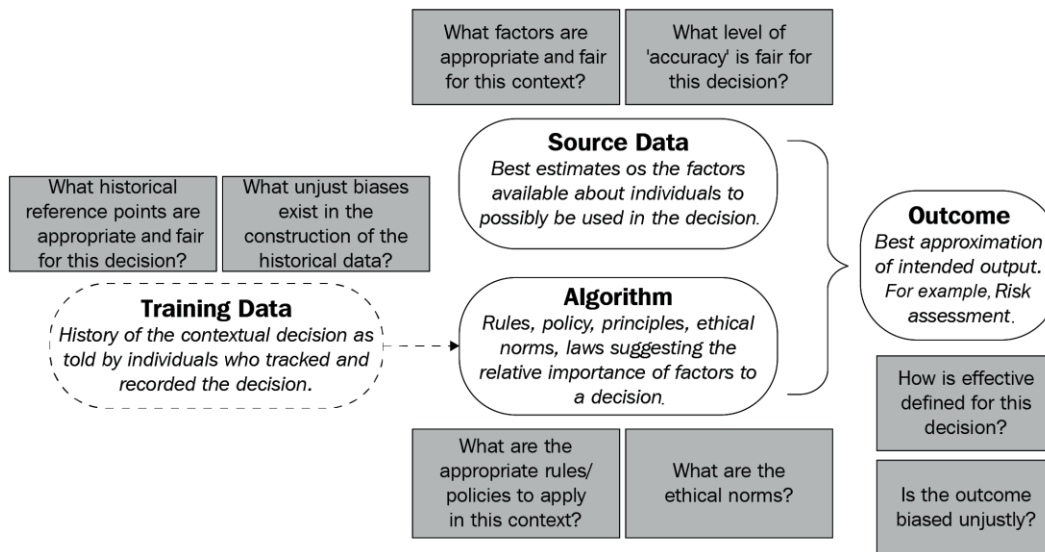
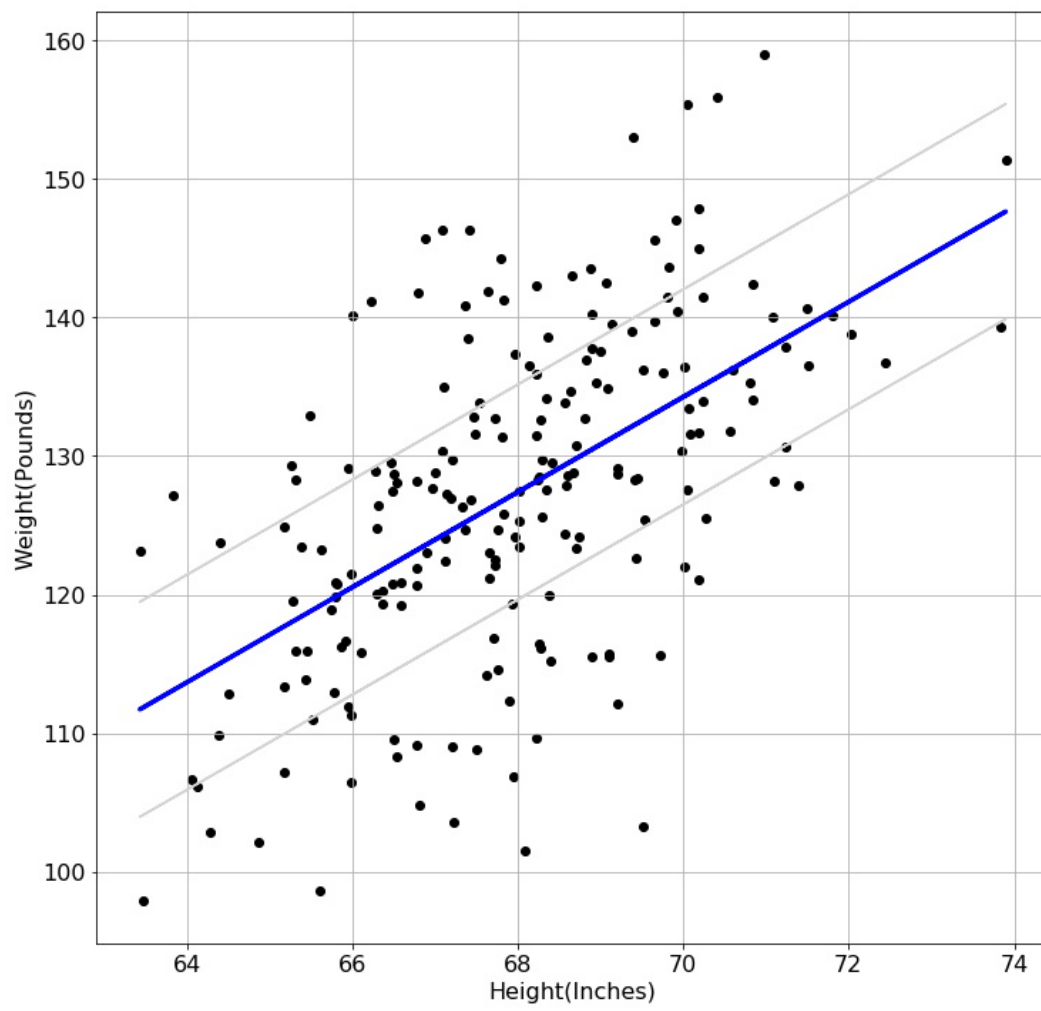
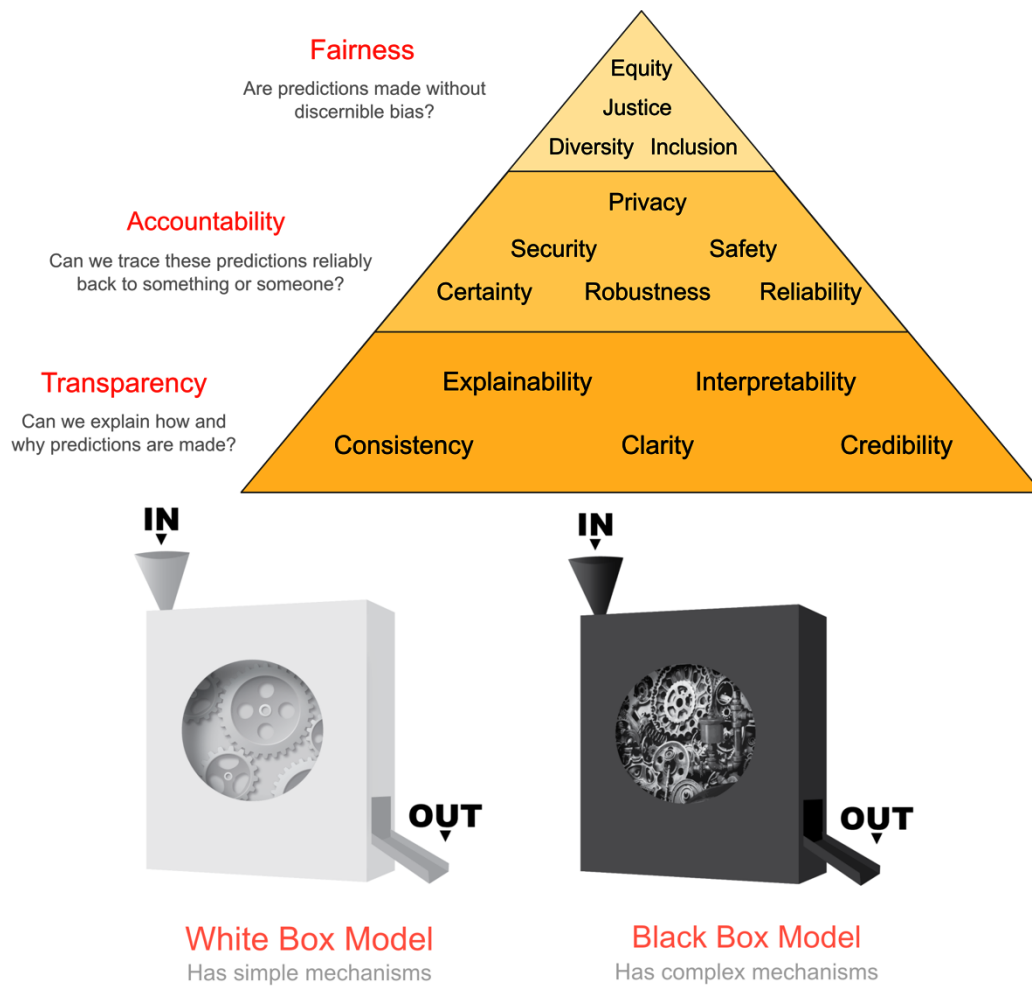
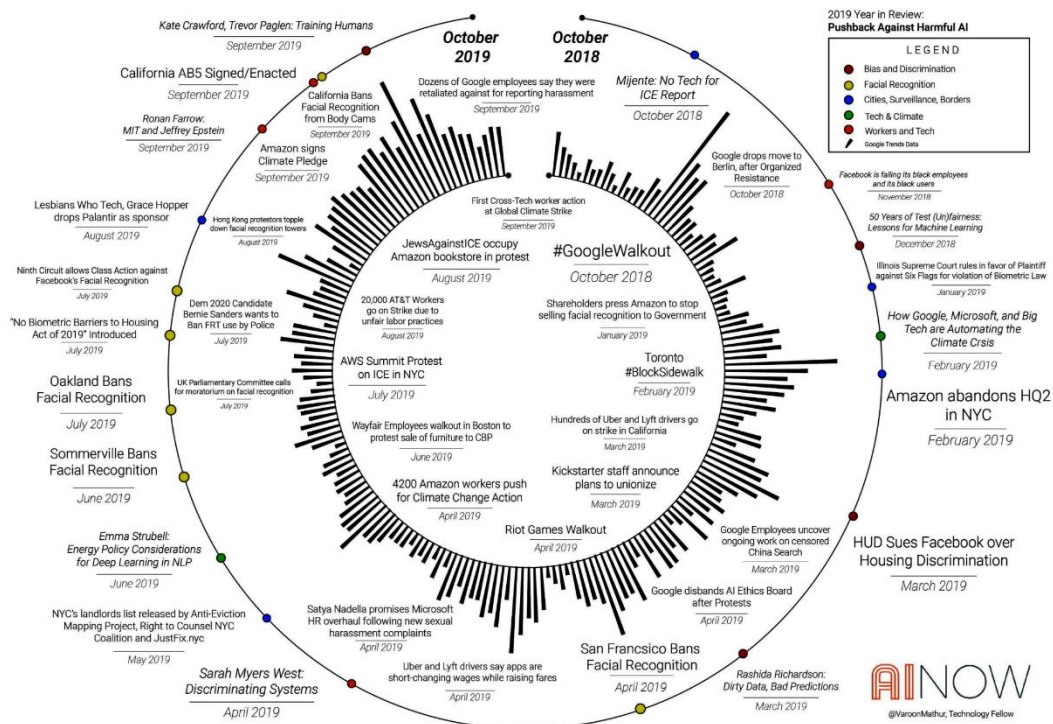


Chapter 1: Interpretation, Interpretability and Explainability; and why does it all matter?

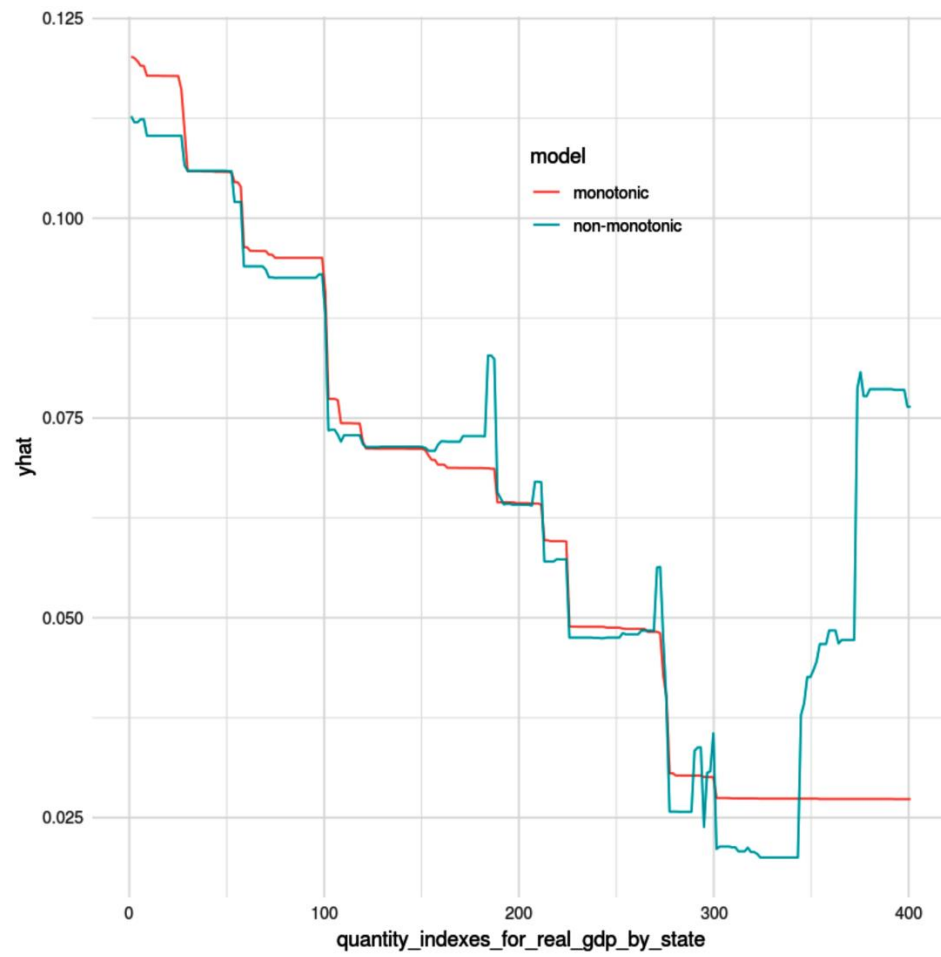




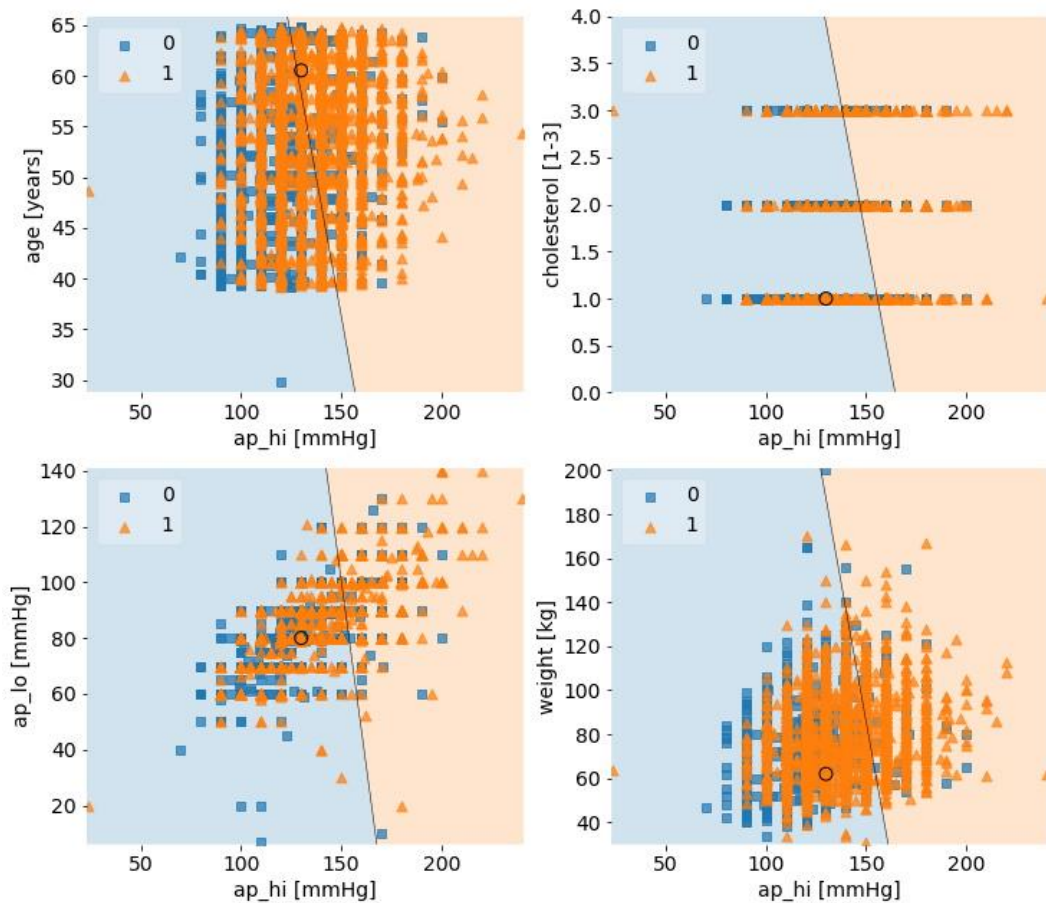


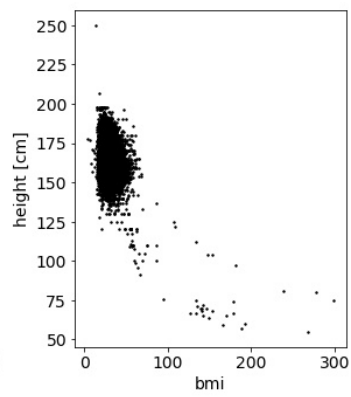
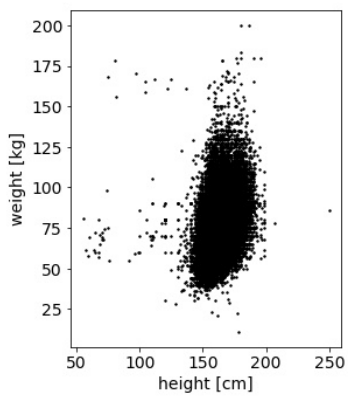
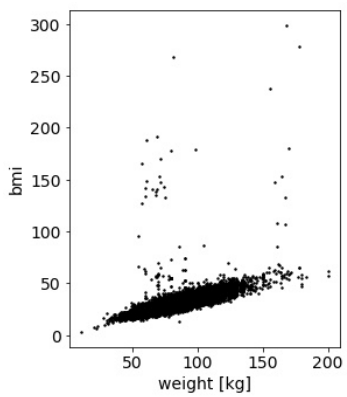
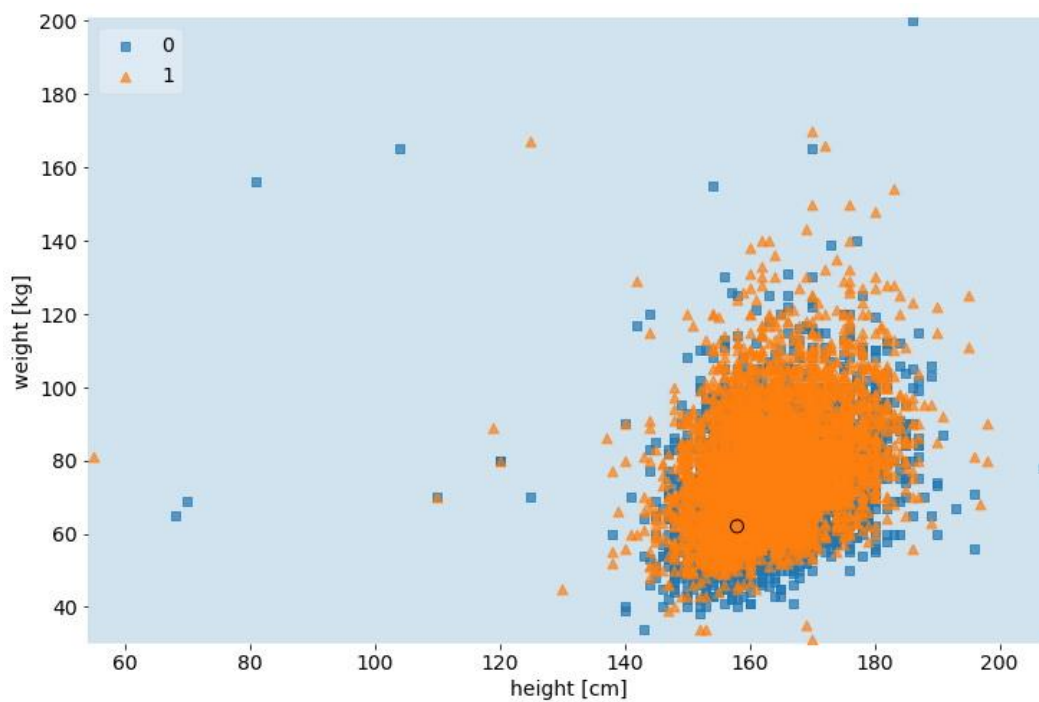


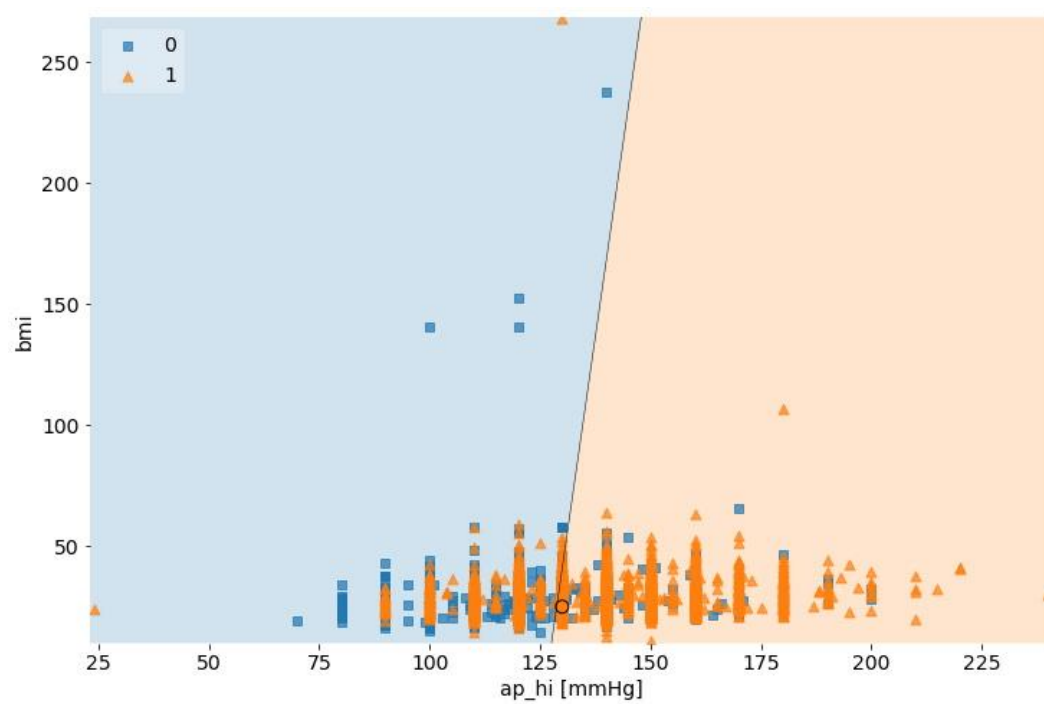
Chapter 2: Key Concepts of Interpretability



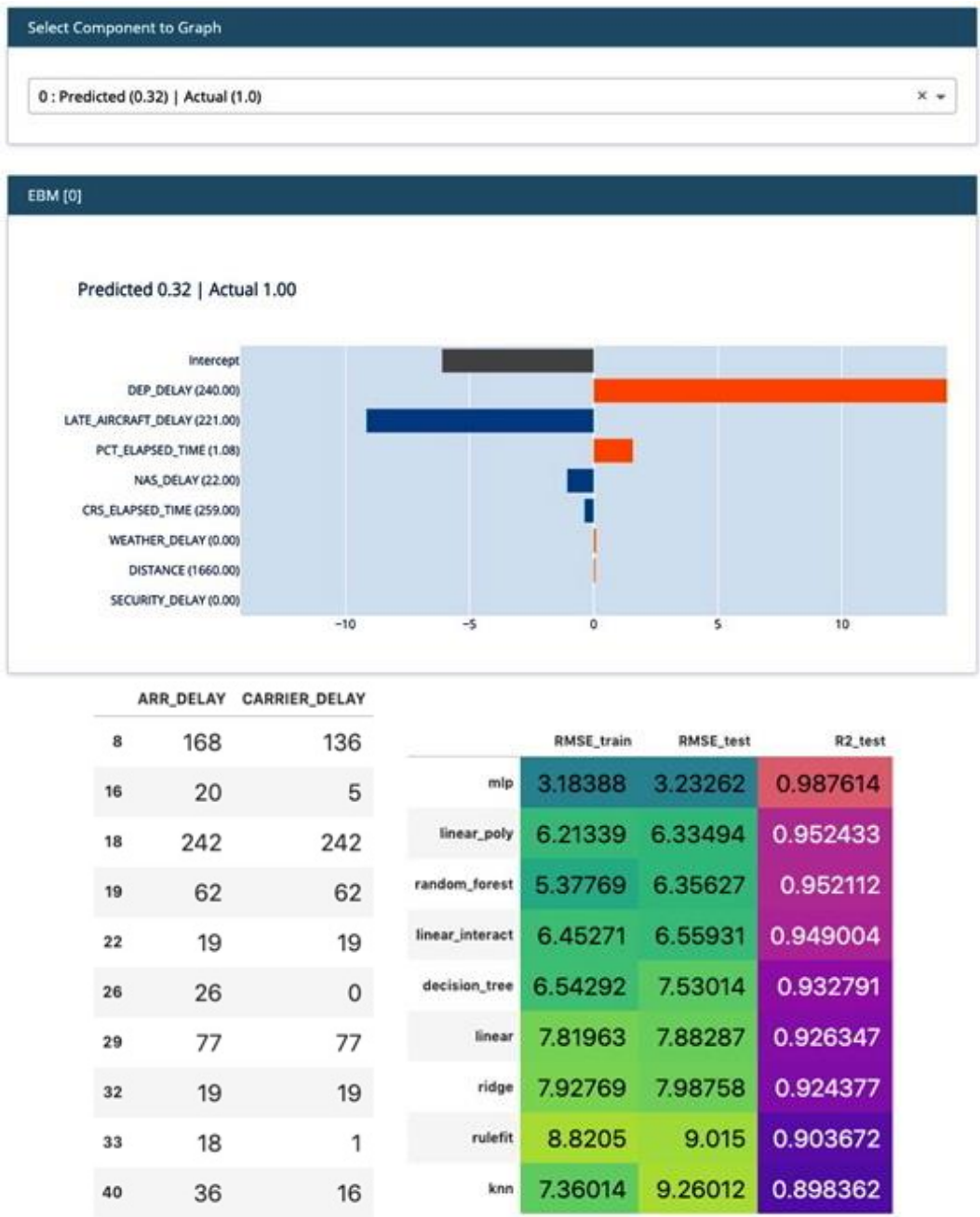
	count	mean	std	min	25%	50%	75%	max
age	70000.0	53.304309	6.755152	29.564122	48.36272	53.945351	58.391742	64.924433
gender	70000.0	1.349571	0.476838	1.000000	1.00000	1.000000	2.000000	2.000000
height	70000.0	164.359229	8.210126	55.000000	159.00000	165.000000	170.000000	250.000000
weight	70000.0	74.205690	14.395757	10.000000	65.00000	72.000000	82.000000	200.000000
ap_hi	70000.0	128.817286	154.011419	-150.000000	120.00000	120.000000	140.000000	16020.000000
ap_lo	70000.0	96.630414	188.472530	-70.000000	80.00000	80.000000	90.000000	11000.000000
cholesterol	70000.0	1.366871	0.680250	1.000000	1.00000	1.000000	2.000000	3.000000
gluc	70000.0	1.226457	0.572270	1.000000	1.00000	1.000000	1.000000	3.000000
smoke	70000.0	0.088129	0.283484	0.000000	0.00000	0.000000	0.000000	1.000000
alco	70000.0	0.053771	0.225568	0.000000	0.00000	0.000000	0.000000	1.000000
active	70000.0	0.803729	0.397179	0.000000	1.00000	1.000000	1.000000	1.000000
cardio	70000.0	0.499700	0.500003	0.000000	0.00000	0.000000	1.000000	1.000000



