

Pyridine Phosphination - Best Model Feature Analysis

Table of Contents

Published Model	1
5th Place with 5 Features.....	1
Model Analysis - Variable Impact.....	7
1. Avg_NPA_SM Removed.....	7
2. fr_Ar_N Removed.....	12
3. fr_aryl_methyl Removed.....	17
4. fr_benzene Removed.....	22
5. fr_halogen Removed	27
Section Results Summary.....	32

Published Model

5th Place with 5 Features

Table:

Full Model Stats - Overall Accuracy and Pseudo-R2

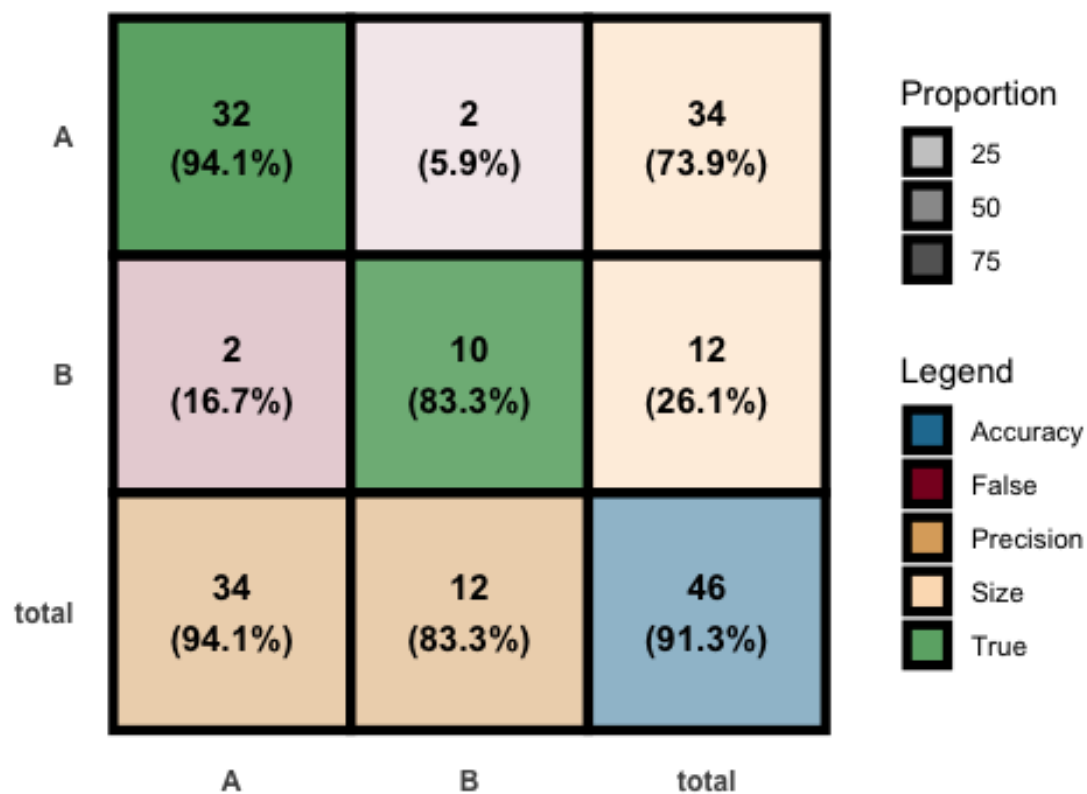
Accuracy	McFadden_R2
91.3%	0.578

Model Coefficients

	x
(Intercept)	38.099048
Avg_NPA_SM	48.412036
fr_Ar_N1	-31.903705
fr_Ar_N2	-60.673746
fr_Ar_N3	-66.149054
fr_aryl_methyl1	-4.592317
fr_aryl_methyl2	-37.006653
fr_benzene1	-1.970290
fr_benzene2	-77.887220
fr_halogen1	1.607589
fr_halogen2	2.929693
fr_halogen3	2.595277
fr_halogen4	-2.087011

