

Univariate Analysis Report

Cancer Model

2024-02-02

Overview

Cancer Model Univariate Analysis Report

These sorted results for the features in this report indicate the average cross-validated test scores for each feature, if it were used as the only predictor in a simple linear model. Keep in mind that these results are based on the average, without considering the standard deviation. This means that the results are not necessarily the best predictors, but they are the best on average, and provide a fine starting point for grouping those predictors that are on average better than others. This means that nothing was done to account for possible sampling variability in the sorted results. This is a limitation of the univariate analysis, and it is important to keep this in mind when interpreting the results. It is also important to consider further that depending on the purpose of the model, the most appropriate features may not be the ones with the highest average test scores, if a different metric is more important.

In particular, this should not be taken as an opinion (actuarial or otherwise) regarding the most appropriate features to use in a model, but it rather provides a starting point for further analysis.

	Accuracy	Precision	Recall	AUC	F1	MCC	Ave.
mean_smoothness	62.3%	62.3%	100.0%	50.0%	76.8%	0.00e+00	58.6%
mean_symmetry	62.3%	62.3%	100.0%	50.0%	76.8%	0.00e+00	58.6%
mean_fractal_dimension	62.3%	62.3%	100.0%	50.0%	76.8%	0.00e+00	58.6%
mean_radius	37.7%	0.00e+00	0.00e+00	50.0%	0.00e+00	0.00e+00	14.6%
mean_perimeter	37.7%	0.00e+00	0.00e+00	50.0%	0.00e+00	0.00e+00	14.6%

This table shows an overview of the results for the variables in this file, representing those whose average test score are ranked between 2 and 6 of the variables passed to the Cancer Model.

Univariate Report

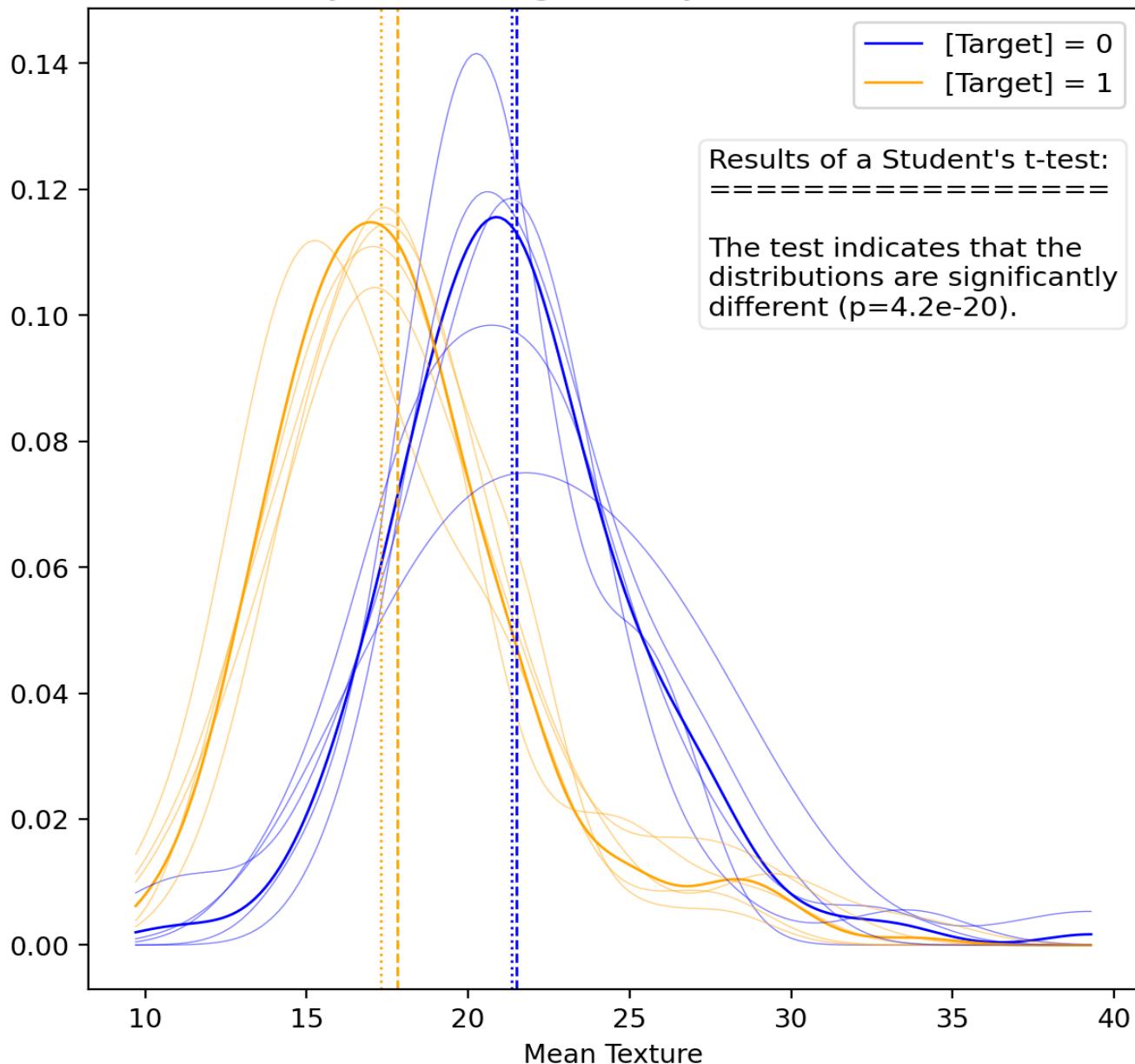
Mean Texture - Results

	Fold-1	Fold-2	Fold-3	Fold-4	Fold-5	Agg. Mean	Agg. SD
Fitted Coef.	1.3e-02	1.7e-02	1.9e-02	2.0e-02	1.5e-02	1.7e-02	2.8e-03
Fitted p-Value	1.7e-02	2.0e-03	6.5e-04	3.6e-04	5.3e-03	6.6e-04	7.0e-03
Fitted Std. Err.	5.3e-03	5.4e-03	5.4e-03	5.5e-03	5.4e-03	4.8e-03	7.3e-05
Conf. Int. Lower	2.2e-03	6.2e-03	7.9e-03	8.8e-03	4.5e-03	7.0e-03	2.6e-03
Conf. Int. Upper	2.3e-02	2.7e-02	2.9e-02	3.0e-02	2.6e-02	2.6e-02	2.9e-03
Train Accuracy	60.8%	63.2%	63.9%	64.3%	62.1%	62.9%	1.4%
Val Accuracy	71.1%	61.5%	58.7%	57.1%	66.3%	62.3%	5.7%
Train AUC	50.0%	50.0%	50.0%	50.0%	50.0%	50.0%	0.0%
Val AUC	50.0%	50.0%	50.0%	50.0%	50.0%	50.0%	0.0%
Train F1	75.6%	77.5%	78.0%	78.3%	76.6%	77.2%	1.1%
Test F1	83.1%	76.1%	74.0%	72.7%	79.7%	76.8%	4.3%
Train Precision	60.8%	63.2%	63.9%	64.3%	62.1%	62.9%	1.4%
Val Precision	71.1%	61.5%	58.7%	57.1%	66.3%	62.3%	5.7%
Train Recall	100.0%	100.0%	100.0%	100.0%	100.0%	100.0%	0.0%
Val Recall	100.0%	100.0%	100.0%	100.0%	100.0%	100.0%	0.0%
Train MCC	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%
Val MCC	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%
Train Log-Loss	14.12	13.25	13.01	12.87	13.68	13.39	0.51
Val Log-Loss	10.41	13.89	14.89	15.45	12.15	13.60	2.07

Univariate Report

Mean Texture - Kernel Density Plot

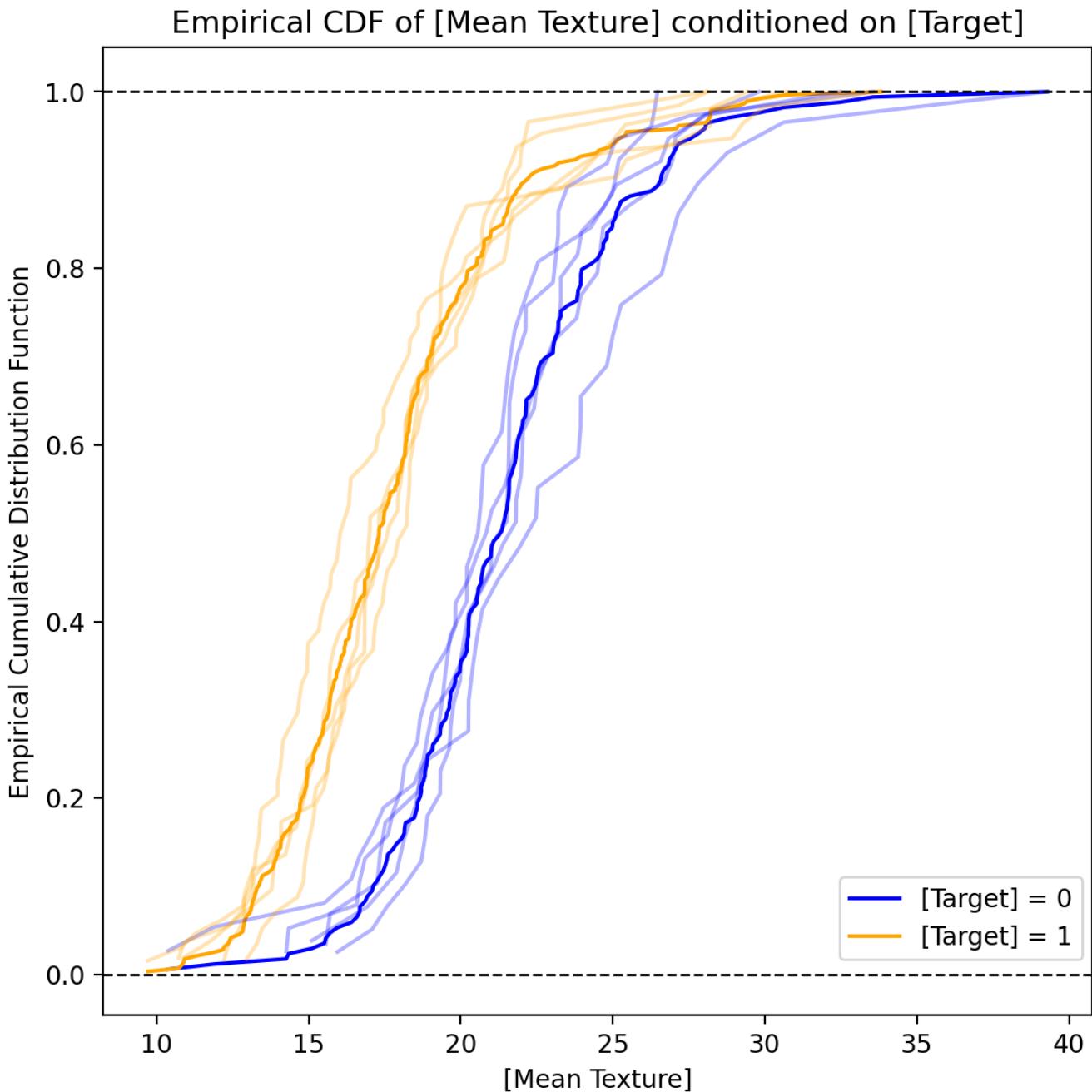
Kernel Density Plot of [Mean Texture] by [Target]
Distributions by level are significantly different at the 95% level.



This plot shows the Gaussian kernel density for each level of the target variable, both in total and for each fold. The x-axis represents the feature variable, and the y-axis represents the density of the target variable. The cross-validation folds are included in slightly washed-out colors to help understand the variability of the data. There are annotations with the results of a t-test for the difference in means between the feature variable at each level of the target variable. The annotations corresponding to the color of the target variable level show the mean/median ratio to help understand differences in skewness between the levels of the target variable.

