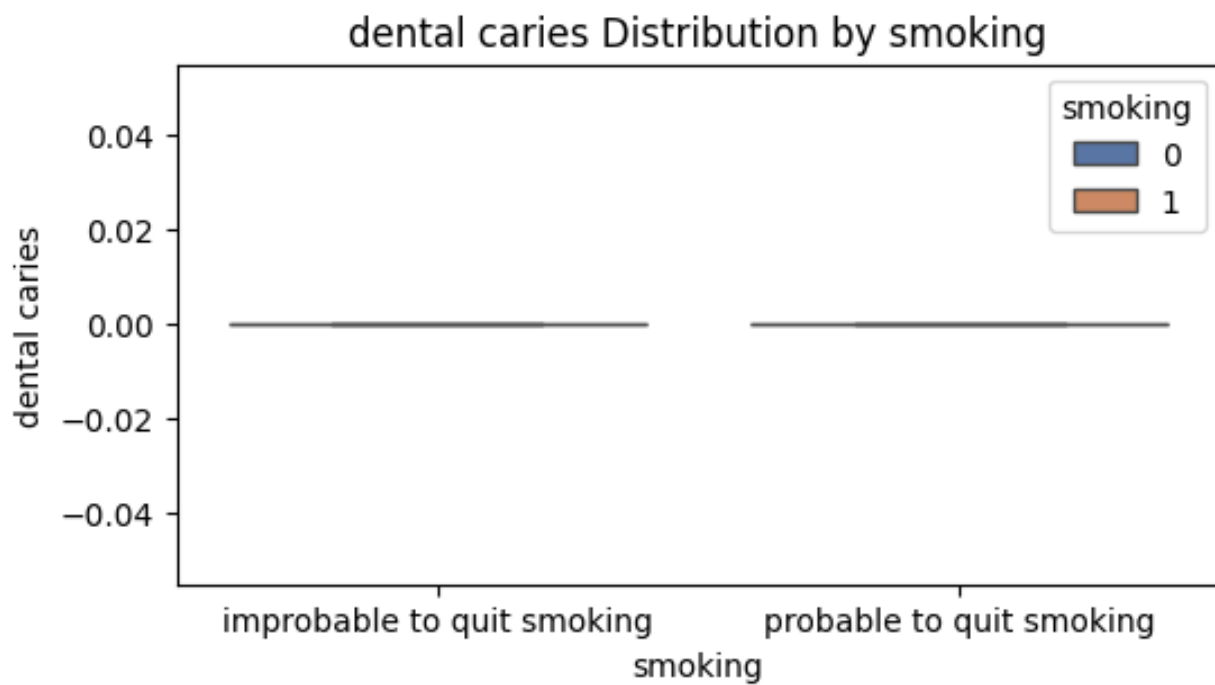


Confusion Matrix for Logistic Regression