	x
fr_halogen5	2.187832
fr_halogen6	-2.311938

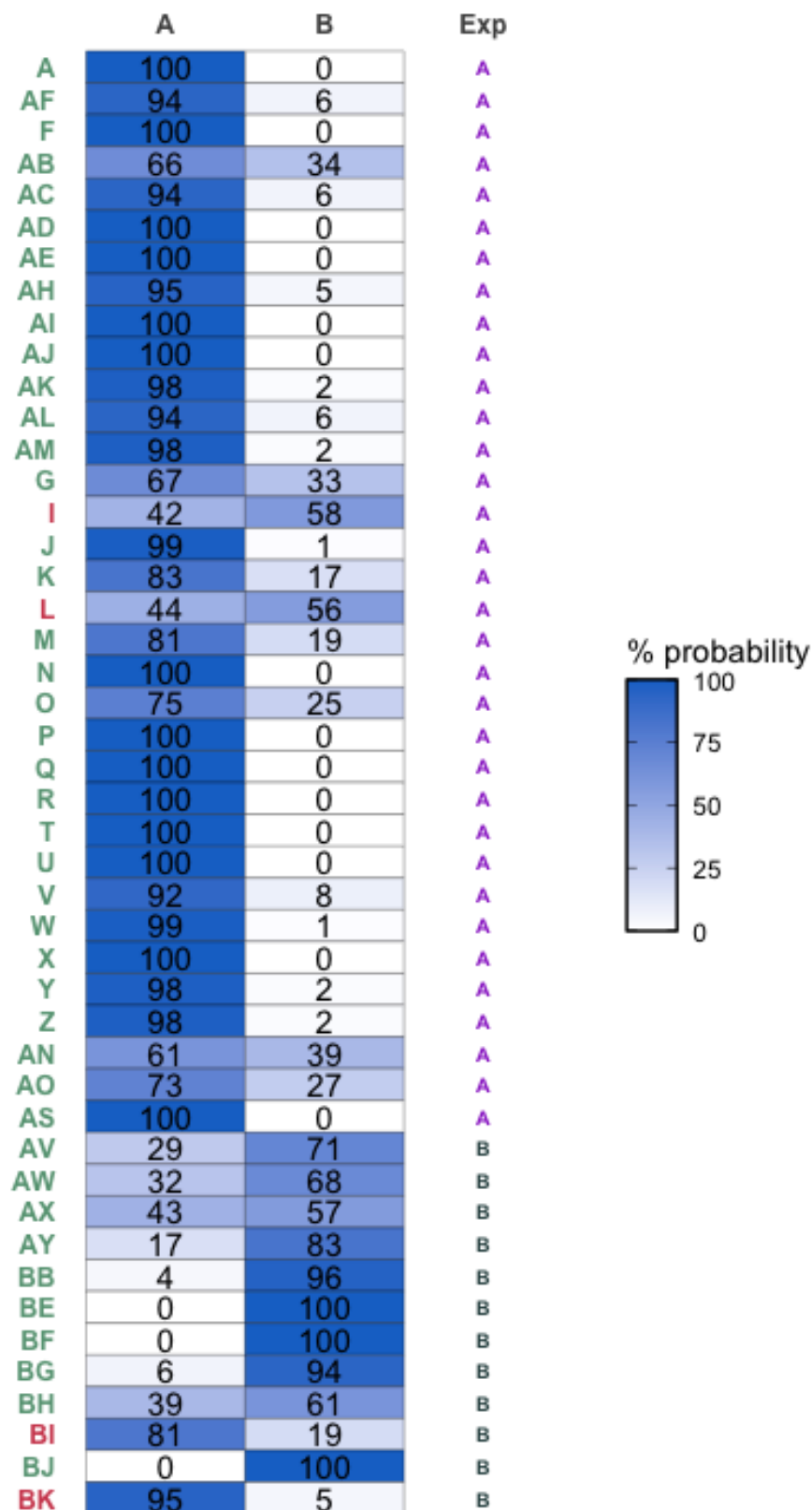
Training Set Confusion Matrix



5. 5th Place - 5 features

Training Set

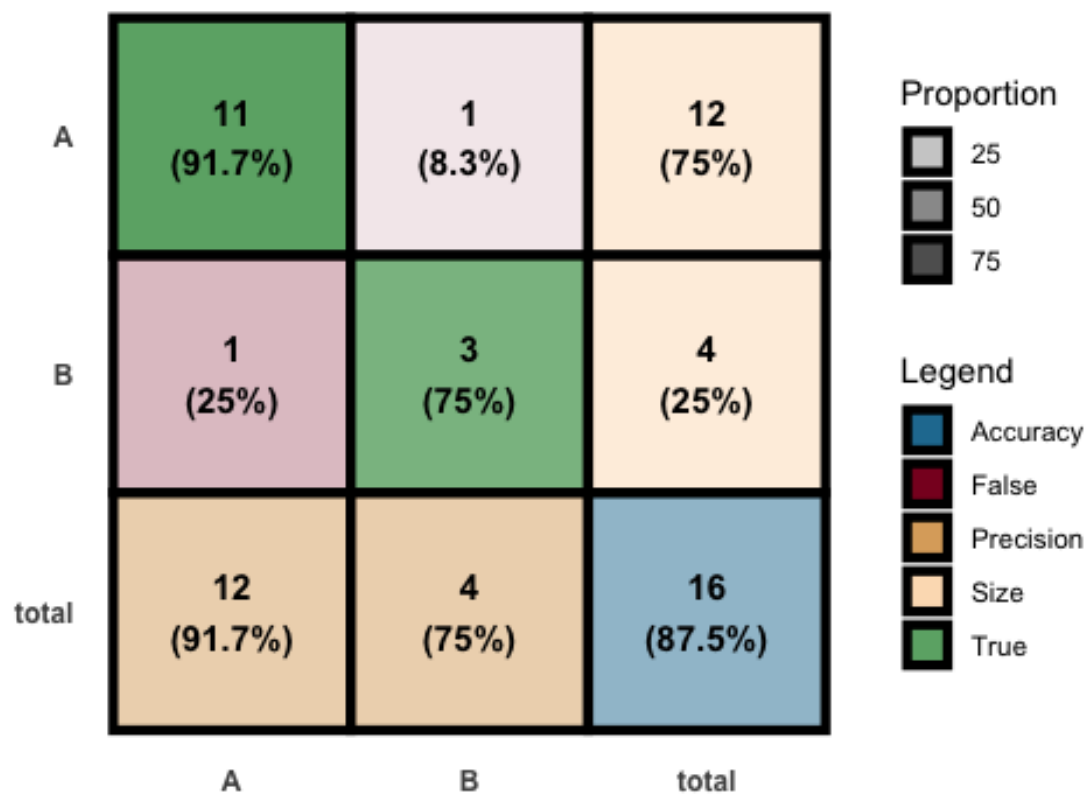
Probability Heatmap



5. 5th Place - 5 features

Test Set

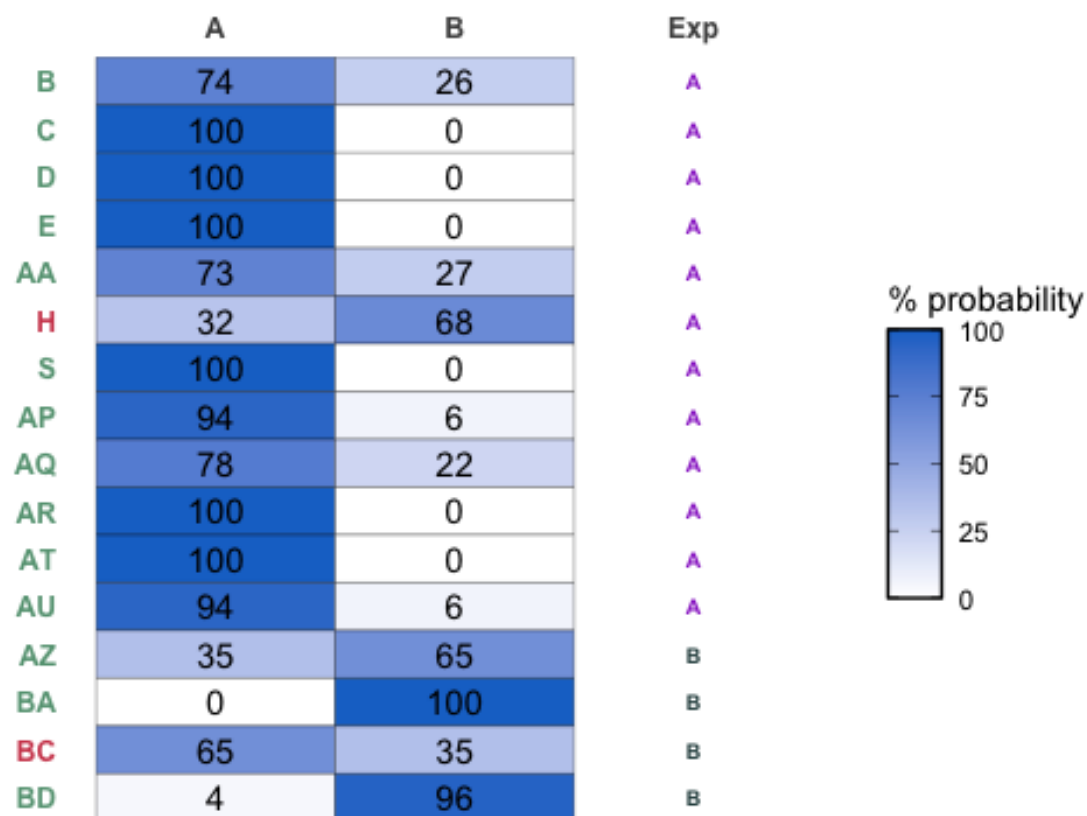
Confusion Matrix



5. 5th Place - 5 features

Test Set

Probability Heatmap



5. 5th Place - 5 features

Model Analysis - Variable Impact

1. Avg_NPA_SM Removed

formula

class ~ fr_Ar_N + fr_aryl_methyl + fr_benzene + fr_halogen

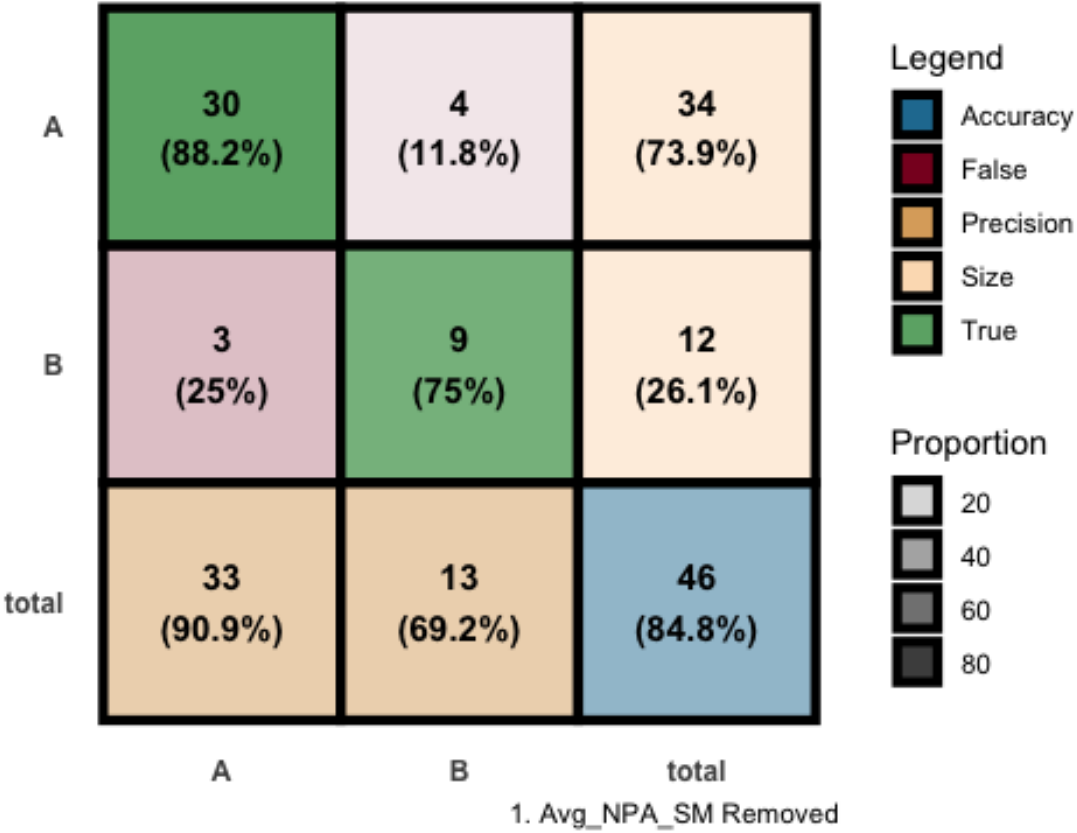
Full Model Stats - Overall Accuracy and Pseudo-R2

Accuracy	McFadden_R2
84.78%	0.52

Model Coefficients

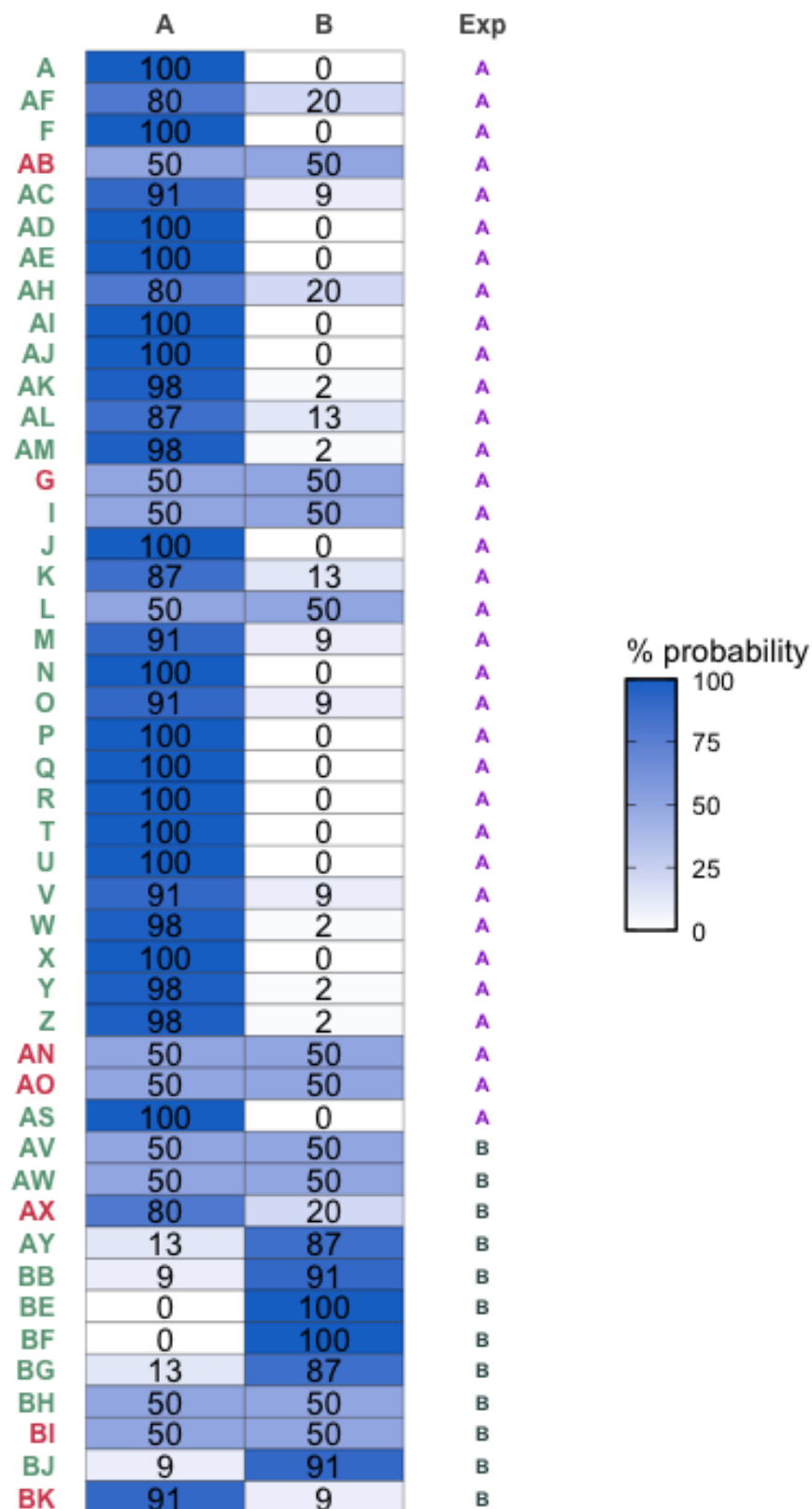
	x
(Intercept)	14.653960
fr_Ar_N1	-14.641780
fr_Ar_N2	-28.886527
fr_Ar_N3	-30.760596
fr_aryl_methyl1	-3.720066
fr_aryl_methyl2	-16.691481
fr_benzene1	-2.364301
fr_benzene2	-35.077472
fr_halogen1	2.342569
fr_halogen2	4.212291
fr_halogen3	1.848104
fr_halogen4	-1.081430
fr_halogen5	2.352174
fr_halogen6	-9.373352