Univariate Report

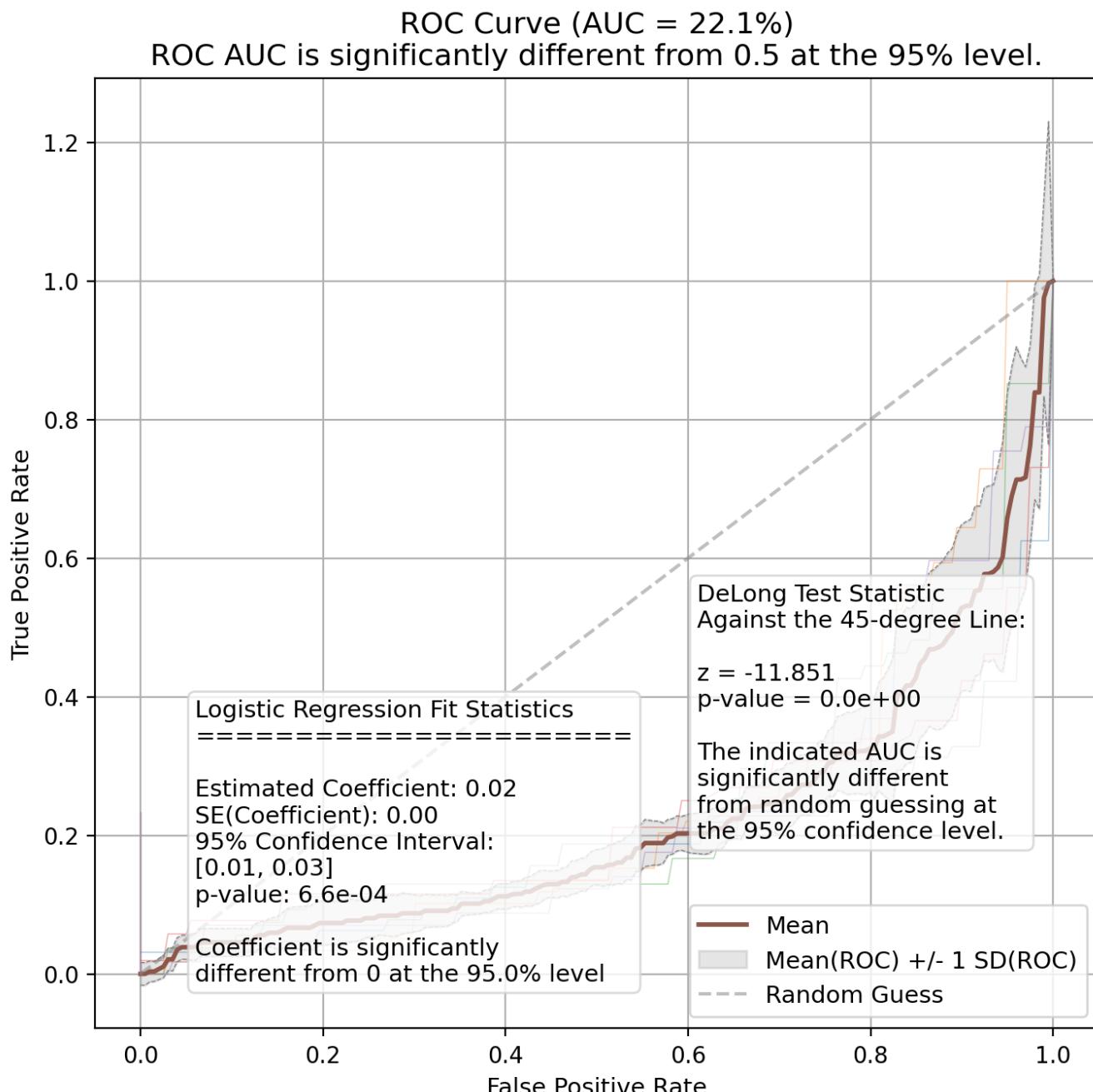
Mean Texture - Empirical CDF Plot



This plot shows the empirical cumulative distribution function for each level of the target variable, both in total and for each fold. The x-axis represents the feature variable, and the y-axis represents the cumulative distribution of the target variable. The cross-validation folds are included in slightly washed-out colors to help understand the variability of the data, and whether or not it is reasonable to assume that the data is drawn from different distributions.

Univariate Report

Mean Texture - ROC Curve



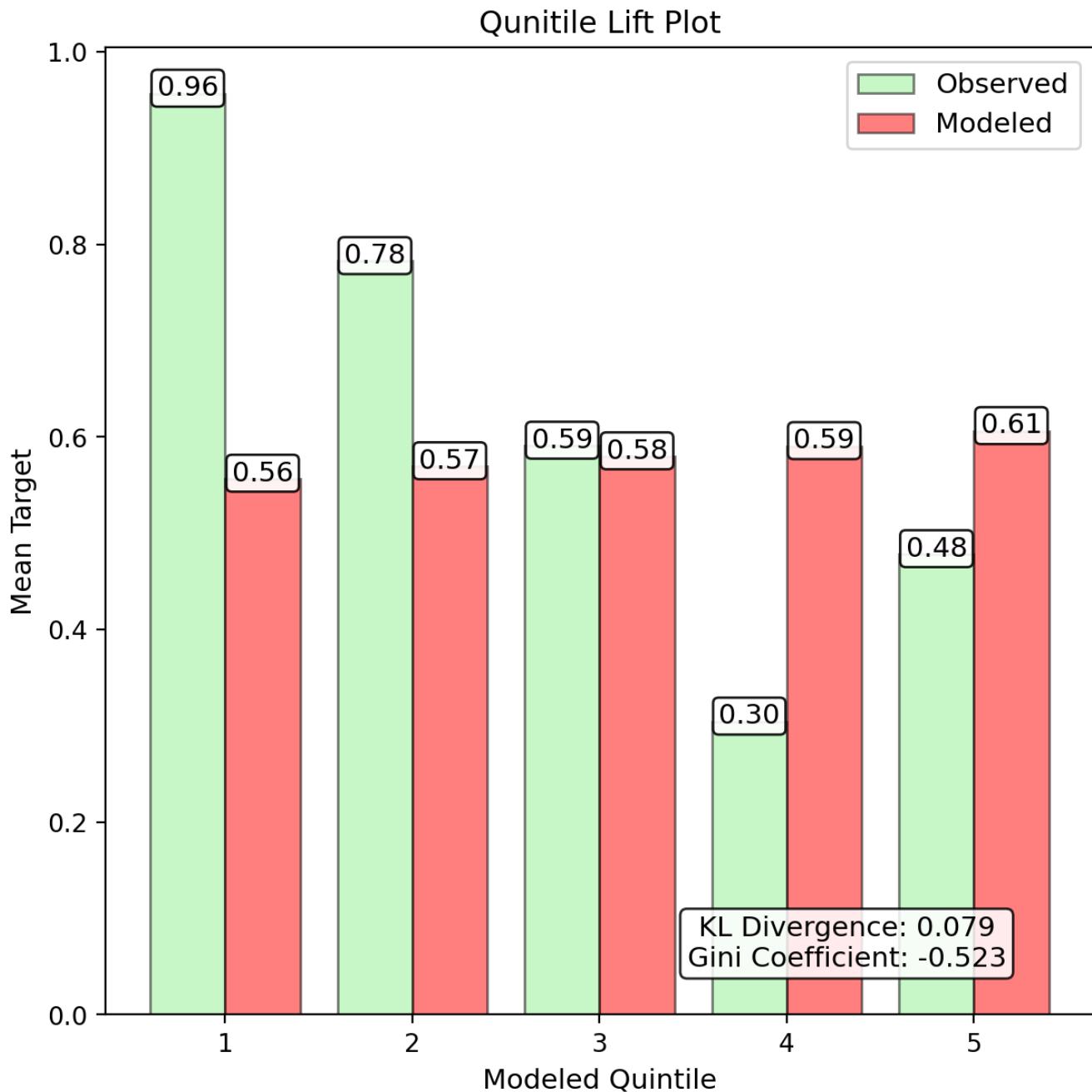
This plot shows the receiver operating characteristic (ROC) curve for the target variable in total and for each fold. The x-axis represents the false positive rate, and the y-axis represents the true positive rate. This is based on a simple Logistic Regression model with no regularization, no intercept, and no other features. Annotations are on the plot to help understand the results of the model, including the coefficient, standard error, and p-value for the feature variable. The cross-validation folds are used to create the grey region around the mean ROC curve to help understand the variability of the data.

Significance of the ROC curve is determined based on the method from DeLong et al. (1988). In brief, the AUC is assumed to be normally distributed, and the z-score is calculated based on the AUC and the standard error. This z-score is compared to a +/- two standard deviations from a standard normal

distribution to get the p-value.

Univariate Report

Mean Texture - Quintile Lift



The quintile lift plot is meant to show the power of the single feature to discriminate between the highest and lowest quintiles of the target variable.