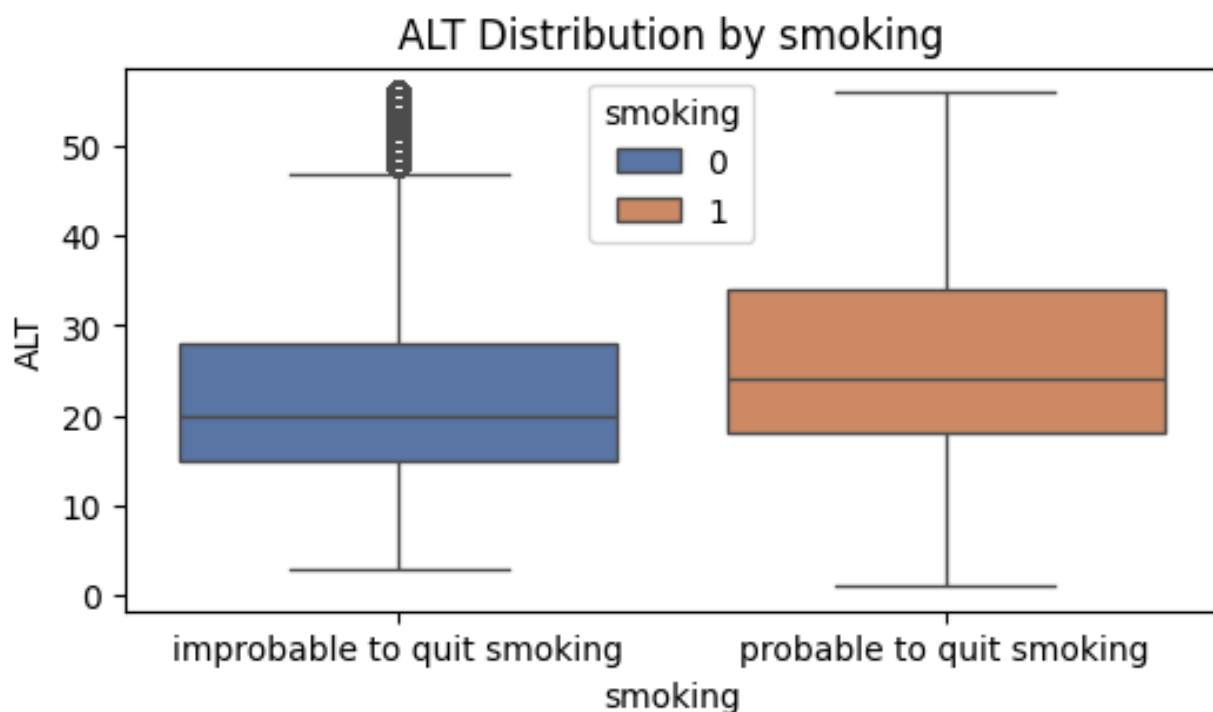
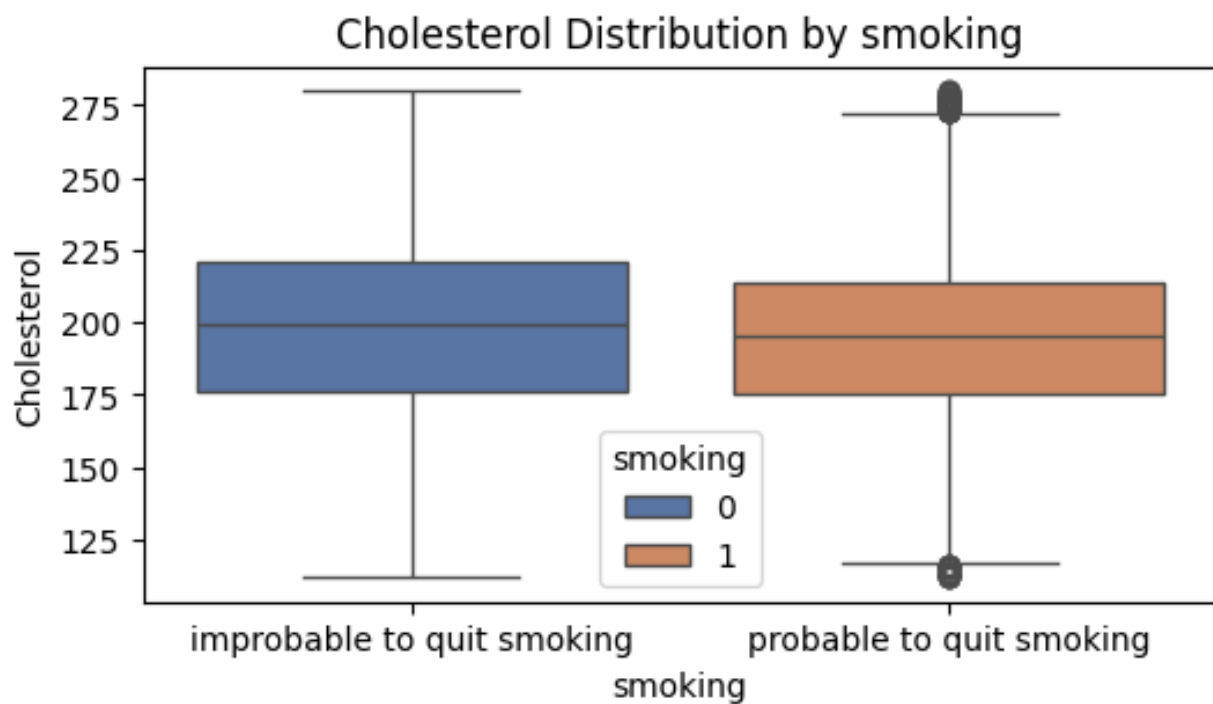