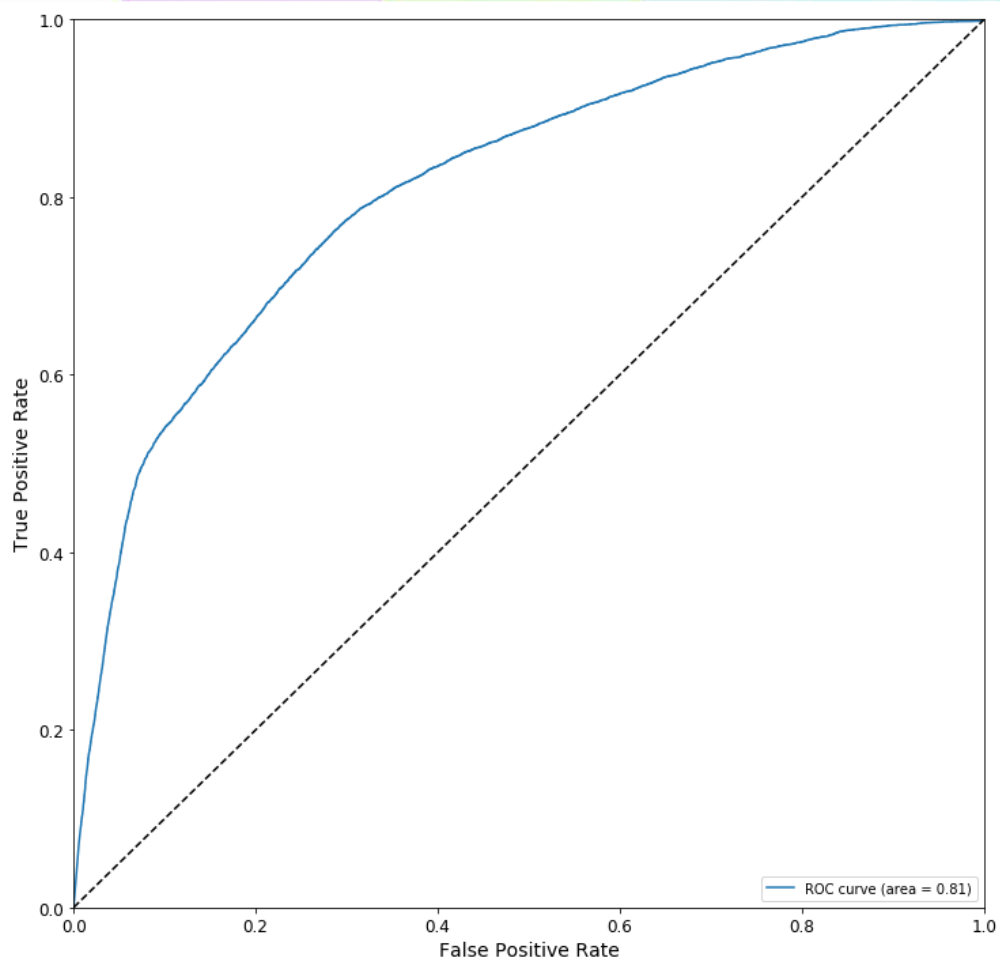


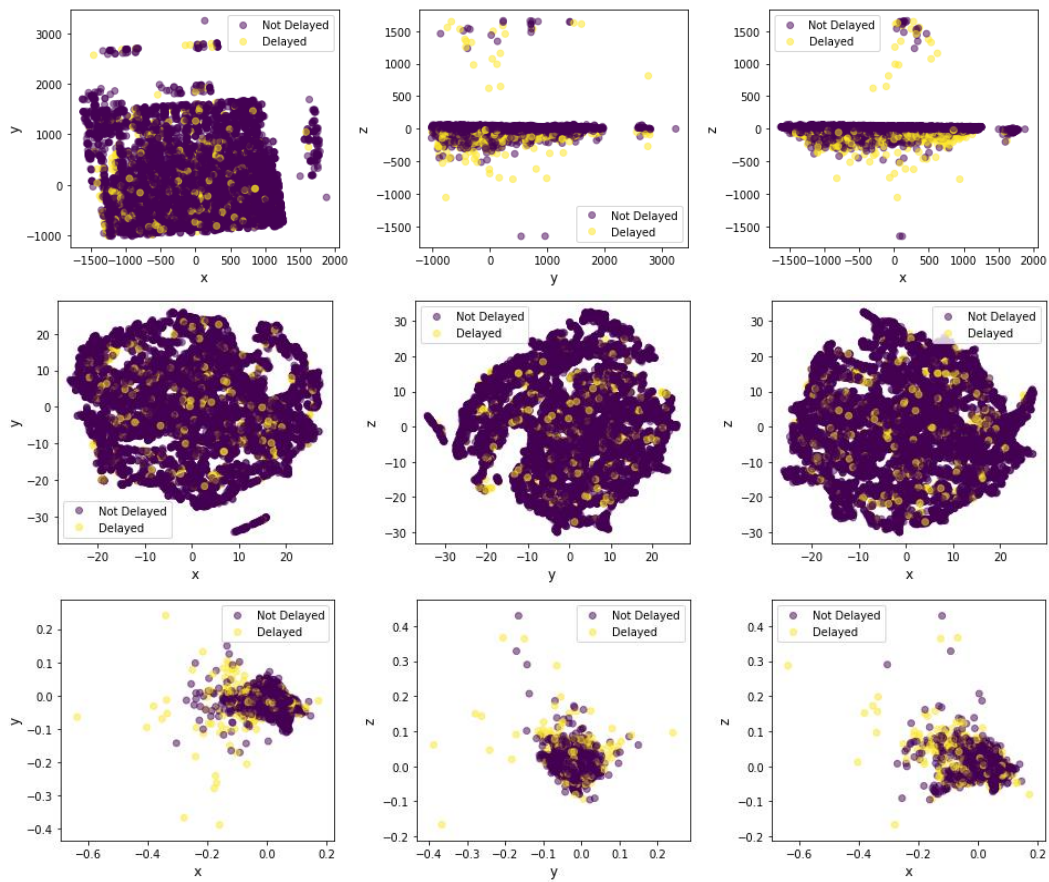


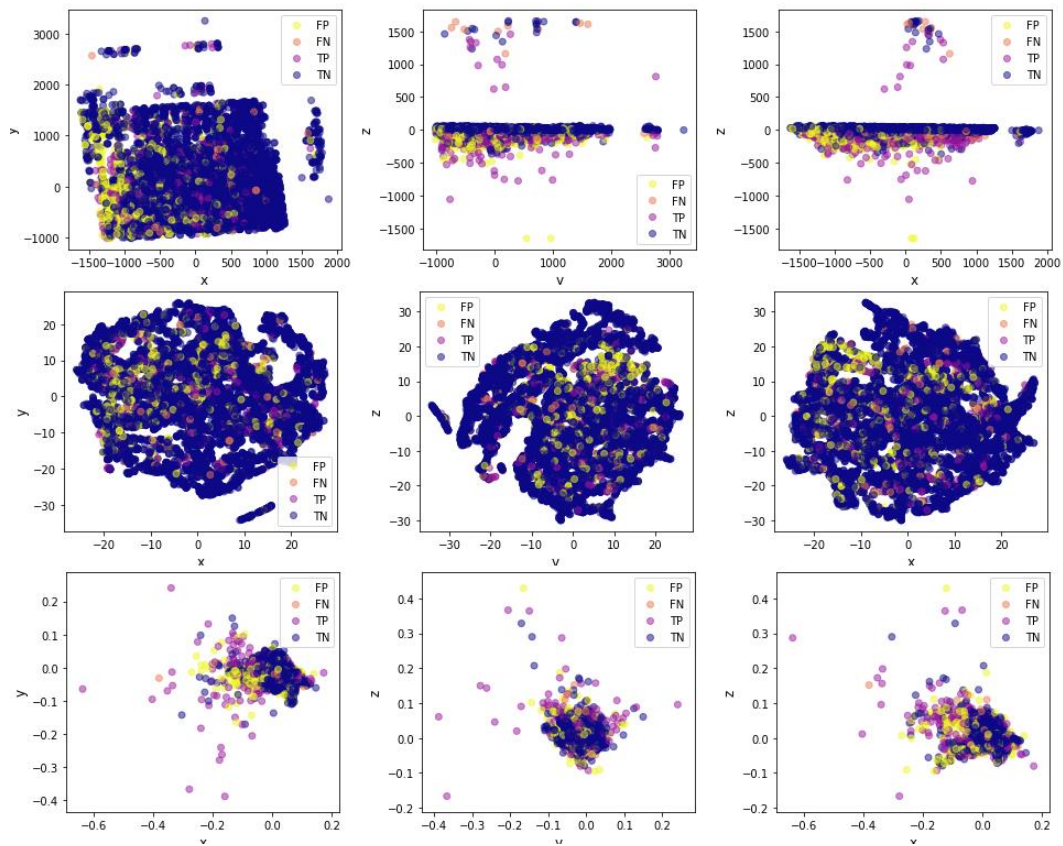
Chapter 3: Interpretation Challenges



	Accuracy_train	Accuracy_test	Recall_train	Recall_test	ROC_AUC_test	F1_test	MCC_test
mip	0.998364	0.99851	0.984826	0.986444	0.999909	0.987819	0.987027
gradient_boosting	0.991725	0.991662	0.89293	0.893851	0.998885	0.929223	0.925619
decision_tree	0.983297	0.982895	0.856969	0.852215	0.994932	0.859182	0.85011
random_forest	0.938783	0.937879	0.997546	0.990559	0.992844	0.661333	0.677145
logistic	0.9786	0.978381	0.743923	0.742677	0.971935	0.807953	0.800067
knn	0.97289	0.965123	0.680667	0.607722	0.948387	0.680906	0.668176
naive_bayes	0.925115	0.925561	0.279126	0.274268	0.811872	0.310922	0.275073
ridge	0.890447	0.891255	0.777002	0.77802	0	0.466998	0.463706







	feature	coef
0	CRS_DEP_TIME	0.00454956
1	DEP_TIME	-0.00525032
2	DEP_DELAY	0.894124
3	TAXI_OUT	0.125274
4	WHEELS_OFF	-0.0006468
5	CRS_ARR_TIME	-0.000369914
6	CRS_ELAPSED_TIME	-0.0126273
7	DISTANCE	0.000676793
8	WEATHER_DELAY	-0.906354
9	NAS_DELAY	-0.674053
10	SECURITY_DELAY	-0.917398
11	LATE_AIRCRAFT_DELAY	-0.929841
12	DEP_AFPH	-0.0152963
13	ARR_AFPH	0.000548174
14	DEP_MONTH	-0.039835
15	DEP_DOW	-0.0182132
16	DEP_RFPH	-0.469474
17	ARR_RFPH	0.373844
18	ORIGIN_HUB	-1.02909
19	DEST_HUB	-0.394899
20	PCT_ELAPSED_TIME	45.0116

OLS Regression Results

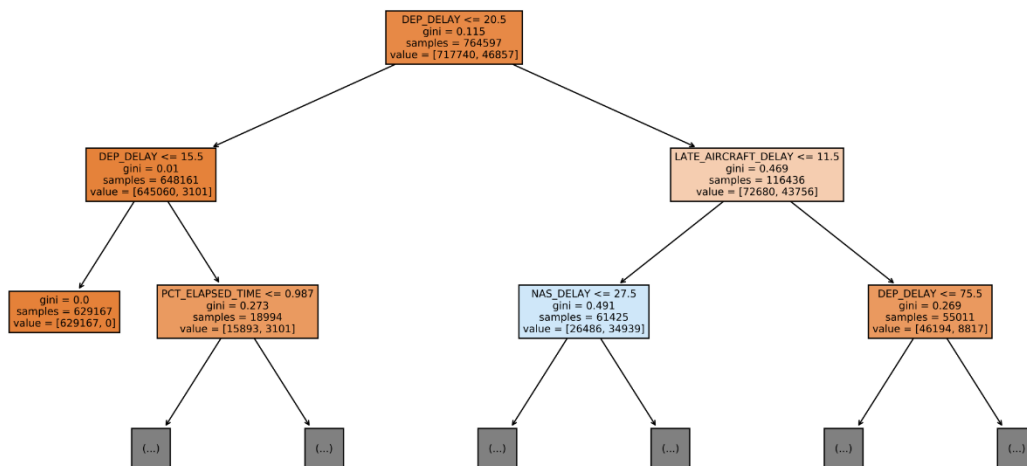
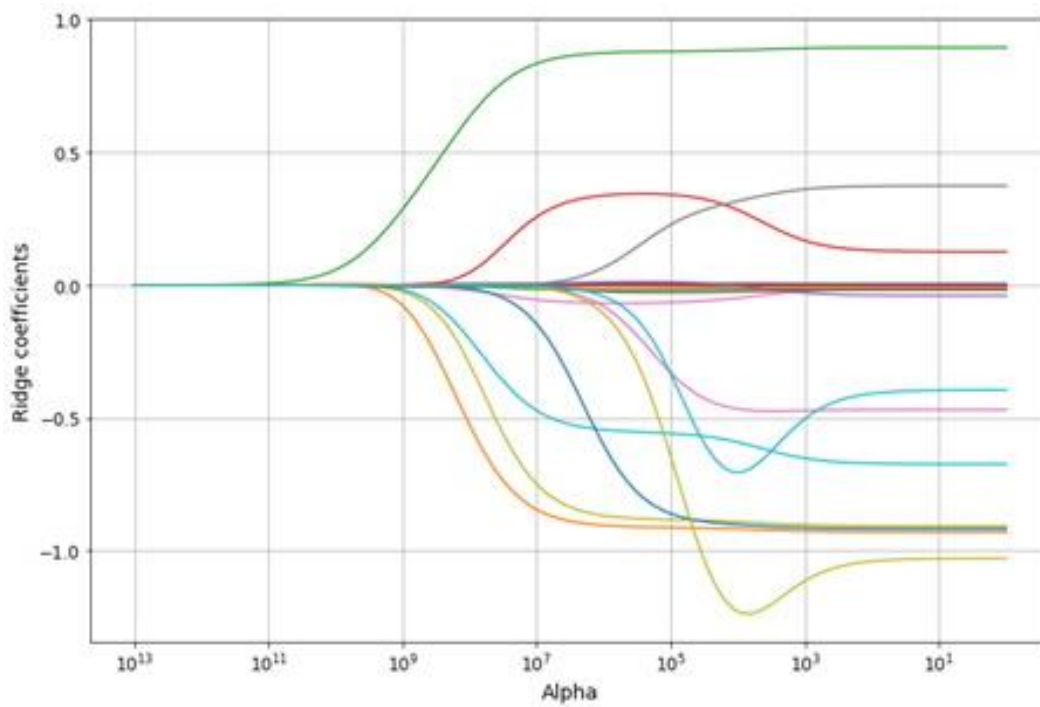
Dep. Variable:	CARRIER_DELAY	R-squared:	0.921
Model:	OLS	Adj. R-squared:	0.921
Method:	Least Squares	F-statistic:	4.251e+05
Date:	Wed, 02 Sep 2020	Prob (F-statistic):	0.00
Time:	13:32:20	Log-Likelihood:	-2.6574e+06
No. Observations:	764597	AIC:	5.315e+06
Df Residuals:	764575	BIC:	5.315e+06
Df Model:	21		
Covariance Type:	nonrobust		

	coef	std err	t	P> t	[0.025	0.975]
const	-37.8618	0.125	-301.763	0.000	-38.108	-37.616
CRS_DEP_TIME	0.0045	7.24e-05	62.872	0.000	0.004	0.005
DEP_TIME	-0.0053	9.19e-05	-57.116	0.000	-0.005	-0.005
DEP_DELAY	0.8941	0.000	2951.056	0.000	0.894	0.895
DEP_AFPH	-0.0153	0.000	-47.725	0.000	-0.016	-0.015
DEP_RFPH	-0.4696	0.017	-27.353	0.000	-0.503	-0.436
TAXI_OUT	0.1253	0.001	104.120	0.000	0.123	0.128
WHEELS_OFF	-0.0006	6.7e-05	-9.646	0.000	-0.001	-0.001
CRS_ELAPSED_TIME	-0.0126	0.001	-19.132	0.000	-0.014	-0.011
PCT_ELAPSED_TIME	45.0113	0.117	384.073	0.000	44.782	45.241
DISTANCE	0.0007	8.02e-05	8.429	0.000	0.001	0.001
CRS_ARR_TIME	-0.0004	2.18e-05	-16.939	0.000	-0.000	-0.000
ARR_AFPH	0.0005	0.000	1.651	0.099	-0.000	0.001
ARR_RFPH	0.3739	0.013	28.386	0.000	0.348	0.400
WEATHER_DELAY	-0.9064	0.001	-995.366	0.000	-0.908	-0.905
NAS_DELAY	-0.6741	0.001	-829.129	0.000	-0.676	-0.672
SECURITY_DELAY	-0.9174	0.005	-167.857	0.000	-0.928	-0.907
LATE_AIRCRAFT_DELAY	-0.9298	0.001	-1827.018	0.000	-0.931	-0.929
DEP_MONTH	-0.0397	0.003	-15.019	0.000	-0.045	-0.034
DEP_DOW	-0.0180	0.004	-4.005	0.000	-0.027	-0.009
ORIGIN_HUB	-1.0291	0.027	-38.589	0.000	-1.081	-0.977
DEST_HUB	-0.3949	0.026	-15.041	0.000	-0.446	-0.343

Omnibus:	211121.387	Durbin-Watson:	2.001
Prob(Omnibus):	0.000	Jarque-Bera (JB):	24359701.834
Skew:	0.098	Prob(JB):	0.00
Kurtosis:	30.651	Cond. No.	5.69e+04

	feature	Coef.	Std.Err.	t	P> t	[0.025	0.975]	t_abs
2	DEP_DELAY	0.894124	0.000302981	2951.09	0	0.89353	0.894717	2951.09
11	LATE_AIRCRAFT_DELAY	-0.929841	0.000508937	-1827.03	0	-0.930839	-0.928844	1827.03
8	WEATHER_DELAY	-0.906354	0.000910567	-995.373	0	-0.908138	-0.904569	995.373
9	NAS_DELAY	-0.674053	0.000812964	-829.13	0	-0.675646	-0.67246	829.13
20	PCT_ELAPSED_TIME	45.0116	0.117195	384.076	0	44.7819	45.2413	384.076
10	SECURITY_DELAY	-0.917398	0.00546544	-167.855	0	-0.928111	-0.906686	167.855
3	TAXI_OUT	0.125274	0.00120321	104.117	0	0.122916	0.127633	104.117
0	CRS_DEP_TIME	0.00454956	7.23674e-05	62.8675	0	0.00440772	0.00469139	62.8675
1	DEP_TIME	-0.00525032	9.19302e-05	-57.1121	0	-0.0054305	-0.00507014	57.1121
12	DEP_AFPH	-0.0152963	0.000320506	-47.7256	0	-0.0159245	-0.0146681	47.7256
18	ORIGIN_HUB	-1.02909	0.0266686	-38.5879	0	-1.08136	-0.976818	38.5879
17	ARR_RFPH	0.373844	0.0131708	28.3844	3.89612e-177	0.34803	0.399658	28.3844
16	DEP_RFPH	-0.469474	0.0171688	-27.3446	1.50325e-164	-0.503124	-0.435824	27.3446
6	CRS_ELAPSED_TIME	-0.0126273	0.000659852	-19.1366	1.3093e-81	-0.0139206	-0.011334	19.1366
5	CRS_ARR_TIME	-0.000369914	2.18388e-05	-16.9384	2.4083e-64	-0.000412717	-0.00032711	16.9384
14	DEP_MONTH	-0.039835	0.00264082	-15.0844	2.08773e-51	-0.045011	-0.0346591	15.0844
19	DEST_HUB	-0.394899	0.0262564	-15.0401	4.07781e-51	-0.44636	-0.343437	15.0401

	feature	coef_linear	coef_ridge
0	CRS_DEP_TIME	0.00454956	0.00501961
1	DEP_TIME	-0.00525032	-0.00441738
2	DEP_DELAY	0.894124	0.894292
3	TAXI_OUT	0.125274	0.125165
4	WHEELS_OFF	-0.0006468	0.000232365
5	CRS_ARR_TIME	-0.000369914	-0.00189765
6	CRS_ELAPSED_TIME	-0.0126273	-0.0125826
7	DISTANCE	0.000676793	0.0021406
8	WEATHER_DELAY	-0.906354	-0.906168
9	NAS_DELAY	-0.674053	-0.67396
10	SECURITY_DELAY	-0.917398	-0.917398
11	LATE_AIRCRAFT_DELAY	-0.929841	-0.929537
12	DEP_AFPH	-0.0152963	-0.0154111
13	ARR_AFPH	0.000548174	0.000532269
14	DEP_MONTH	-0.039835	-0.0398301
15	DEP_DOW	-0.0182132	-0.018213
16	DEP_RFPH	-0.469474	-0.469473
17	ARR_RFPH	0.373844	0.373847
18	ORIGIN_HUB	-1.02909	-1.02909
19	DEST_HUB	-0.394899	-0.394898
20	PCT_ELAPSED_TIME	45.0116	45.0116



	feature	importance
2	DEP_DELAY	0.527482
11	LATE_AIRCRAFT_DELAY	0.199153
20	PCT_ELAPSED_TIME	0.105381
8	WEATHER_DELAY	0.101649
9	NAS_DELAY	0.0627577
10	SECURITY_DELAY	0.00199756
7	DISTANCE	0.000993382
6	CRS_ELAPSED_TIME	0.000280958
3	TAXI_OUT	0.000238682
4	WHEELS_OFF	3.46469e-05
12	DEP_AFPH	3.10537e-05
5	CRS_ARR_TIME	0
1	DEP_TIME	0
13	ARR_AFPH	0
14	DEP_MONTH	0
15	DEP_DOW	0
16	DEP_RFPH	0
17	ARR_RFPH	0
18	ORIGIN_HUB	0
19	DEST_HUB	0
0	CRS_DEP_TIME	0

	rule	type	coef	support	importance
129	LATE_AIRCRAFT_DELAY <= 222.5 & WEATHER_DELAY <= 166.0 & DEP_DELAY > 344.0	rule	207.246	0.0016835	8.49625
80	DEP_DELAY > 477.5 & LATE_AIRCRAFT_DELAY <= 333.5	rule	170.948	0.00112233	5.72377
53	WEATHER_DELAY > 255.0 & DEP_DELAY > 490.5	rule	-333.579	0.000187056	4.56188
11	LATE_AIRCRAFT_DELAY	linear	-0.383065	1	4.48841
2	DEP_DELAY	linear	0.162592	1	4.25384
46	LATE_AIRCRAFT_DELAY <= 198.0 & DEP_DELAY <= 788.0 & DEP_DELAY > 341.5	rule	-95.8115	0.00149645	3.70359
57	DEP_DELAY > 1206.0	rule	254.29	0.000187056	3.47755
84	DEP_DELAY > 300.0 & DEP_DELAY > 576.5 & LATE_AIRCRAFT_DELAY <= 158.5	rule	121.199	0.000748223	3.31401
64	DEP_DELAY > 880.5	rule	102.969	0.000748223	2.81552
147	DEP_DELAY <= 37.5 & DEP_DELAY <= 370.5	rule	-9.13357	0.898429	2.7591
52	LATE_AIRCRAFT_DELAY <= 19.5 & DEP_DELAY <= 849.0 & DEP_DELAY > 66.5 & NAS_DELAY > 43.5	rule	-41.4699	0.00430228	2.71422
63	WEATHER_DELAY <= 61.0 & DEP_DELAY <= 849.0 & LATE_AIRCRAFT_DELAY <= 19.5 & DEP_DELAY > 270.0 & NAS_DELAY <= 43.5 & DEP_DELAY > 66.5	rule	99.0067	0.000748223	2.70718
153	WEATHER_DELAY <= 61.0 & DEP_DELAY <= 849.0 & LATE_AIRCRAFT_DELAY <= 19.5 & NAS_DELAY <= 43.5 & DEP_DELAY > 109.0 & DEP_DELAY > 66.5 & DEP_DELAY <= 270.0	rule	29.733	0.00598578	2.29348
169	WEATHER_DELAY > 61.0 & DEP_DELAY <= 849.0 & LATE_AIRCRAFT_DELAY <= 19.5 & NAS_DELAY <= 43.5 & DEP_DELAY > 66.5	rule	-45.9107	0.00224467	2.17271
162	DEP_DELAY > 117.0 & WEATHER_DELAY <= 10.0 & DEP_DELAY <= 225.0 & LATE_AIRCRAFT_DELAY <= 56.5 & DEP_DELAY <= 459.0 & DEP_DELAY > 68.5 & NAS_DELAY <= 66.0	rule	28.4973	0.00467639	1.9442
38	LATE_AIRCRAFT_DELAY <= 32.5 & NAS_DELAY <= 40.5 & DEP_DELAY <= 491.5 & DEP_DELAY > 57.5 & DEP_DELAY <= 245.5 & WEATHER_DELAY <= 20.0	rule	12.1724	0.0226337	1.81044
51	DEP_DELAY <= 20.5 & DEP_DELAY <= 68.5 & DEP_DELAY <= 459.0	rule	-4.56733	0.846053	1.64834

White Box?	Model Class	Properties that Increase Interpretability					Task		Performance Rank	
		Expl.	Linear	Monotone	Non-Interactive	Regul.	Regr.	Classif.	Regr.	Classif.
✓	Linear Regression	✓	✓	✓	✓	✓	✓	✗	6	
✓	Regularized Regression	✓	✓	✓	✓	✓	✓	✓	7	8
✓	Logistic Regression	✓	✓	✓	✓	✓	✗	✓		5
✓	Gaussian Naive Bayes	✓	✓	✓	✓	✓	✗	✓		7
✓	Polynomial Regression	✓	✓	✓	✓	✓	✓	✓	2	
✓	RuleFit	✓	✓	✓	✓	✓	✓	✓	8	
✓	Decision Tree	✓	✓	✓	✓	✓	✓	✓	5	3
✓	k-Nearest Neighbors	✓	✓	✓	✓	✓	✓	✓	9	6
✗	Random Forest	✓	✓	✓	✓	✓	✓	✓	3	4
✗	Gradient Boosted Trees	✓	✓	✓	✓	✓	✓	✓		2
✗	Multi-layer Perceptron	✓	✓	✓	✓	✓	✓	✓	1	

	White Box	Glass Box	Black Box
Interpretability	High	Mid-High	Low
Predictive Performance	Mid	High	High
Execution Speed Performance	High	Low	Mid

Select Component to Graph

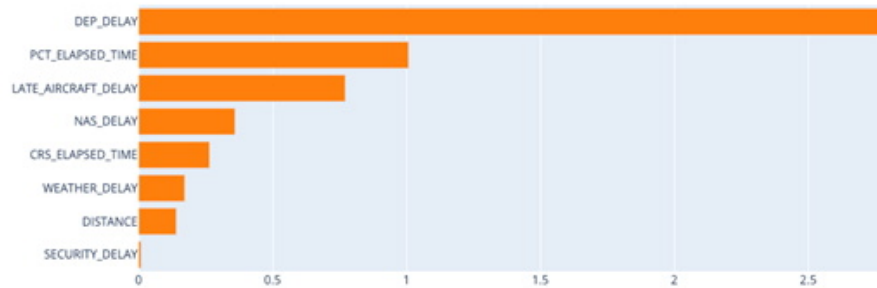
Summary

X

ExplainableBoostingClassifier_0 (Overall)

🔍 + 📄 📊 📉 📈 📉 📈 📉 📈 📉 📈 📉

Overall Importance:
Mean Absolute Score



Select Component to Graph

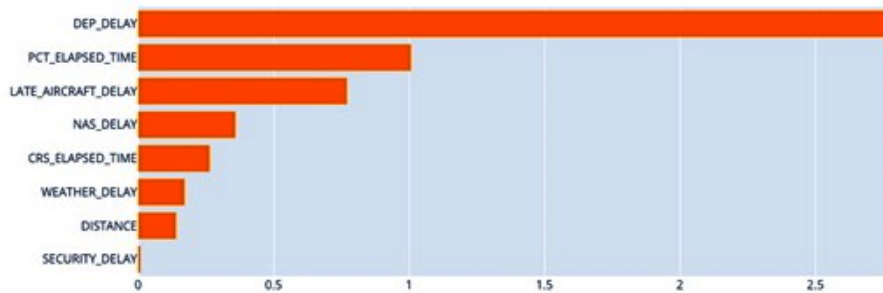
Summary

X

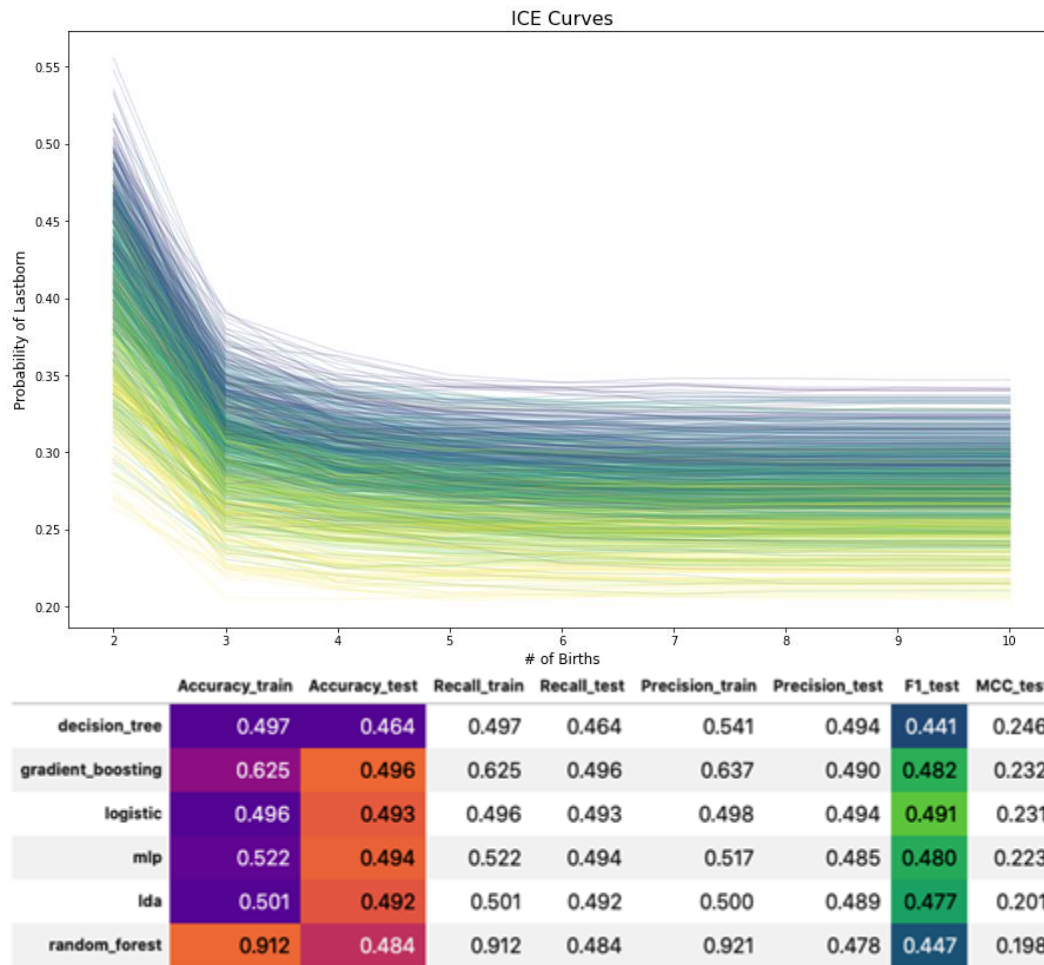
ExplainableBoostingClassifier_0 (Overall)