Univariate Report

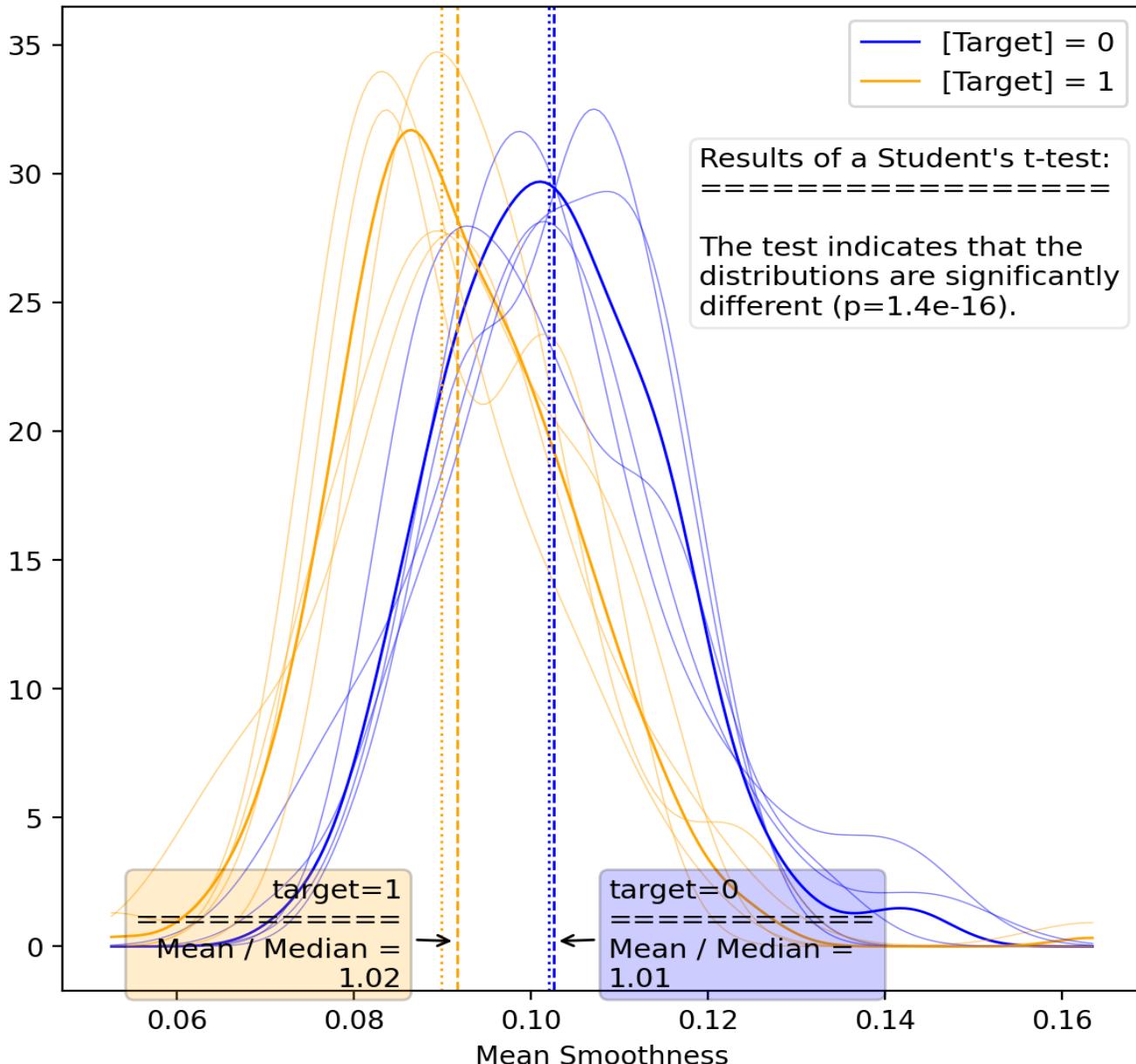
Mean Smoothness - Results

	Fold-1	Fold-2	Fold-3	Fold-4	Fold-5	Agg. Mean	Agg. SD
Fitted Coef.	3.23	4.34	4.58	4.74	3.93	4.16	0.60
Fitted p-Value	3.2e-03	1.1e-04	4.2e-05	2.1e-05	3.2e-04	2.7e-05	1.4e-03
Fitted Std. Err.	1.10	1.12	1.12	1.11	1.09	0.99	0.01
Conf. Int. Lower	1.09	2.14	2.39	2.56	1.79	2.22	0.58
Conf. Int. Upper	5.38	6.54	6.78	6.92	6.07	6.10	0.62
Train Accuracy	60.8%	63.2%	63.9%	64.3%	62.1%	62.9%	1.4%
Val Accuracy	71.1%	61.5%	58.7%	57.1%	66.3%	62.3%	5.7%
Train AUC	50.0%	50.0%	50.0%	50.0%	50.0%	50.0%	0.0%
Val AUC	50.0%	50.0%	50.0%	50.0%	50.0%	50.0%	0.0%
Train F1	75.6%	77.5%	78.0%	78.3%	76.6%	77.2%	1.1%
Test F1	83.1%	76.1%	74.0%	72.7%	79.7%	76.8%	4.3%
Train Precision	60.8%	63.2%	63.9%	64.3%	62.1%	62.9%	1.4%
Val Precision	71.1%	61.5%	58.7%	57.1%	66.3%	62.3%	5.7%
Train Recall	100.0%	100.0%	100.0%	100.0%	100.0%	100.0%	0.0%
Val Recall	100.0%	100.0%	100.0%	100.0%	100.0%	100.0%	0.0%
Train MCC	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%
Val MCC	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%
Train Log-Loss	14.12	13.25	13.01	12.87	13.68	13.39	0.51
Val Log-Loss	10.41	13.89	14.89	15.45	12.15	13.60	2.07

Univariate Report

Mean Smoothness - Kernel Density Plot

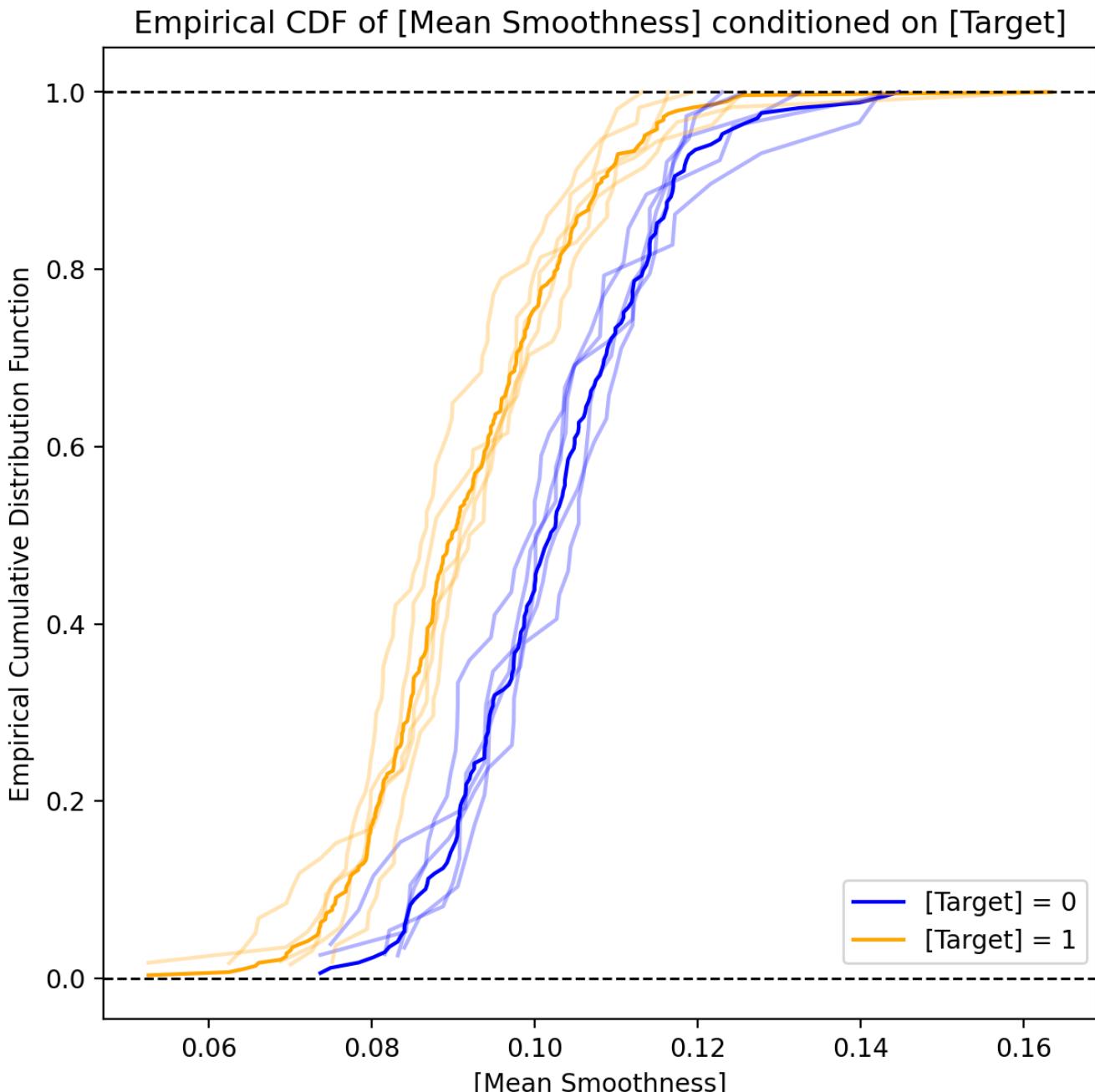
Kernel Density Plot of [Mean Smoothness] by [Target]
Distributions by level are significantly different at the 95% level.



This plot shows the Gaussian kernel density for each level of the target variable, both in total and for each fold. The x-axis represents the feature variable, and the y-axis represents the density of the target variable. The cross-validation folds are included in slightly washed-out colors to help understand the variability of the data. There are annotations with the results of a t-test for the difference in means between the feature variable at each level of the target variable. The annotations corresponding to the color of the target variable level show the mean/median ratio to help understand differences in skewness between the levels of the target variable.

Univariate Report

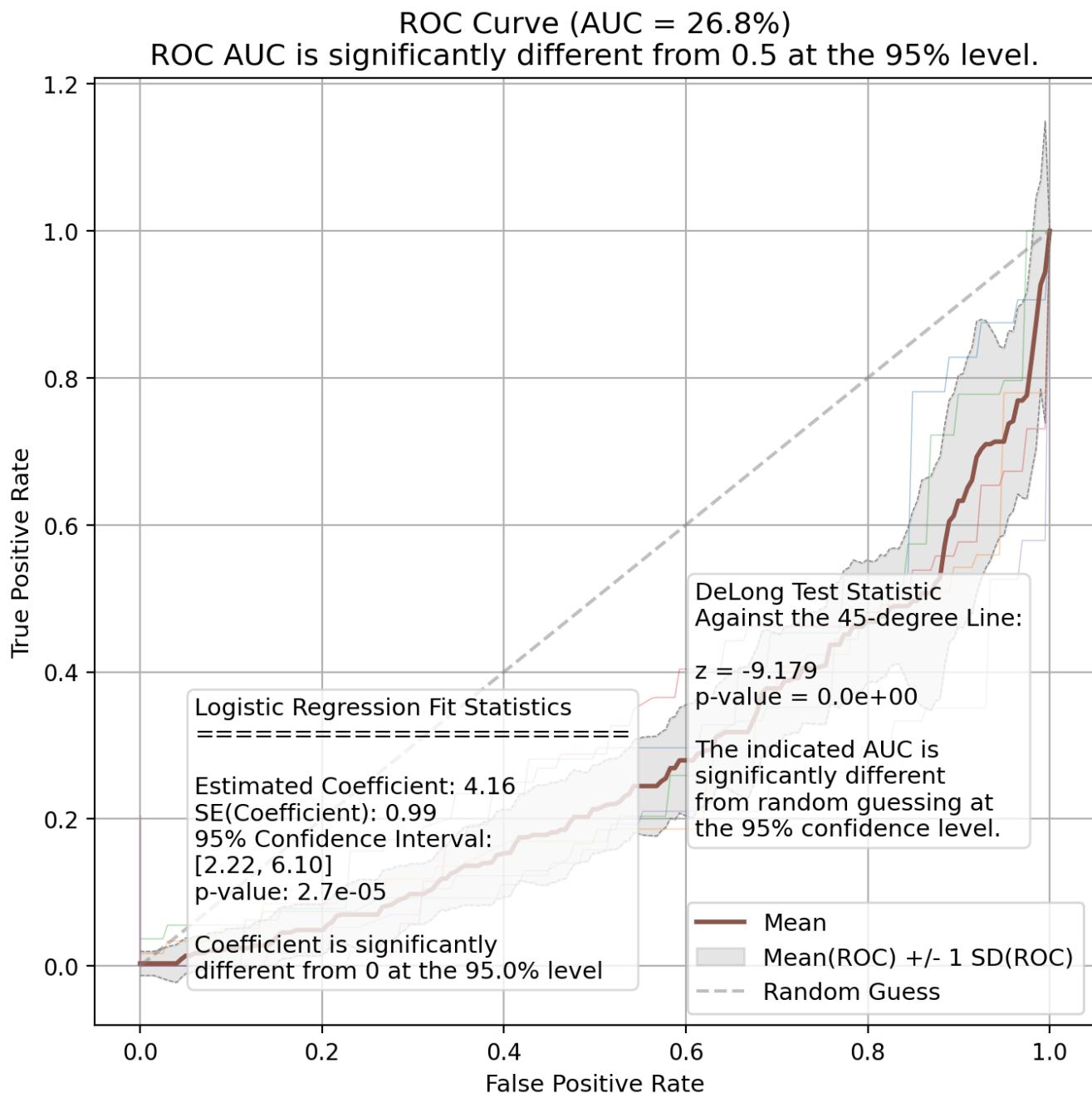
Mean Smoothness - Empirical CDF Plot



This plot shows the empirical cumulative distribution function for each level of the target variable, both in total and for each fold. The x-axis represents the feature variable, and the y-axis represents the cumulative distribution of the target variable. The cross-validation folds are included in slightly washed-out colors to help understand the variability of the data, and whether or not it is reasonable to assume that the data is drawn from different distributions.