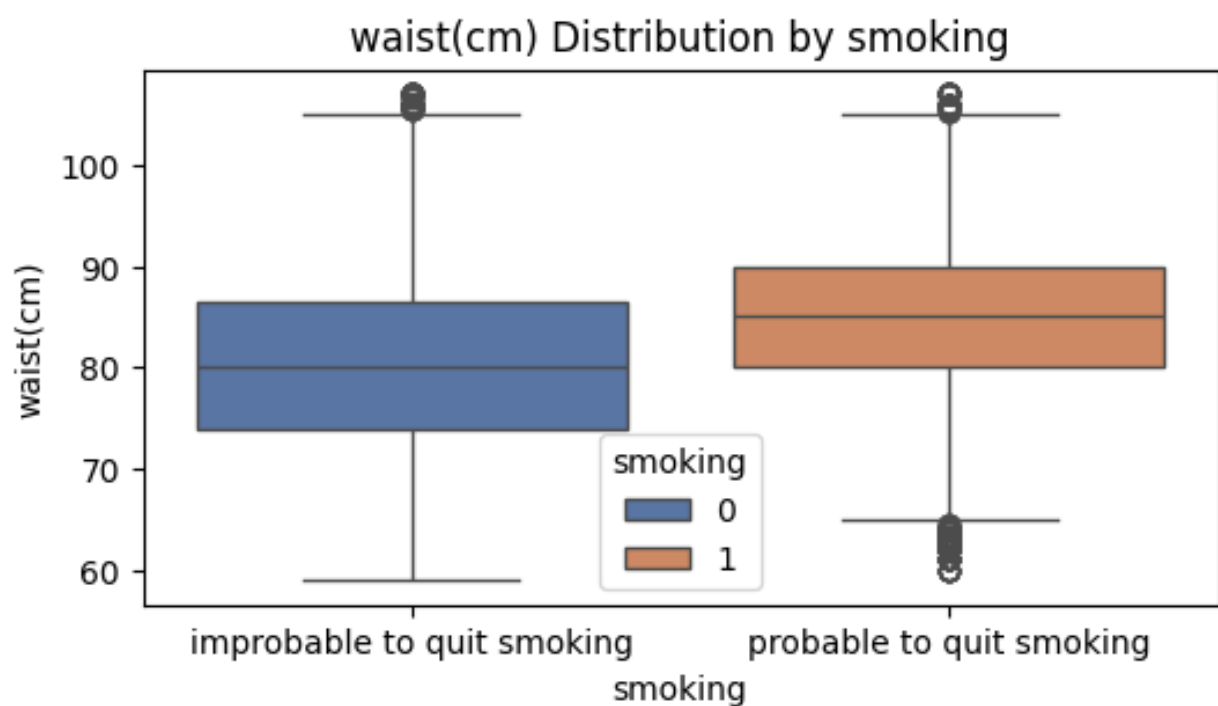
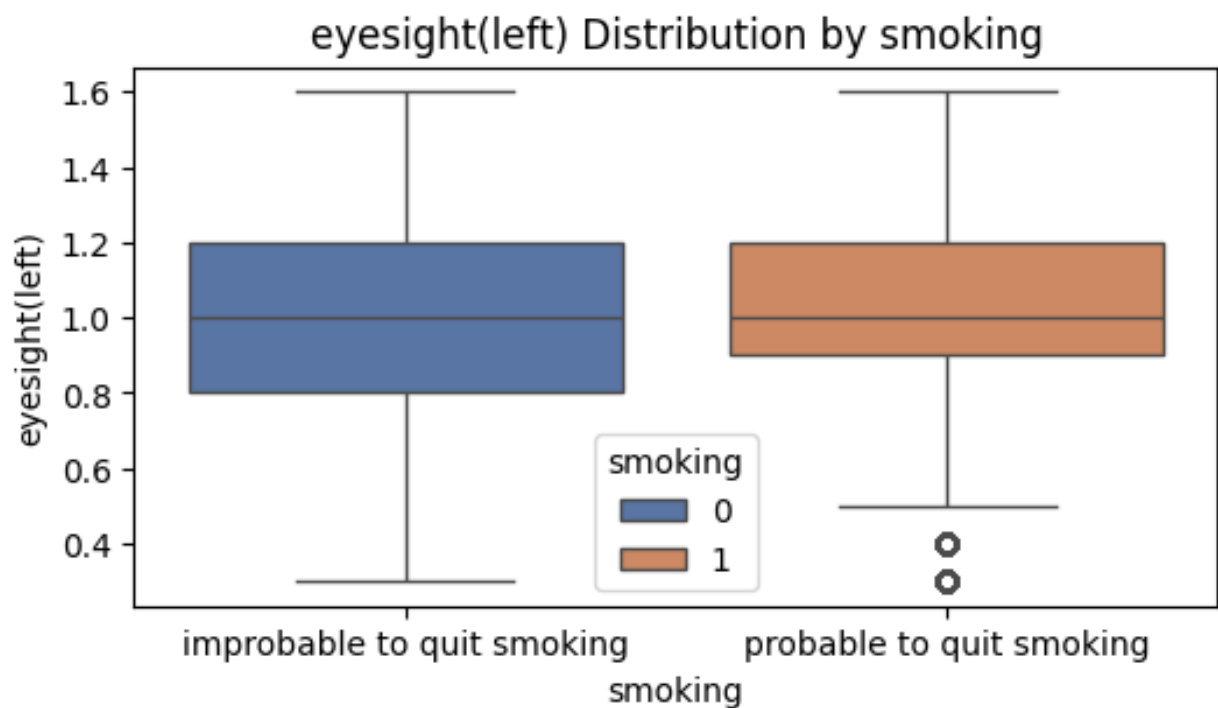
Confusion Matrix for Logistic Regression