🔍 + 📄 📊 📉 📈 📉 📈 📉 📈 📉 📈 📉

Overall Importance:
Mean Absolute Score



Chapter 4: Fundamentals of Feature Importance and Impact



	name	dt_imp	dt_rank	gb_imp	gb_rank	rf_imp	rf_rank	avg_rank
28	birthn	0.851533	1	0.371305	1	0.198748	1	1
82	testelapse	0.0137081	3	0.0335579	2	0.0275725	2	2.33333
26	age	0.00667898	7	0.030532	3	0.0248301	3	4.33333
0	Q1	0.0253401	2	0.0236222	6	0.0159306	6	4.66667
81	introelapse	0.00505607	9	0.0297233	4	0.0224896	5	6
12	Q13	0.0080825	4	0.014516	7	0.0113429	8	6.33333
	:	:	:	:	:	:	:	:
90	country_GB	0	91	0.000755431	91	0.00194744	90	90.6667
92	country_NZ	0	93	0.00103713	90	0.000736748	91	91.3333
84	gender_undefined	0	87	0.000316311	94	0.000302447	94	91.6667
91	country_IE	0	92	0.000596172	92	0.000499432	93	92.3333
	name	first_coef_norm	middle_coef_norm	last_coef_norm	coef_weighted_avg			
28	birthn	-0.412945	1.3538	-0.0132044	0.499431			
26	age	0.0552764	-0.0265002	-0.149019	0.0804694			
0	Q1	0.110523	0.0224566	-0.00631052	0.0540604			
12	Q13	0.0793163	-0.0382582	-0.000743793	0.0427518			
15	Q16	0.0604051	-0.0542339	-0.000581668	0.0385124			
19	Q20	-0.0609848	0.0508594	0.0015853	0.0382996			
39	EST1	0.0498431	-0.0622704	0.00411372	0.0371717			
3	Q4	0.044028	-0.0576418	-0.000594055	0.0324218			
59	CSN1	0.0316447	-0.0699186	0.00127486	0.0303448			
	:	:	:	:	:			
90	country_GB	-2.218e-05	-0.00138172	4.47252e-05	0.000352116			
91	country_IE	-0.00014727	0.000314903	7.59173e-07	0.000136968			
92	country_NZ	7.70629e-05	-0.000394417	5.52904e-06	0.000127852			
87	gender_other	7.0736e-05	0.000394679	-8.67944e-06	0.000126324			
84	gender_undefined	5.65254e-05	-9.94834e-05	-3.80303e-07	4.75334e-05			

	name	first_coef_norm	middle_coef_norm	last_coef_norm	coef_weighted_avg
28	birthn	-0.315215	1.00305	-0.307128	0.475483
0	Q1	0.0899109	-0.0122606	-0.102456	0.0757905
12	Q13	0.0564803	-0.0337293	-0.0462968	0.0476102
51	AGR3	-0.0392475	-0.00558213	0.0523123	0.0357299
15	Q16	0.0395618	-0.0363935	-0.0235674	0.0333487
6	Q7	-0.00407858	0.0644172	-0.0396745	0.0305362
24	Q25	-0.0350918	0.0343628	0.01946	0.0295807
16	Q17	0.034915	-0.00978912	-0.036297	0.0294317
77	OPN9	-0.0326552	0.0447374	0.00925253	0.0275268
33	EXT5	-0.000968246	0.0064997	-0.00331607	0.0030811
85	gender_male	0.00272646	-0.00239968	-0.00169943	0.00229827
81	introelapse	-0.00168655	0.00484121	-0.00127851	0.00229491
36	EXT8	-0.000826513	-0.00324433	0.00327175	0.00223465
7	Q8	-0.000763955	0.0039999	-0.00183324	0.0018961

	name	first_coef_norm	middle_coef_norm	last_coef_norm	coef_weighted_avg
28	birthn	-0.315215	1.00305	-0.307128	0.475483
0	Q1	0.0899109	-0.0122606	-0.102456	0.0757905
12	Q13	0.0564803	-0.0337293	-0.0462968	0.0476102
51	AGR3	-0.0392475	-0.00558213	0.0523123	0.0357299
15	Q16	0.0395618	-0.0363935	-0.0235674	0.0333487
6	Q7	-0.00407858	0.0644172	-0.0396745	0.0305362
24	Q25	-0.0350918	0.0343628	0.01946	0.0295807
16	Q17	0.034915	-0.00978912	-0.036297	0.0294317
77	OPN9	-0.0326552	0.0447374	0.00925253	0.0275268

33	EXT5	-0.000968246	0.0064997	-0.00331607	0.0030811
85	gender_male	0.00272646	-0.00239968	-0.00169943	0.00229827
81	introelapse	-0.00168655	0.00484121	-0.00127851	0.00229491
36	EXT8	-0.000826513	-0.00324433	0.00327175	0.00223465
7	Q8	-0.000763955	0.0039999	-0.00183324	0.0018961

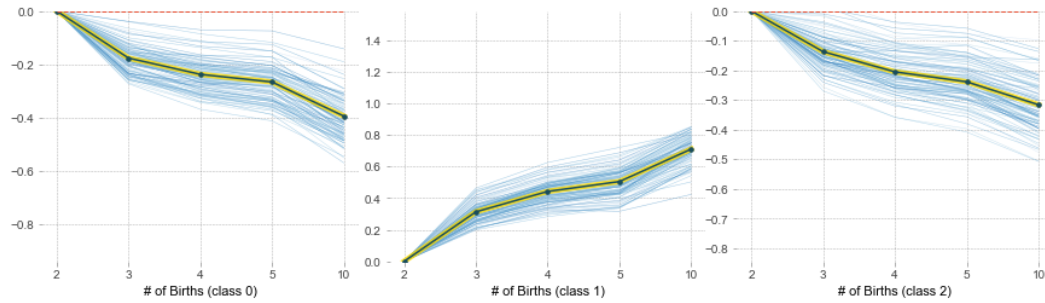
	name	dt_imp	gb_imp	rf_imp	log_imp	lda_imp	mlp_imp	avg_imp
28	birthn	0.1385	0.10735	0.07604	0.11818	0.08199	0.11172	0.10563
0	Q1	0.00832	0.00688	0.00428	0.00509	0.01103	0.0093	0.00749
26	age	0.00107	0.00327	0.00496	0.00713	-0.00122	0.00183	0.00284
12	Q13	0.00098	-0.00252	-6e-05	0.00428	0.00235	0.00499	0.00167
3	Q4	0	0.00274	0.00163	0.00178	0.0006	0.00214	0.00148
16	Q17	0.00119	-0.00201	0.00255	0.00122	0.00179	0.00273	0.00124
51	AGR3	0.00032	-7e-05	-0.00156	0.00109	0.00339	0.0039	0.00118
24	Q25	0	-6e-05	-0.00087	0.00106	0.00112	0.00465	0.00098
30	EXT2	0	0.00073	0.00161	0.00075	-0.00076	0.00348	0.00097
	:	:	:	:	:	:	:	:
69	OPN1	0.00015	-0.00035	-0.00175	-0.00062	-0.00018	0.00088	-0.00031
21	Q22	0	-0.00279	-0.00025	0.00019	-0.00207	0.00242	-0.00042
79	source	0	-0.00048	6e-05	0.00094	-0.00135	-0.0017	-0.00042
25	Q26	0.00126	-0.00144	-0.00216	0.00015	-0.00211	0.0007	-0.0006
22	Q23	0	-0.00169	-0.00012	-7e-05	-0.0017	-0.00028	-0.00064

Accuracy_test

decision_tree	0.325639
gradient_boosting	0.388483
random_forest	0.40958
logistic	0.383053
lda	0.409265
mlp	0.37977

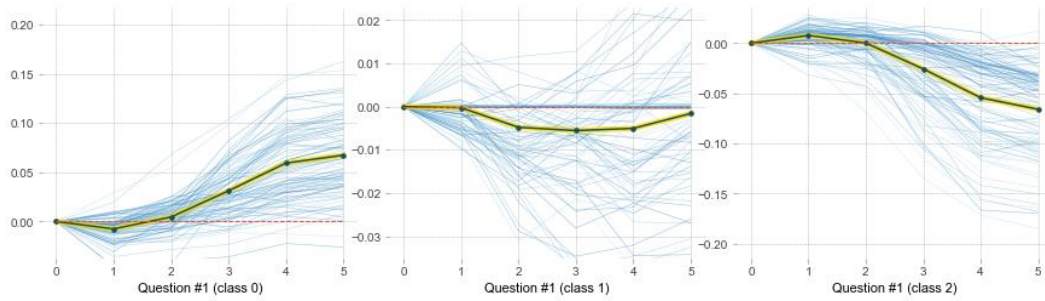
PDP for feature "# of Births"

Number of unique grid points: 5



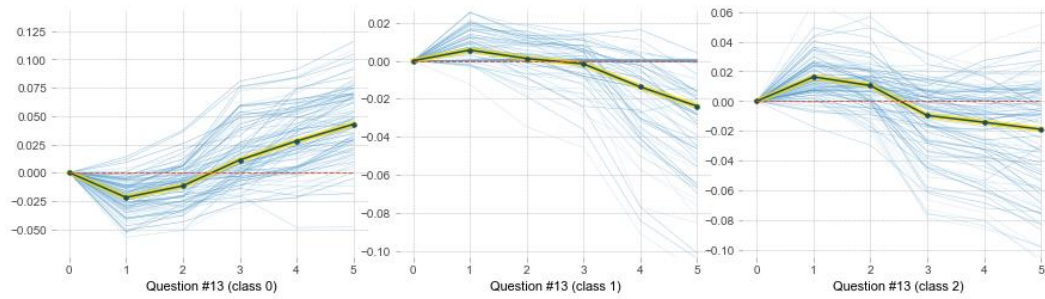
PDP for feature "Question #1"

Number of unique grid points: 6



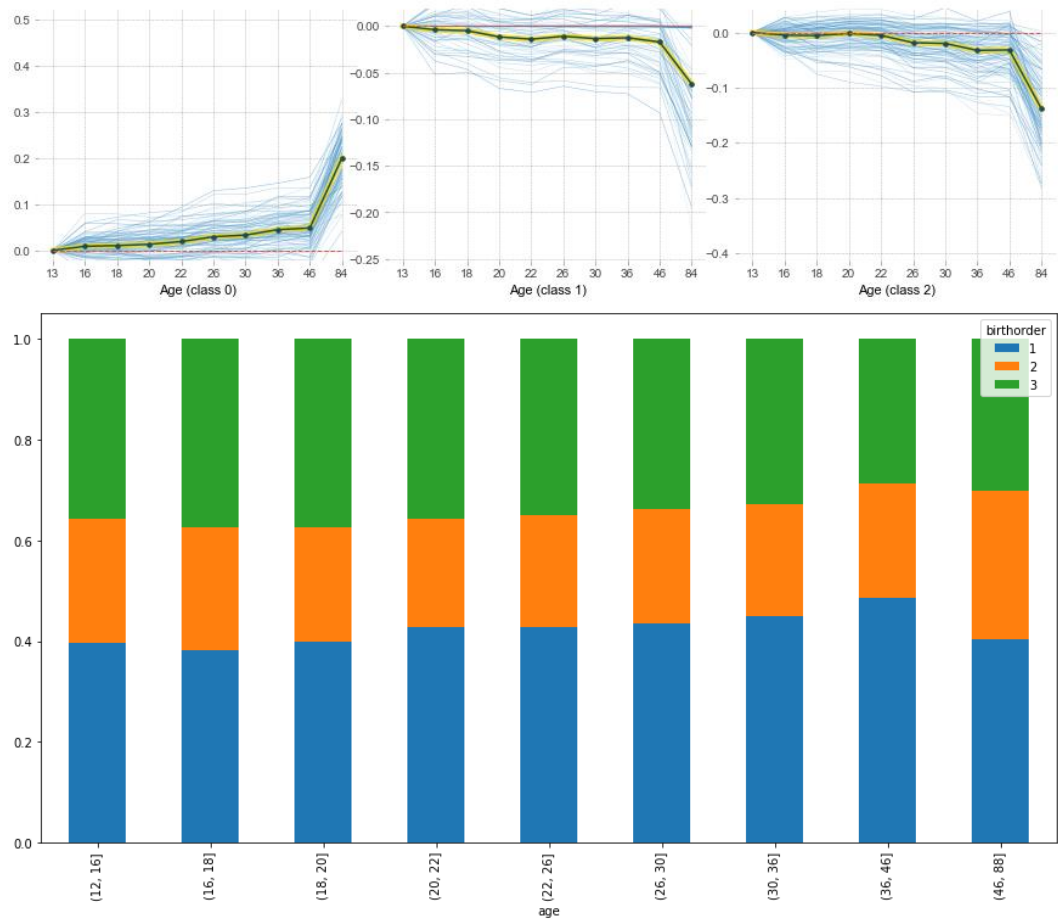
PDP for feature "Question #13"

Number of unique grid points: 6



PDP for feature "Age"

Number of unique grid points: 10



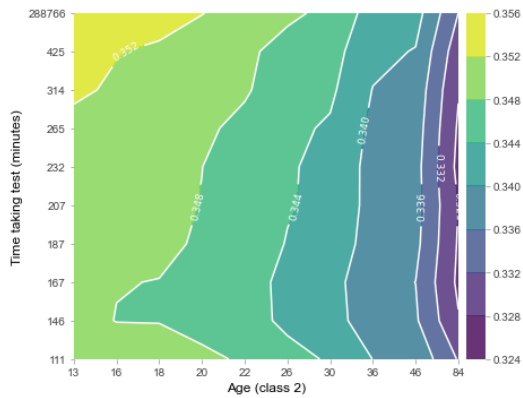
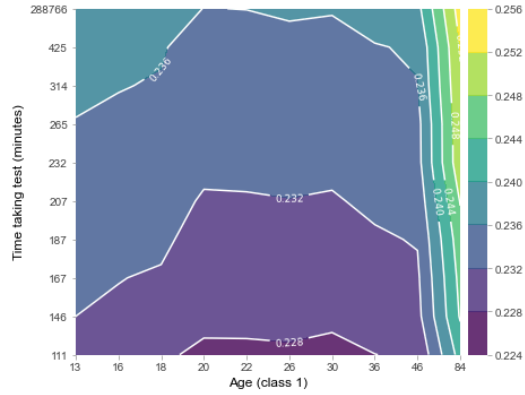
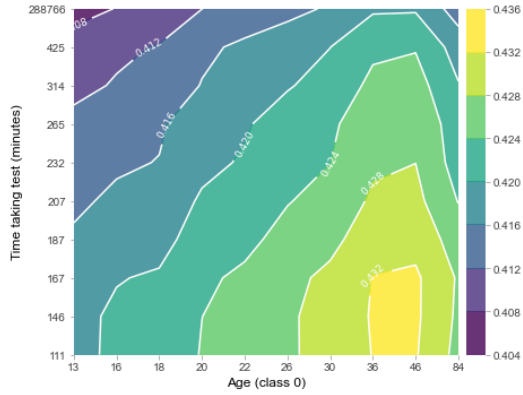
PDP interact for "# of Births" and "Question #1"

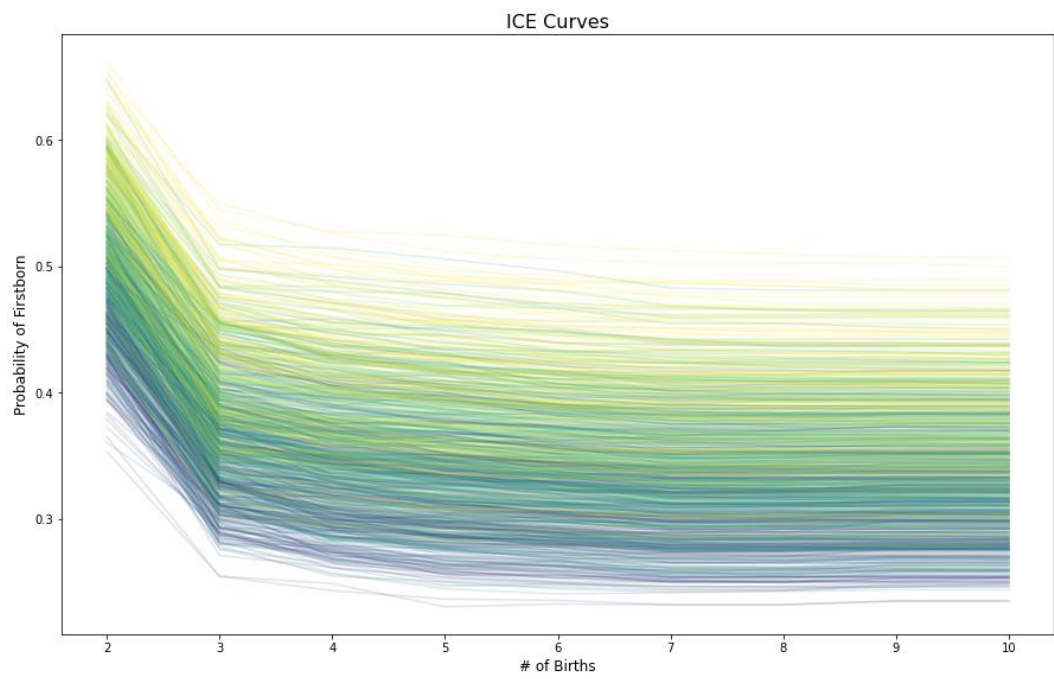
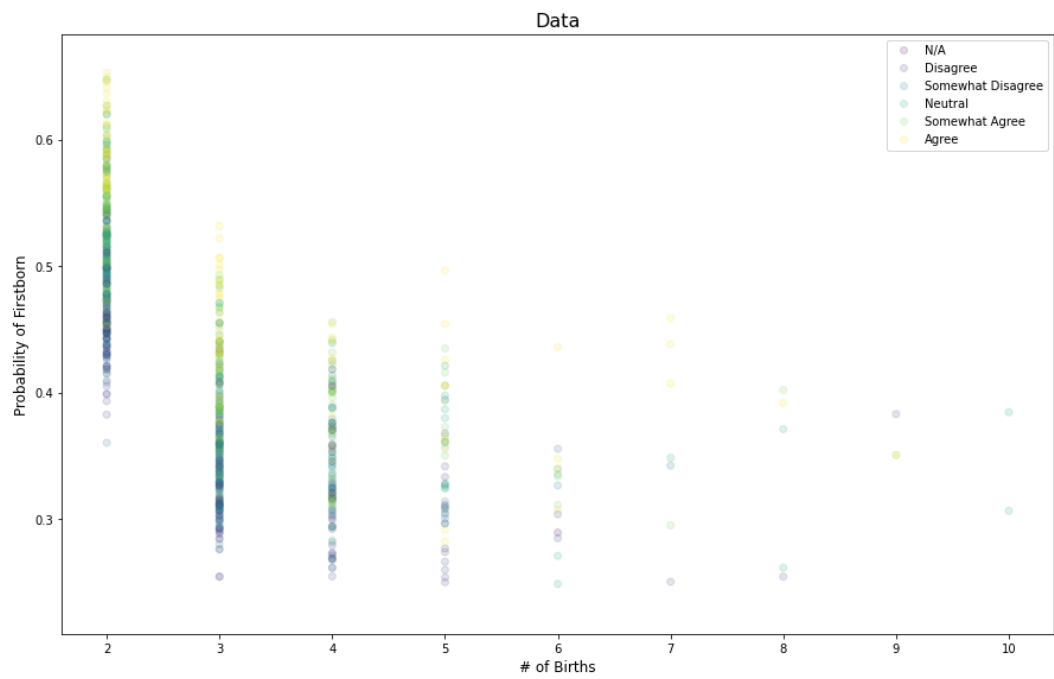
Number of unique grid points: (# of Births: 5, Question #1: 6)

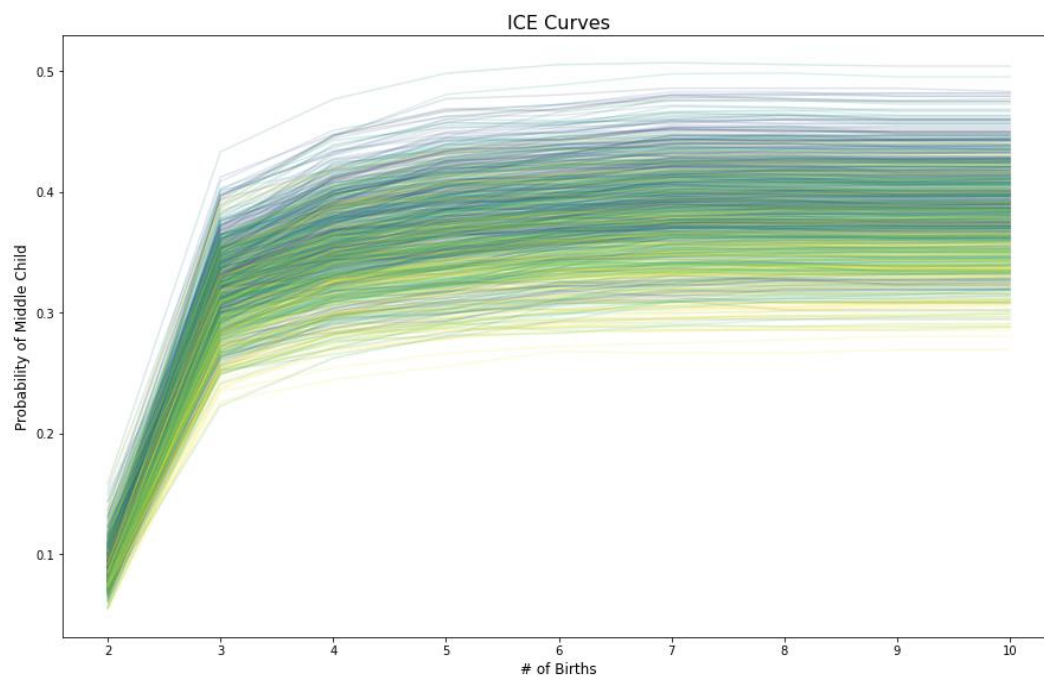


PDP interact for "Age" and "Time taking test (minutes)"

Number of unique grid points: (Age: 10, Time taking test (minutes): 10)

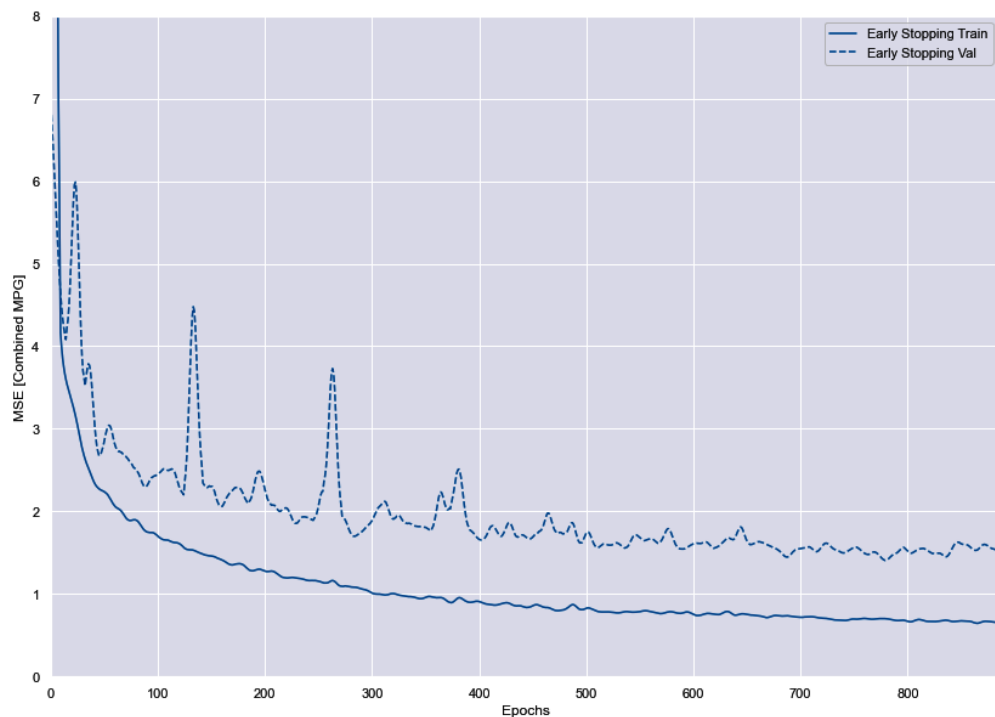


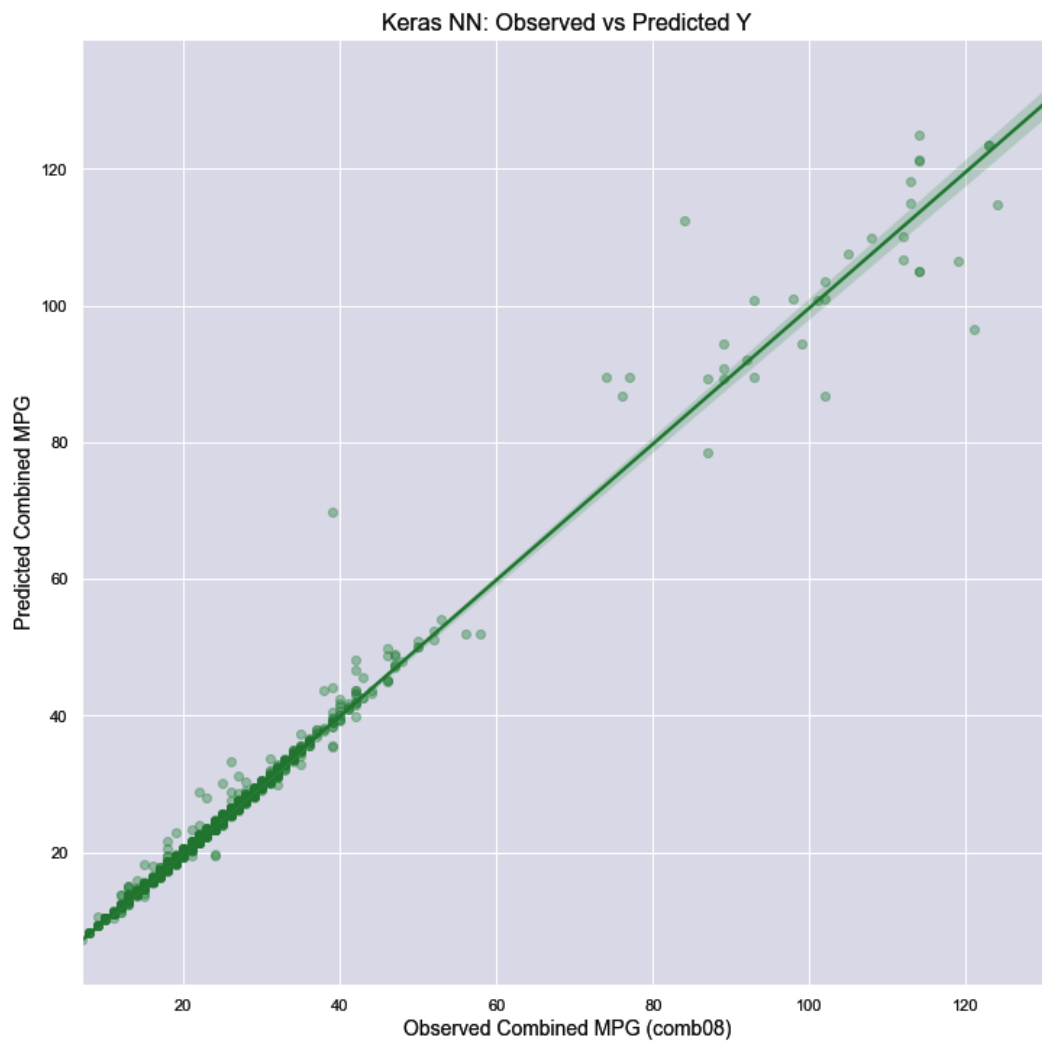


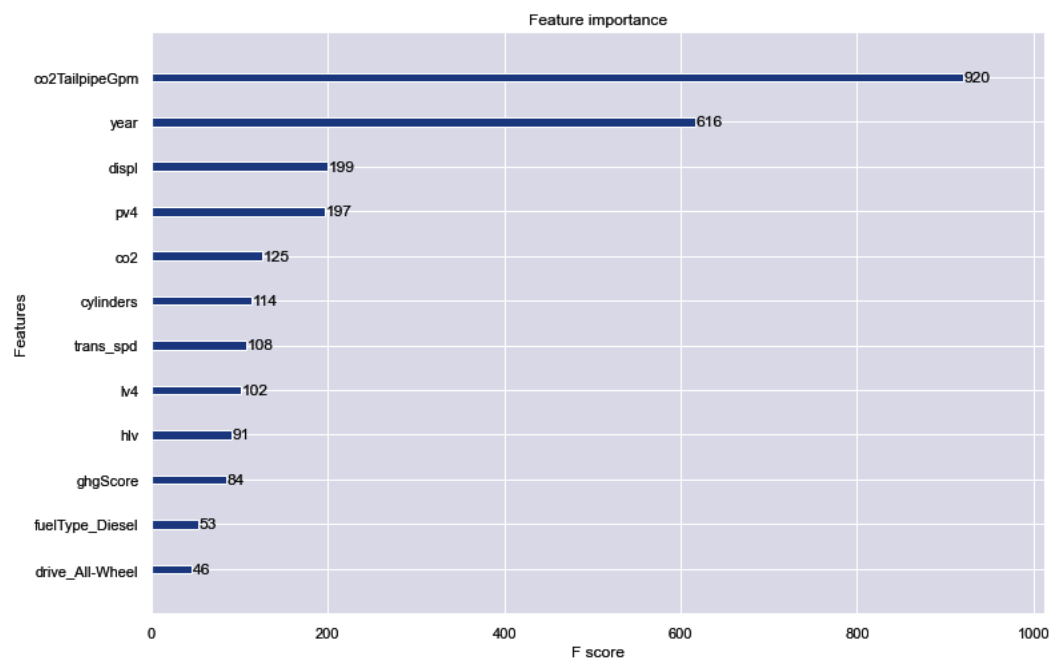
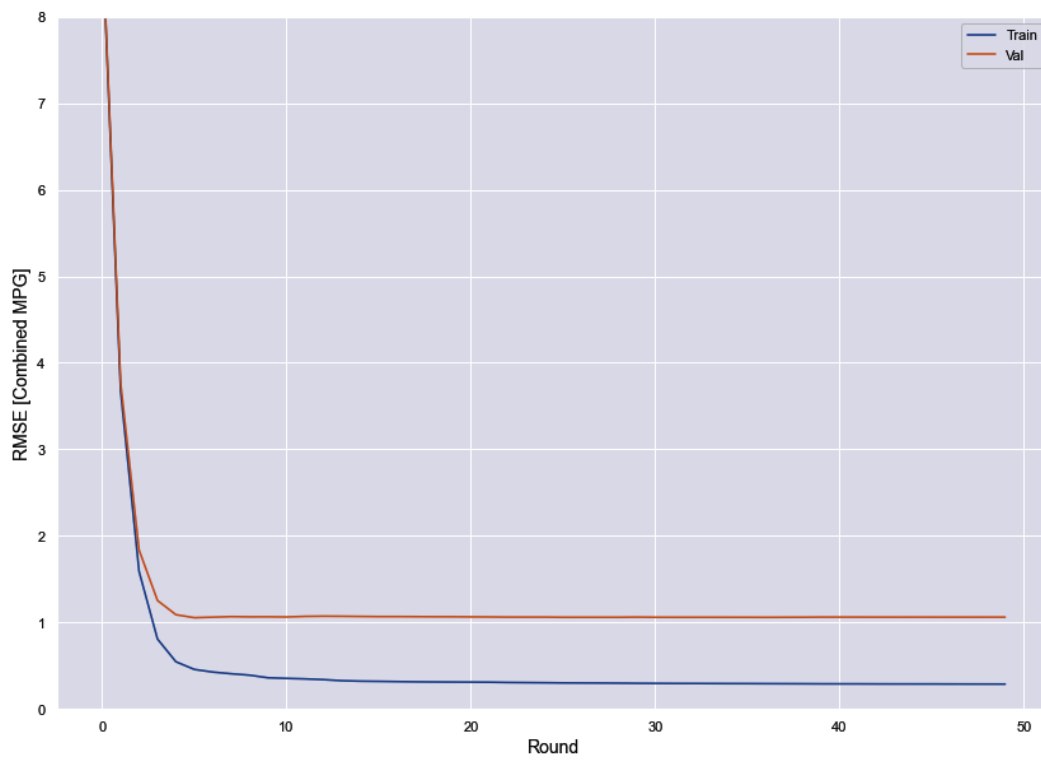