Univariate Report

Mean Smoothness - ROC Curve



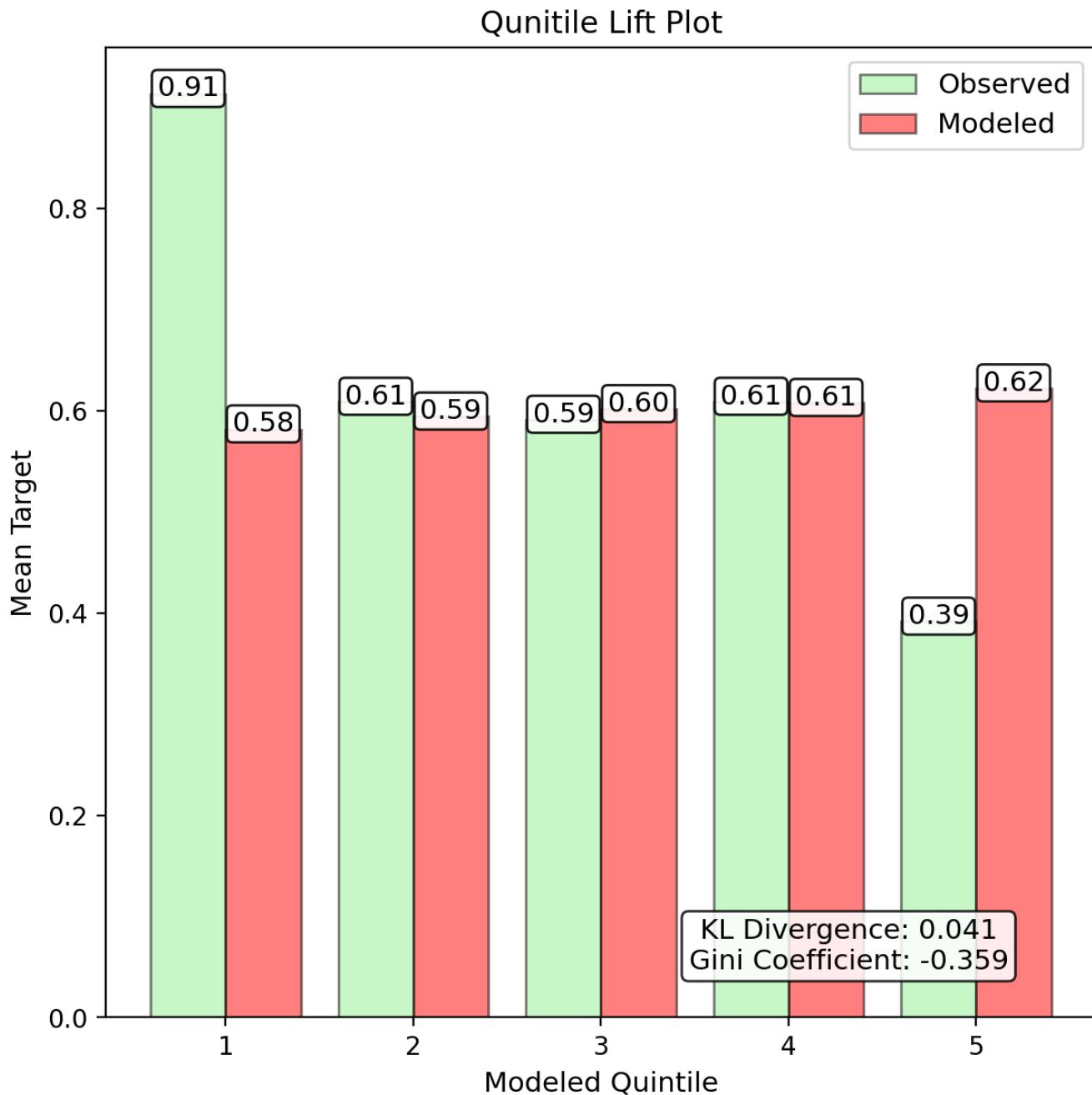
This plot shows the receiver operating characteristic (ROC) curve for the target variable in total and for each fold. The x-axis represents the false positive rate, and the y-axis represents the true positive rate. This is based on a simple Logistic Regression model with no regularization, no intercept, and no other features. Annotations are on the plot to help understand the results of the model, including the coefficient, standard error, and p-value for the feature variable. The cross-validation folds are used to create the grey region around the mean ROC curve to help understand the variability of the data.

Significance of the ROC curve is determined based on the method from DeLong et al. (1988). In brief, the AUC is assumed to be normally distributed, and the z-score is calculated based on the AUC and the standard error. This z-score is compared to a +/- two standard deviations from a standard normal

distribution to get the p-value.

Univariate Report

Mean Smoothness - Quintile Lift



The quintile lift plot is meant to show the power of the single feature to discriminate between the highest and lowest quintiles of the target variable.

Univariate Report

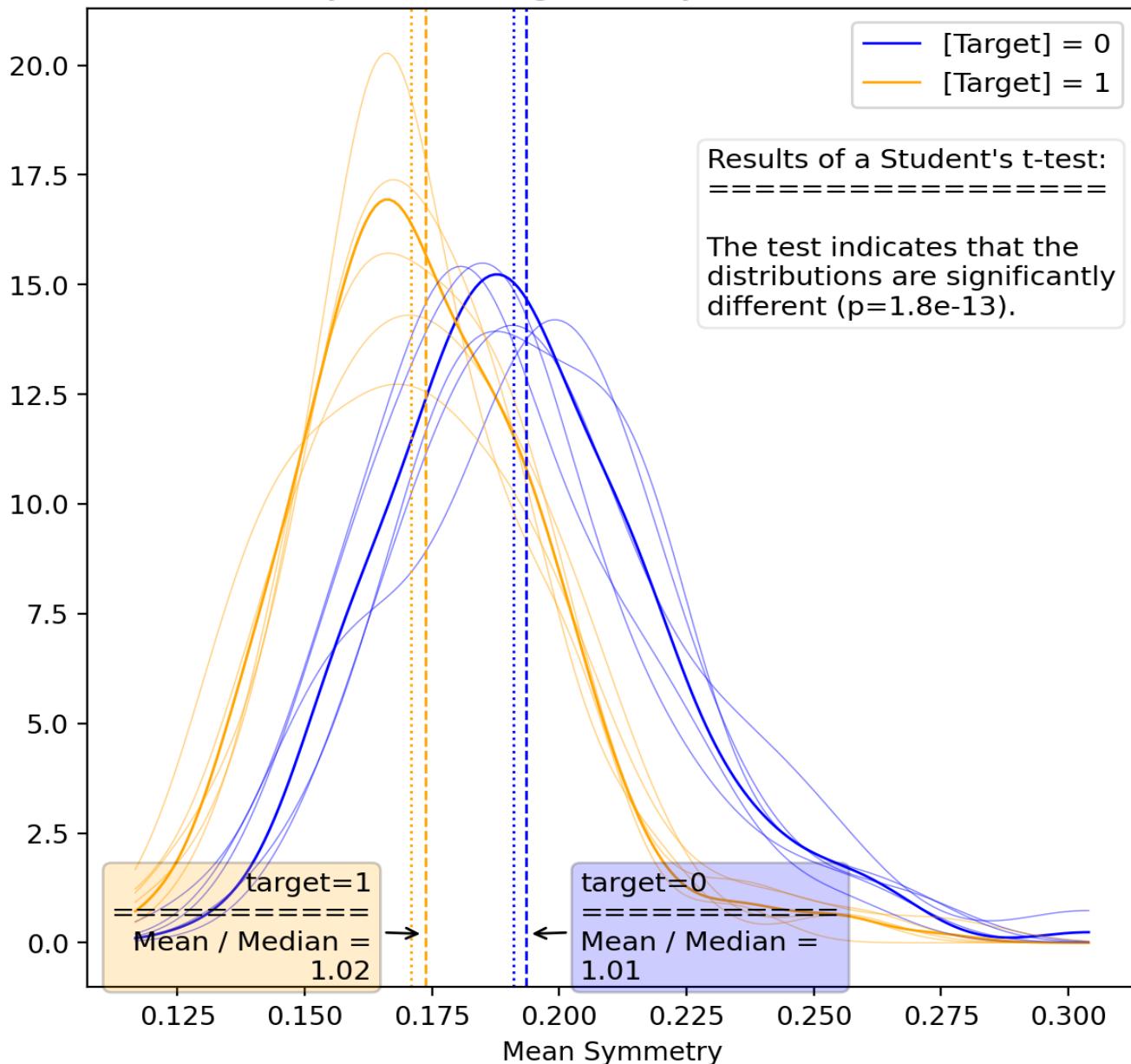
Mean Symmetry - Results

	Fold-1	Fold-2	Fold-3	Fold-4	Fold-5	Agg. Mean	Agg. SD
Fitted Coef.	1.77	2.37	2.48	2.49	1.98	2.21	0.33
Fitted p-Value	2.3e-03	6.6e-05	2.7e-05	2.2e-05	6.0e-04	2.4e-05	9.6e-04
Fitted Std. Err.	0.578	0.594	0.591	0.586	0.578	0.523	0.007
Conf. Int. Lower	0.63	1.21	1.32	1.34	0.85	1.19	0.31
Conf. Int. Upper	2.90	3.53	3.64	3.64	3.12	3.24	0.34
Train Accuracy	60.8%	63.2%	63.9%	64.3%	62.1%	62.9%	1.4%
Val Accuracy	71.1%	61.5%	58.7%	57.1%	66.3%	62.3%	5.7%
Train AUC	50.0%	50.0%	50.0%	50.0%	50.0%	50.0%	0.0%
Val AUC	50.0%	50.0%	50.0%	50.0%	50.0%	50.0%	0.0%
Train F1	75.6%	77.5%	78.0%	78.3%	76.6%	77.2%	1.1%
Test F1	83.1%	76.1%	74.0%	72.7%	79.7%	76.8%	4.3%
Train Precision	60.8%	63.2%	63.9%	64.3%	62.1%	62.9%	1.4%
Val Precision	71.1%	61.5%	58.7%	57.1%	66.3%	62.3%	5.7%
Train Recall	100.0%	100.0%	100.0%	100.0%	100.0%	100.0%	0.0%
Val Recall	100.0%	100.0%	100.0%	100.0%	100.0%	100.0%	0.0%
Train MCC	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%
Val MCC	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%
Train Log-Loss	14.12	13.25	13.01	12.87	13.68	13.39	0.51
Val Log-Loss	10.41	13.89	14.89	15.45	12.15	13.60	2.07

Univariate Report

Mean Symmetry - Kernel Density Plot

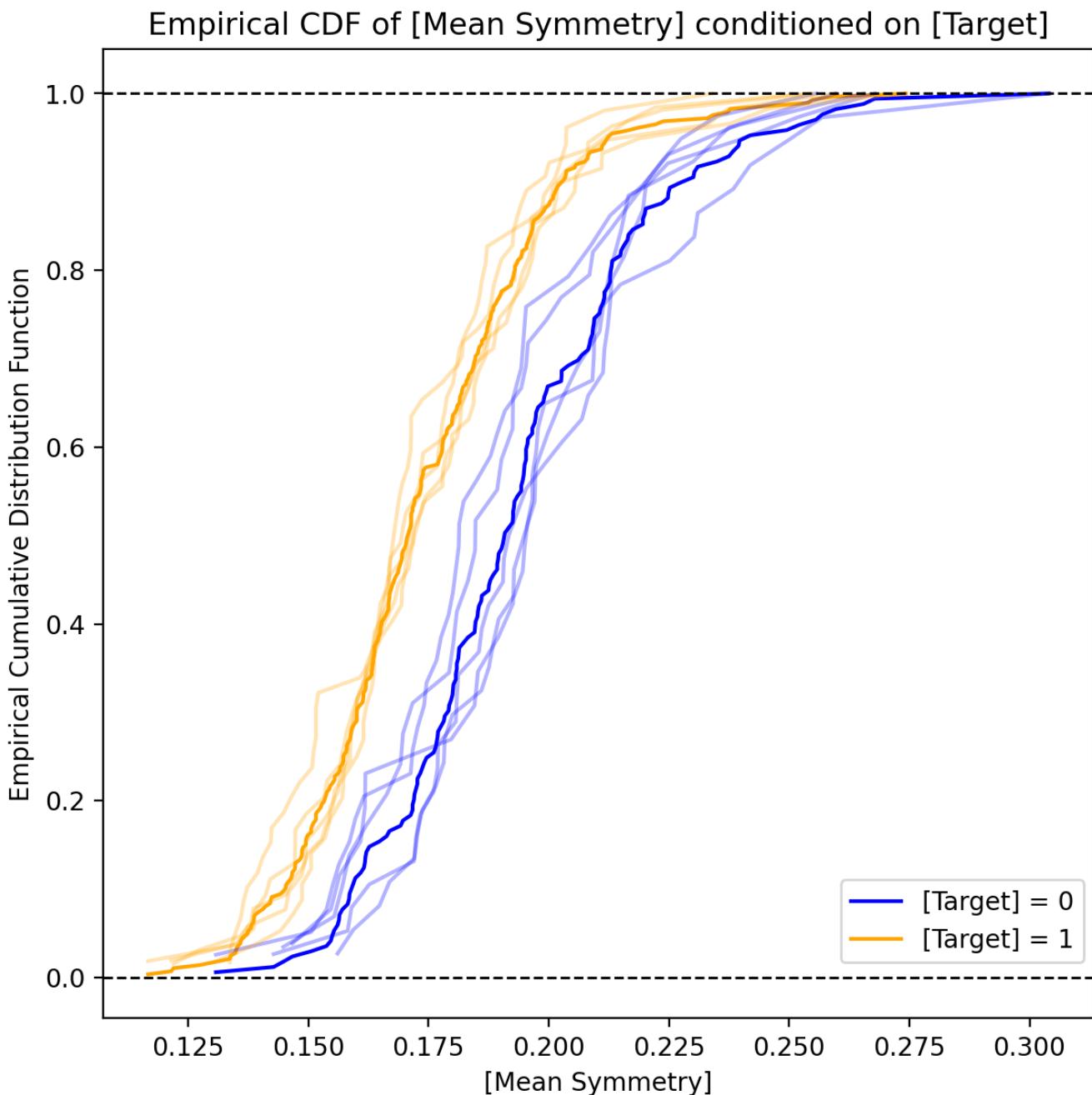
Kernel Density Plot of [Mean Symmetry] by [Target]
Distributions by level are significantly different at the 95% level.