Training Set
Confusion Matrix



Training Set

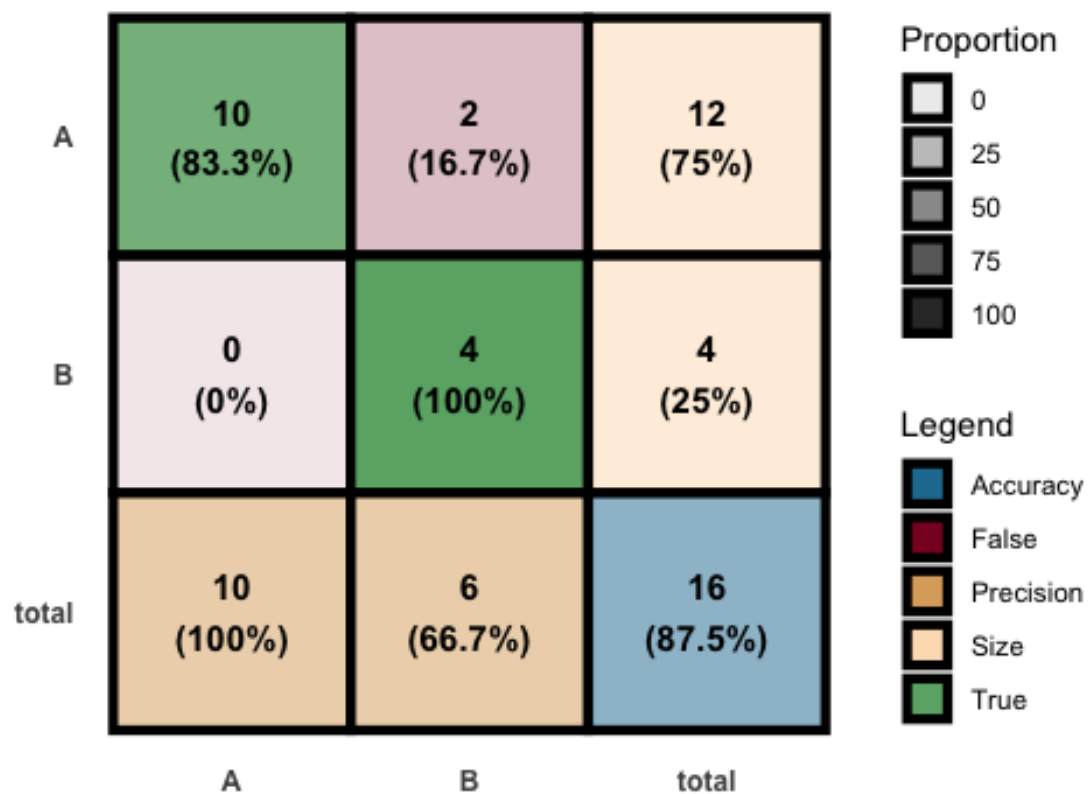
Probability Heatmap



1. Avg_NPA_SM Removed

Test Set

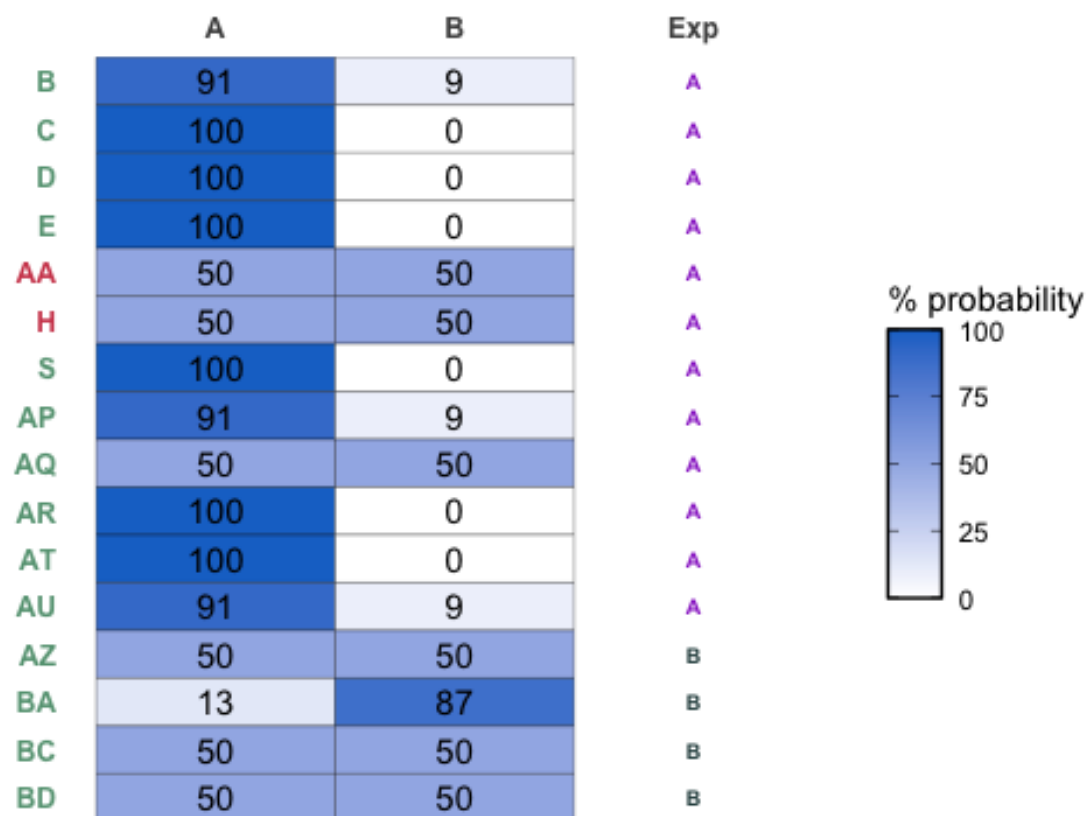
Confusion Matrix



1. Avg_NPA_SM Removed

Test Set

Probability Heatmap



1. Avg_NPA_SM Removed

2. fr_Ar_N Removed

formula

class ~ Avg_NPA_SM + fr_aryl_methyl + fr_benzene + fr_halogen

Full Model Stats - Overall Accuracy and Pseudo-R2

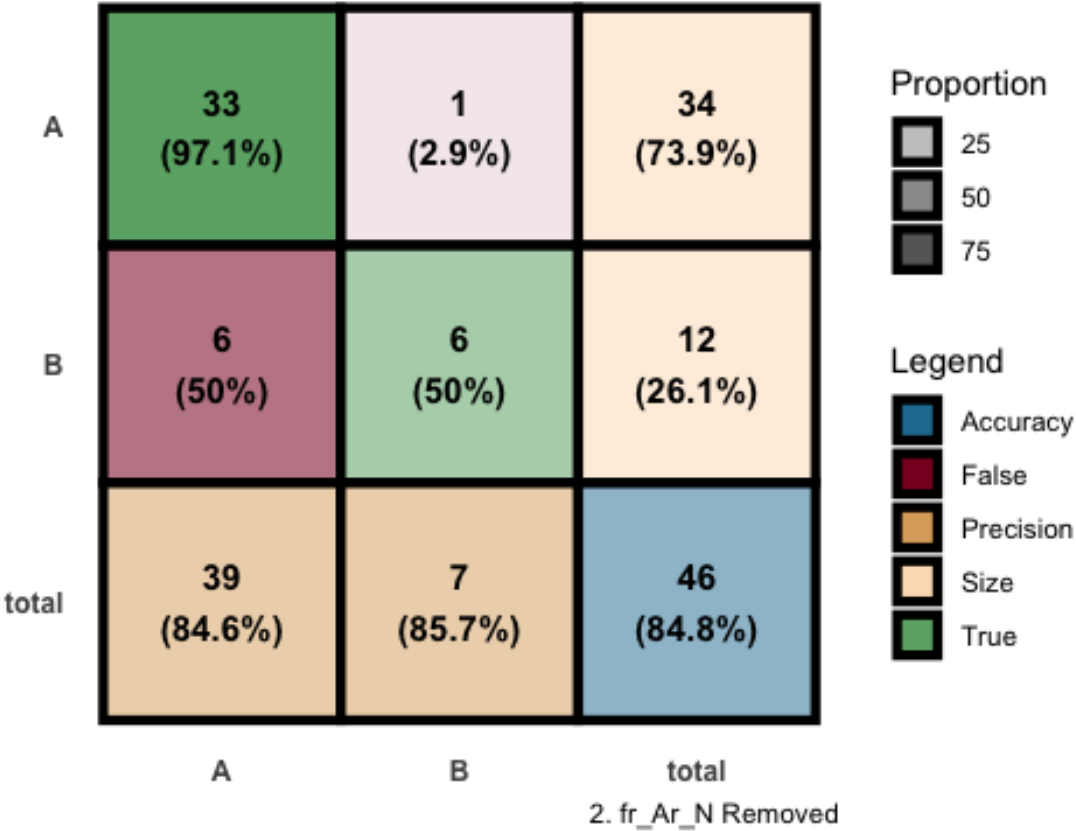
Accuracy	McFadden_R2
----------	-------------