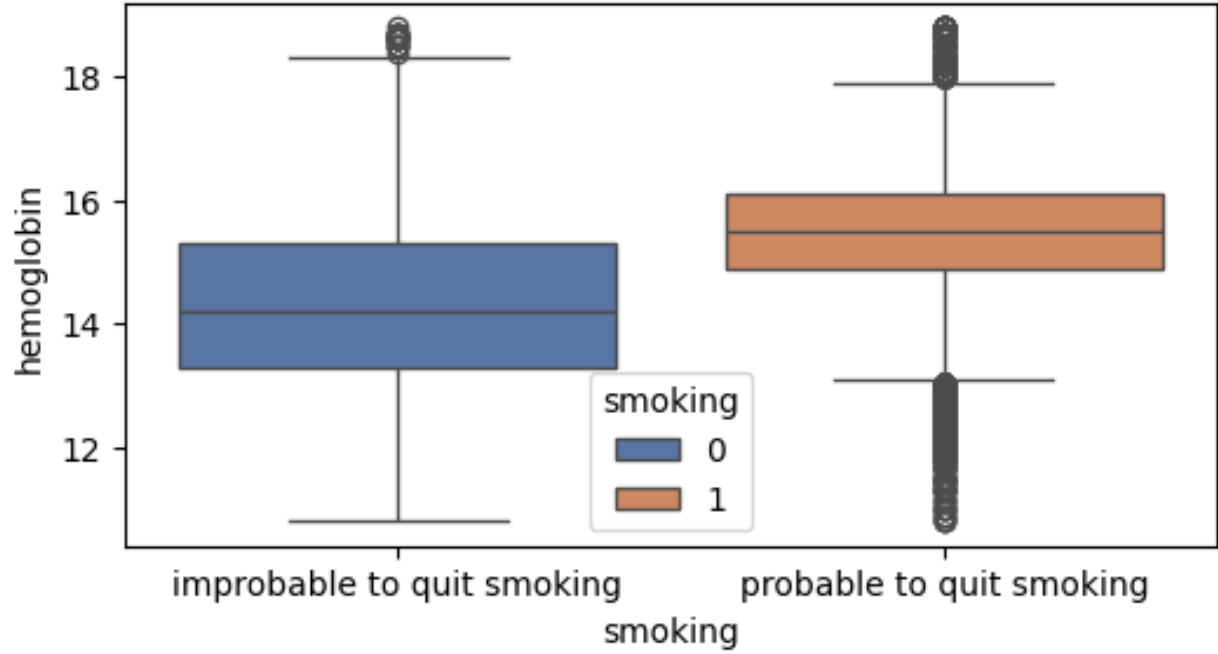


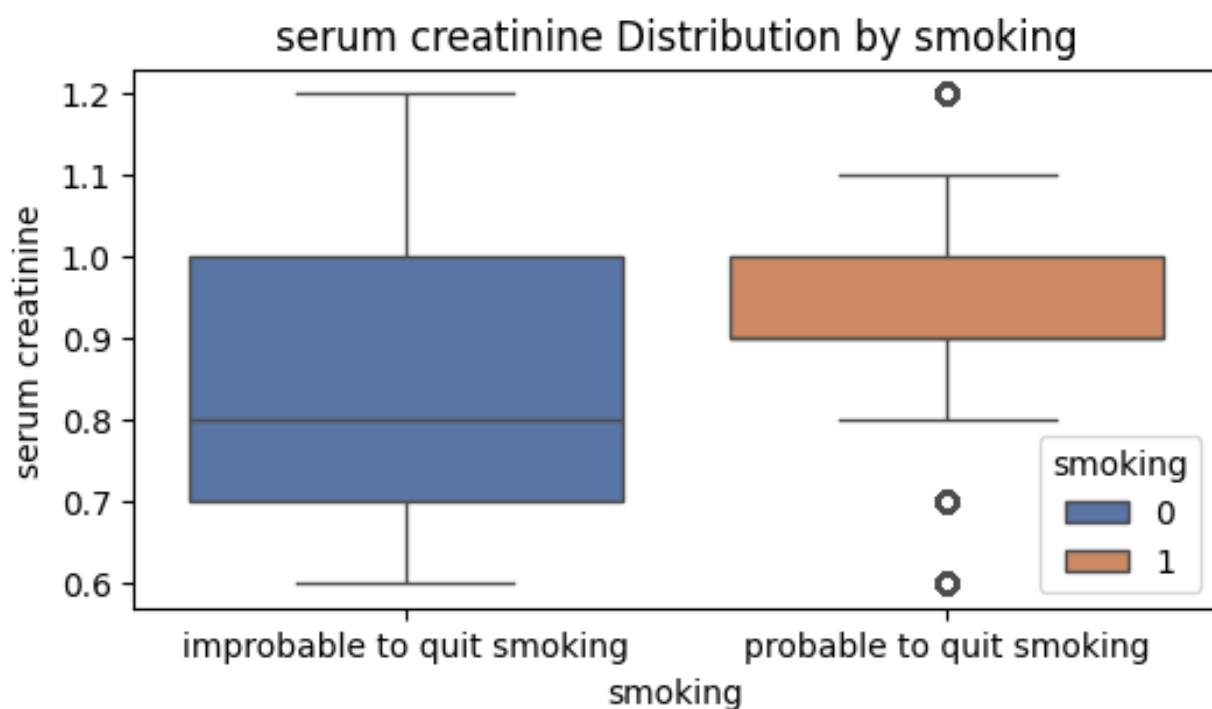