This plot shows the Gaussian kernel density for each level of the target variable, both in total and for each fold. The x-axis represents the feature variable, and the y-axis represents the density of the target variable. The cross-validation folds are included in slightly washed-out colors to help understand the variability of the data. There are annotations with the results of a t-test for the difference in means between the feature variable at each level of the target variable. The annotations corresponding to the color of the target variable level show the mean/median ratio to help understand differences in skewness between the levels of the target variable.

Univariate Report

Mean Symmetry - Empirical CDF Plot



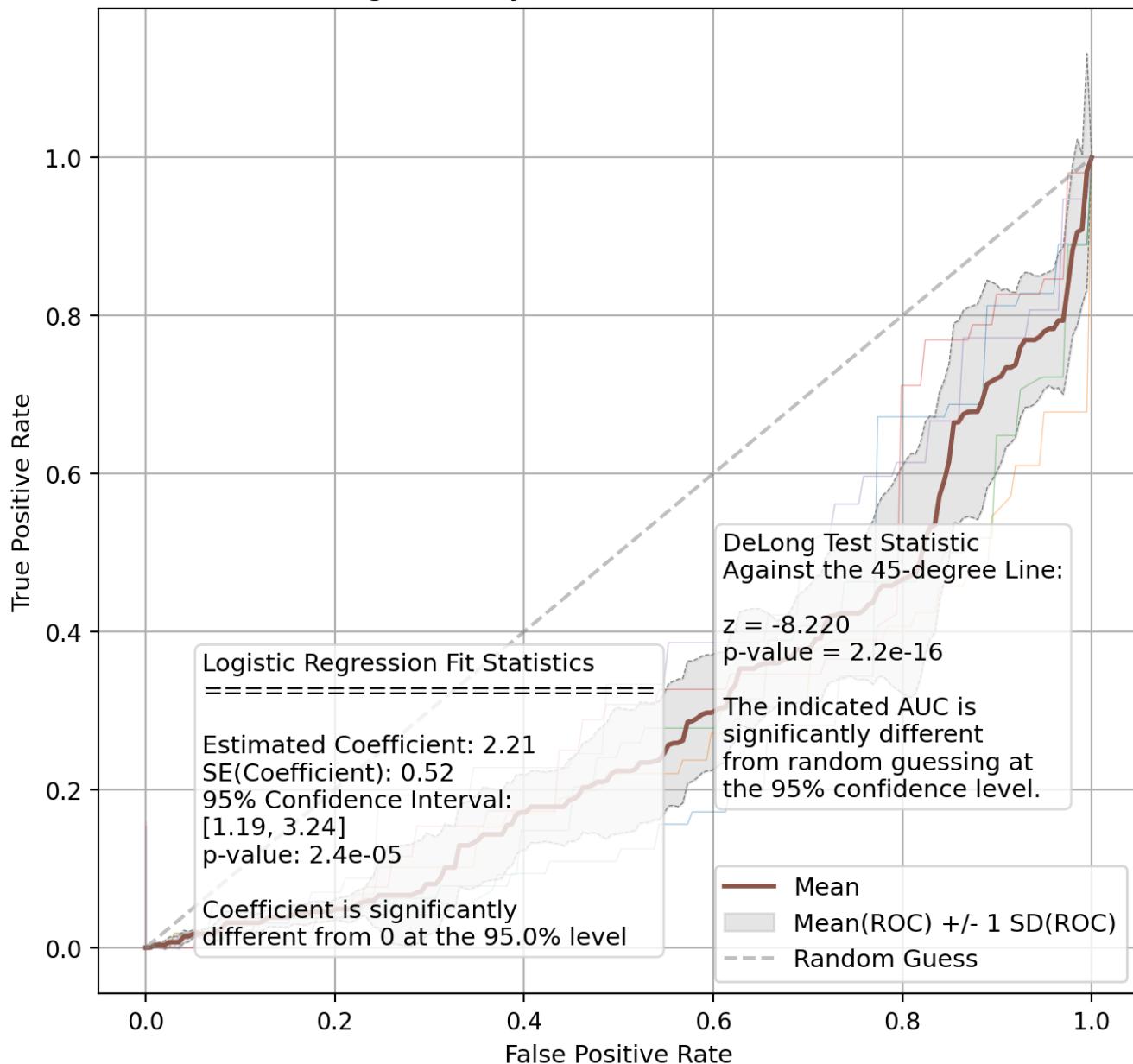
This plot shows the empirical cumulative distribution function for each level of the target variable, both in total and for each fold. The x-axis represents the feature variable, and the y-axis represents the cumulative distribution of the target variable. The cross-validation folds are included in slightly washed-out colors to help understand the variability of the data, and whether or not it is reasonable to assume that the data is drawn from different distributions.

Univariate Report

Mean Symmetry - ROC Curve

ROC Curve (AUC = 28.8%)

ROC AUC is significantly different from 0.5 at the 95% level.



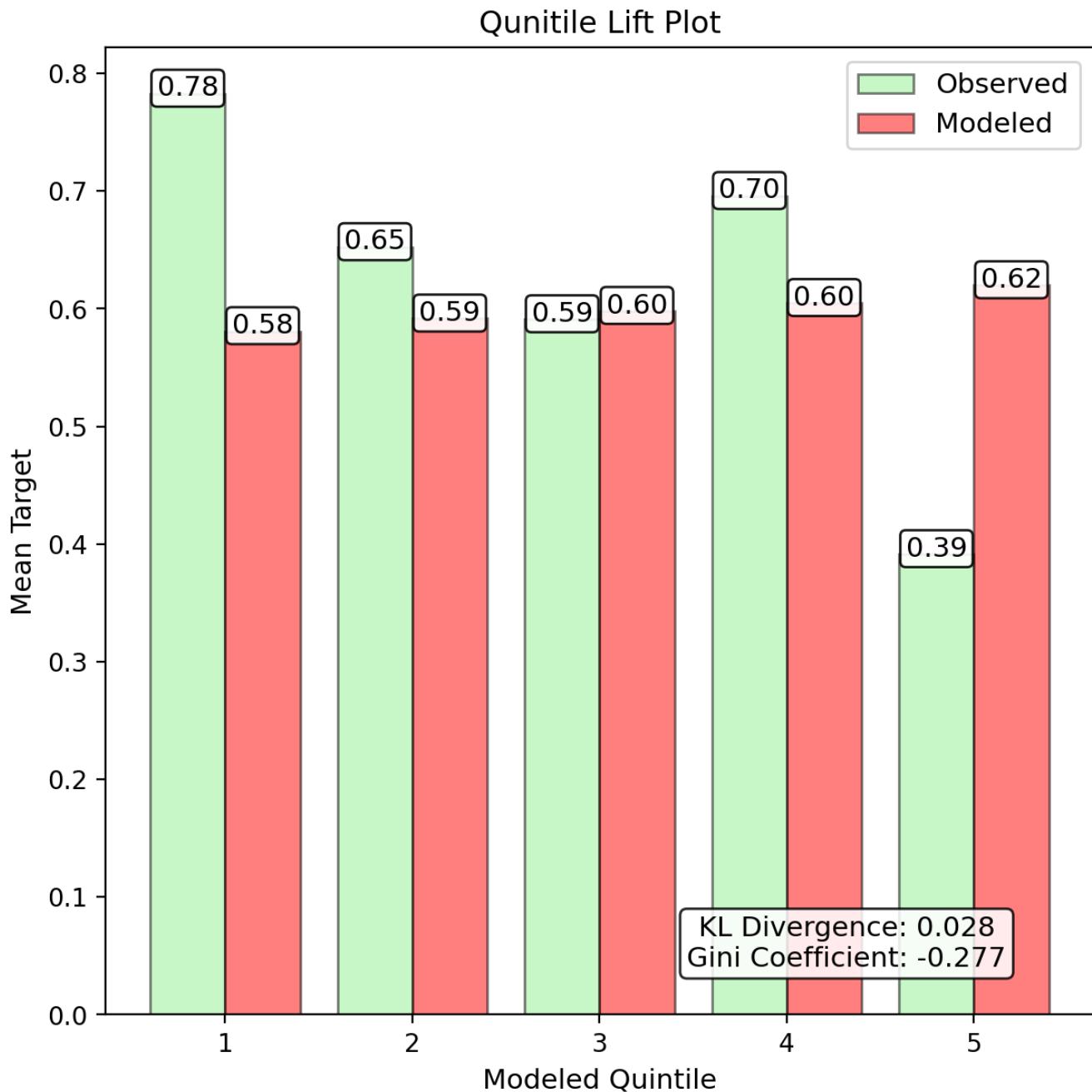
This plot shows the receiver operating characteristic (ROC) curve for the target variable in total and for each fold. The x-axis represents the false positive rate, and the y-axis represents the true positive rate. This is based on a simple Logistic Regression model with no regularization, no intercept, and no other features. Annotations are on the plot to help understand the results of the model, including the coefficient, standard error, and p-value for the feature variable. The cross-validation folds are used to create the grey region around the mean ROC curve to help understand the variability of the data.

Significance of the ROC curve is determined based on the method from DeLong et al. (1988). In brief, the AUC is assumed to be normally distributed, and the z-score is calculated based on the AUC and the standard error. This z-score is compared to a +/- two standard deviations from a standard normal

distribution to get the p-value.

Univariate Report

Mean Symmetry - Quintile Lift



The quintile lift plot is meant to show the power of the single feature to discriminate between the highest and lowest quintiles of the target variable.