Chapter 5: Global Model-Agnostic Interpretation Methods

	rule	type	coef	support	importance
4	co2TailpipeGpm	linear	-0.034222	1.000000	3.865460
90	fuelType_Electricity > 0.5	rule	18.393562	0.006206	1.444470
104	co2TailpipeGpm <= 367.0 & fuelType_Electricity > 0.5 & ghgScore > 4.5	rule	20.776127	0.003546	1.235005
148	atvType_EV > 0.5 & pv4 > 42.0	rule	15.512440	0.003546	0.922113
193	eng_dscr_PFI > 0.5 & displ <= 0.30000001192092896	rule	12.607181	0.005319	0.917024
14	fuelType_Diesel	linear	2.967849	1.000000	0.488102
95	hpv > 45.0 & co2TailpipeGpm <= 28.5	rule	7.200693	0.003546	0.428034
127	co2TailpipeGpm > 124.0 & trans_spd > 0.5 & atvType_Other > 0.5 & cylinders > 1.5 & co2TailpipeGpm <= 408.5 & co2TailpipeGpm <= 302.5	rule	-7.168212	0.003546	0.426103
215	co2TailpipeGpm <= 250.53571319580078 & co2TailpipeGpm <= 320.5 & co2TailpipeGpm <= 410.5 & co2TailpipeGpm > 40.5	rule	3.457322	0.015071	0.421223
146	VClass_Small_Sport_Utility_Vehicle_2WD <= 0.5 & co2TailpipeGpm <= 45.5	rule	7.574405	0.002660	0.390100









EXPLAINER

- TreeExplainer
- DeepExplainer
- GradientExplainer
- LinearExplainer
- KernelExplainer
- SamplingExplainer
- PermutationExplainer
- PartitionExplainer
- AdditiveExplainer*

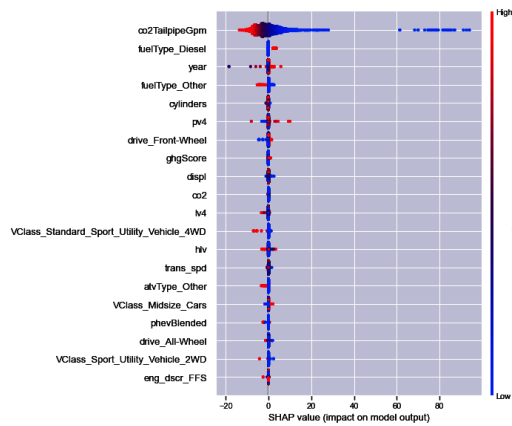
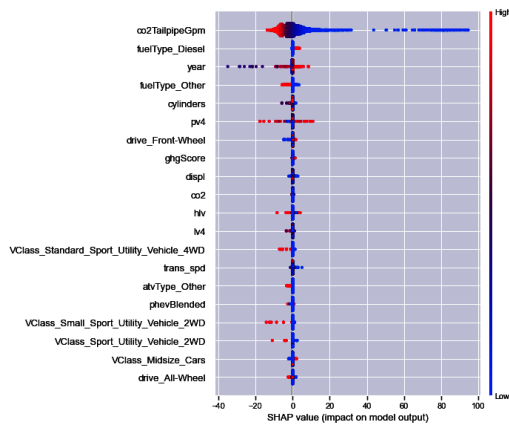
METHOD "UNIFIED"

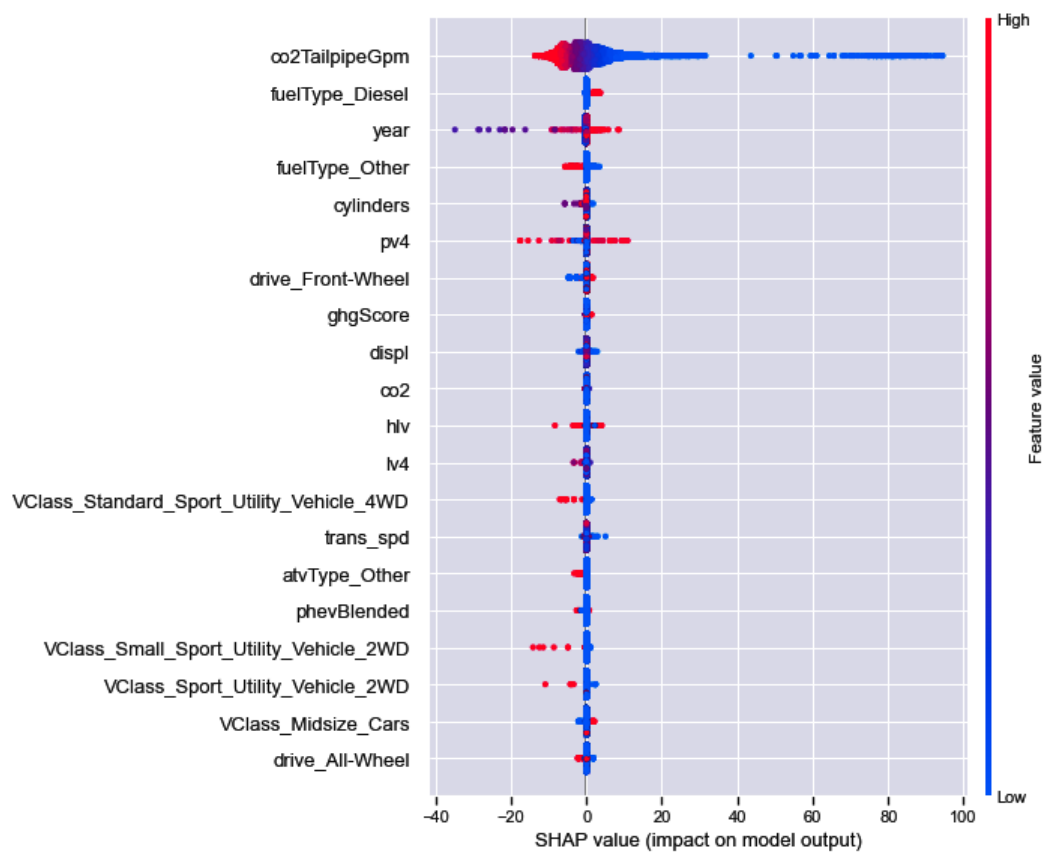
- TreeSHAP
(Lundberg et al. 2019)
- DeepLIFT
(Shrikumar et al. 2017)
- Integrated Gradients
(Sundararajan et al. 2017)
- Shapely Regression Values
(Lipovetsky & Conklin 2001)
- LIME
(Ribeiro et al. 2016)
- Shapely Sampling Values
(Strumbelj & Kononenko 2013)

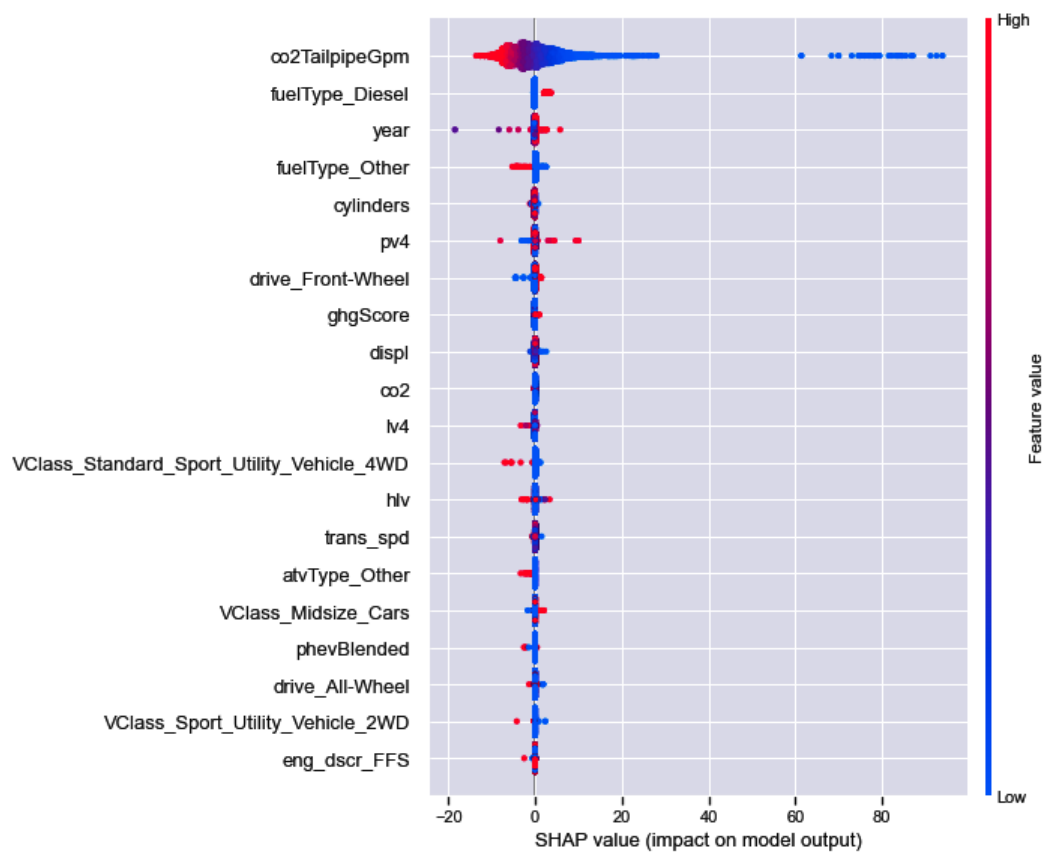
COMPATIBILITY WITH...

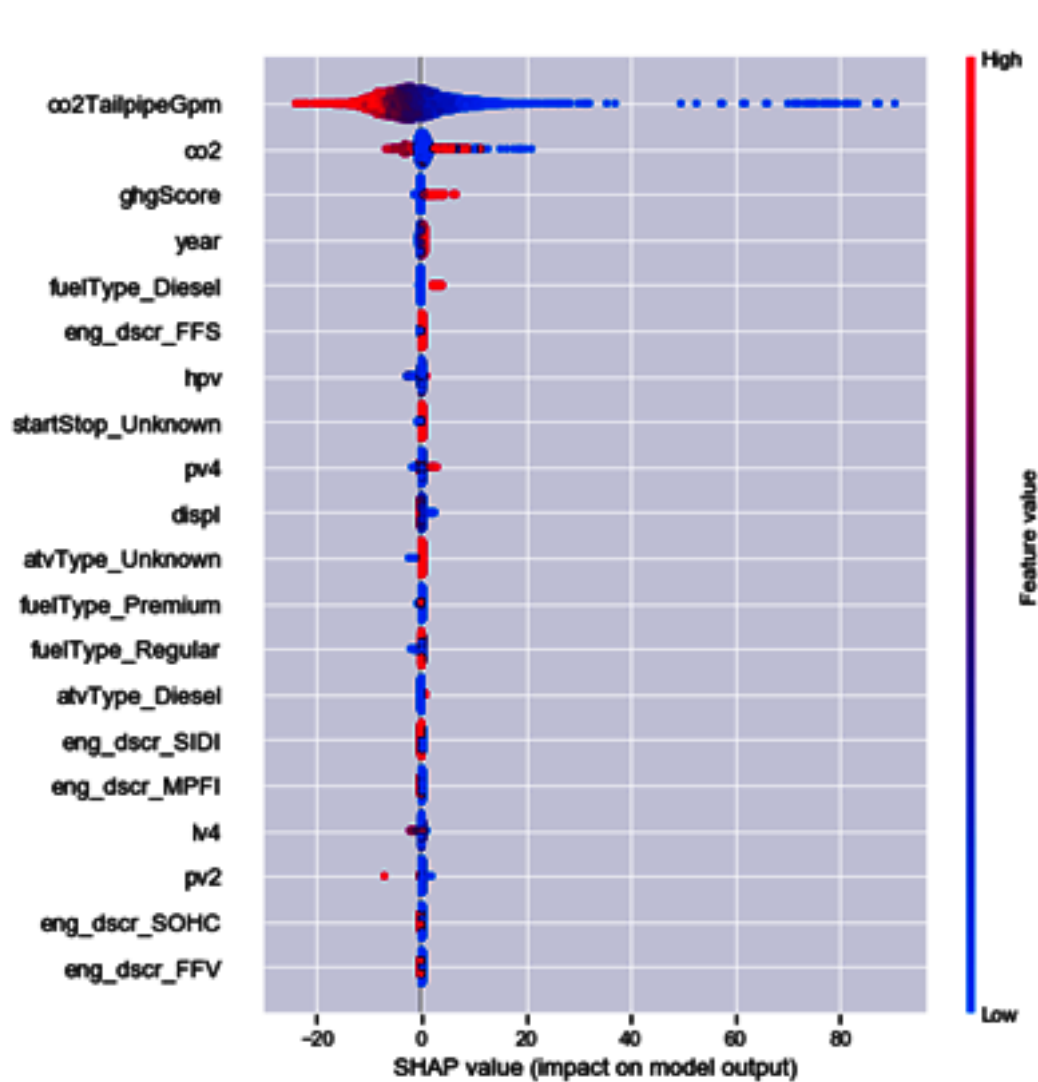
- XGBoost, LightGBM, CatBoost, Pyspark, sklearn.tree.*, sklearn.ensemble.*
- tf.keras.Model, torch.nn.Module
- tf.keras.Model, torch.nn.Module
- sklearn.linear_model.*

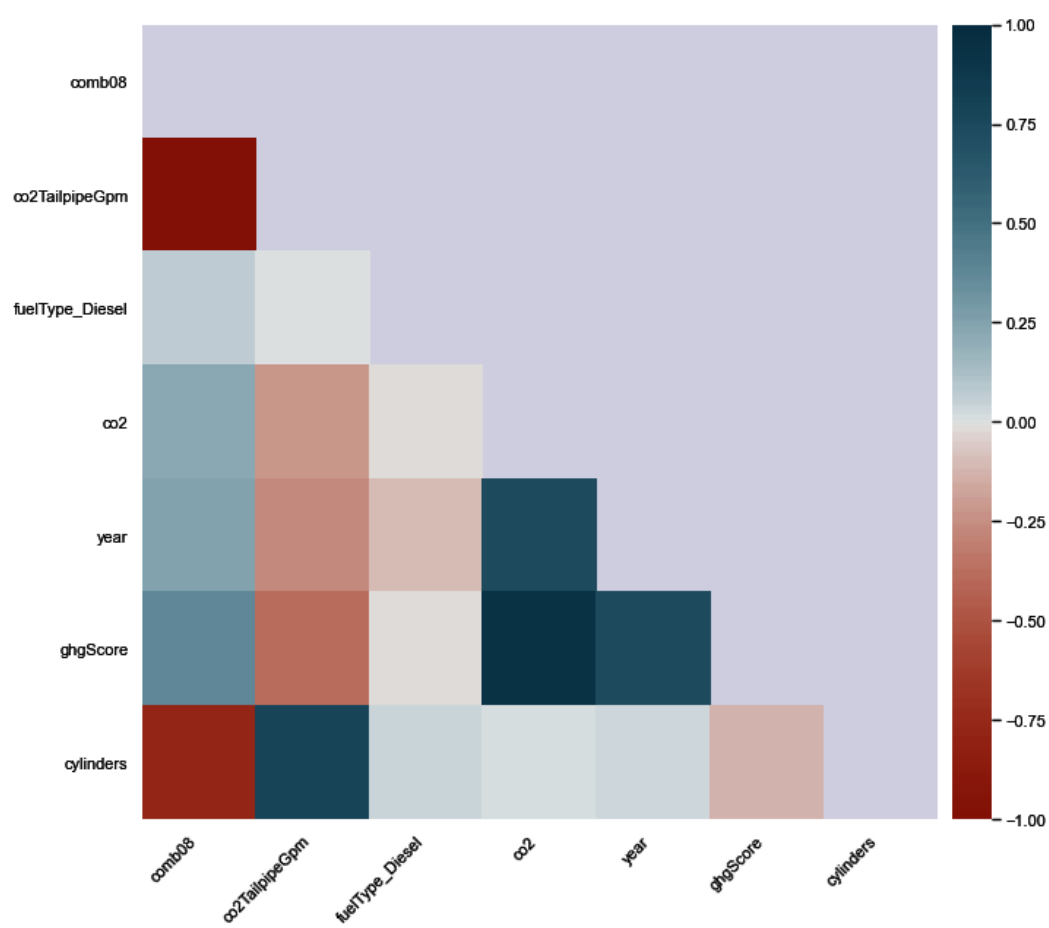
Model-Agnostic

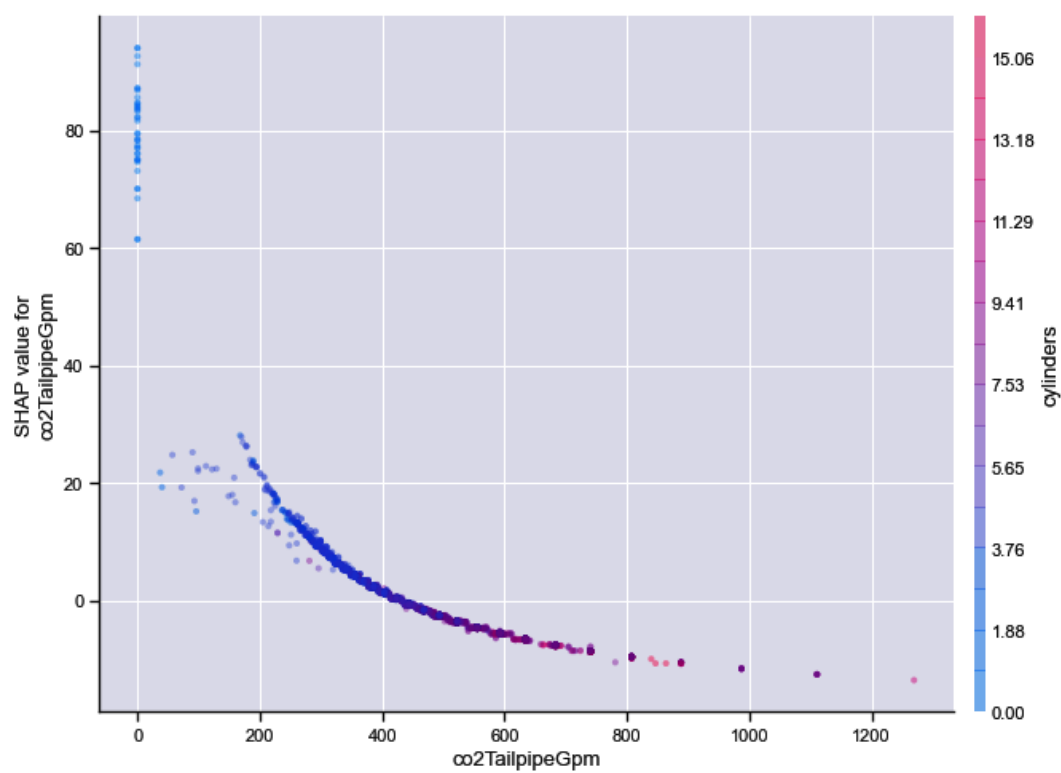




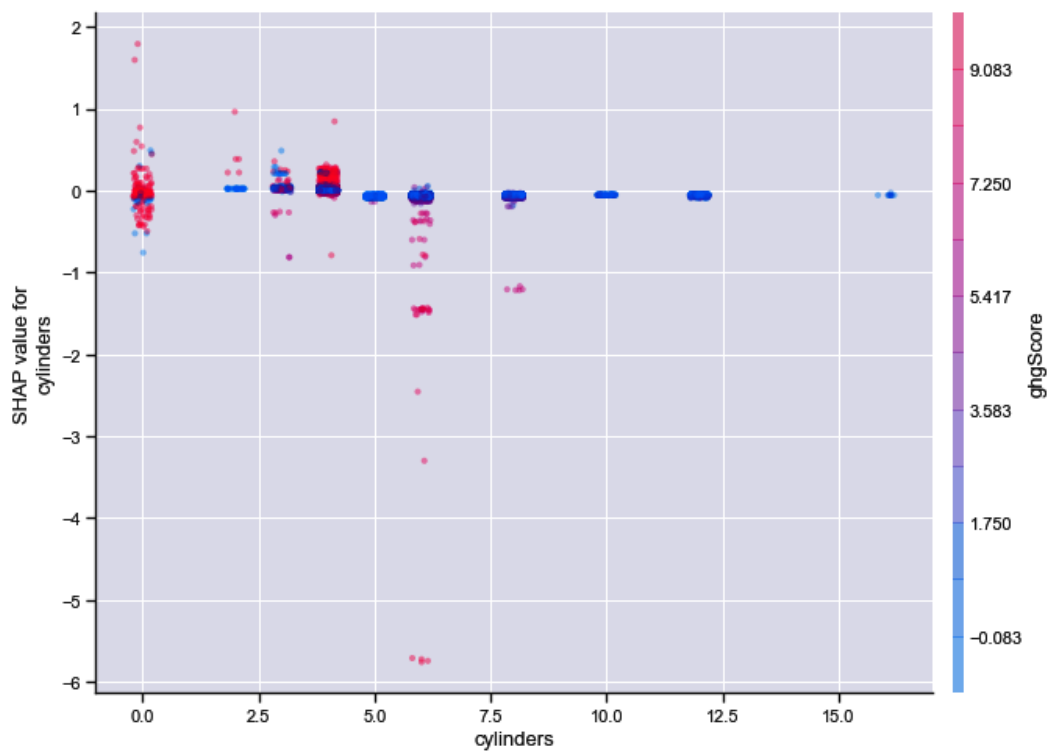




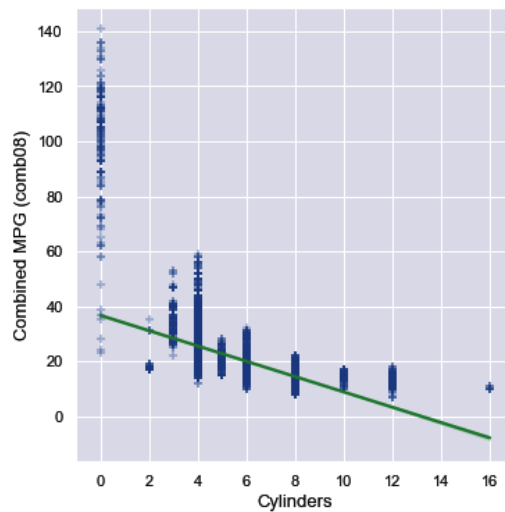
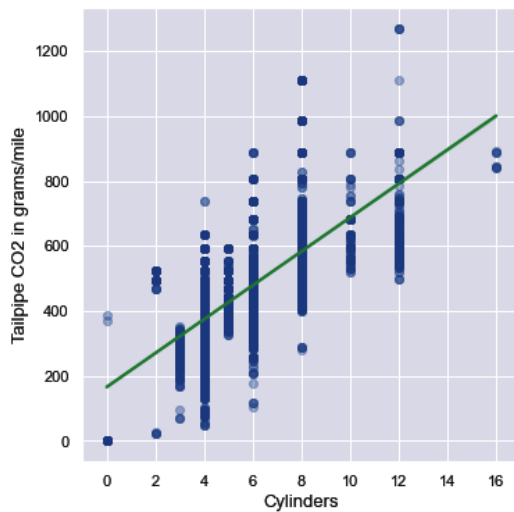


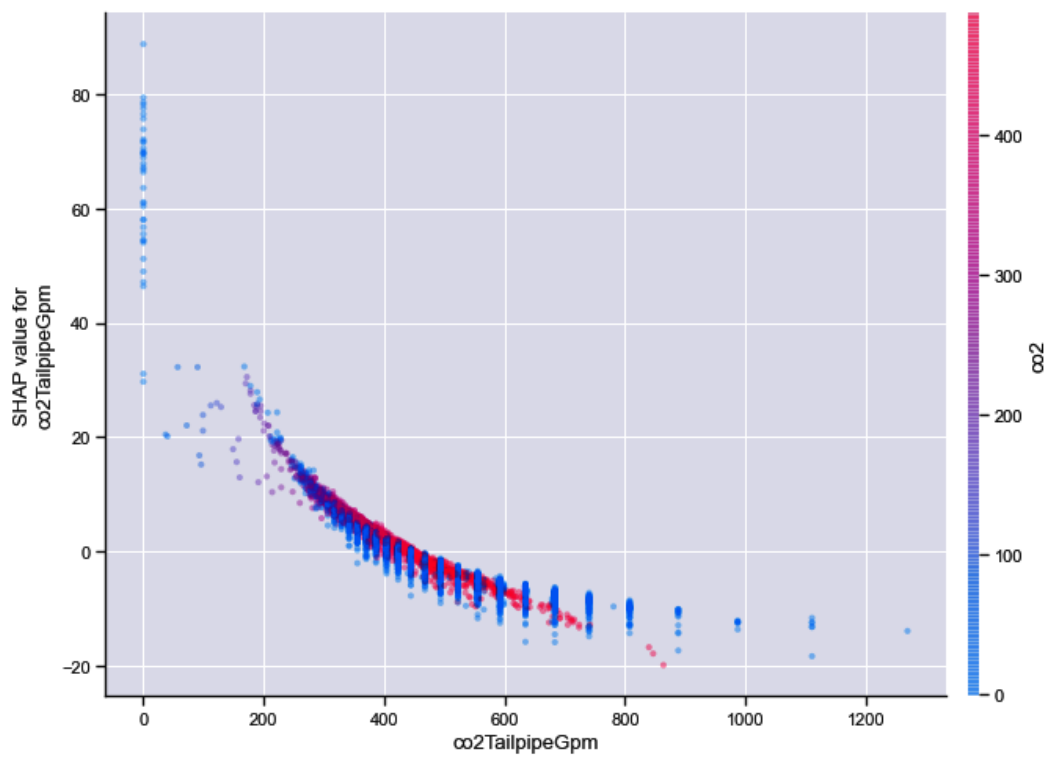


spearman cylinders→co2TailpipeGpm corr: 0.787 p-val:
0.0000

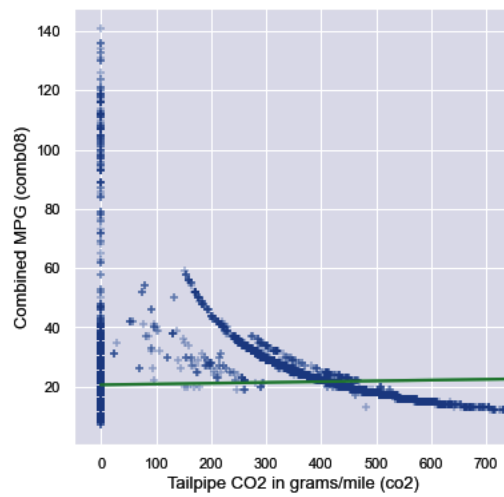
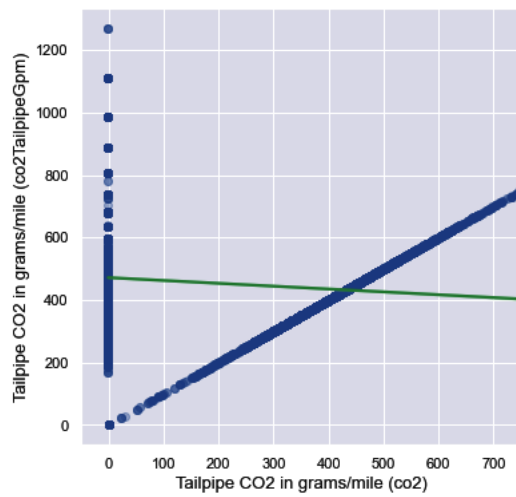