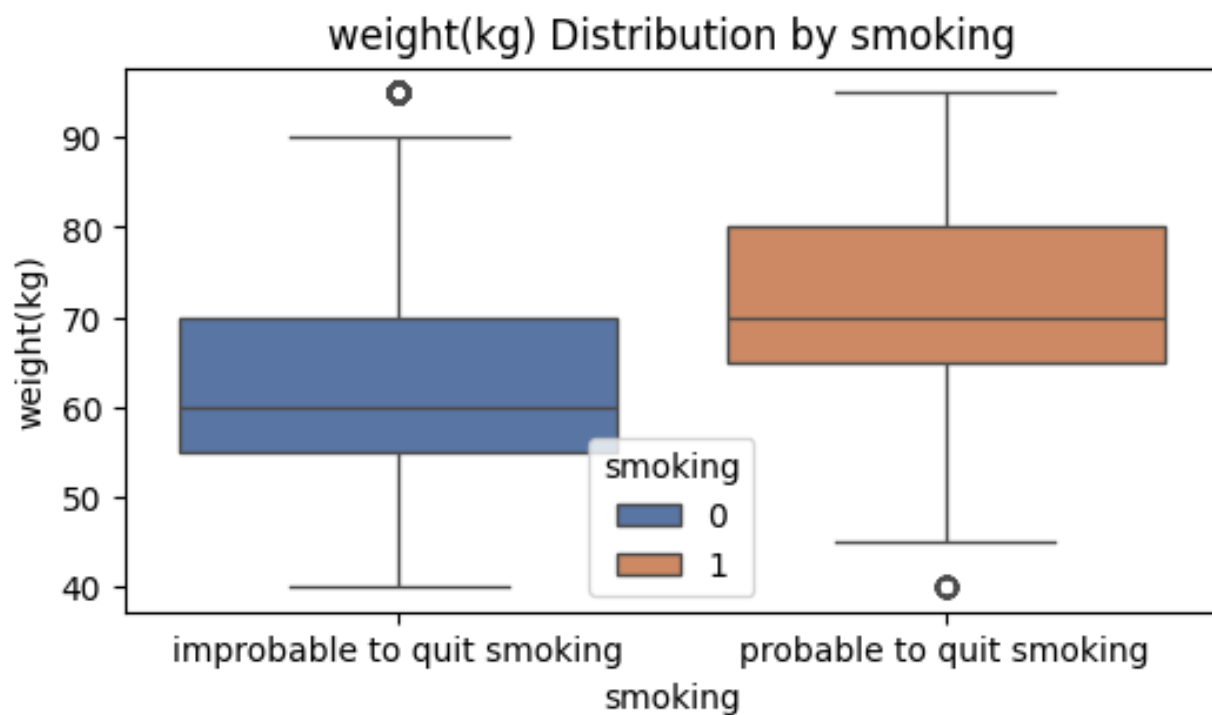
- This again shows no relation to smoking