Univariate Report

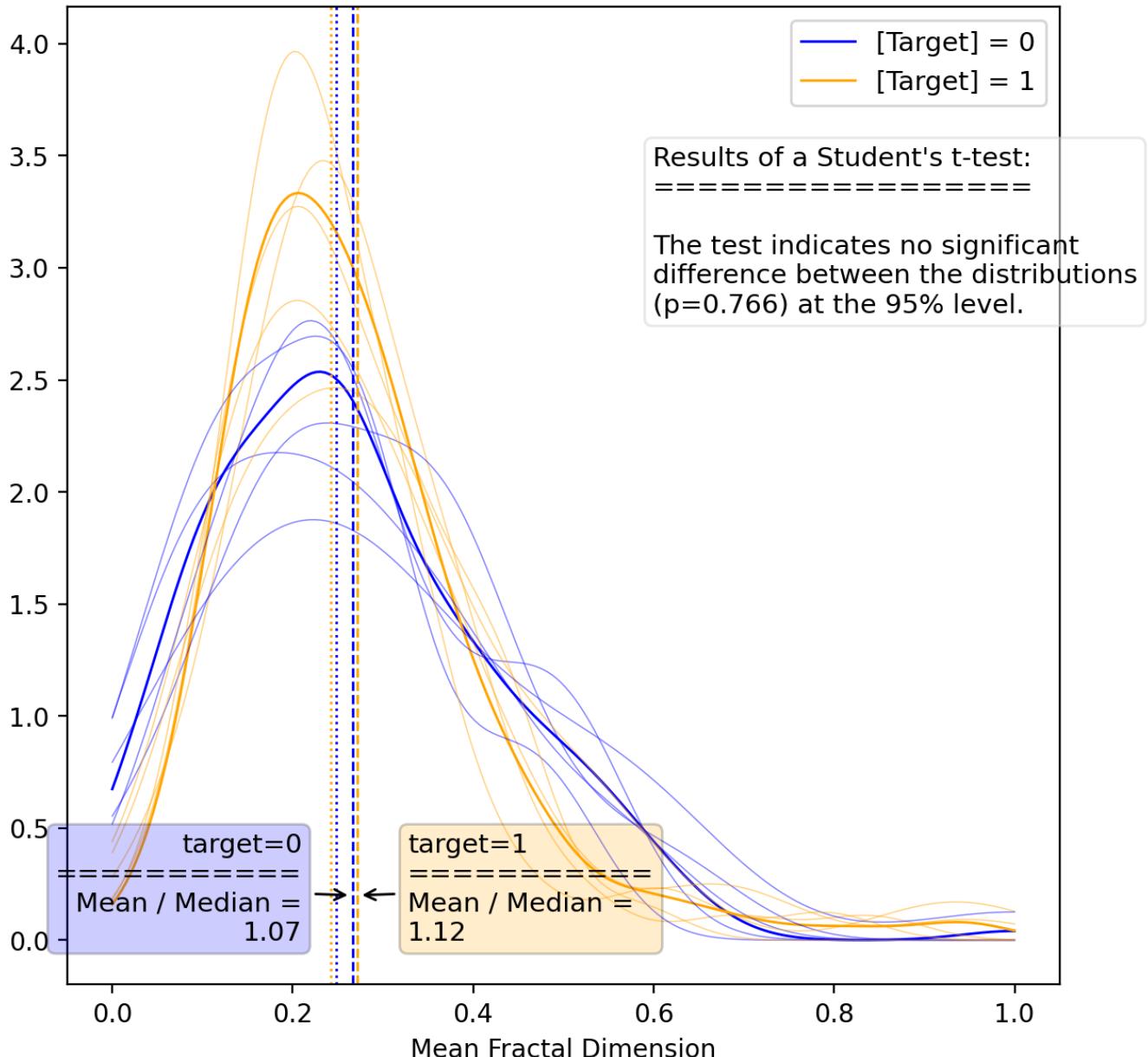
Mean Fractal Dimension - Results

	Fold-1	Fold-2	Fold-3	Fold-4	Fold-5	Agg. Mean	Agg. SD
Fitted Coef.	1.24	1.65	1.63	1.55	1.61	1.53	0.17
Fitted p-Value	4.2e-04	9.4e-06	1.2e-05	1.2e-05	9.6e-06	2.3e-06	1.8e-04
Fitted Std. Err.	0.351	0.371	0.372	0.355	0.363	0.324	0.009
Conf. Int. Lower	0.552	0.918	0.897	0.859	0.895	0.896	0.154
Conf. Int. Upper	1.93	2.37	2.36	2.25	2.32	2.17	0.18
Train Accuracy	61.1%	63.5%	64.2%	64.3%	62.3%	63.1%	1.4%
Val Accuracy	71.1%	61.5%	58.7%	58.2%	66.3%	62.3%	5.5%
Train AUC	50.3%	50.4%	50.4%	50.0%	50.4%	50.3%	0.2%
Val AUC	50.0%	50.0%	50.0%	51.3%	50.0%	50.0%	0.6%
Train F1	75.8%	77.6%	78.1%	78.3%	76.7%	77.3%	1.0%
Test F1	83.1%	76.1%	74.0%	73.2%	79.7%	76.8%	4.1%
Train Precision	61.0%	63.4%	64.1%	64.3%	62.2%	63.0%	1.4%
Val Precision	71.1%	61.5%	58.7%	57.8%	66.3%	62.3%	5.6%
Train Recall	100.0%	100.0%	100.0%	100.0%	100.0%	100.0%	0.0%
Val Recall	100.0%	100.0%	100.0%	100.0%	100.0%	100.0%	0.0%
Train MCC	6.5%	6.9%	7.0%	0.0%	6.7%	6.1%	3.0%
Val MCC	0.0%	0.0%	0.0%	12.2%	0.0%	0.0%	5.4%
Train Log-Loss	14.02	13.15	12.91	12.87	13.58	13.31	0.49
Val Log-Loss	10.41	13.89	14.89	15.05	12.15	13.60	1.97

Univariate Report

Mean Fractal Dimension - Kernel Density Plot

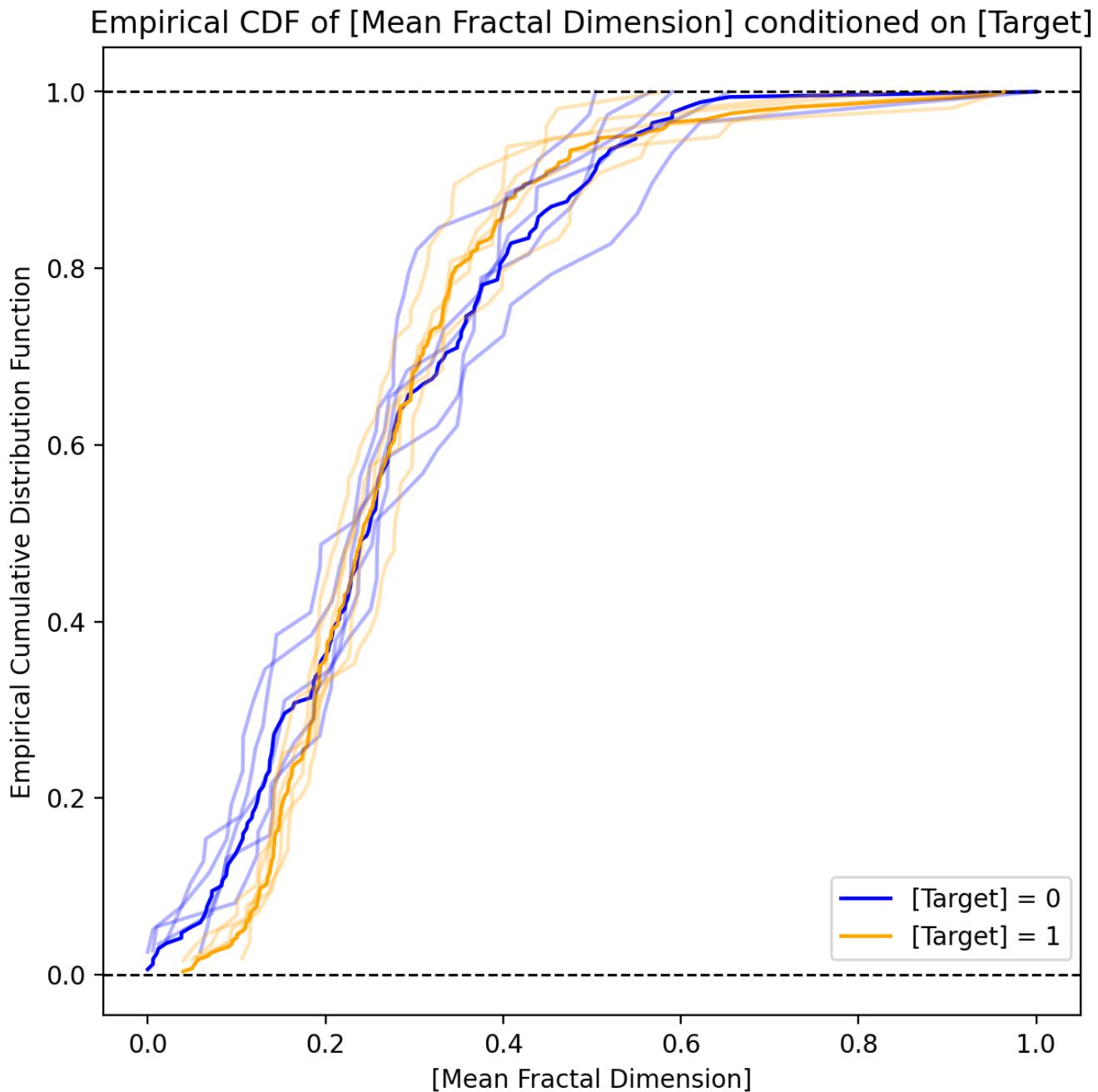
Kernel Density Plot of [Mean Fractal Dimension] by [Target]
Distributions by level are not significantly different at the 95% level.



This plot shows the Gaussian kernel density for each level of the target variable, both in total and for each fold. The x-axis represents the feature variable, and the y-axis represents the density of the target variable. The cross-validation folds are included in slightly washed-out colors to help understand the variability of the data. There are annotations with the results of a t-test for the difference in means between the feature variable at each level of the target variable. The annotations corresponding to the color of the target variable level show the mean/median ratio to help understand differences in skewness between the levels of the target variable.

Univariate Report

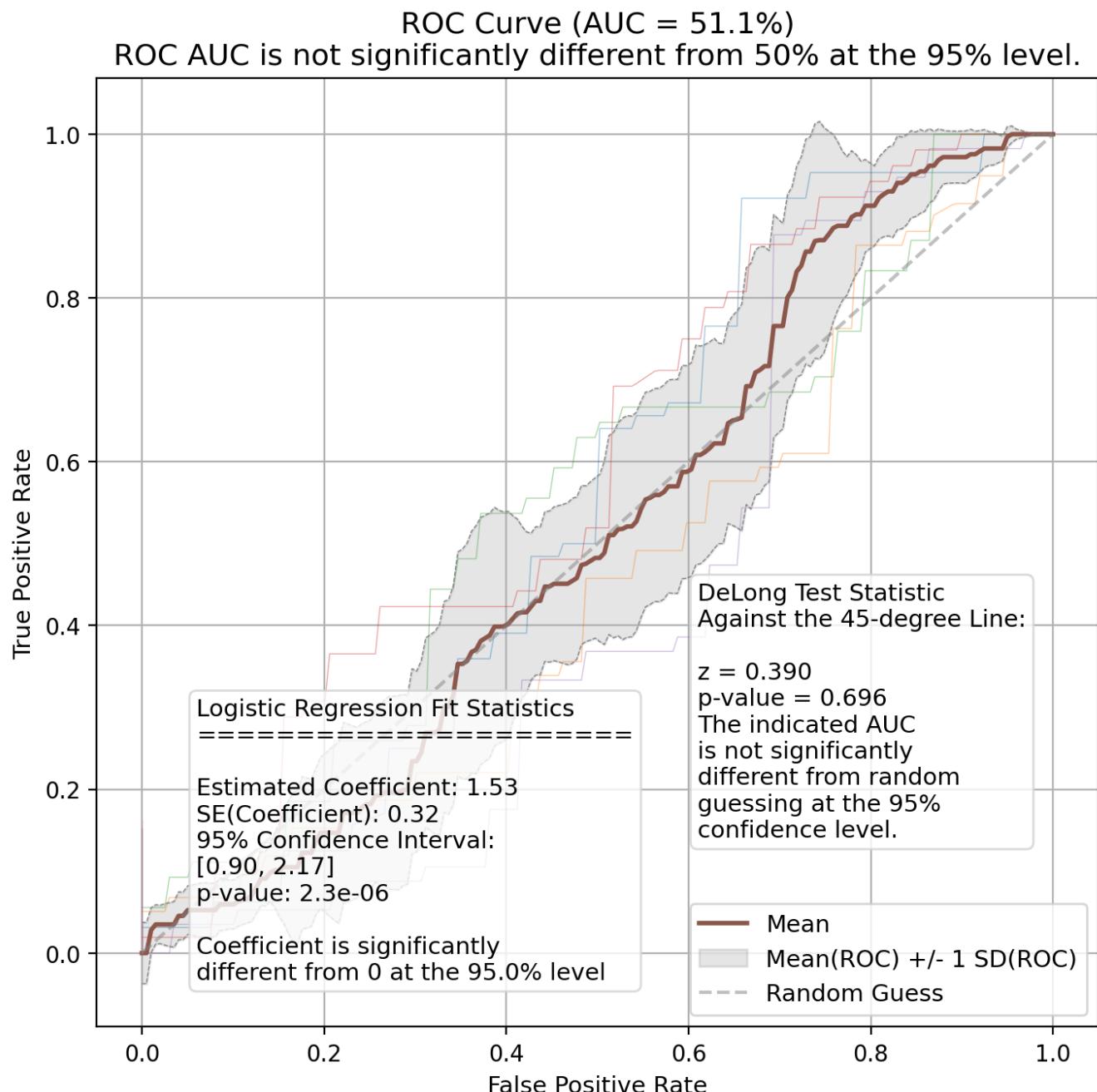
Mean Fractal Dimension - Empirical CDF Plot



This plot shows the empirical cumulative distribution function for each level of the target variable, both in total and for each fold. The x-axis represents the feature variable, and the y-axis represents the cumulative distribution of the target variable. The cross-validation folds are included in slightly washed-out colors to help understand the variability of the data, and whether or not it is reasonable to assume that the data is drawn from different distributions.