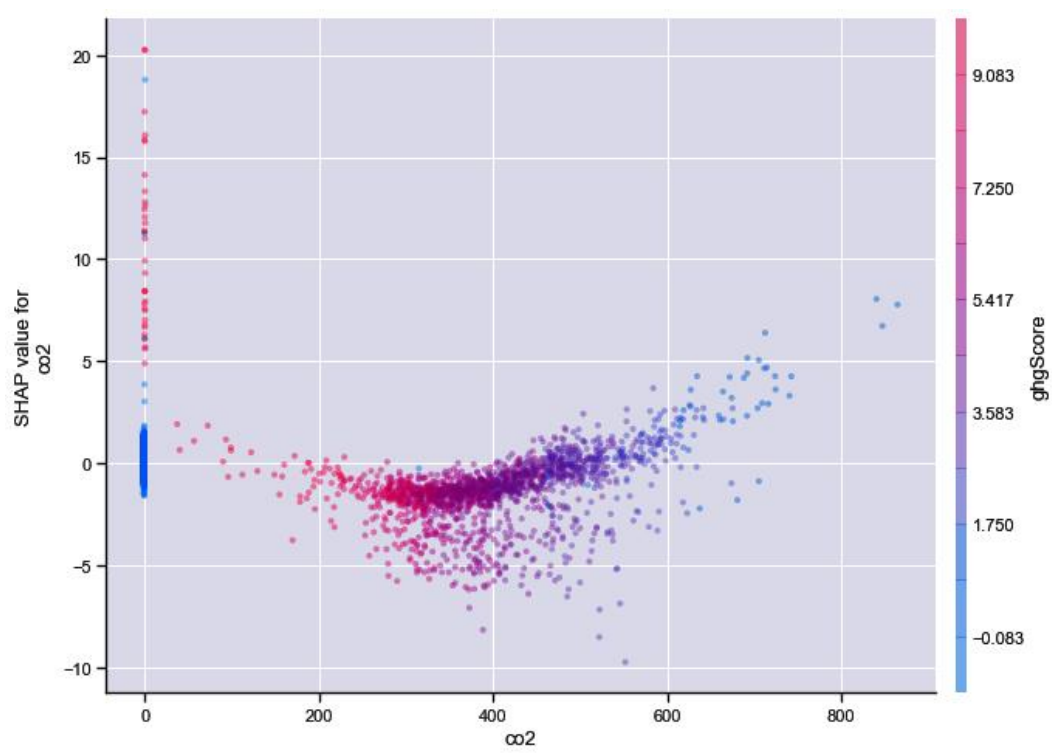
spearman ghgScore→cylinders corr: -0.117 p-val:
0.0000



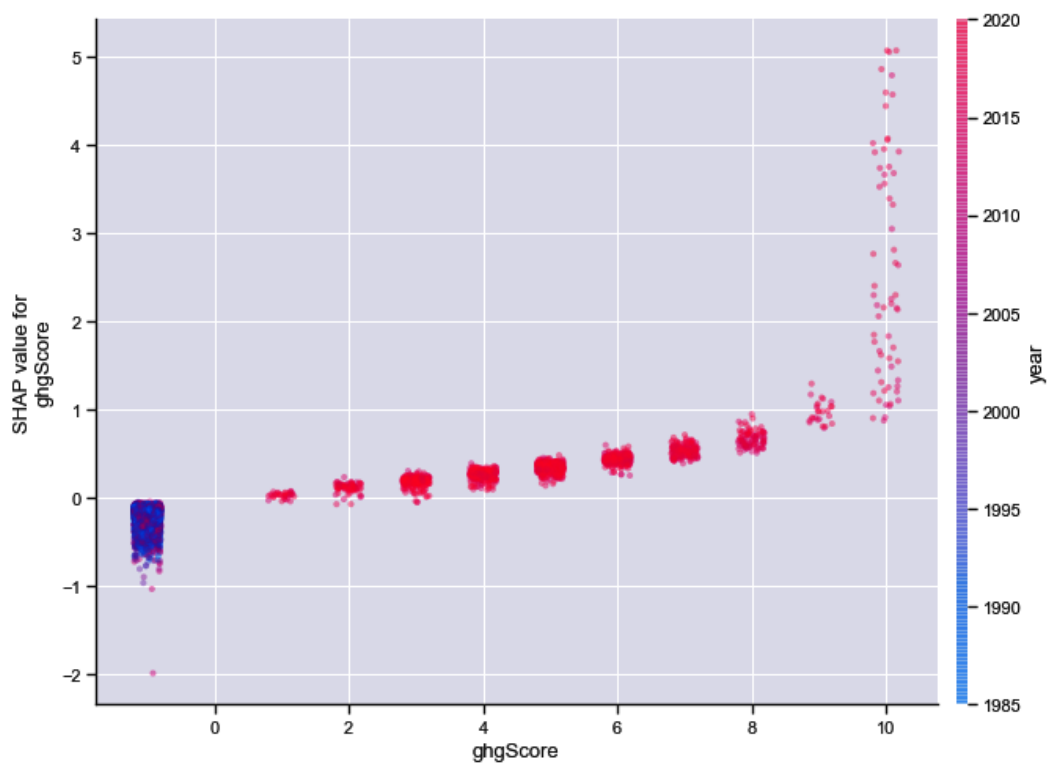


spearman co2→co2TailpipeGpm corr: -0.222 p-val:
0.0000

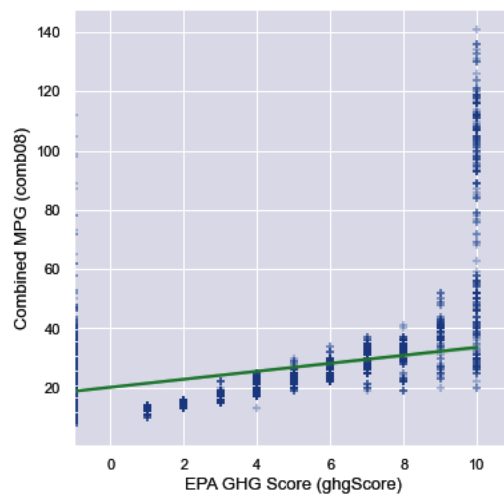
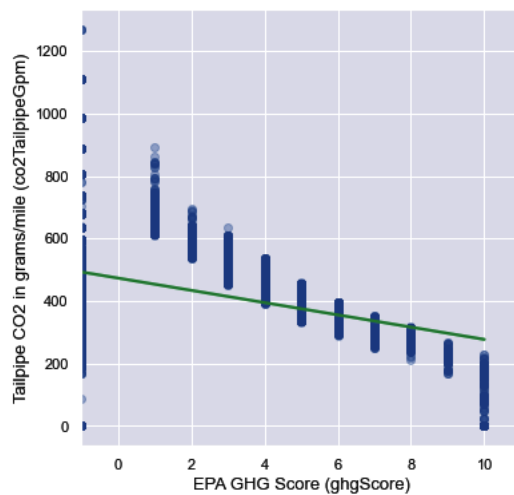


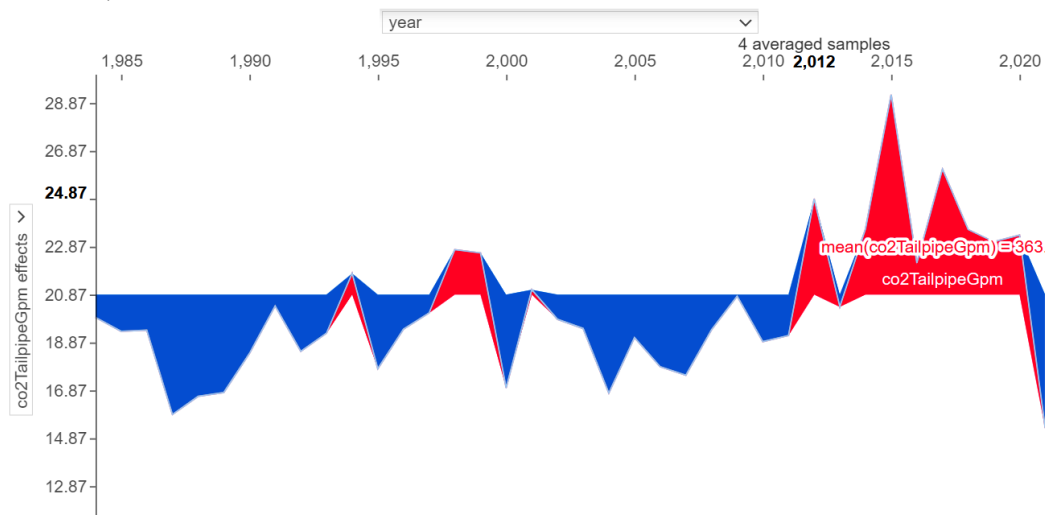
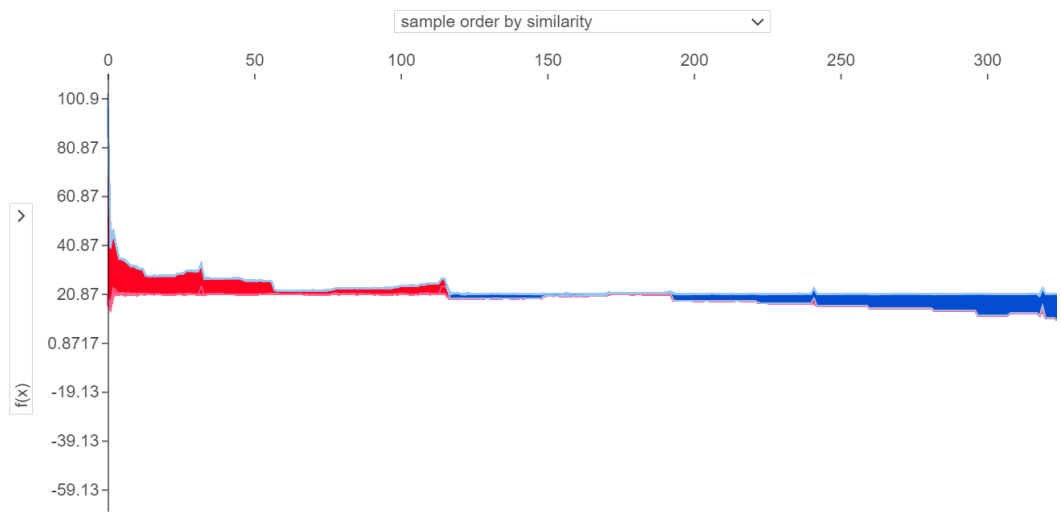


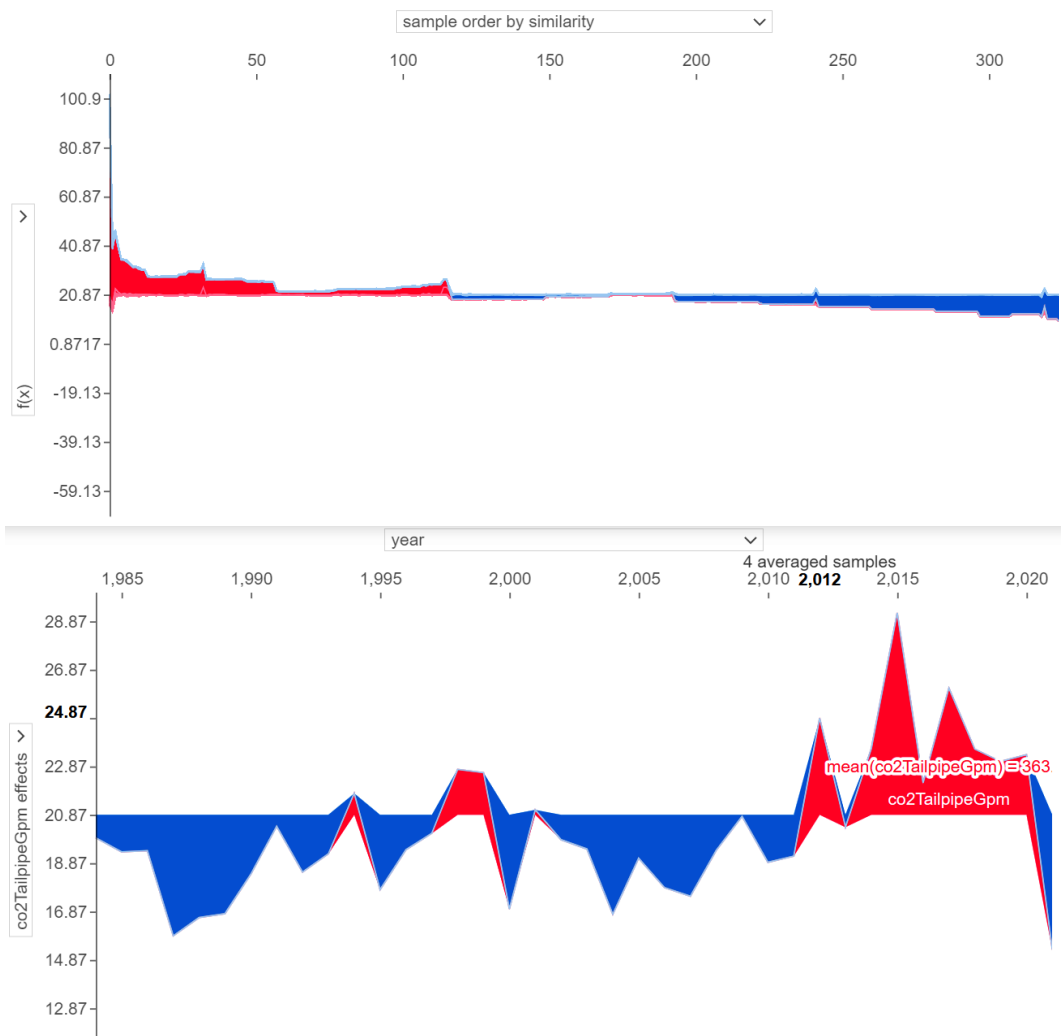
spearman ghgScore→co2 corr: 0.942 p-val: 0.0000

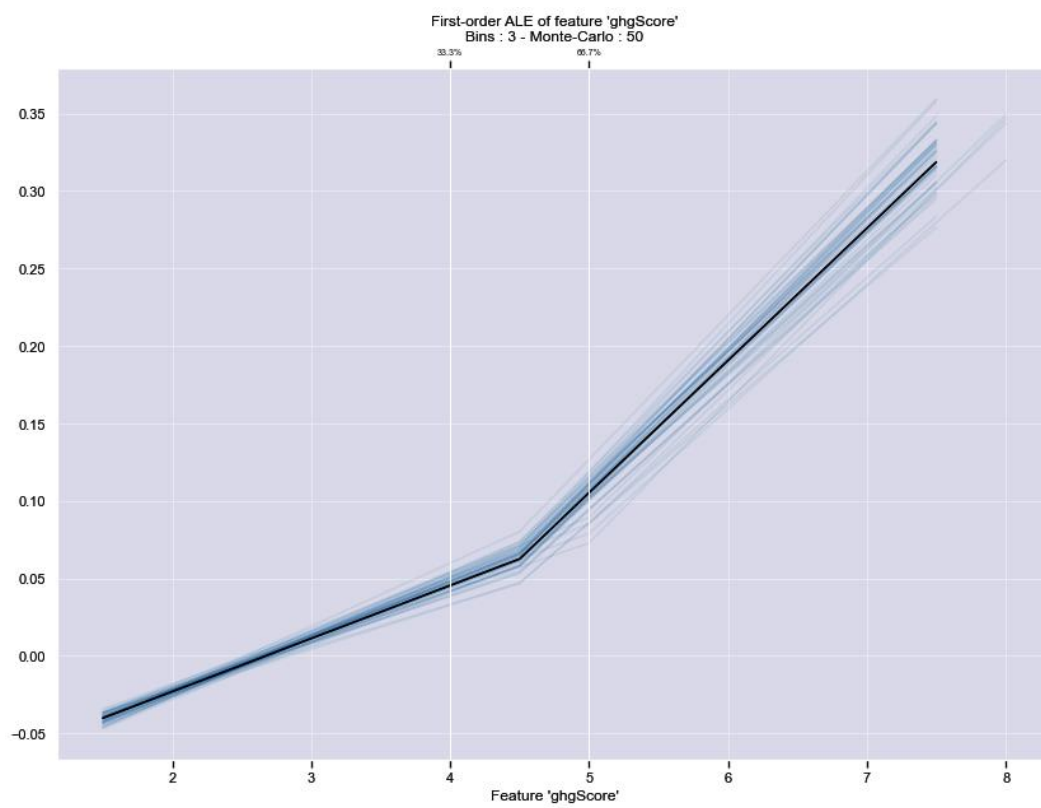


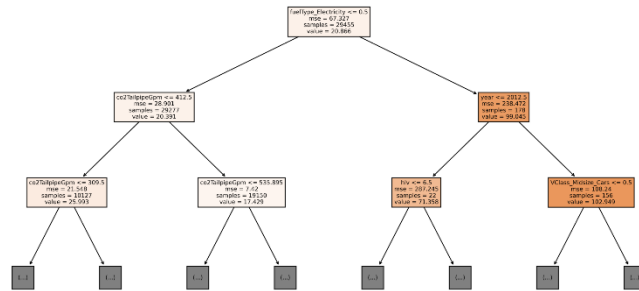
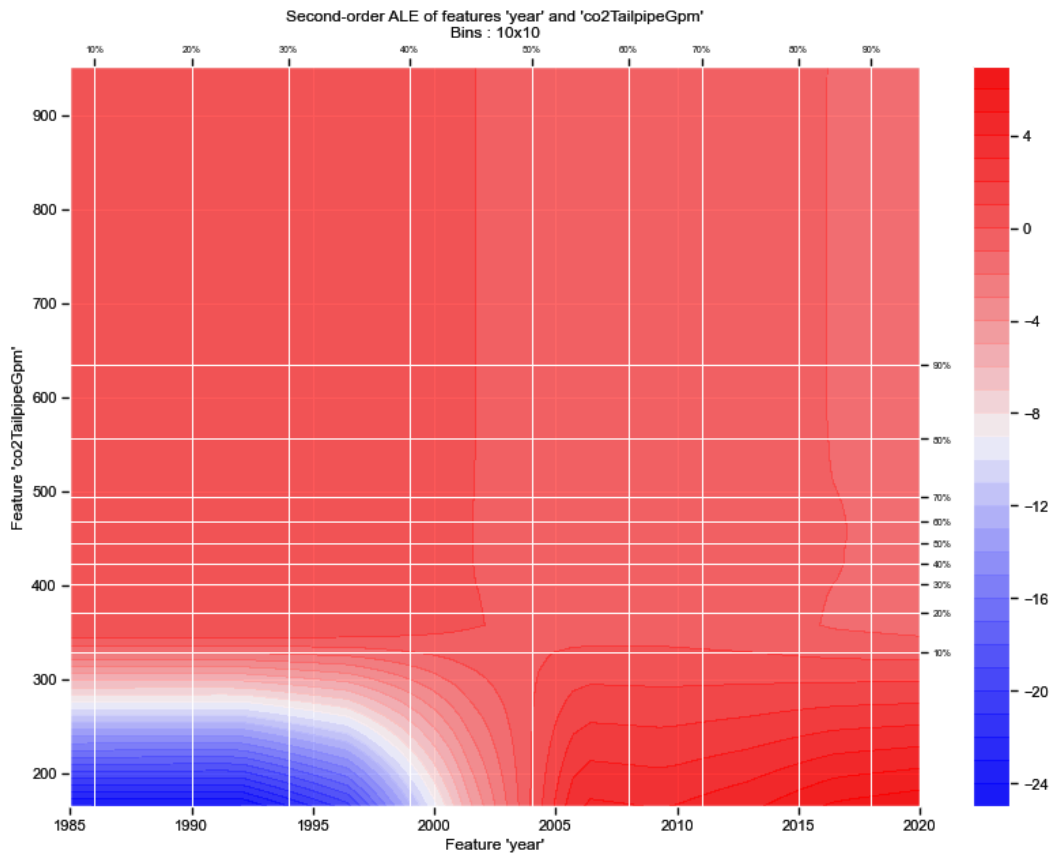
spearman ghgScore→year corr: 0.744 p-val: 0.0000