84.78%	0.328
--------	-------

Model Coefficients

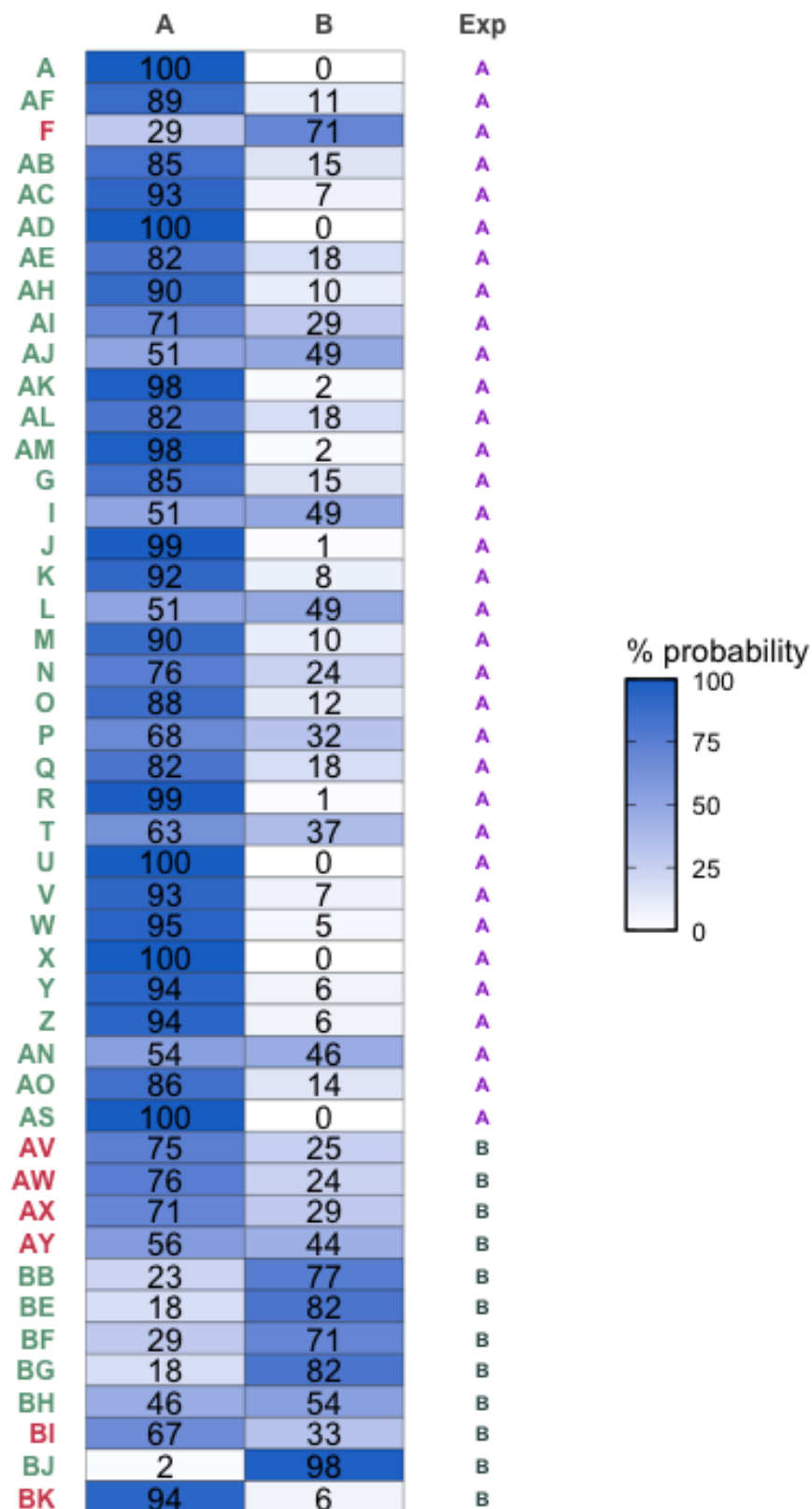
	x
(Intercept)	0.9188853
Avg_NPA_SM	18.4857033
fr_aryl_methyl1	-2.7592175
fr_aryl_methyl2	-12.9215911
fr_benzene1	-0.8918341
fr_benzene2	-32.6158024
fr_halogen1	2.0595415
fr_halogen2	3.2075153
fr_halogen3	1.5695469
fr_halogen4	-8.4512972
fr_halogen5	2.4217225
fr_halogen6	-4.8502196