Univariate Report

Mean Fractal Dimension - ROC Curve



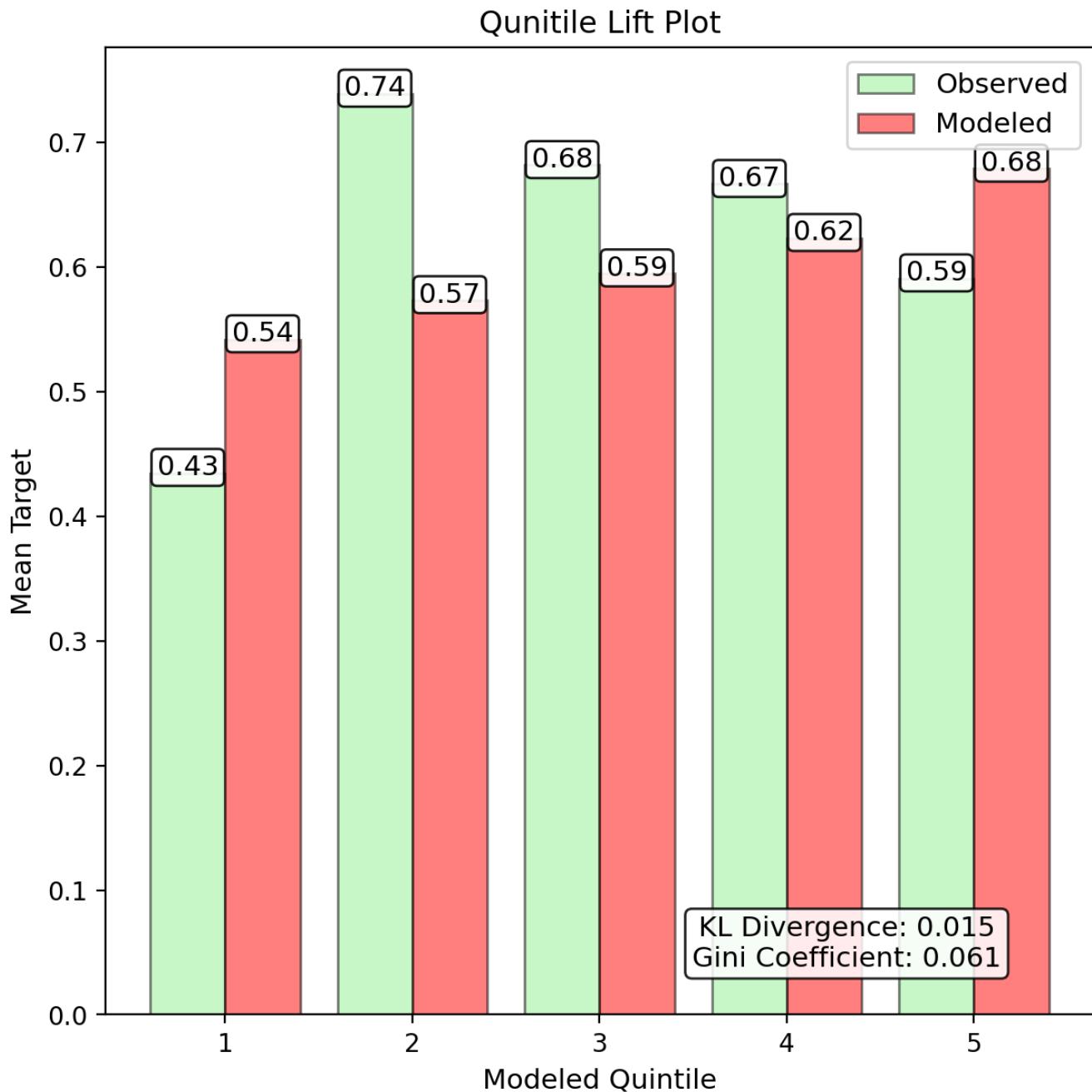
This plot shows the receiver operating characteristic (ROC) curve for the target variable in total and for each fold. The x-axis represents the false positive rate, and the y-axis represents the true positive rate. This is based on a simple Logistic Regression model with no regularization, no intercept, and no other features. Annotations are on the plot to help understand the results of the model, including the coefficient, standard error, and p-value for the feature variable. The cross-validation folds are used to create the grey region around the mean ROC curve to help understand the variability of the data.

Significance of the ROC curve is determined based on the method from DeLong et al. (1988). In brief, the AUC is assumed to be normally distributed, and the z-score is calculated based on the AUC and the standard error. This z-score is compared to a +/- two standard deviations from a standard normal

distribution to get the p-value.

Univariate Report

Mean Fractal Dimension - Quintile Lift



The quintile lift plot is meant to show the power of the single feature to discriminate between the highest and lowest quintiles of the target variable.

Univariate Report

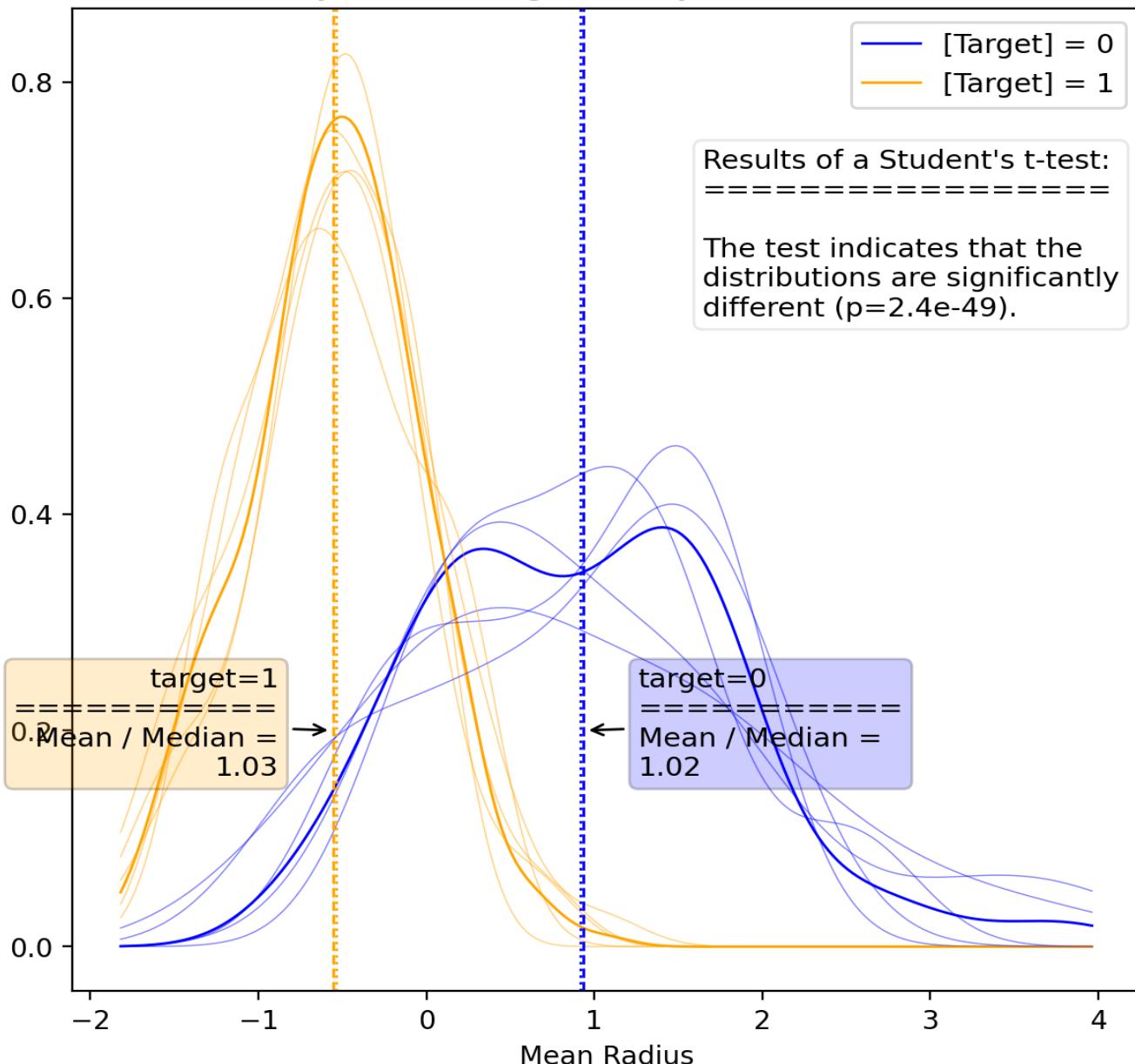
Mean Radius - Results

	Fold-1	Fold-2	Fold-3	Fold-4	Fold-5	Agg. Mean	Agg. SD
Fitted Coef.	-3.5e+00	-3.0e+00	-3.3e+00	-3.1e+00	-3.5e+00	-3.3e+00	2.3e-01
Fitted p-Value	7.6e-21	2.4e-21	5.9e-21	3.8e-21	7.1e-21	5.8e-26	2.2e-21
Fitted Std. Err.	0.376	0.313	0.355	0.332	0.369	0.310	0.026
Conf. Int. Lower	-4.3e+00	-3.6e+00	-4.0e+00	-3.8e+00	-4.2e+00	-3.9e+00	2.8e-01
Conf. Int. Upper	-2.8e+00	-2.4e+00	-2.6e+00	-2.5e+00	-2.7e+00	-2.7e+00	1.8e-01
Train Accuracy	86.6%	84.1%	87.6%	84.9%	86.4%	85.9%	1.4%
Val Accuracy	83.3%	92.7%	79.3%	90.1%	83.7%	37.7%	5.4%
Train AUC	86.4%	83.8%	87.3%	84.5%	86.3%	85.7%	1.5%
Val AUC	82.6%	92.6%	79.3%	90.1%	82.6%	50.0%	5.6%
Train F1	88.8%	87.1%	90.1%	88.0%	88.8%	88.6%	1.1%
Test F1	87.8%	94.0%	81.9%	91.3%	87.5%	0.0%	4.6%
Train Precision	90.2%	89.4%	91.9%	90.1%	90.9%	90.5%	1.0%
Val Precision	91.5%	94.8%	84.3%	92.2%	89.1%	0.0%	4.0%
Train Recall	87.4%	85.0%	88.4%	85.9%	86.9%	86.7%	1.3%
Val Recall	84.4%	93.2%	79.6%	90.4%	86.0%	0.0%	5.3%
Train MCC	72.1%	66.6%	73.6%	67.8%	71.7%	70.4%	3.0%
Val MCC	62.1%	84.7%	58.0%	79.9%	64.3%	0.0%	11.7%
Train Log-Loss	4.84	5.72	4.47	5.45	4.88	5.07	0.50
Val Log-Loss	6.01	2.63	7.44	3.56	5.87	22.45	1.96