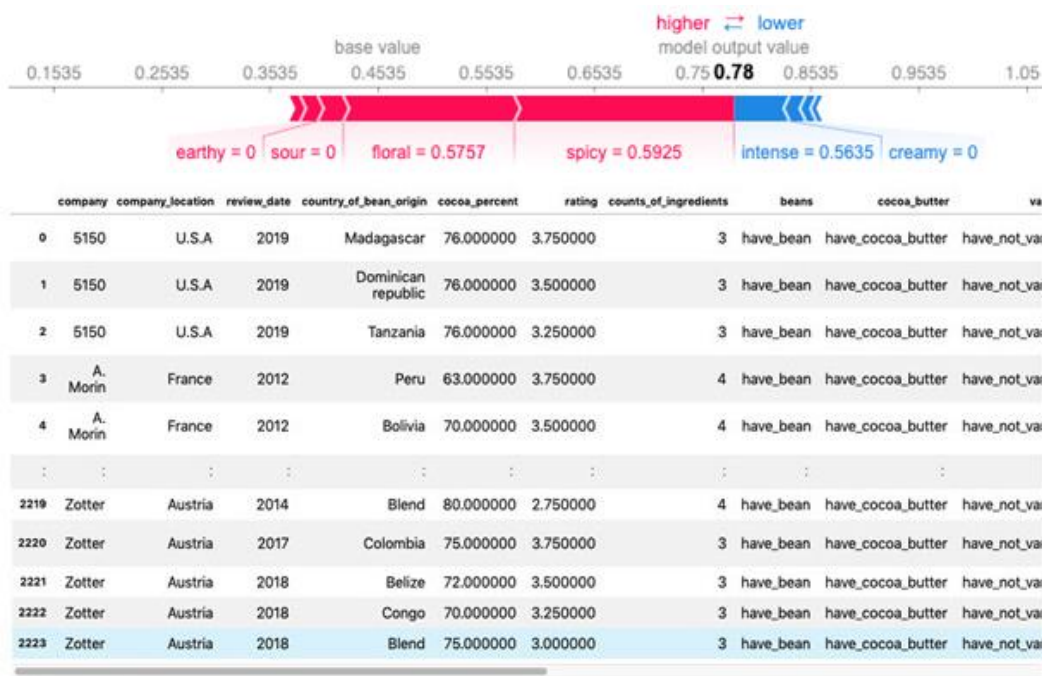






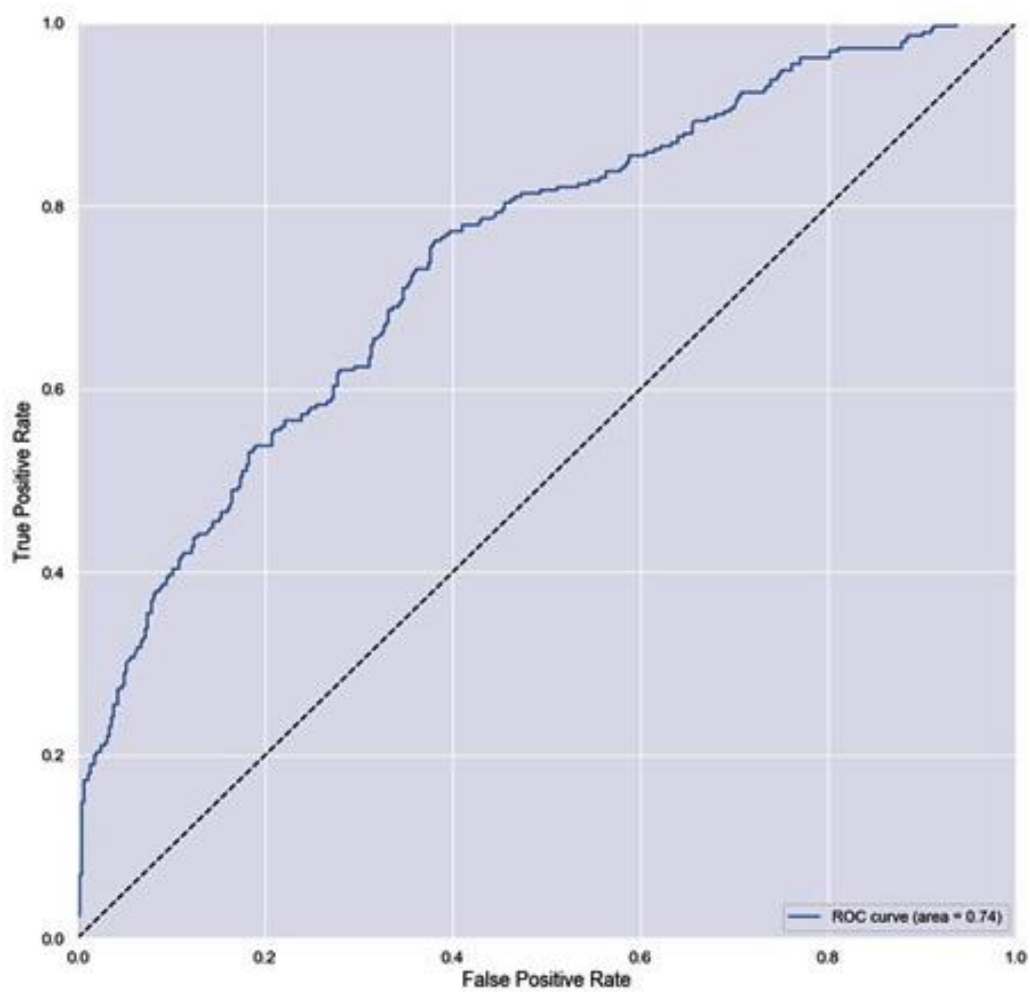


Chapter 6: Local Model-Agnostic Interpretation Methods

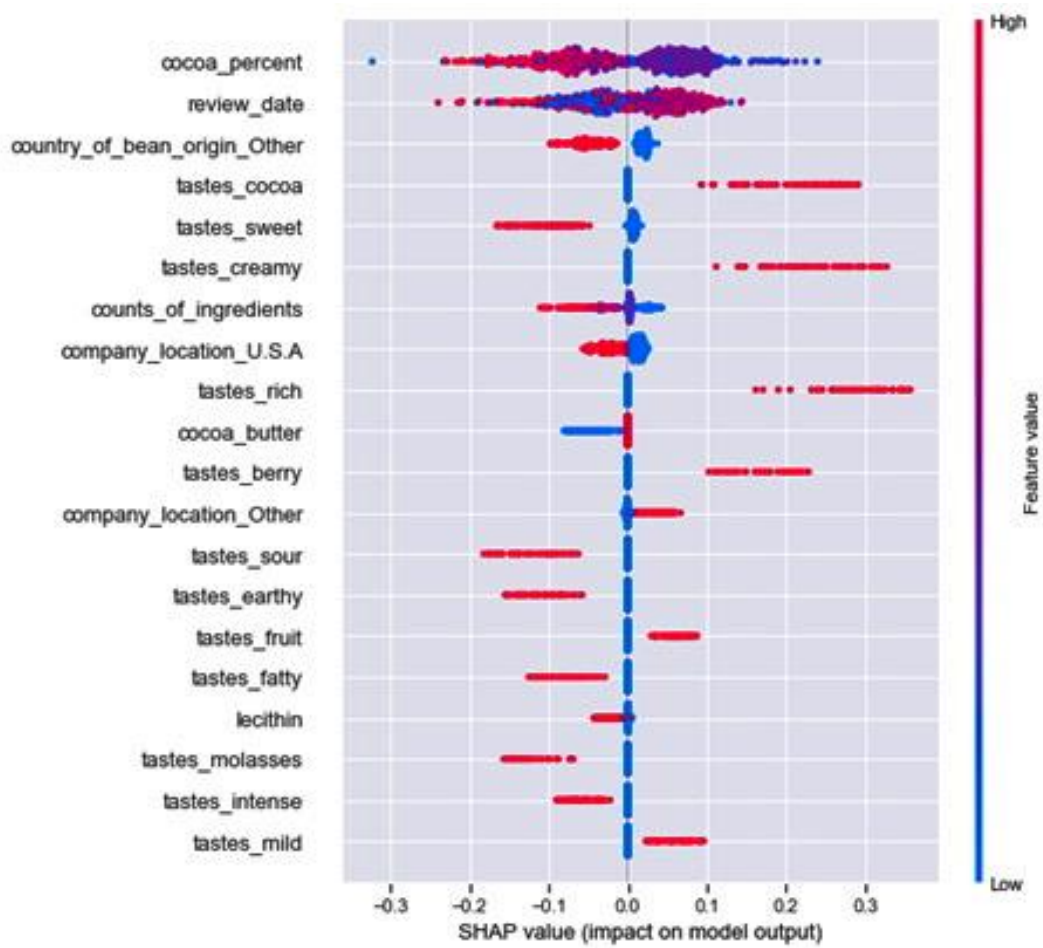


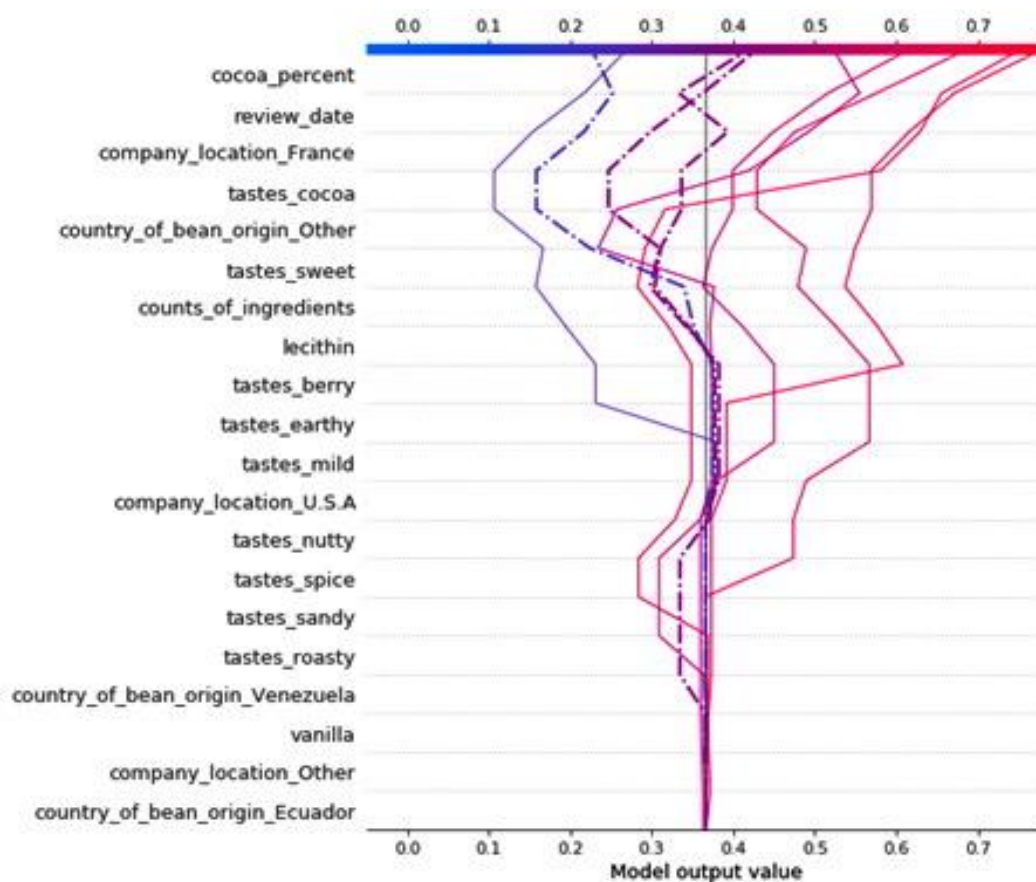
2224 rows x 14 columns

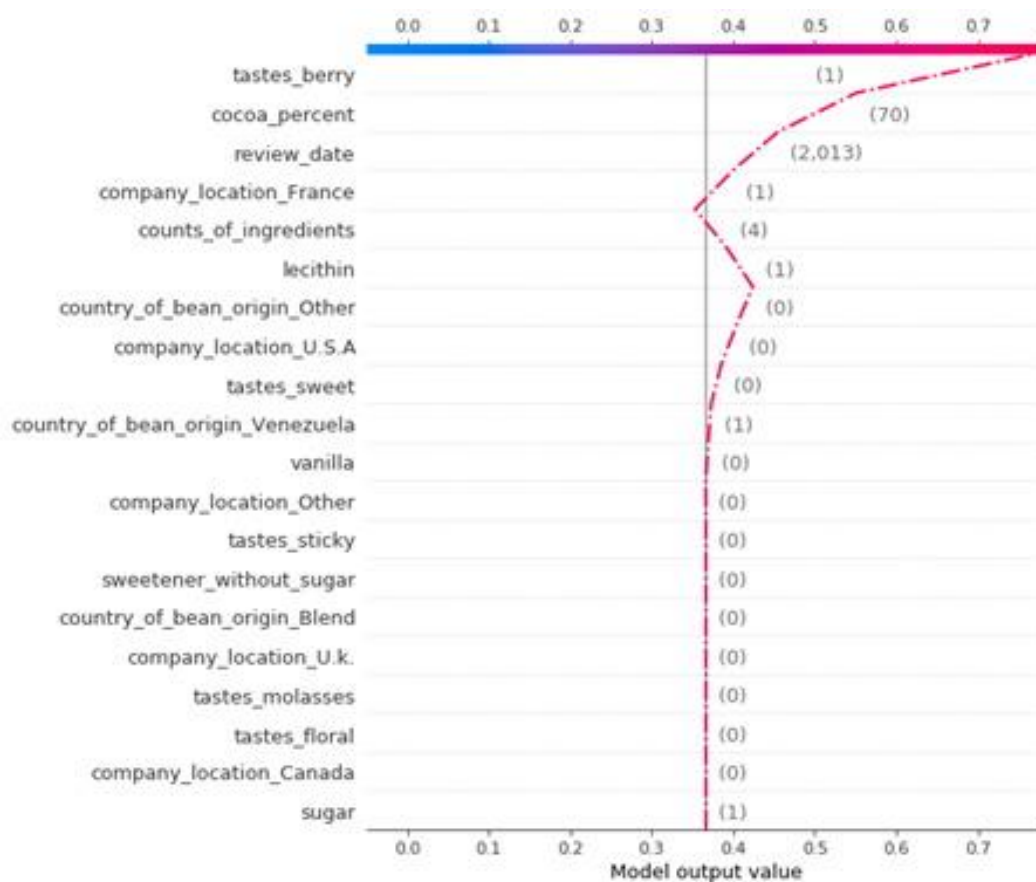
	first_taste	second_taste	third_taste	fourth_taste
80	oily	vegetal	nutty	cocoa
81	oily	vanilla	melon	cocoa
82	rich	sour	mild smoke	nan
83	fruity	sour	nan	nan
84	high roast	high astringent	nan	nan
85	smokey	savory	nan	nan
86	sandy	roasty	nutty	nan
87	roasty	brownie	nutty	nan
88	red wine	rich cocoa	long	nan
89	creamy	fruit	cocoa	nan



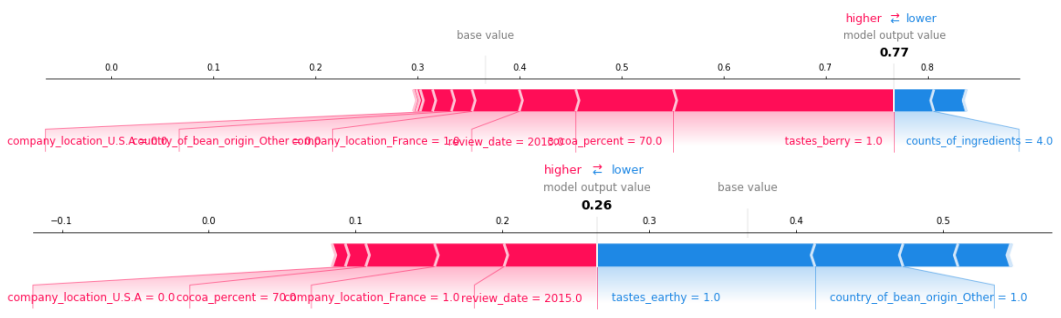
Accuracy_train:	0.7315	Accuracy_test:	0.6962	
Precision_test:	0.6772	Recall_test:	0.4414	
ROC-AUC_test:	0.7449	F1_test:	0.5344	MCC_test:
	0.3399			

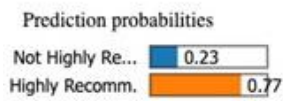




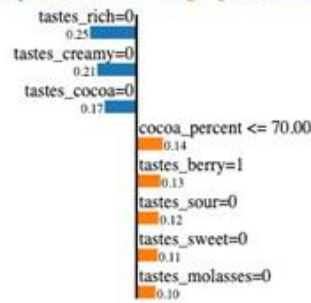


	5	24
rating	4.00	2.75
y	1.00	0.00
y_pred	1.00	0.00
review_date	2013.00	2015.00
cocoa_percent	70.00	70.00
counts_of_ingredients	4.00	4.00
cocoa_butter	1.00	1.00
vanilla	0.00	0.00
lecithin	1.00	1.00
salt	0.00	0.00
sugar	1.00	1.00
sweetener_without_sugar	0.00	0.00
company_location_Canada	0.00	0.00
company_location_France	1.00	1.00
:	:	:
country_of_bean_origin_Other	0.00	1.00
country_of_bean_origin_Peru	0.00	0.00
country_of_bean_origin_Venezuela	1.00	0.00
:	:	:
tastes_earthy	0.00	1.00
:	:	:
tastes_berry	1.00	0.00
tastes_vanilla	0.00	0.00
tastes_creamy	0.00	0.00

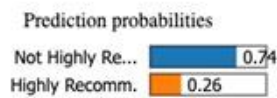




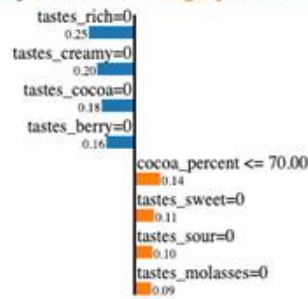
Not Highly Recomm. Highly Recomm.



Feature	Value
tastes_rich=0	True
tastes_creamy=0	True
tastes_cocoa=0	True
cocoa_percent	70.00
tastes_berry=1	True
tastes_sour=0	True
tastes_sweet=0	True
tastes_molasses=0	True



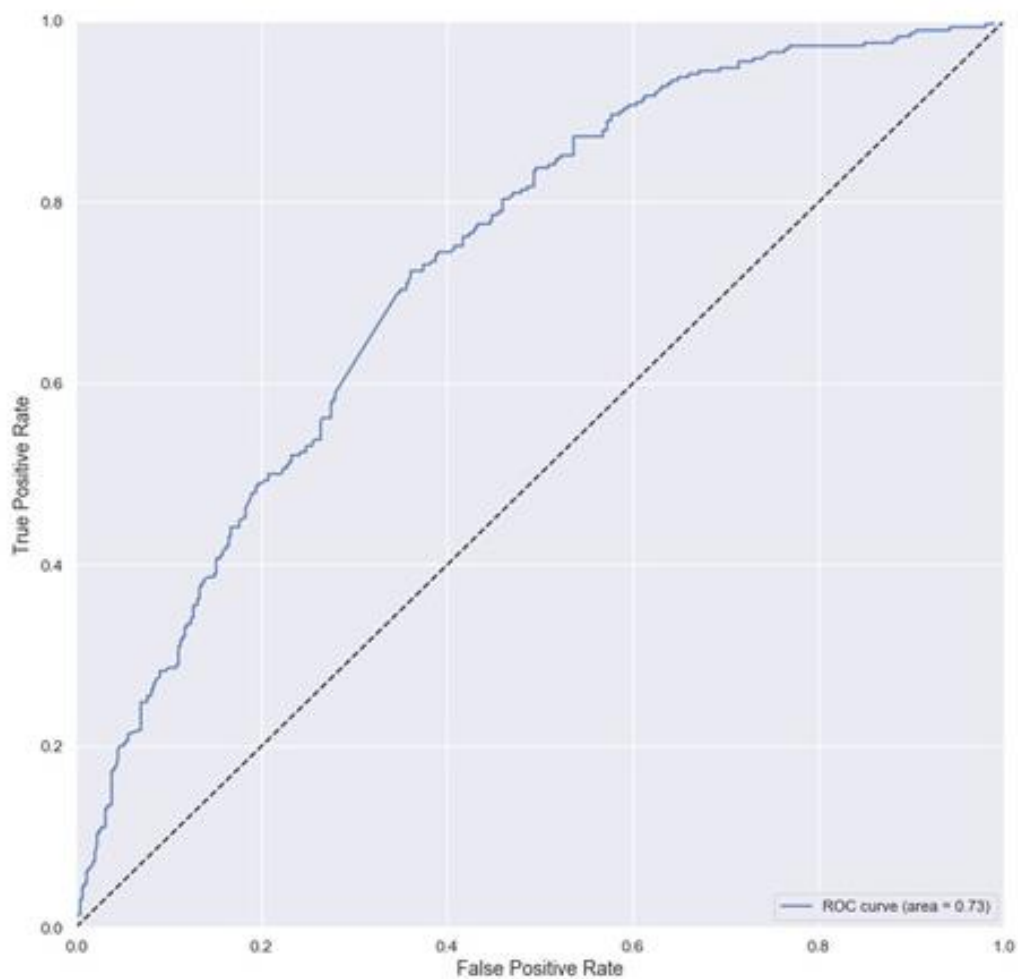
Not Highly Recomm. Highly Recomm.



Feature	Value
tastes_rich=0	True
tastes_creamy=0	True
tastes_cocoa=0	True
tastes_berry=0	True
cocoa_percent	70.00
tastes_sweet=0	True
tastes_sour=0	True
tastes_molasses=0	True

	taste	tf-idf
305	raspberry	0.59
259	nut	0.49
265	oily	0.46
64	caramel	0.45
274	papaya	0.00
:	:	:
135	edge	0.00
134	easy	0.00
133	easter	0.00
415	yogurt	0.00

416 rows × 2 columns

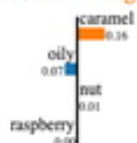


Accuracy_train:	0.7953	Accuracy_test:	0.6798	
Precision_test:	0.6233	Recall_test:	0.4793	
ROC-AUC_test:	0.7328	F1_test:	0.5419	MCC_test:
	0.3084			

Prediction probabilities

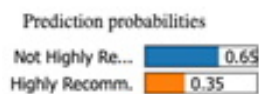
Not Highly Re...	0.43
Highly Recomm.	0.57

Not Highly Recomm. Highly Recomm.

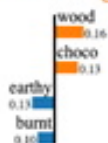


Text with highlighted words

oily nut **caramel** raspberry

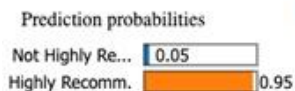


Not Highly Recomm. Highly Recomm.

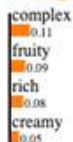


Text with highlighted words

burnt wood earthy choco

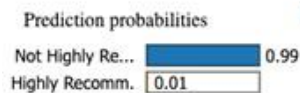


Not Highly Recomm. Highly Recomm.



Text with highlighted words

creamy rich complex fruity

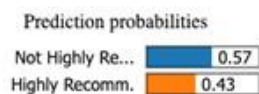


Not Highly Recomm. Highly Recomm.

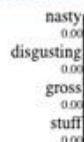


Text with highlighted words

sour bitter roasty molasses

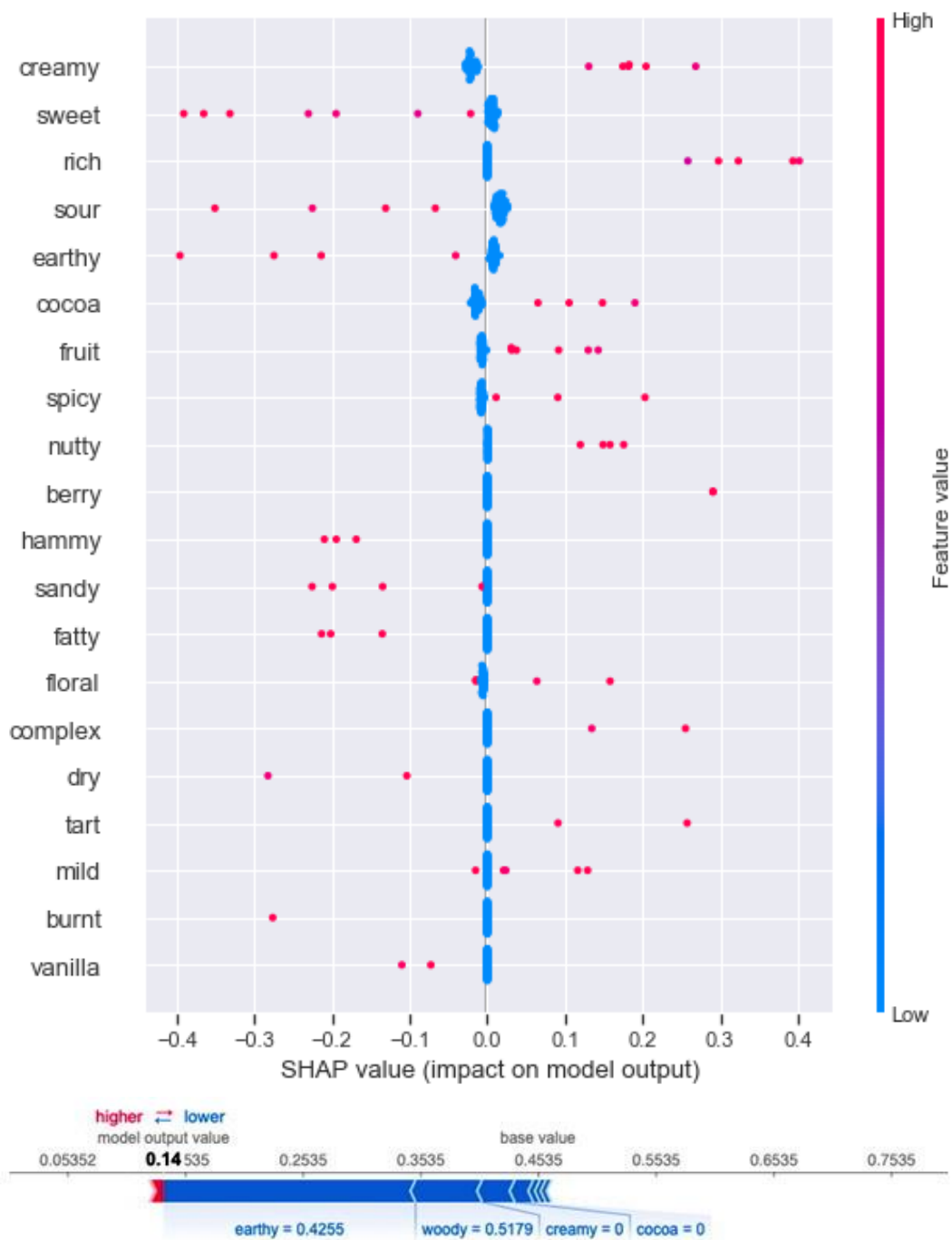


Not Highly Recomm. Highly Recomm.



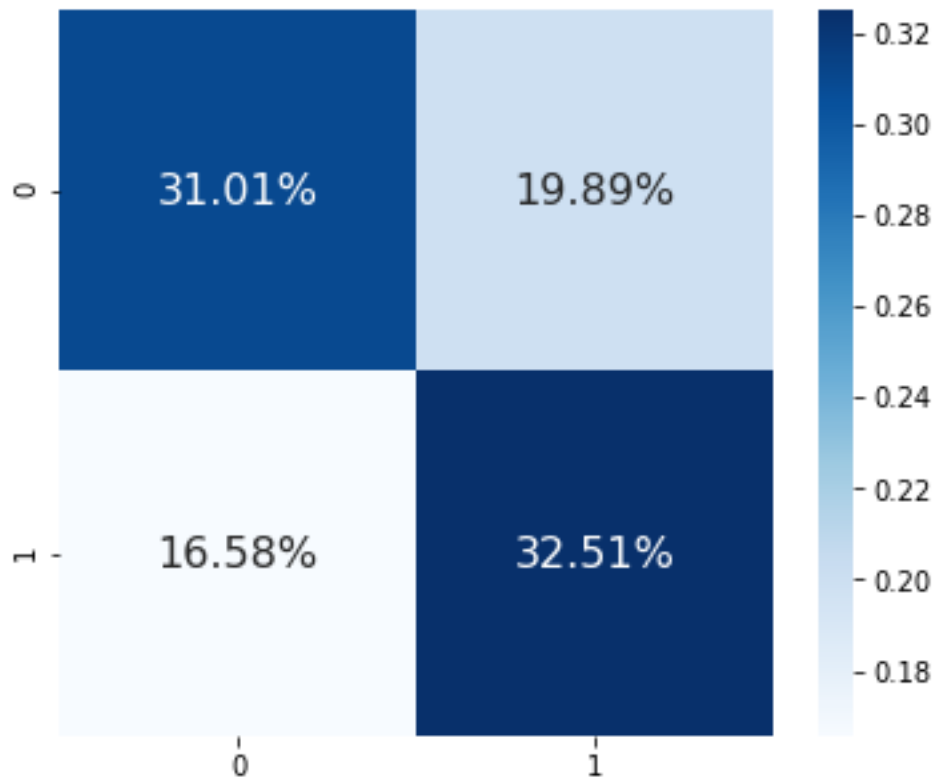
Text with highlighted words

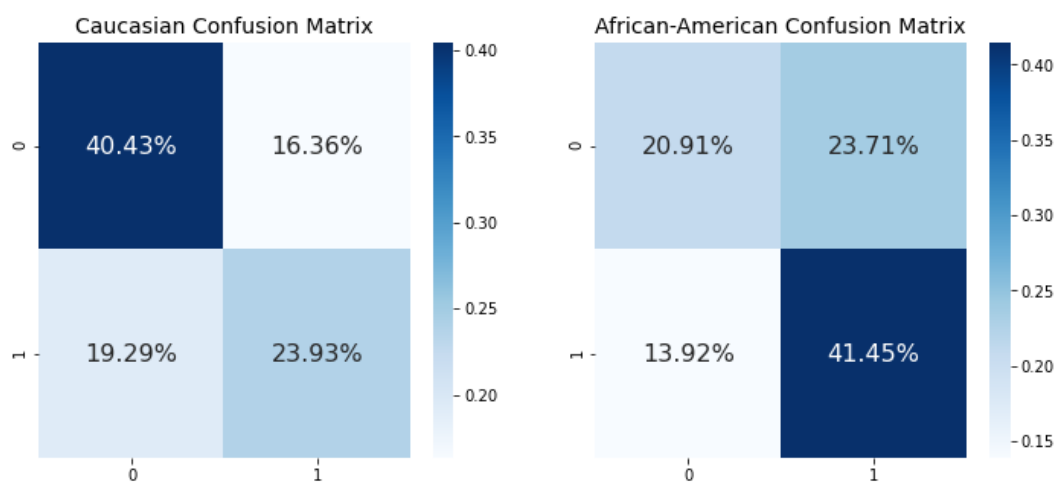
nasty disgusting gross stuff



Chapter 7: Anchor and Counterfactual Explanations

	Feature	x	PN	PN-x	PP
0	age	23	23.000000	0.000000	0.000000
4	priors_count	2	2.000000	0.000000	0.000000
5	sex_Female	0	0.397589	0.397589	0.000000
6	sex_Male	1	1.000000	0.000000	0.000000
7	race_African-American	1	0.457206	-0.542794	0.000000
16	c_charge_degree_(F7)	1	1.000000	0.000000	0.000000

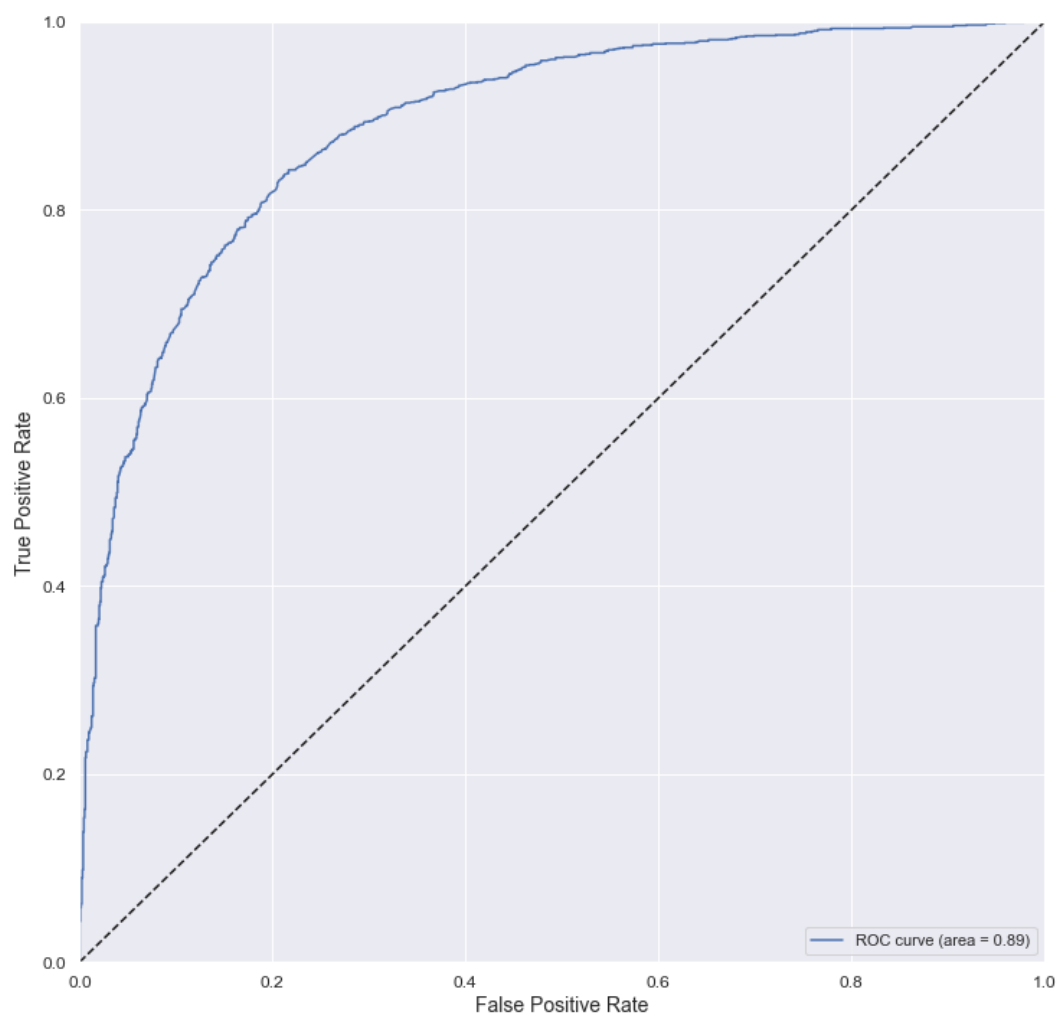




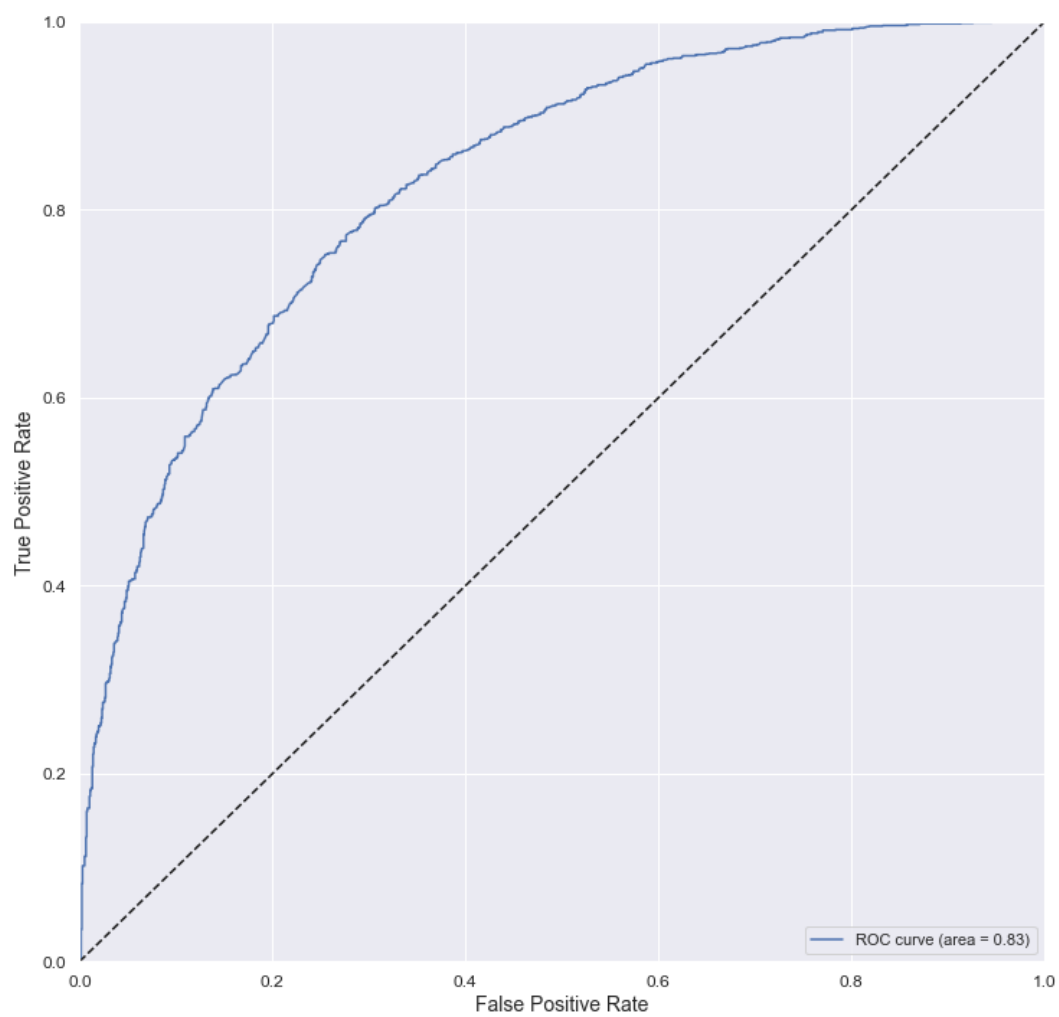
African-American FPR: 53.1%

Caucasian FPR: 28.8%

Ratio FPRs: 1.84 x

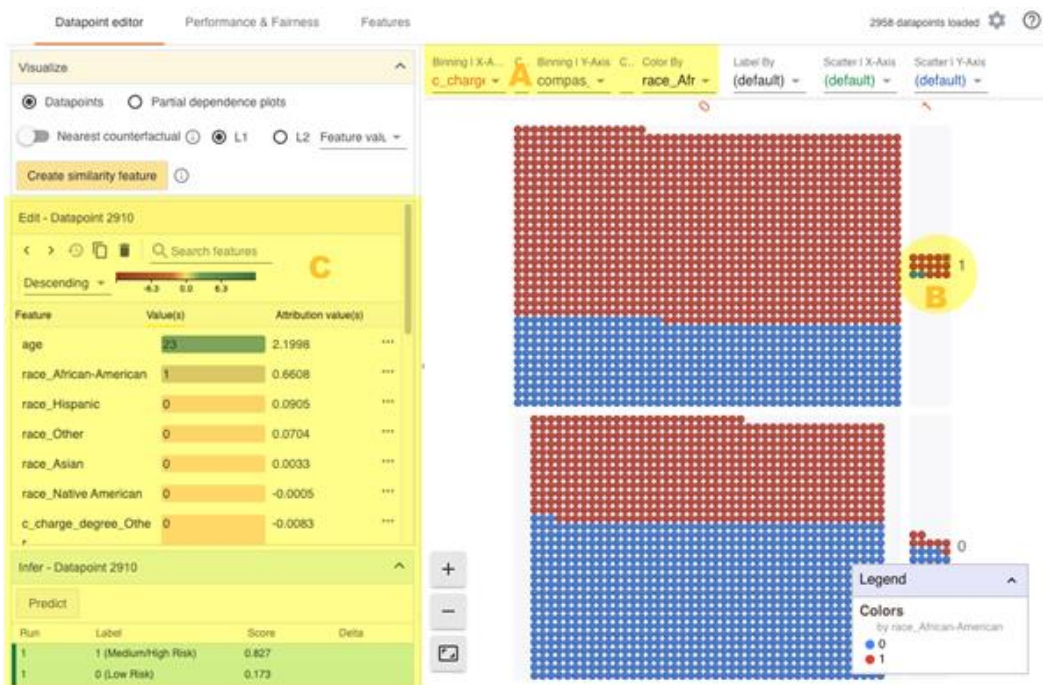


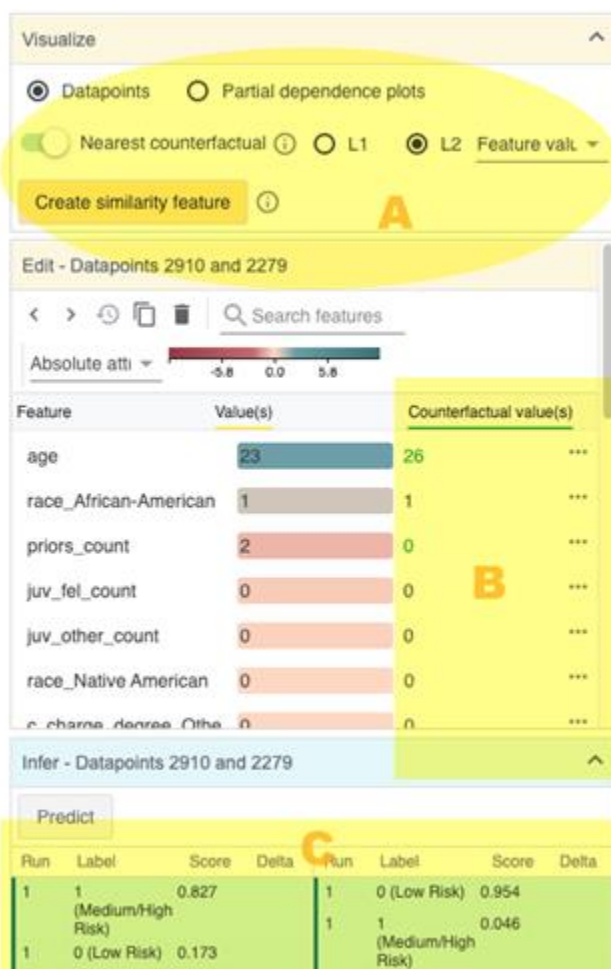
Accuracy_train:	0.8790	Accuracy_test:	0.8087	
Precision_test:	0.8277	Recall_test:	0.8110	
ROC-AUC_test:	0.8927	F1_test:	0.8193	MCC_test:
	0.6162			