Training Set
Confusion Matrix



Training Set

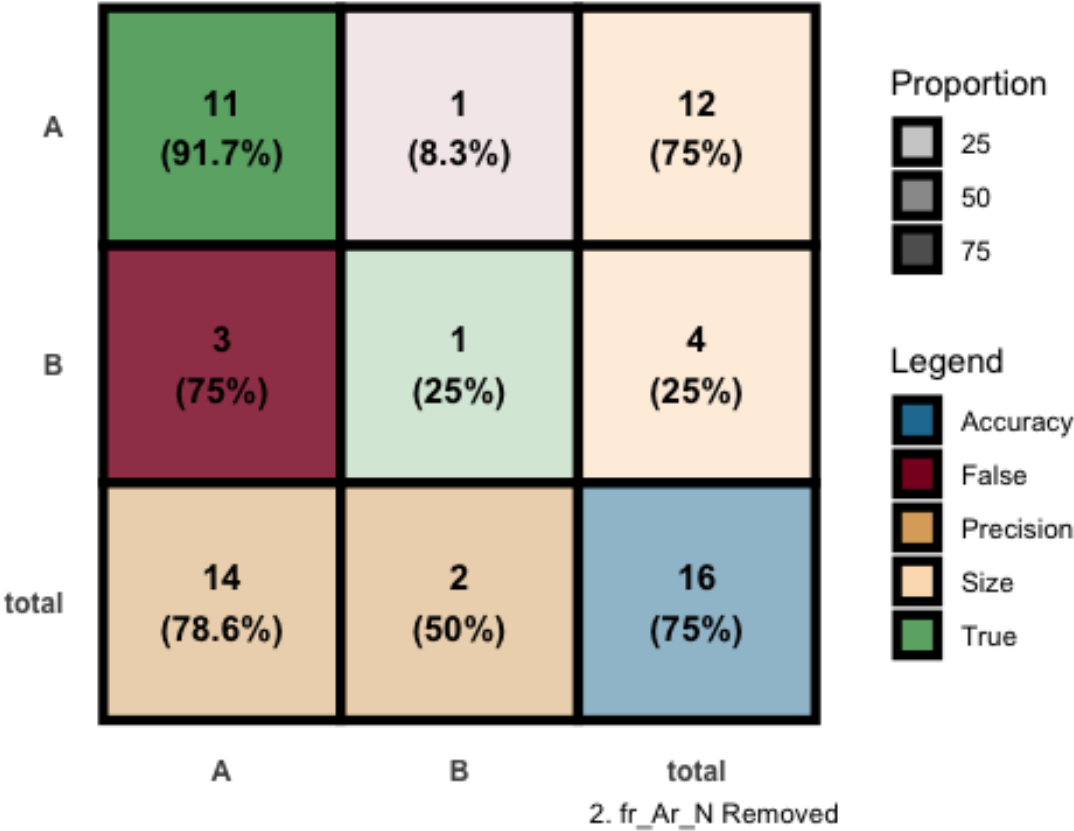
Probability Heatmap



2. fr_Ar_N Removed

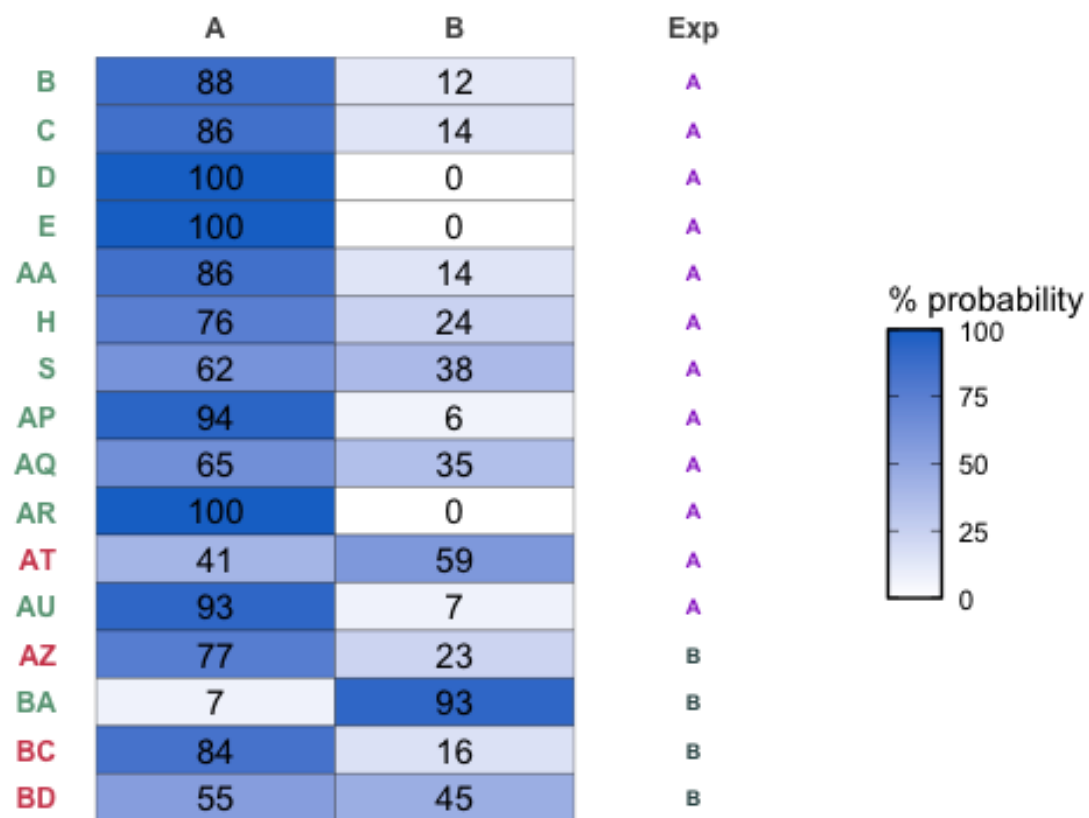
Test Set

Confusion Matrix



Test Set

Probability Heatmap



2. fr_Ar_N Removed

3. fr_aryl_methyl Removed

formula

class ~ Avg_NPA_SM + fr_Ar_N + fr_benzene + fr_halogen

Full Model Stats - Overall Accuracy and Pseudo-R2

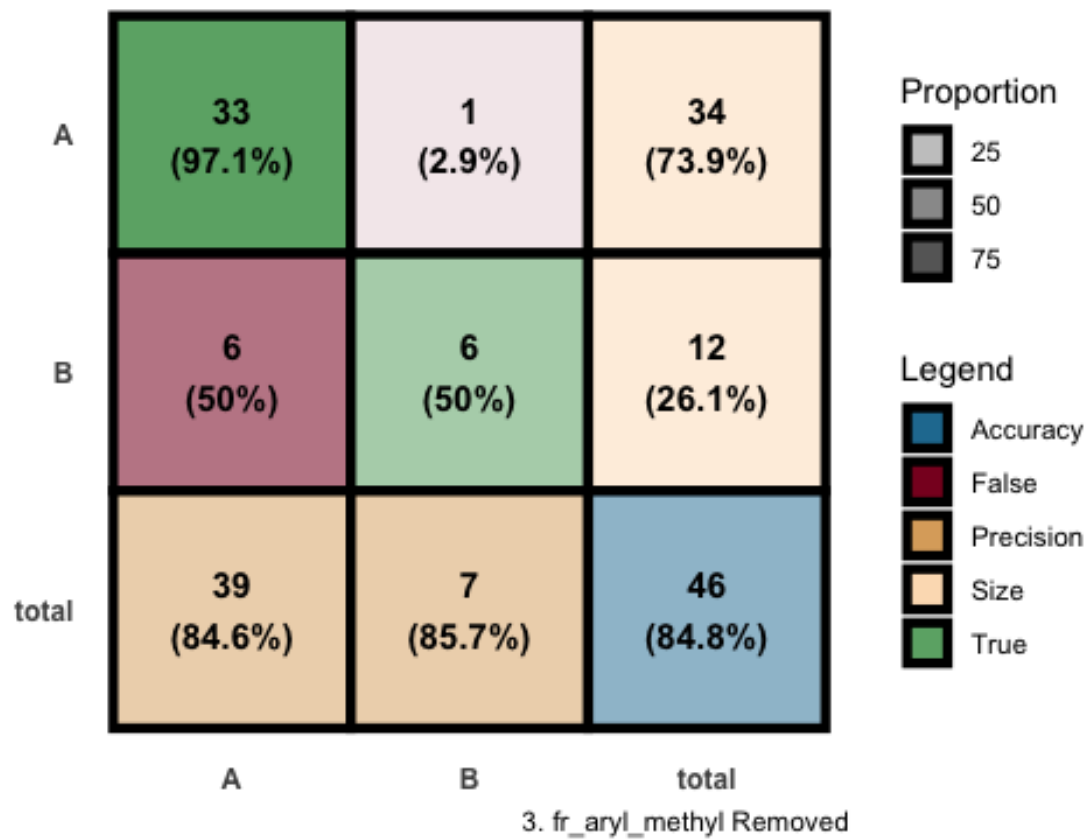
Accuracy	McFadden_R2
----------	-------------

84.78%	0.36
--------	------

Model Coefficients

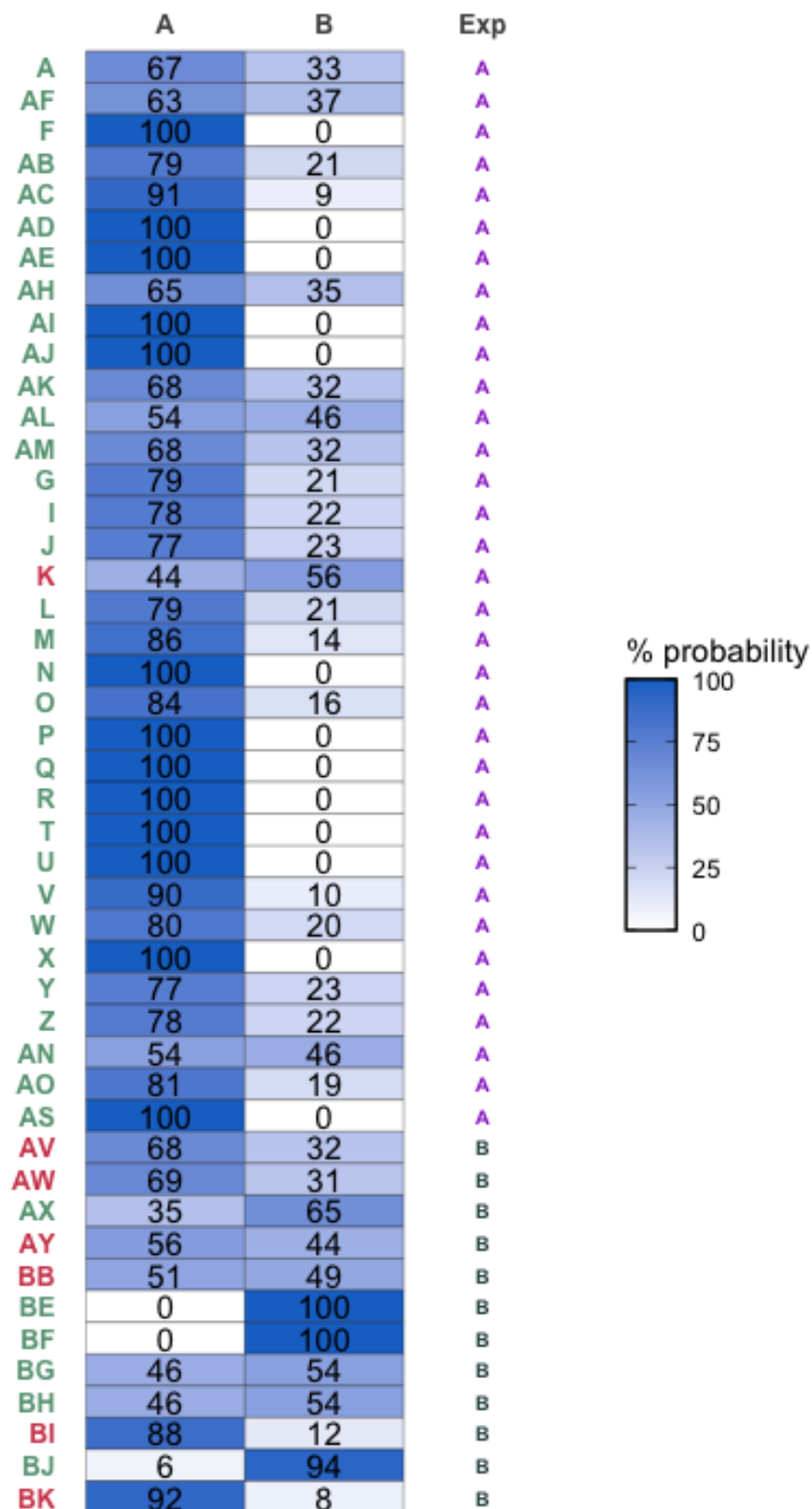
	x
(Intercept)	18.3982999
Avg_NPA_SM	18.5844081
fr_Ar_N1	-17.1017267
fr_Ar_N2	-29.0614501
fr_Ar_N3	-32.9558398
fr_benzene1	-0.9195207
fr_benzene2	-27.8437909
fr_halogen1	0.4496324
fr_halogen2	1.4880482
fr_halogen3	1.1965587
fr_halogen4	-1.7308296
fr_halogen5	2.0846872
fr_halogen6	-9.0103971