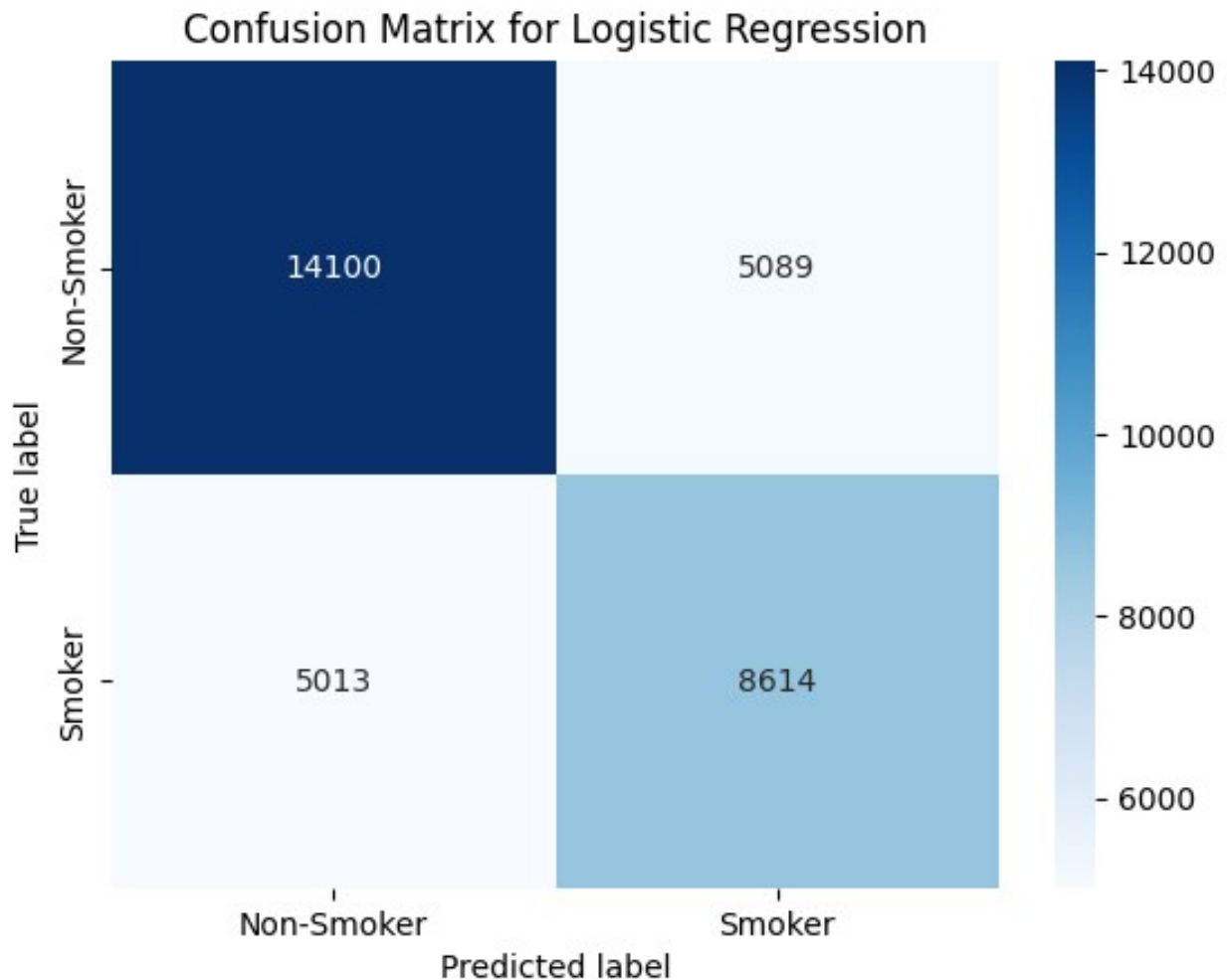


hemoglobin Distribution by smoking





Bivariate:



from analyzing the data relative to the label 'smoking' we were able to identify certain features that correlate the most to the label.

- we note that some features have somewhat positive correlations with one another.
- positive correlations with smoking: [waist, hemoglobin, weight, serum creatinine]
- positive correlations with non-smoking: (ALT, weight) (ALT, hemoglobin) (ALT, waist)
- (waist, serum creatinine) (waist, weight) (waist, hemoglobin)
- (weight, serum creatinine) (weight, hemoglobin)



- Here is a collective graph of the relationships between each two features

Multivariate: Using linear regression (label is binary)