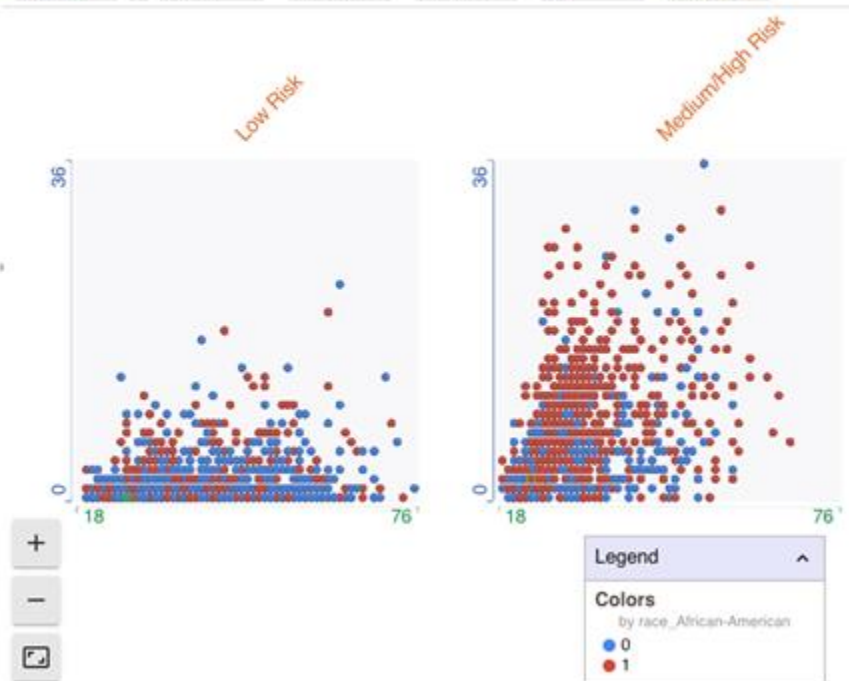
Accuracy_train:	0.7527	Accuracy_test:	0.7471	
Precision_test:	0.7653	Recall_test:	0.7604	
ROC-AUC_test:	0.8320	F1_test:	0.7628	MCC_test:
0.4920				

	10127	2726	5231
y	0	0	1
y_pred	0	0	1
age	24	23	23
:	:	:	:
priors_count	2	2	2
sex_Female	0	0	0
sex_Male	1	1	1
race_African-American	0	0	1
race_Asian	0	0	0
race_Caucasian	1	0	0
race_Hispanic	0	1	0
:	:	:	:
c_charge_degree_(F3)	0	1	0
c_charge_degree_(F7)	0	0	1
c_charge_degree_(M1)	1	0	0
:	:	:	:





Binning | X-Axis C... Binning | Y-Axis Color By Label By Scatter | X-Axis Scatter | Y-Axis
Inference (none) race_Afri (default) age priors_cc



Edit - Datapoint 2910

< > ↺ 📄 🗑️ 🔍 Search features

Absolute attl ▾

-5.8 0.0 5.8

Feature	Value(s)	Attribu
priors_count	1	-2.4239 ***
age	23	1.5873 ***
race_African-American	1	0.8297 ***
race_Caucasian	0	0.3332 ***
juv_fel_count	0	-0.2306 ***

Infer - Datapoint 2910

Predict

Run	Label	Score	Delta
2	0 (Low Risk)	0.665	↑ 0.492183
2	1 (Medium/High Risk)	0.335	↓ -0.492183
1	1 (Medium/High Risk)	0.827	
1	0 (Low Risk)	0.173	

Edit - Datapoint 2910

< > ↺ 📄 🗑️ 🔍 Search features

Absolute attl ▾

-5.8 0.0 5.8

Feature	Value(s)	Attribu
priors_count	2	-1.1525 ***
age	25	1.0107 ***
race_African-American	1	0.8278 ***

Infer - Datapoint 2910

Predict

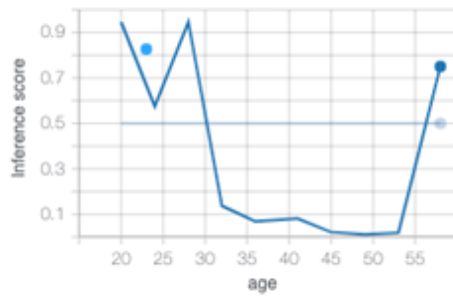
Run	Label	Score	Delta
3	0 (Low Risk)	0.508	↓ -0.157111
3	1 (Medium/High Risk)	0.492	↑ 0.157111
2	0 (Low Risk)	0.665	↑ 0.492183
2	1 (Medium/High Risk)	0.335	↓ -0.492183
1	1 (Medium/High Risk)	0.827	
1	0 (Low Risk)	0.173	

Partial Dependence Plots ⓘ

Sort by variation

☐ Global partial dependence plots

▼ age



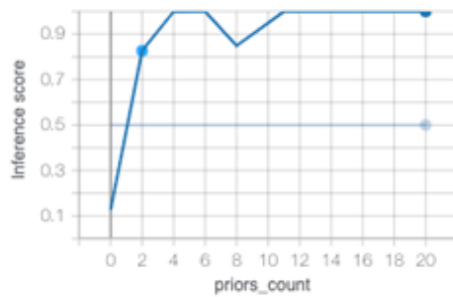
Set range of values to visualize

20

-

58

▼ priors_count

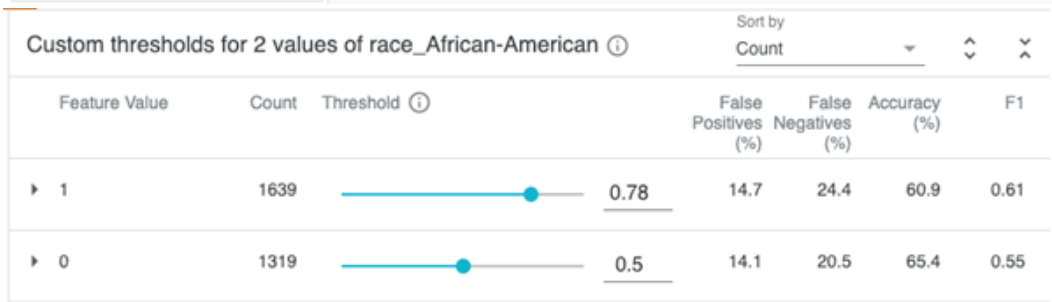
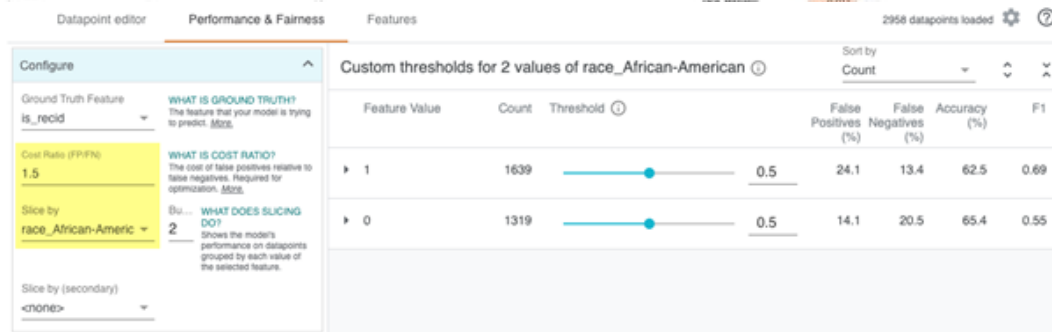
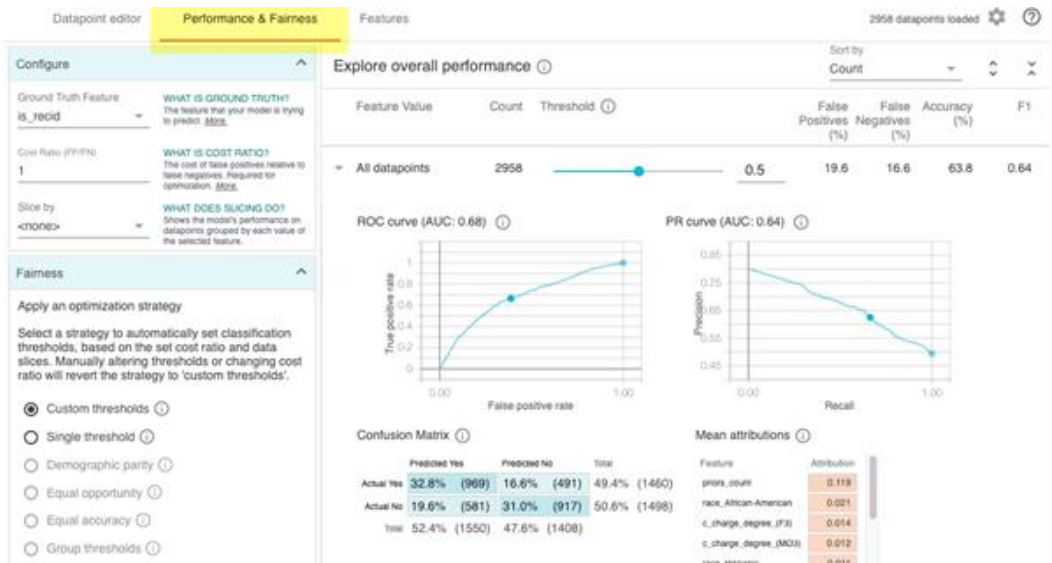


Set range of values to visualize

0

-

20



Chapter 8: Visualizing Convolutional Neural Networks



Apple Golden



Banana



Apple Granny Smith



Clementine



Apple Red



Grapefruit Pink



Avocado



Mango Red



Nectarine



Onion Red



Onion White



Orange



Peach



Pear



Pomegranate



Tomato



Apple Golden



Apple Granny Smith



Apple Red



Avocado



Banana



Clementine



Grapefruit Pink



Mango Red



Nectarine



Onion Red



Onion White



Orange



Peach



Pear

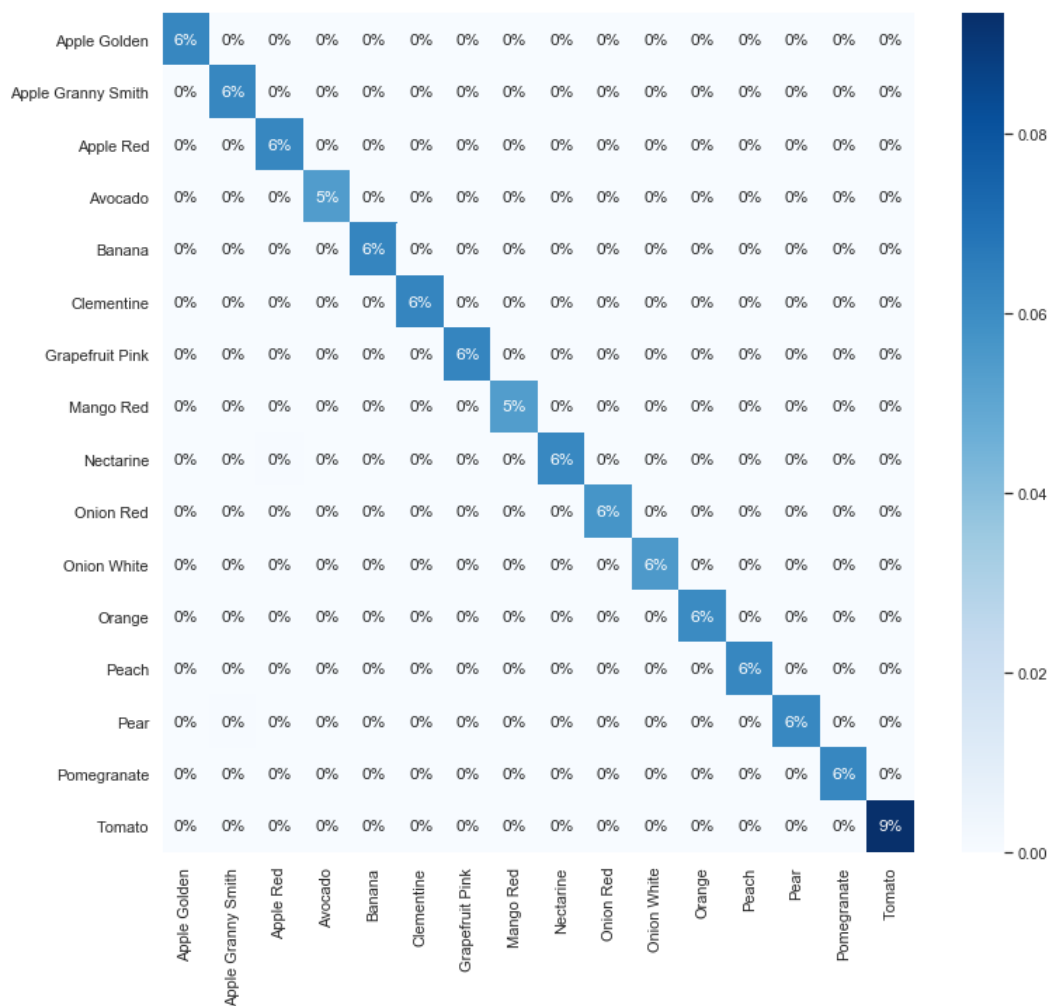


Pomegranate

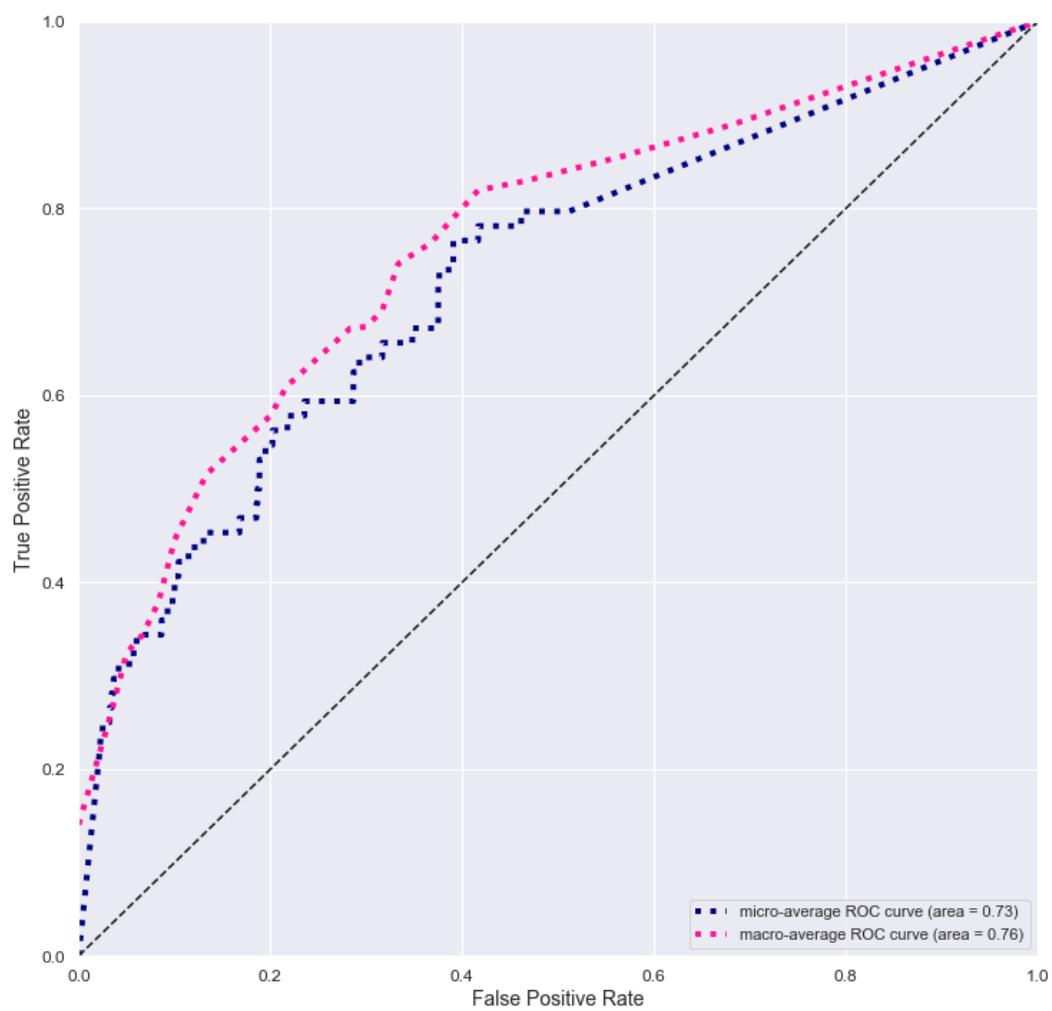


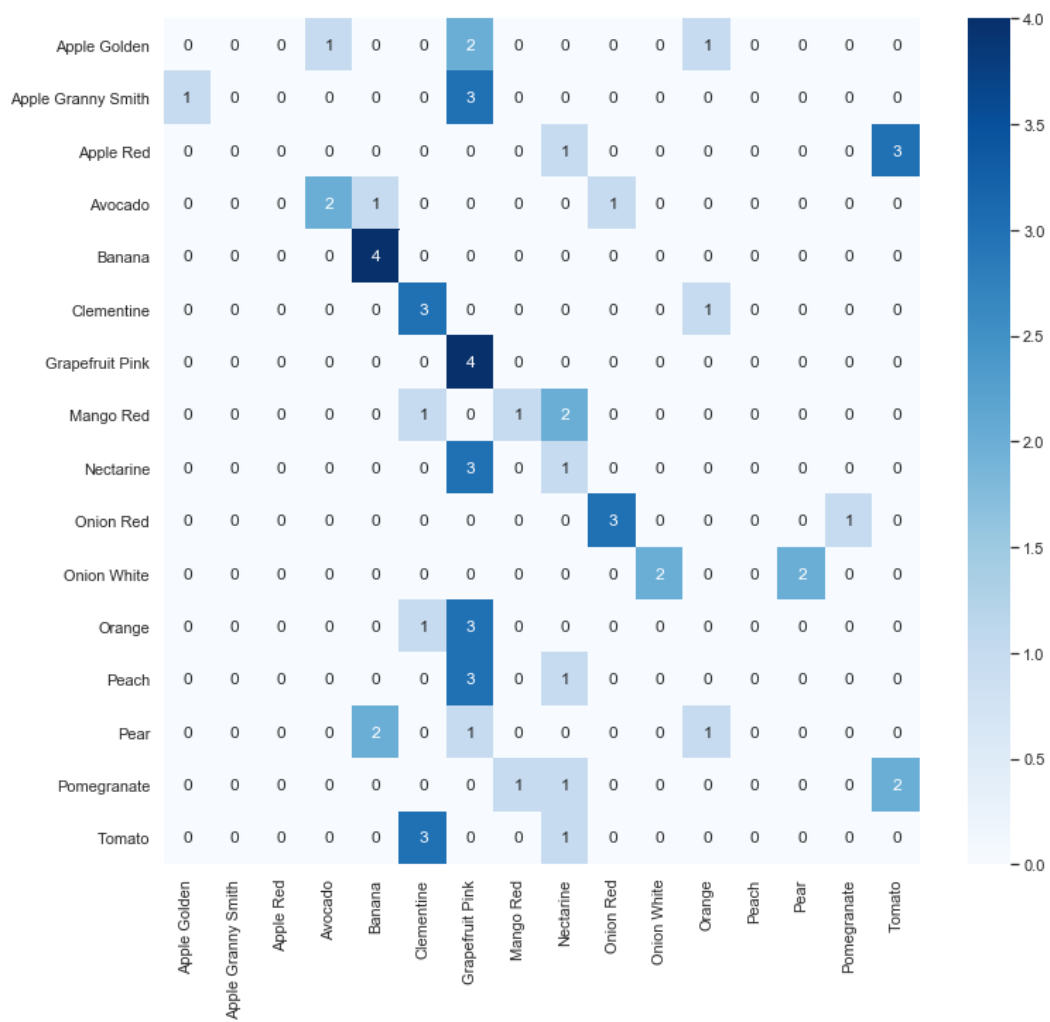
Tomato





	precision	recall	f1-score	support
Apple Golden	1.000	1.000	1.000	164
Apple Granny Smith	0.994	1.000	0.997	164
Apple Red	0.994	1.000	0.997	164
Avocado	1.000	1.000	1.000	143
Banana	1.000	1.000	1.000	166
Clementine	1.000	1.000	1.000	166
Grapefruit Pink	1.000	1.000	1.000	166
Mango Red	1.000	1.000	1.000	142
Nectarine	1.000	0.994	0.997	164
Onion Red	1.000	1.000	1.000	150
Onion White	1.000	1.000	1.000	146
Orange	1.000	1.000	1.000	160
Peach	1.000	1.000	1.000	164
Pear	1.000	0.994	0.997	164
Pomegranate	1.000	1.000	1.000	164
Tomato	1.000	1.000	1.000	246
accuracy			0.999	2633
macro avg	0.999	0.999	0.999	2633
weighted avg	0.999	0.999	0.999	2633

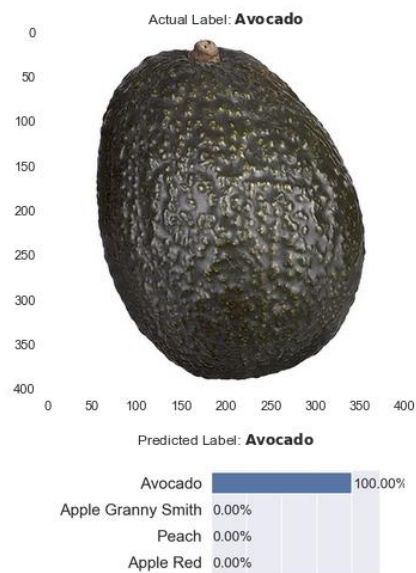
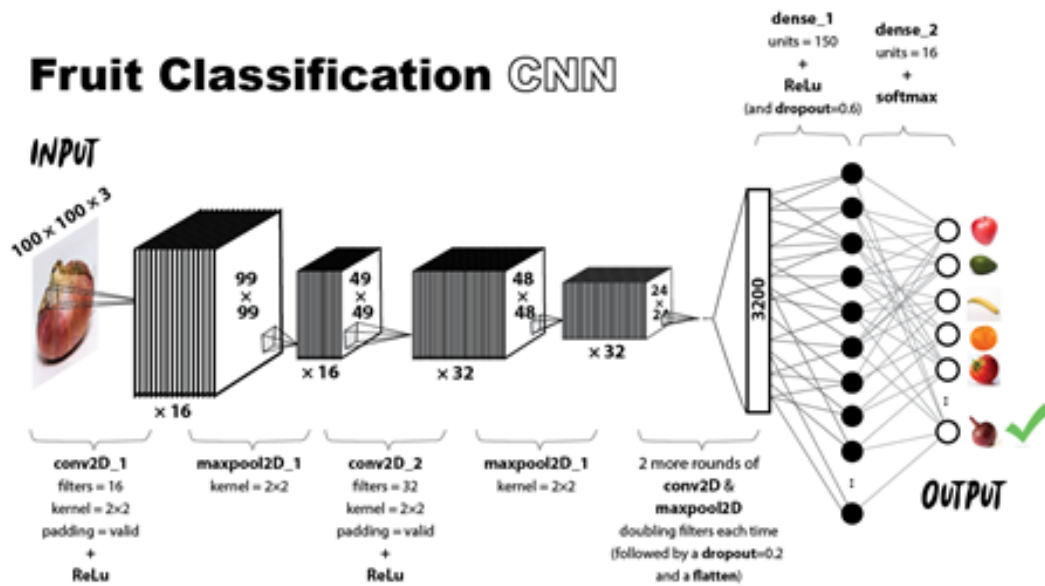




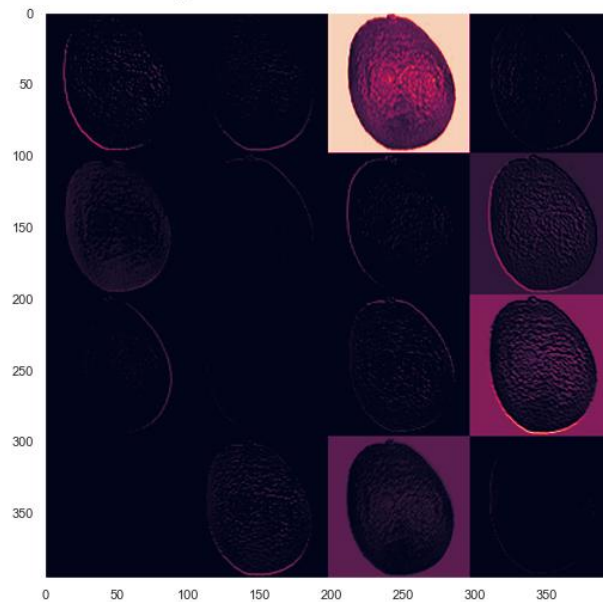
	precision	recall	f1-score	support
Apple Golden	0.000	0.000	0.000	4
Apple Granny Smith	0.000	0.000	0.000	4
Apple Red	0.000	0.000	0.000	4
Avocado	0.667	0.500	0.571	4
Banana	0.571	1.000	0.727	4
Clementine	0.375	0.750	0.500	4
Grapefruit Pink	0.211	1.000	0.348	4
Mango Red	0.500	0.250	0.333	4
Nectarine	0.143	0.250	0.182	4
Onion Red	0.750	0.750	0.750	4
Onion White	1.000	0.500	0.667	4
Orange	0.000	0.000	0.000	4
Peach	0.000	0.000	0.000	4
Pear	0.000	0.000	0.000	4
Pomegranate	0.000	0.000	0.000	4
Tomato	0.000	0.000	0.000	4
accuracy			0.312	64
macro avg	0.264	0.312	0.255	64
weighted avg	0.264	0.312	0.255	64