Training Set Confusion Matrix



Training Set

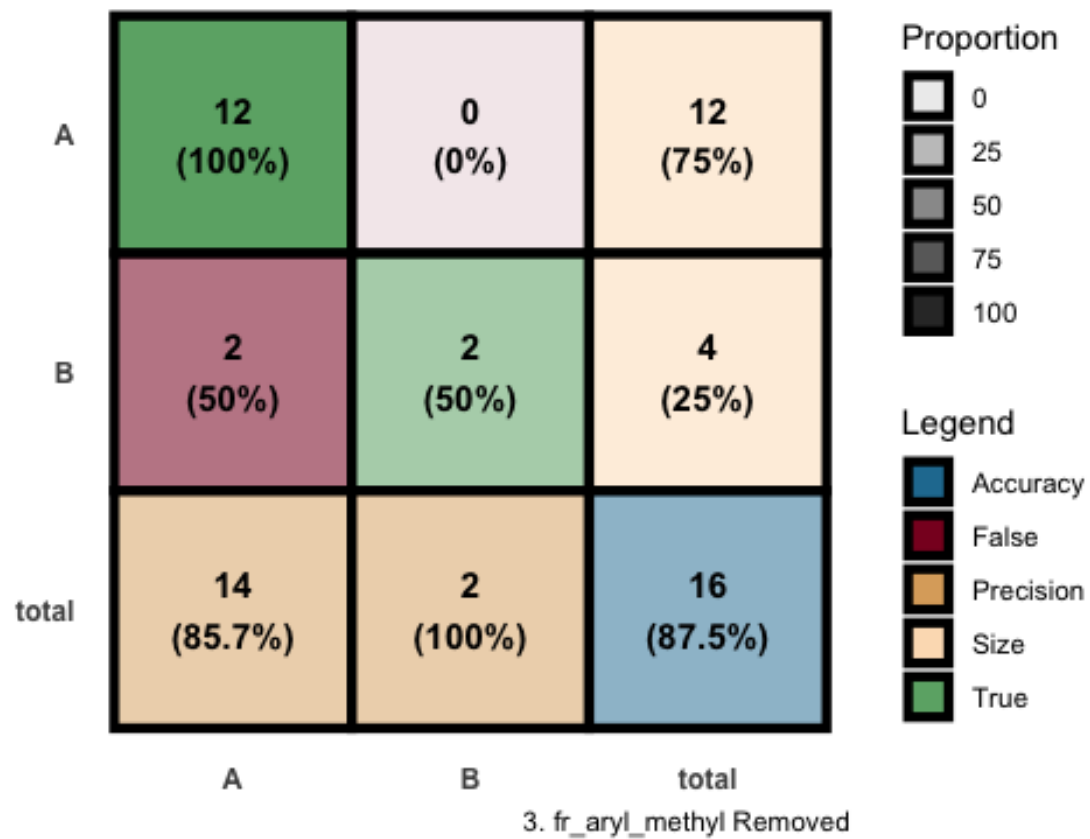
Probability Heatmap



3. fr_aryl_methyl Removed

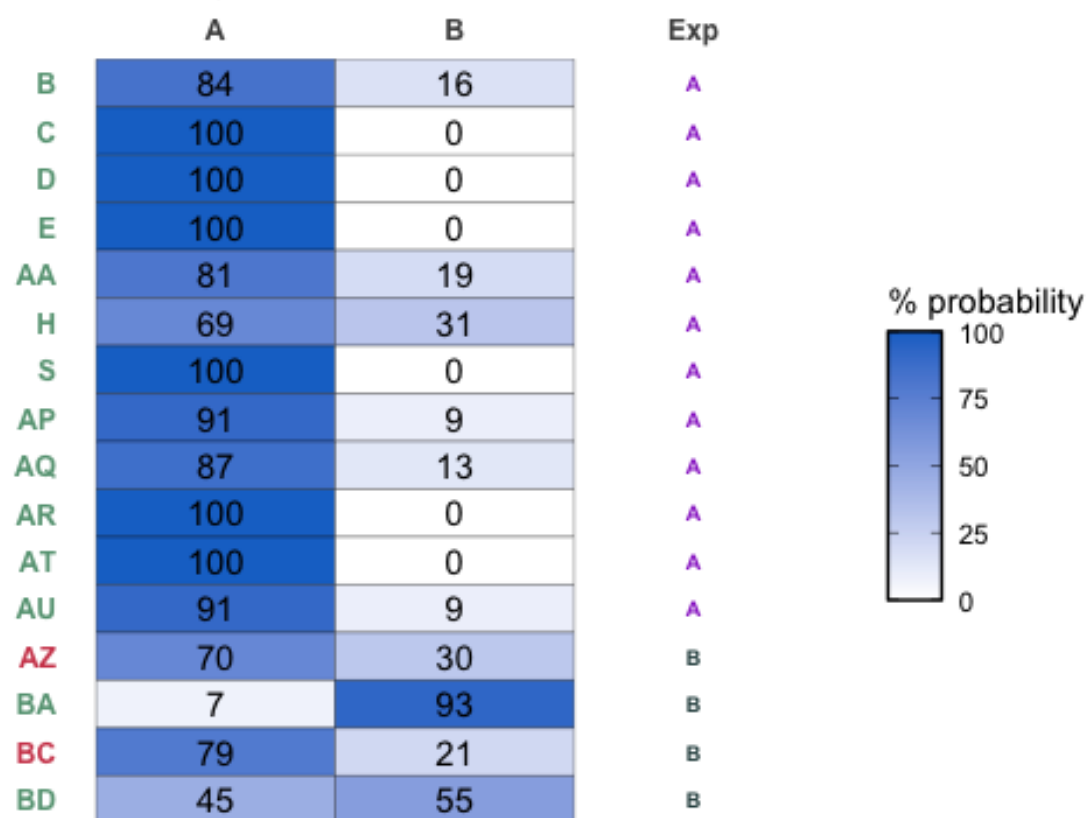
Test Set

Confusion Matrix



Test Set

Probability Heatmap



3. fr_aryl_methyl Removed

4. fr_benzene Removed

formula

class ~ Avg_NPA_SM + fr_Ar_N + fr_aryl_methyl + fr_halogen

Full Model Stats - Overall Accuracy and Pseudo-R2

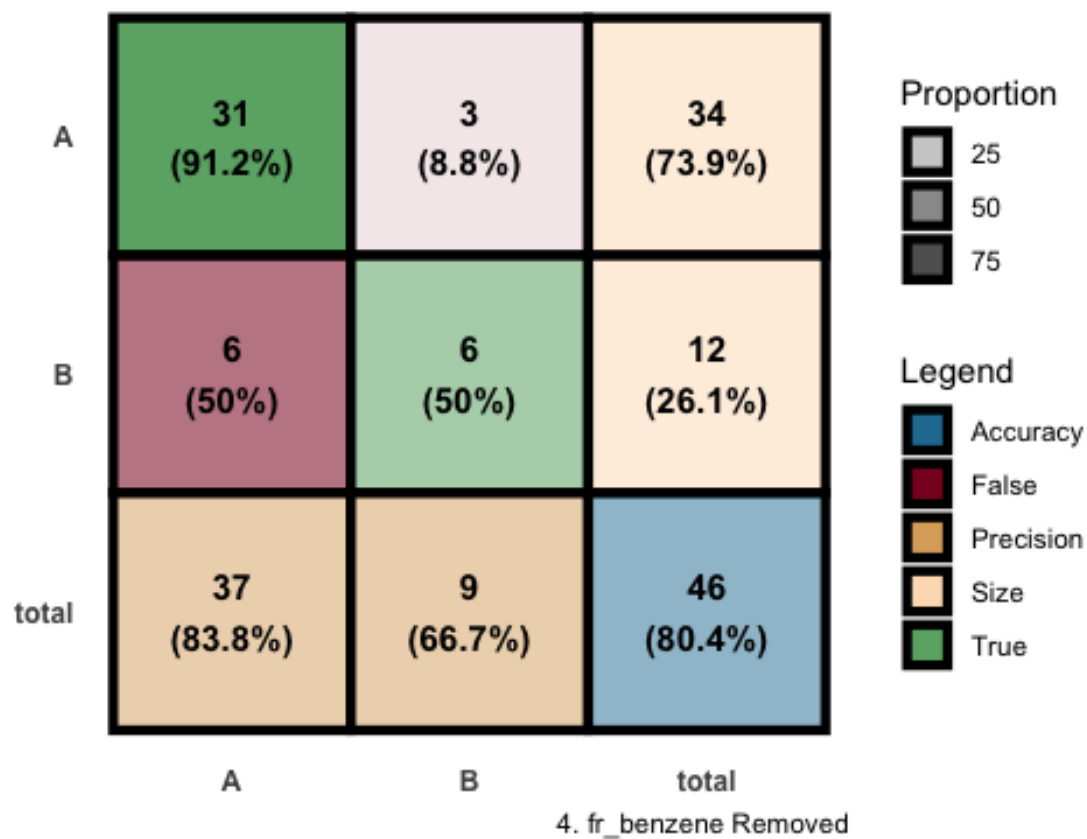
Accuracy	McFadden_R2
----------	-------------

80.43%	0.404
--------	-------

Model Coefficients

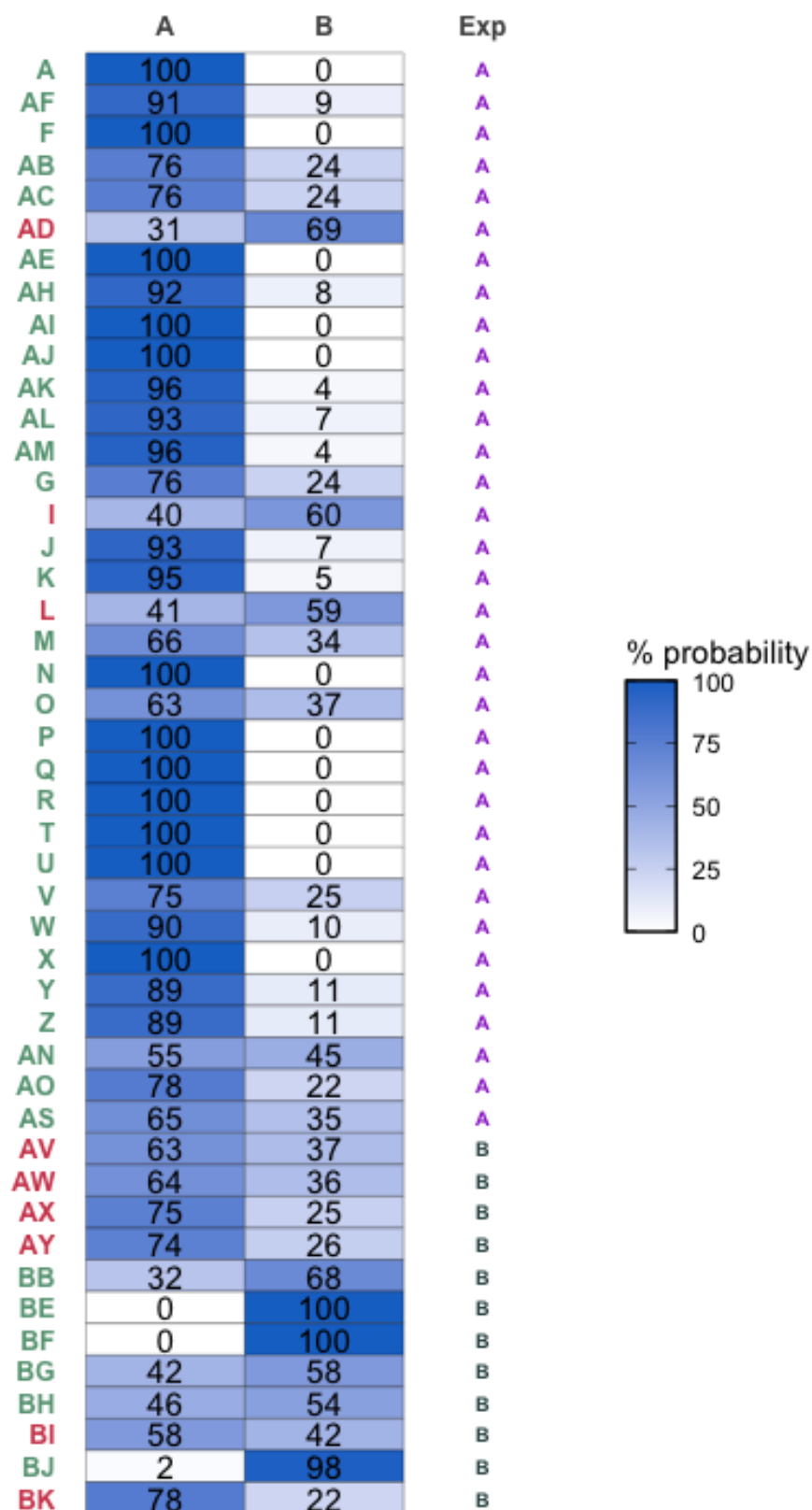
	x
(Intercept)	15.3364807
Avg_NPA_SM	19.3241422
fr_Ar_N1	-13.7390649
fr_Ar_N2	-26.0021302
fr_Ar_N3	-34.6453926
fr_aryl_methyl1	-2.5523230
fr_aryl_methyl2	-17.8132254
fr_halogen1	1.0049517
fr_halogen2	0.4974903
fr_halogen3	0.2314174
fr_halogen4	-2.3972048
fr_halogen5	0.9617090
fr_halogen6	-8.0304964