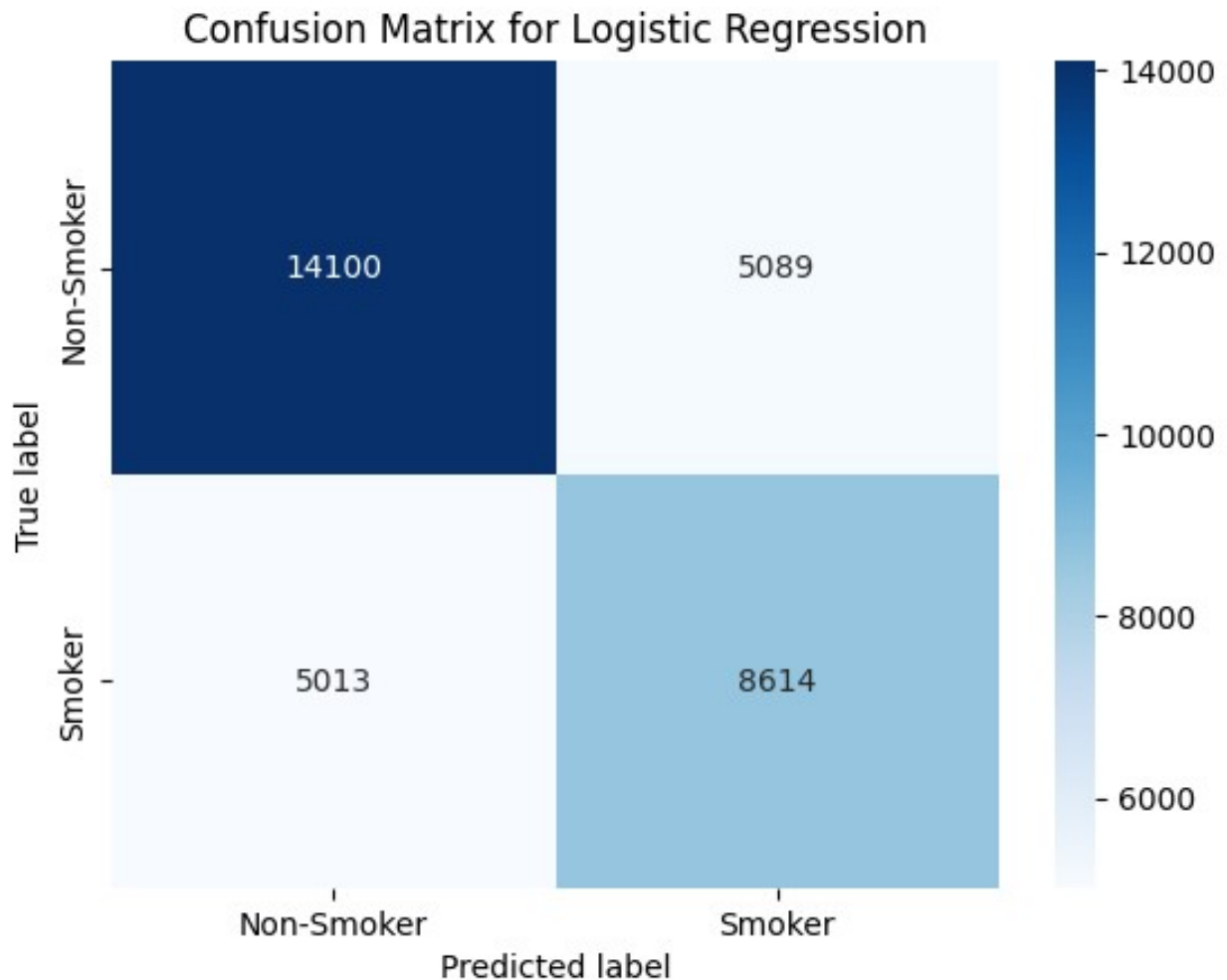


Principal Component Analysis (PCA) is a dimensionality reduction technique that transforms high-dimensional data into a lower-dimensional space while preserving as much variance as possible, making it useful for reducing complexity and mitigating multicollinearity.



multivariate analysis showed that hearing left, hearing right , and dental carries contribute very little. we might drop them for dimensionality reduction.

note that hearing left, hearing right , and dental carries contribute very little. we might drop them for dimensionality reduction. now, we attempt to retrain the model using only seven features(excluding hearing left, hearing right , and dental carries) with high correlation to smoking.



2. Feature Engineering

we concluded that we could reduce the dimensionality only to four features while maintaining the same information with regard to the label. we then engineered a new feature that slightly improved the model accuracy.

3.Ensemble Methods

we created a class for each of the ensemble methods

4. Hyperparametr tuning

Bagging

Bagging, the short form for bootstrap aggregating, is mainly applied in classification and regression. It increases the accuracy of models through decision trees, which reduces variance to a large extent. The

reduction of variance increases accuracy, eliminating overfitting, which is a challenge to many predictive models.

Boosting

Boosting is an ensemble technique that learns from previous predictor mistakes to make better predictions in the future. The technique combines several weak base learners to form one strong learner, thus significantly improving the predictability of models. Boosting works by arranging weak learners in a sequence, such that weak learners learn from the next learner in the sequence to create better predictive models.