Univariate Report

Mean Radius - Kernel Density Plot

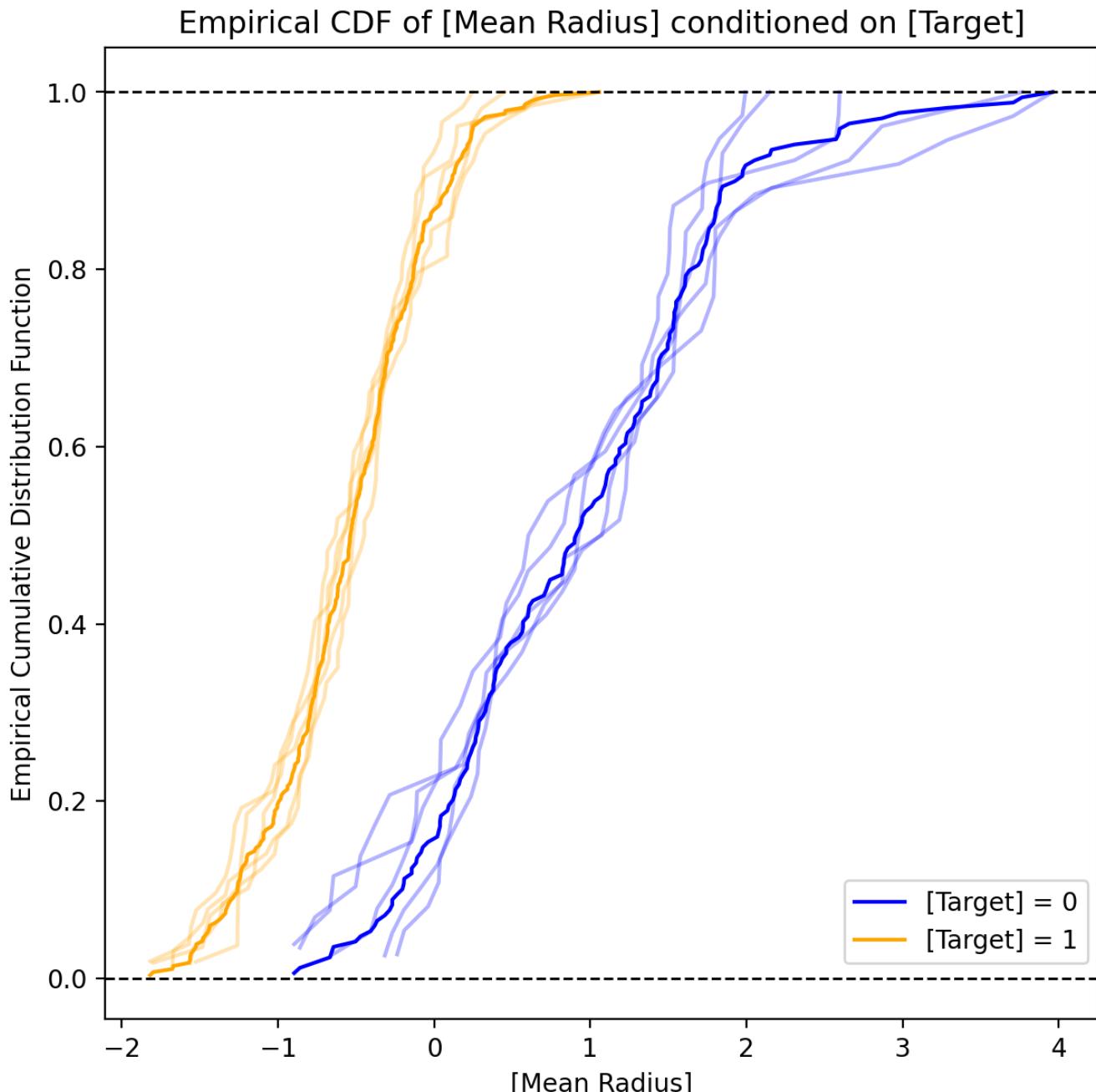
Kernel Density Plot of [Mean Radius] by [Target]
Distributions by level are significantly different at the 95% level.



This plot shows the Gaussian kernel density for each level of the target variable, both in total and for each fold. The x-axis represents the feature variable, and the y-axis represents the density of the target variable. The cross-validation folds are included in slightly washed-out colors to help understand the variability of the data. There are annotations with the results of a t-test for the difference in means between the feature variable at each level of the target variable. The annotations corresponding to the color of the target variable level show the mean/median ratio to help understand differences in skewness between the levels of the target variable.

Univariate Report

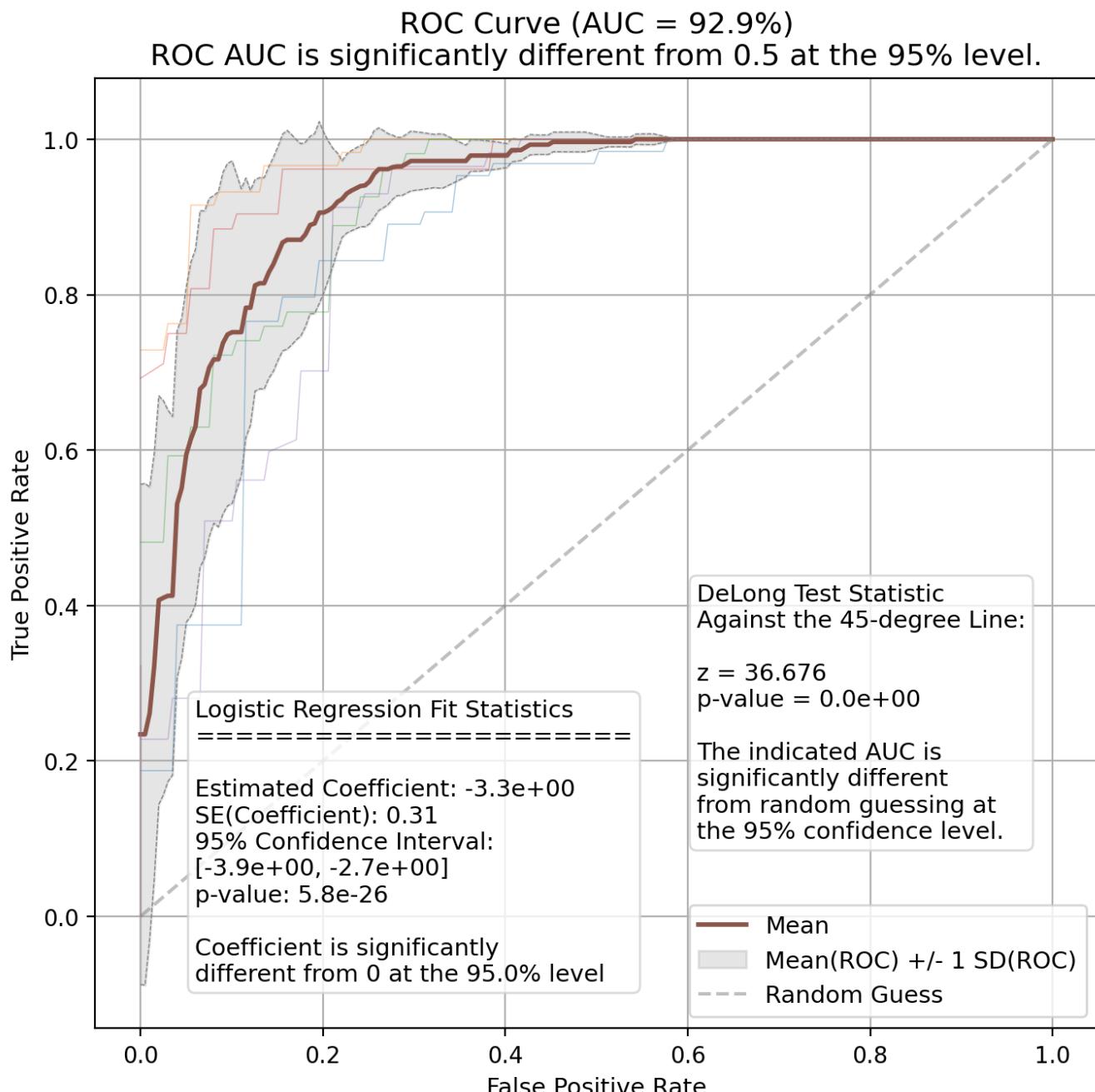
Mean Radius - Empirical CDF Plot



This plot shows the empirical cumulative distribution function for each level of the target variable, both in total and for each fold. The x-axis represents the feature variable, and the y-axis represents the cumulative distribution of the target variable. The cross-validation folds are included in slightly washed-out colors to help understand the variability of the data, and whether or not it is reasonable to assume that the data is drawn from different distributions.

Univariate Report

Mean Radius - ROC Curve



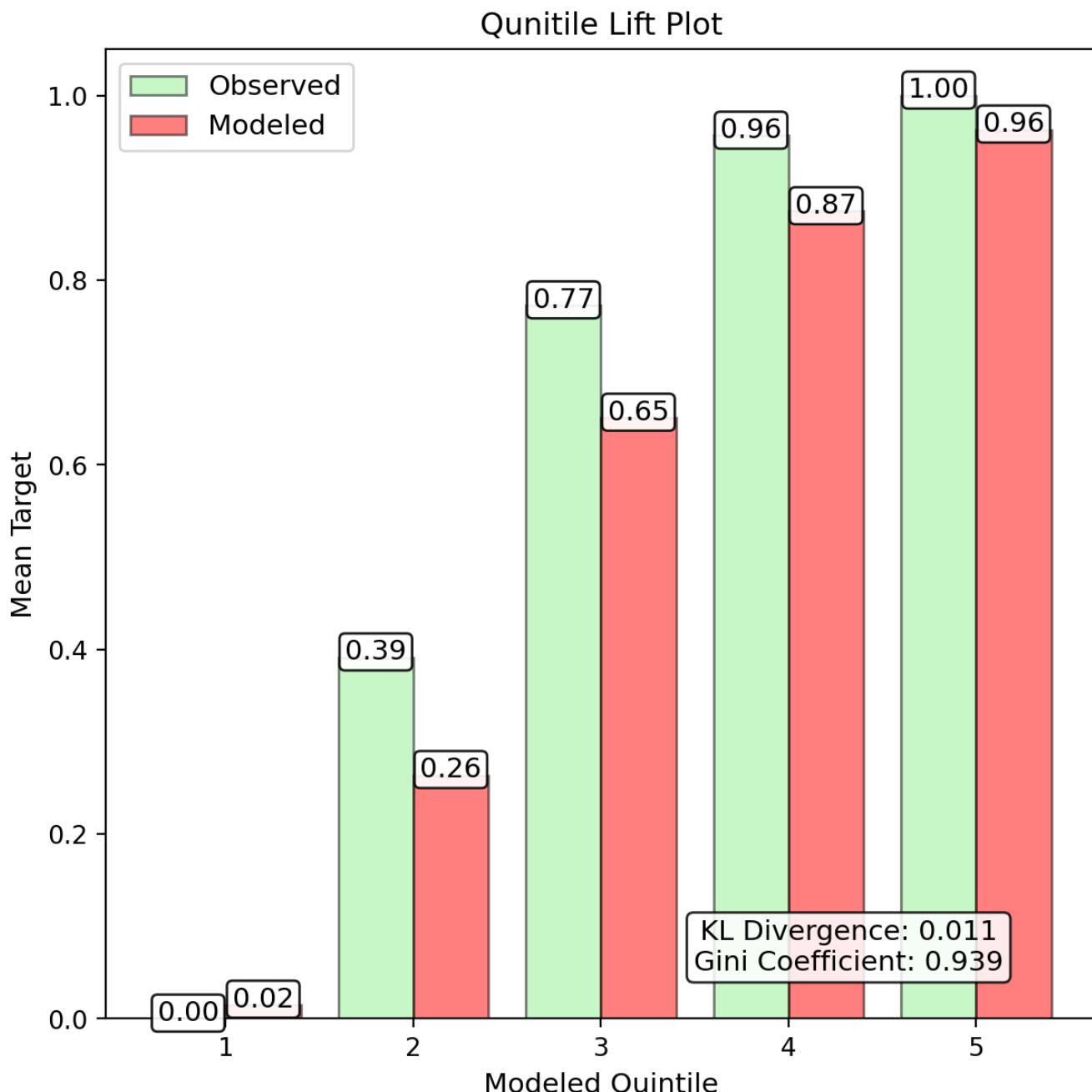
This plot shows the receiver operating characteristic (ROC) curve for the target variable in total and for each fold. The x-axis represents the false positive rate, and the y-axis represents the true positive rate. This is based on a simple Logistic Regression model with no regularization, no intercept, and no other features. Annotations are on the plot to help understand the results of the model, including the coefficient, standard error, and p-value for the feature variable. The cross-validation folds are used to create the grey region around the mean ROC curve to help understand the variability of the data.

Significance of the ROC curve is determined based on the method from DeLong et al. (1988). In brief, the AUC is assumed to be normally distributed, and the z-score is calculated based on the AUC and the standard error. This z-score is compared to a +/- two standard deviations from a standard normal

distribution to get the p-value.

Univariate Report

Mean Radius - Quintile Lift



The quintile lift plot is meant to show the power of the single feature to discriminate between the highest and lowest quintiles of the target variable.