Training Set Confusion Matrix



Training Set

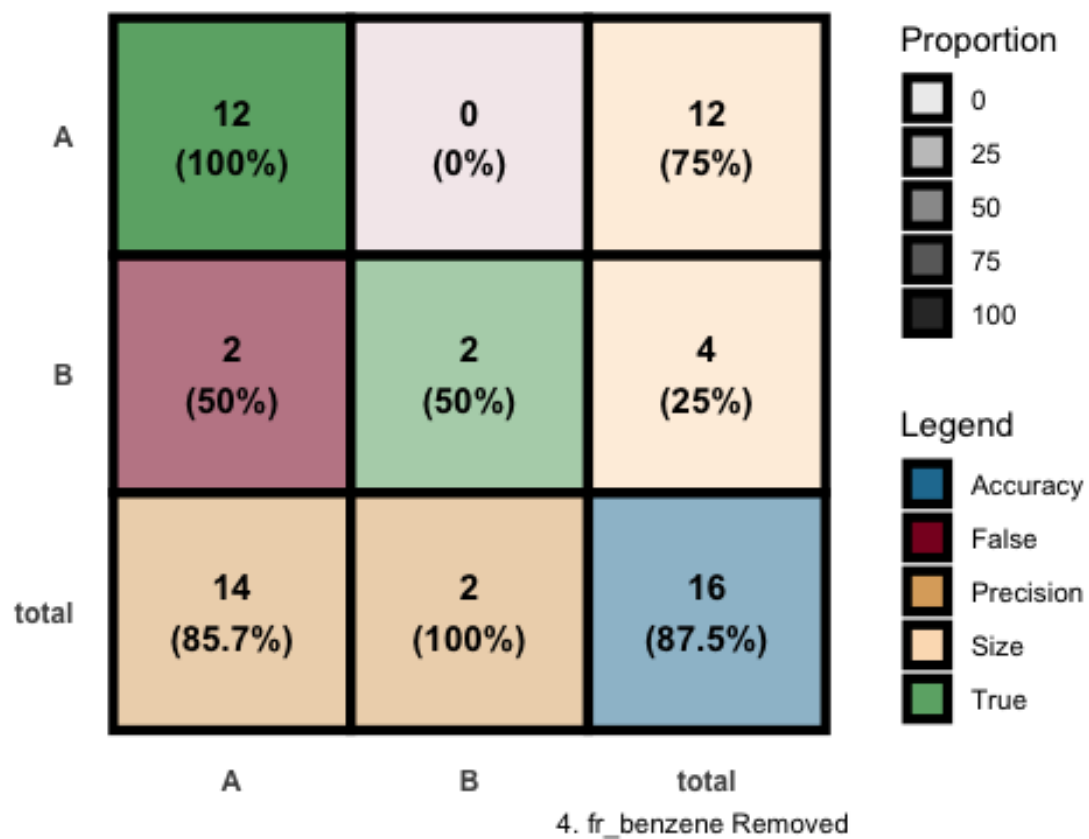
Probability Heatmap



4. fr_benzene Removed

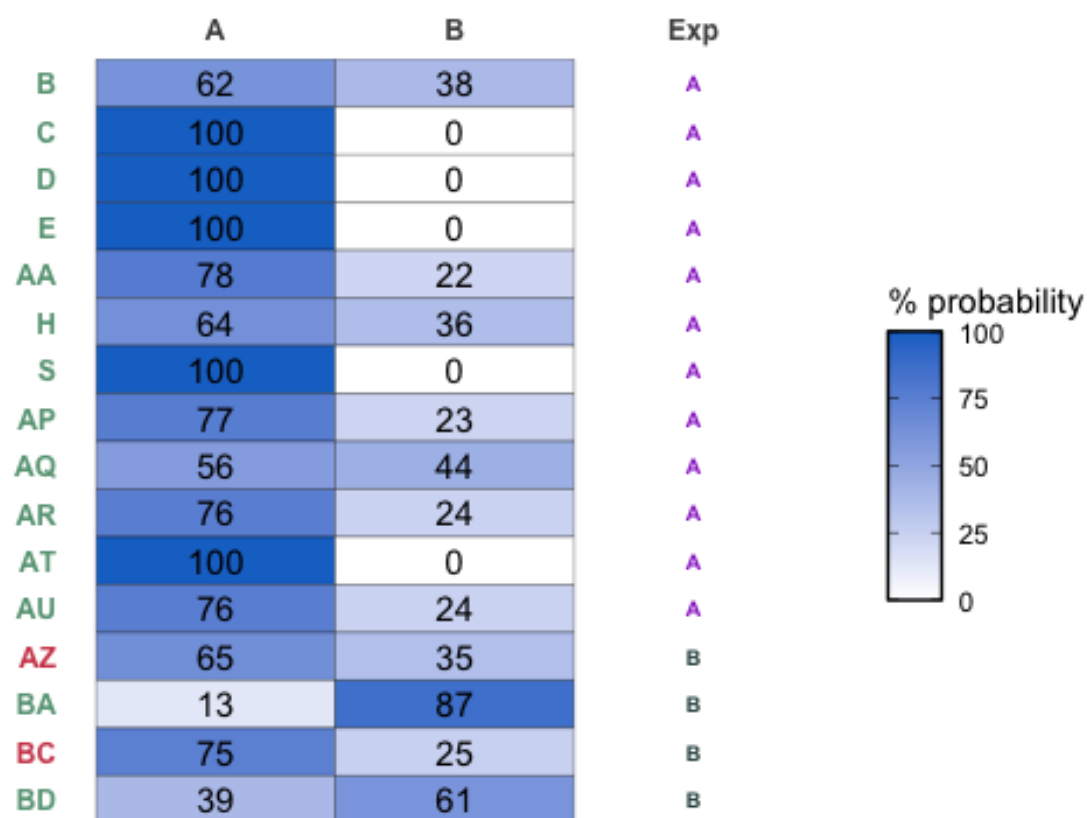
Test Set

Confusion Matrix



Test Set

Probability Heatmap



4. fr_benzene Removed

5. fr_halogen Removed

formula

class ~ Avg_NPA_SM + fr_Ar_N + fr_aryl_methyl + fr_benzene

Full Model Stats - Overall Accuracy and Pseudo-R2

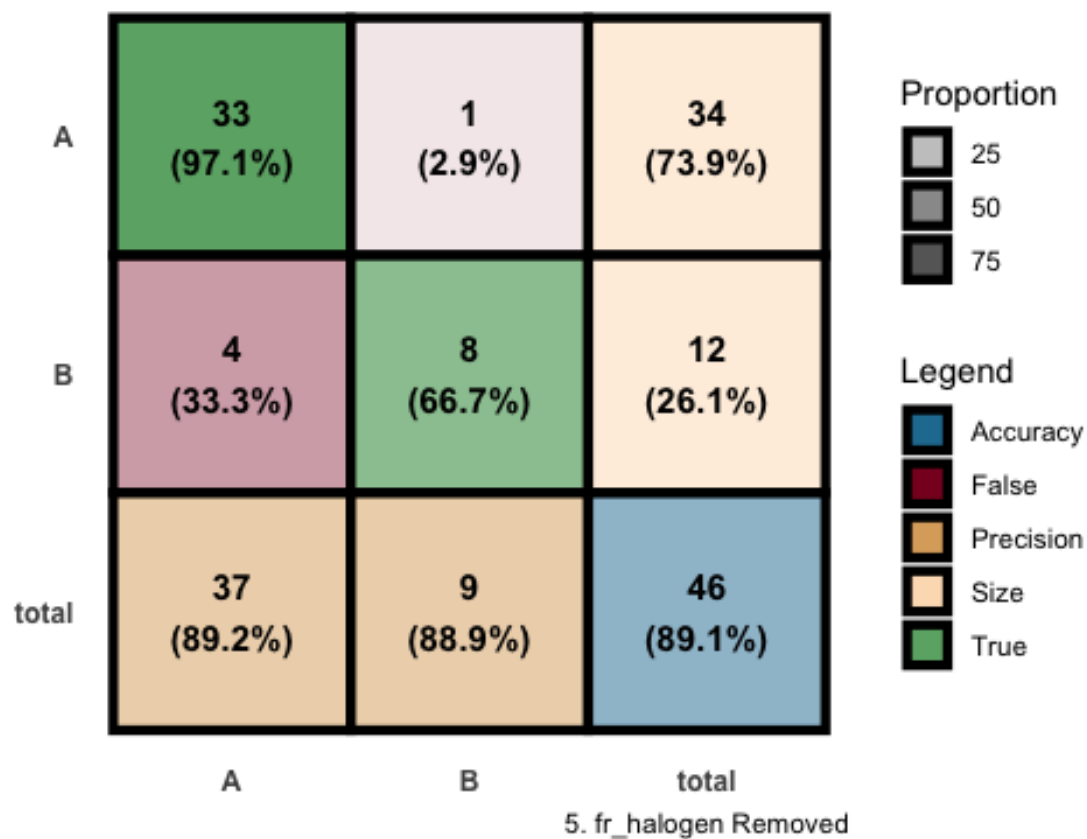
Accuracy	McFadden_R2
----------	-------------

89.13%	0.512
--------	-------

Model Coefficients

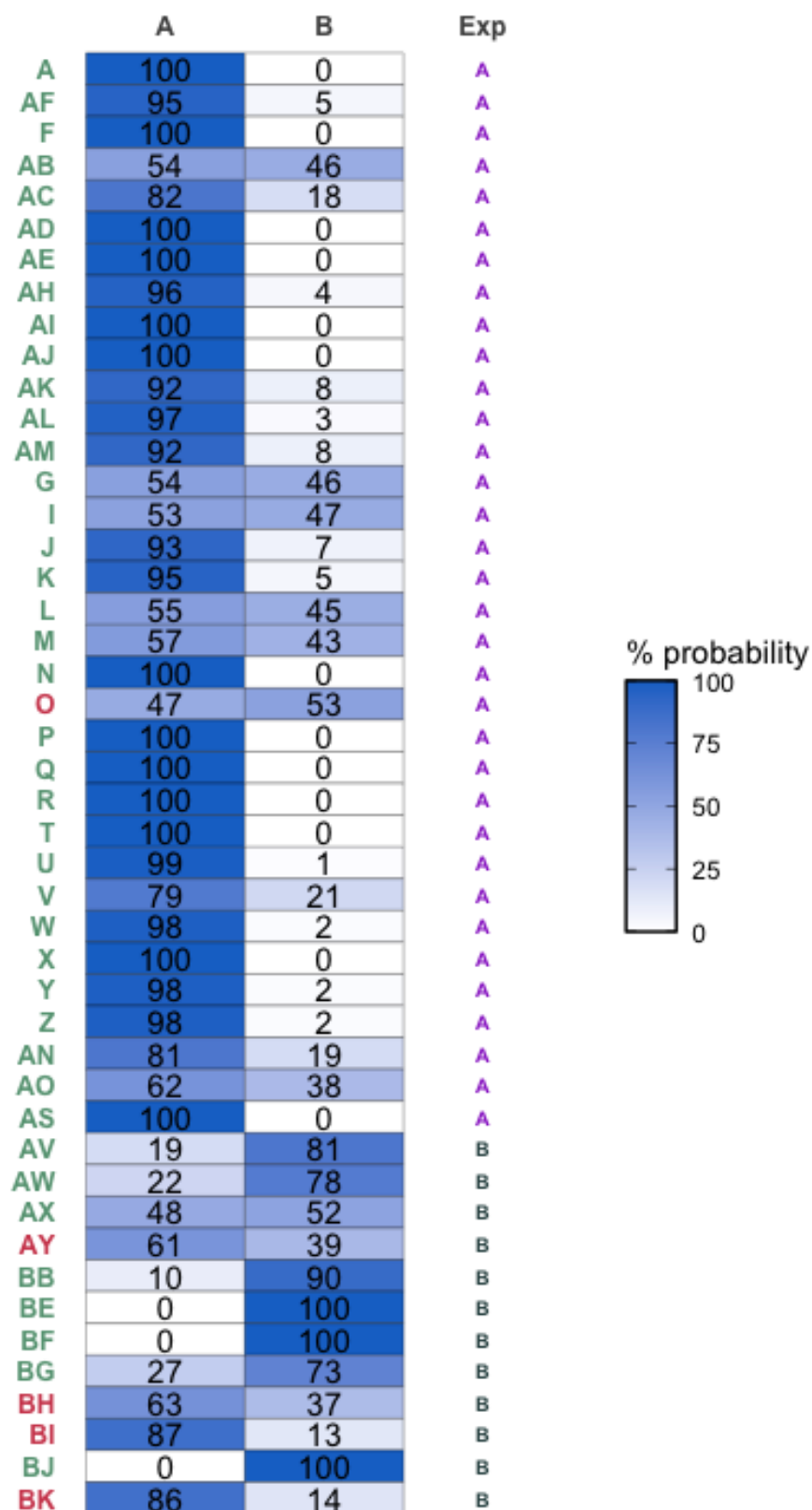
	x
(Intercept)	43.957001
Avg_NPA_SM	49.387661
fr_Ar_N1	-37.097092
fr_Ar_N2	-69.144375
fr_Ar_N3	-69.553894
fr_aryl_methyl1	-3.779621
fr_aryl_methyl2	-51.951799
fr_benzene1	-1.330087
fr_benzene2	-68.569342

Training Set Confusion Matrix



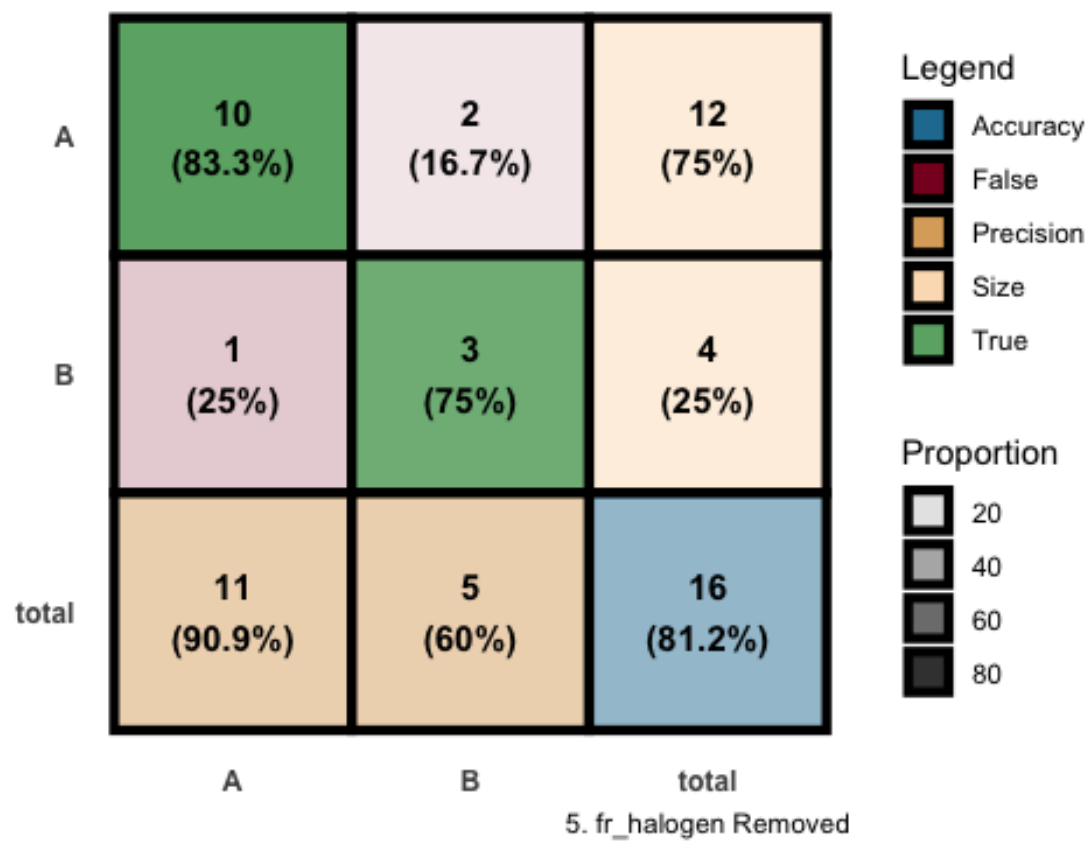
Training Set

Probability Heatmap



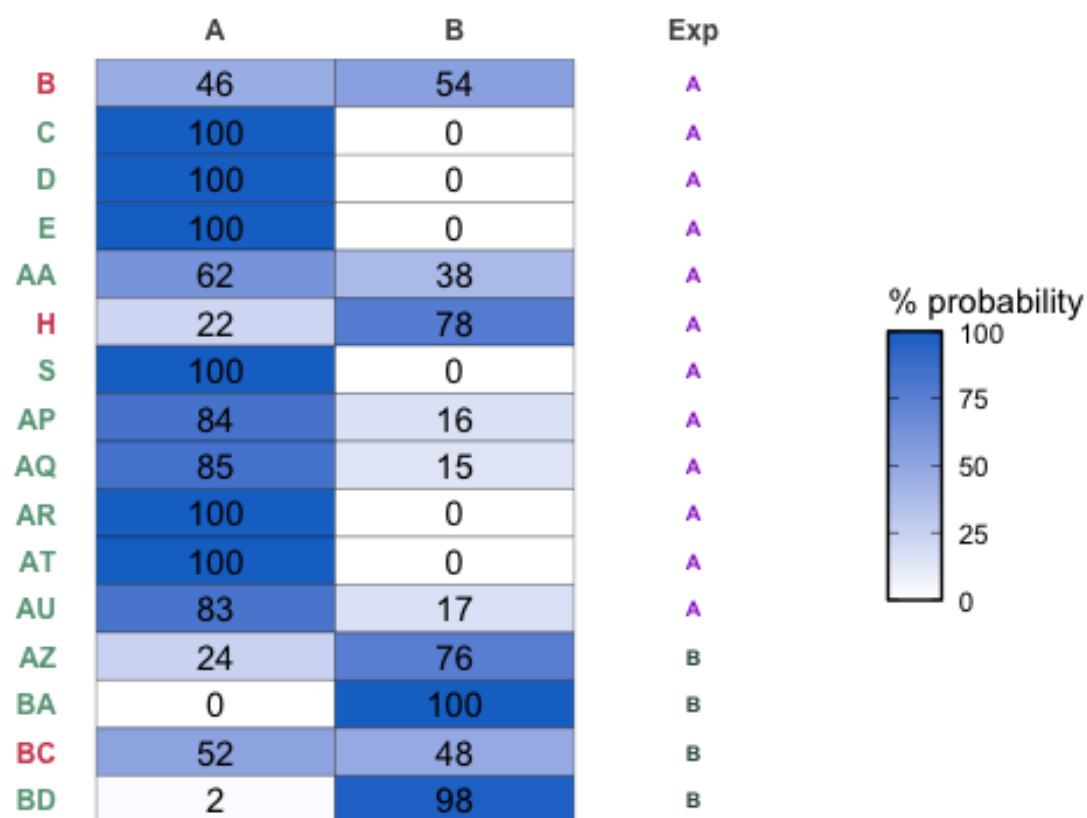
Test Set

Confusion Matrix



Test Set

Probability Heatmap



5. fr_halogen Removed

Section Results Summary

	Training Accuracy	Test Accuracy
5. 5th Place with 5 features	91.30	87.5
1. Avg_NPA_SM Removed	84.78	87.5
2. fr_Ar_N Removed	84.78	75.0
3. fr_aryl_methyl Removed	84.78	87.5
4. fr_benzene Removed	80.43	87.5
5. fr_halogen Removed	89.13	81.2