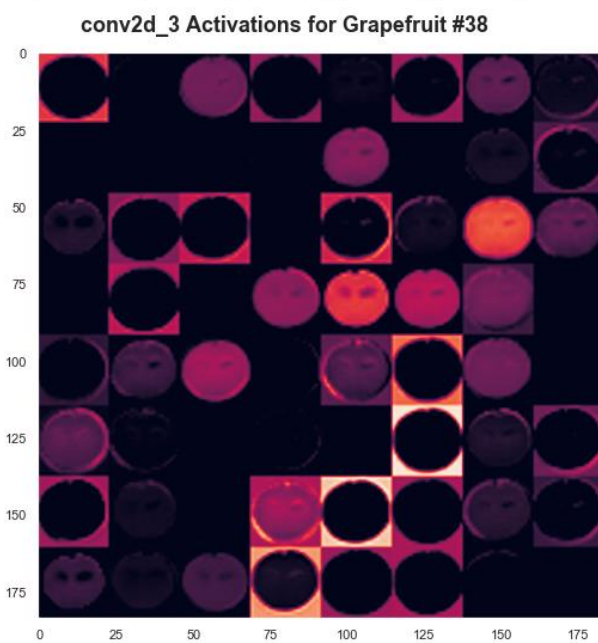
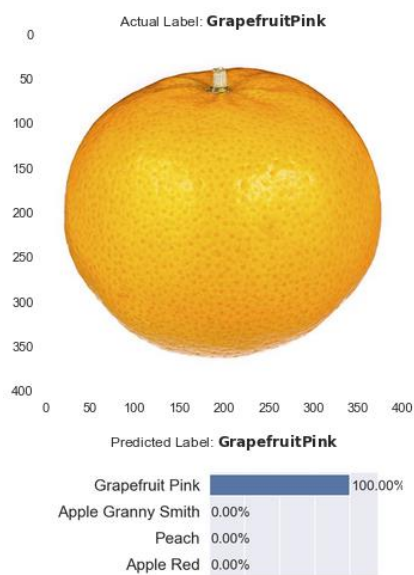
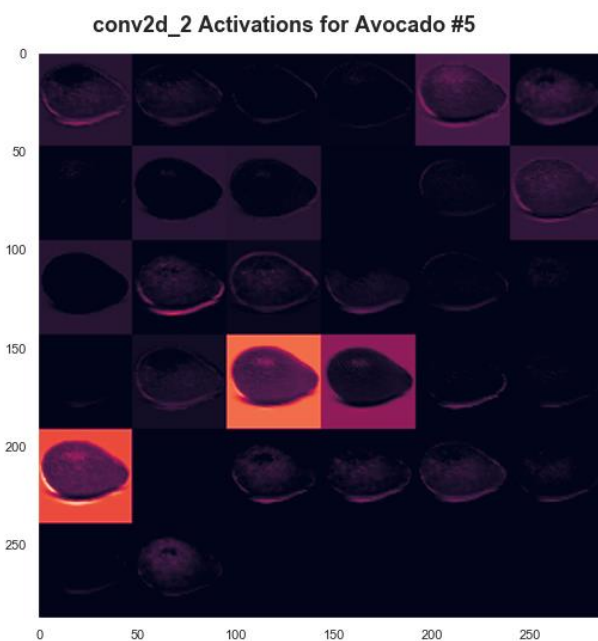
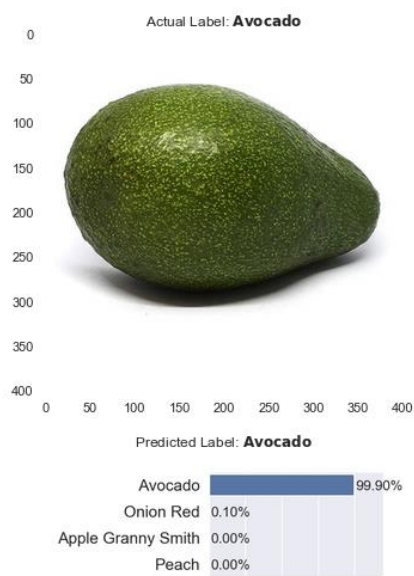
	y_true	y_pred	Grapefruit Pink	Clementine	Banana	Nectarine	Tomato	Onion Red	Avocado	Orange
0	Pear	Banana			100.0					
1	Pear	Banana			100.0					
2	Pear	Orange	5.2							94.8
3	Pear	Grapefruit Pink	96.5		3.5					
4	Avocado	Onion Red						100.0		
5	Avocado	Avocado						0.1	99.9	
6	Avocado	Banana			100.0					
7	Avocado	Avocado							100.0	
8	Pomegranate	Tomato					100.0			
:	:	:	:	:	:	:	:	:	:	:
16	Apple Golden	Avocado							99.8	0.2
17	Apple Golden	Grapefruit Pink	100.0							
18	Apple Golden	Orange								100.0
19	Apple Golden	Grapefruit Pink	100.0							
20	Nectarine	Grapefruit Pink	100.0							
21	Nectarine	Nectarine				100.0				
22	Nectarine	Grapefruit Pink	100.0							
23	Nectarine	Grapefruit Pink	100.0							
24	Clementine	Clementine		100.0						
25	Clementine	Clementine	1.9	98.1						
26	Clementine	Clementine		100.0						
27	Clementine	Orange								100.0
28	Onion White	Pear								
29	Onion White	Pear								
30	Onion White	Onion White								
31	Onion White	Onion White	8.2		0.6					
32	Apple Granny Smith	Grapefruit Pink	100.0							
33	Apple Granny Smith	Apple Golden								
34	Apple Granny Smith	Grapefruit Pink	100.0							
35	Apple Granny Smith	Grapefruit Pink	92.7		7.2				0.1	
36	Grapefruit Pink	Grapefruit Pink	100.0							
37	Grapefruit Pink	Grapefruit Pink	100.0							
:	:	:	:	:	:	:	:	:	:	:
56	Orange	Clementine	38.7	61.3						
57	Orange	Grapefruit Pink	100.0							
58	Orange	Grapefruit Pink	51.8	48.2						
59	Orange	Grapefruit Pink	91.6	8.4						
60	Peach	Grapefruit Pink	100.0							
61	Peach	Grapefruit Pink	100.0							
62	Peach	Nectarine		0.1		99.9				
63	Peach	Grapefruit Pink	87.7				11.9			

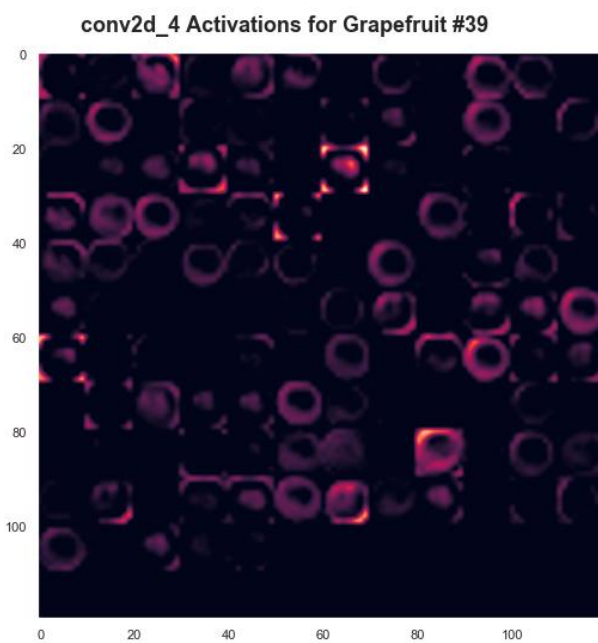
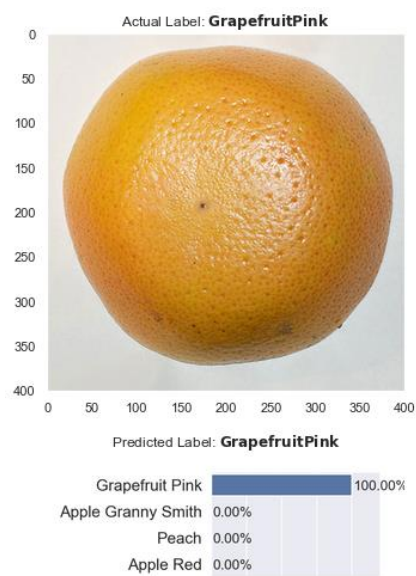
Fruit Classification CNN



conv2d_1 Activations for Avocado #7

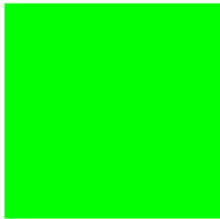






conv2d_1 Layer

Filter #0



Filter #1

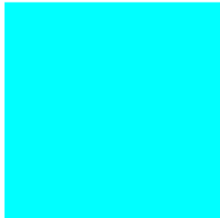


Filter #2



Filter #3

Filter #4



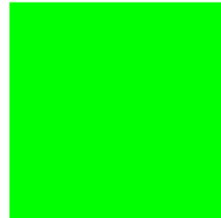
Filter #5



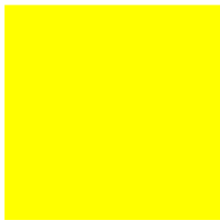
Filter #6



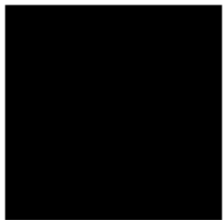
Filter #7



Filter #8



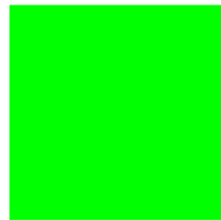
Filter #9



Filter #10



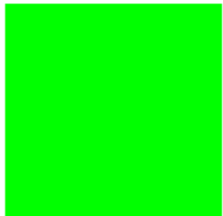
Filter #11



Filter #12



Filter #13



Filter #14

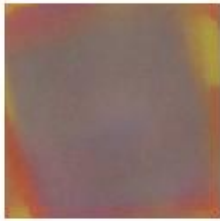


Filter #15

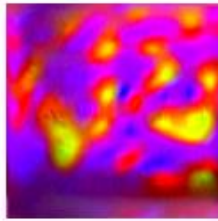


conv2d_4 Layer

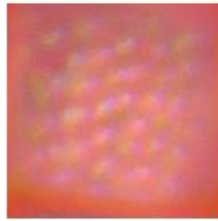
Filter #92



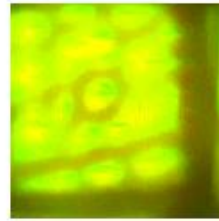
Filter #104



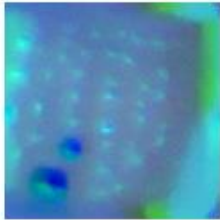
Filter #26



Filter #38



Filter #47



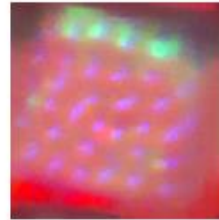
Filter #25



Filter #95



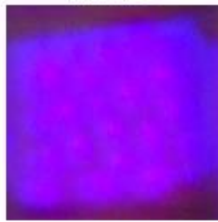
Filter #112



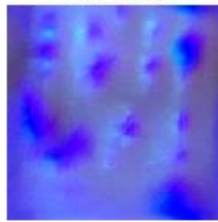
Filter #64



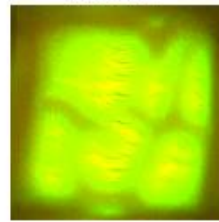
Filter #33



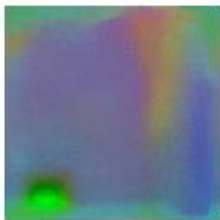
Filter #99



Filter #13



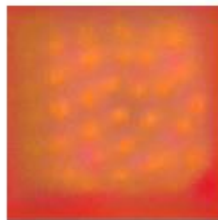
Filter #24



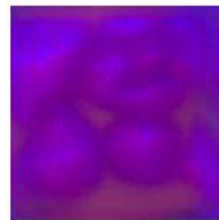
Filter #75

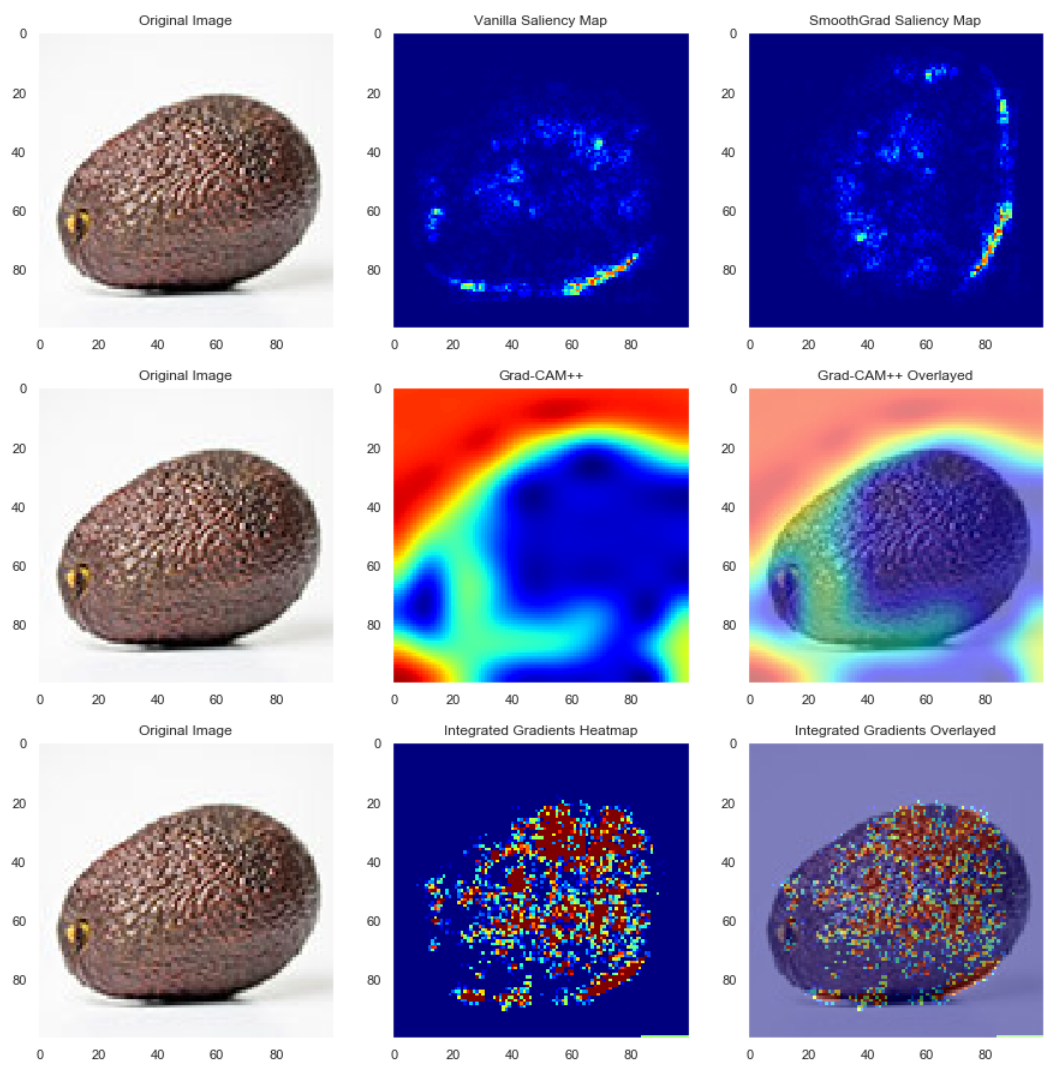


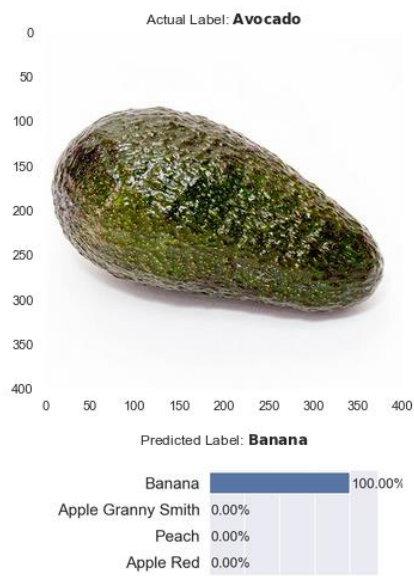
Filter #101



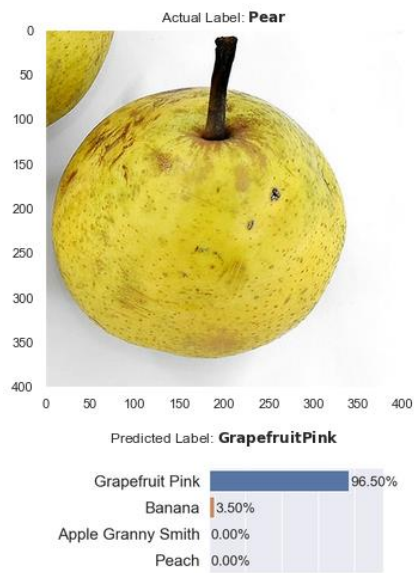
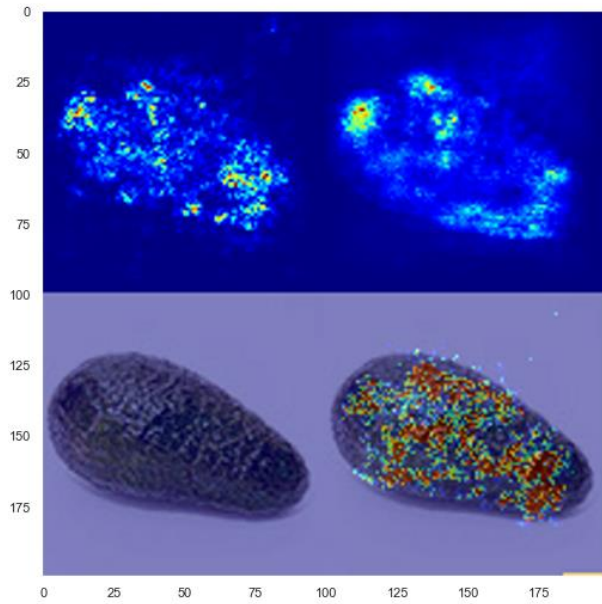
Filter #22



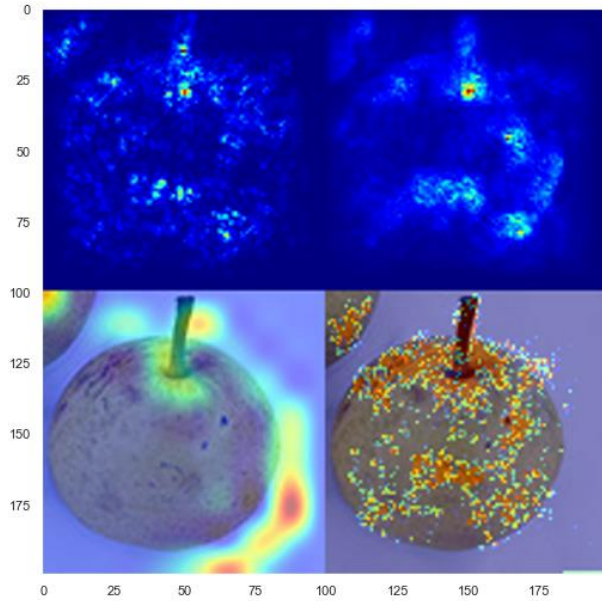


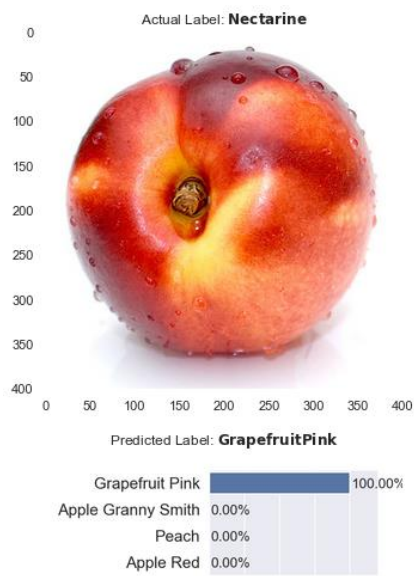


Gradient-Based Attributions for Misclassification #2

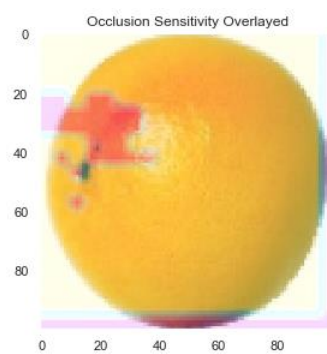
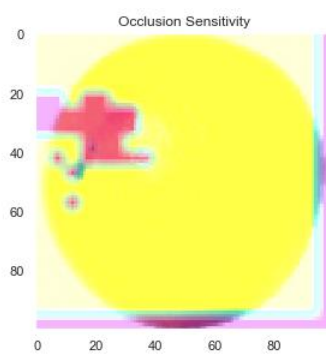
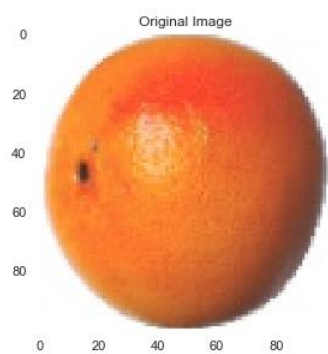
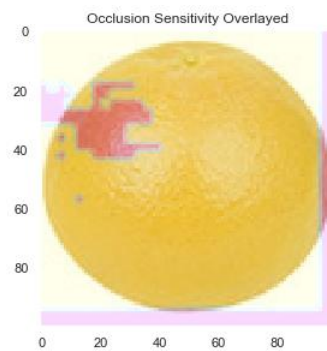
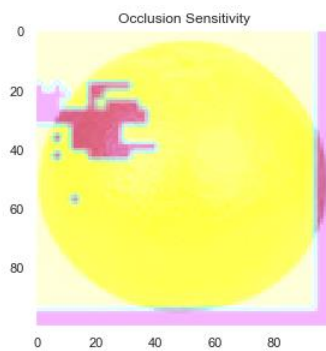
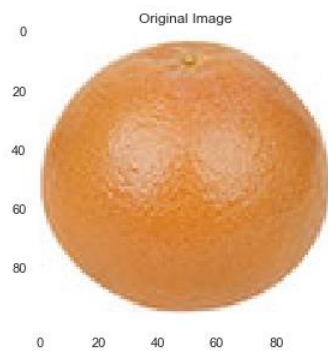
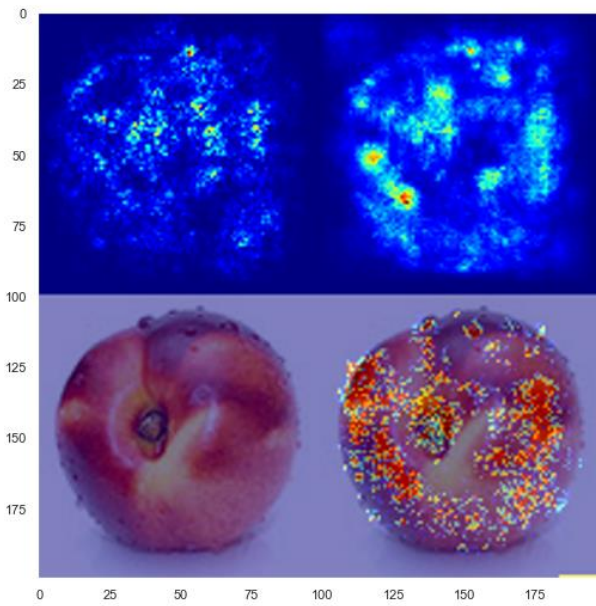


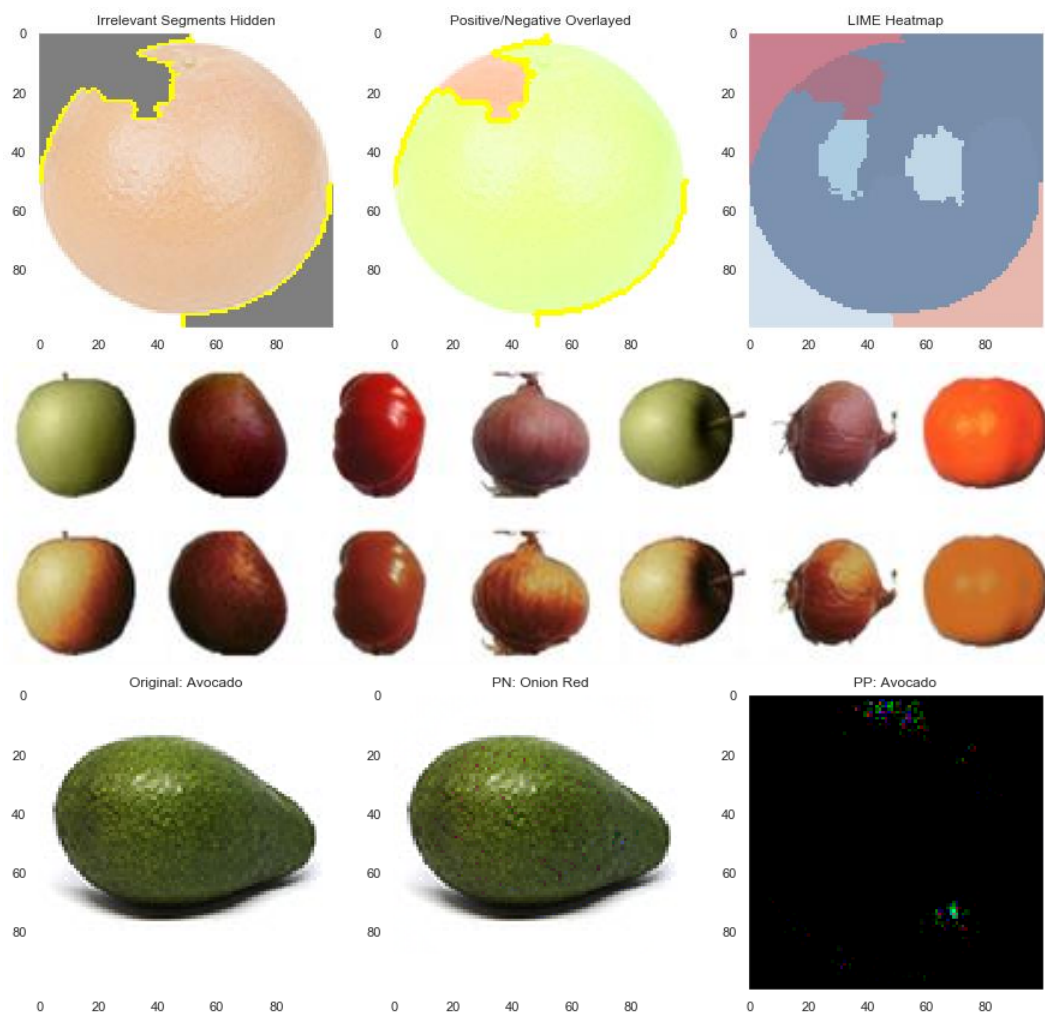
Gradient-Based Attributions for Misclassification #3

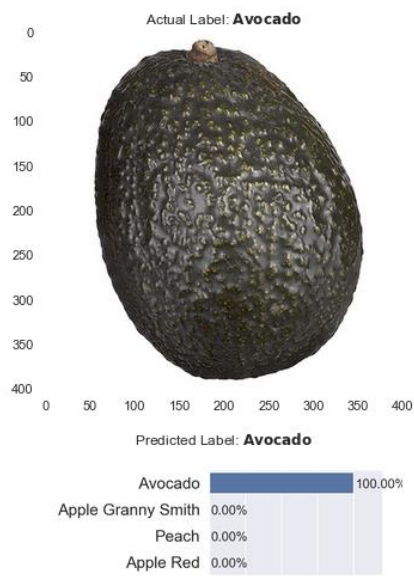




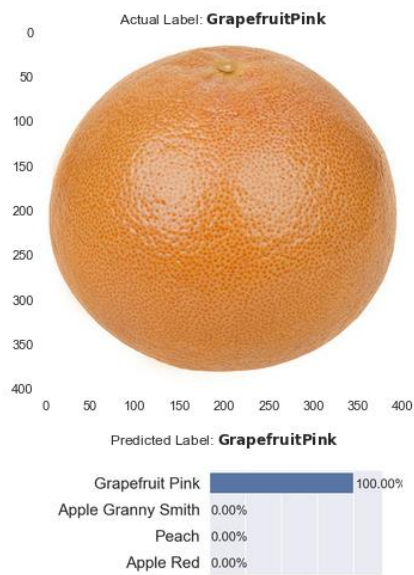
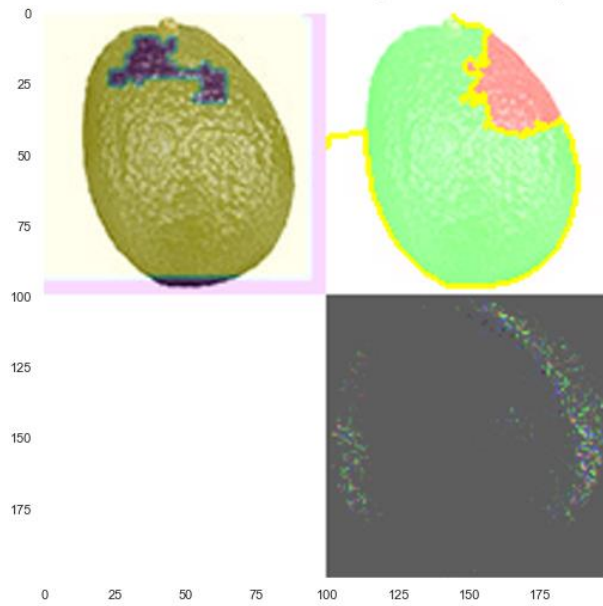
Gradient-Based Attributions for Misclassification #7



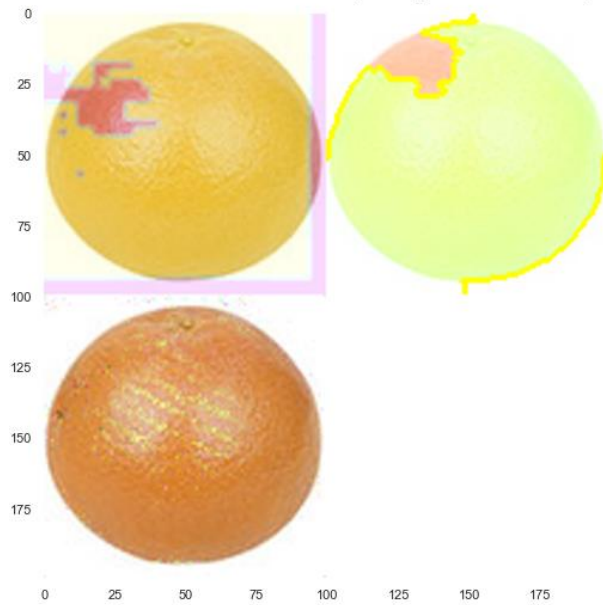


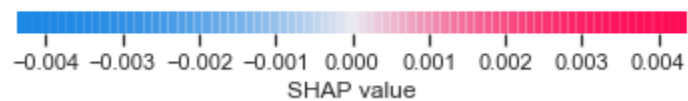


Perturbation-Based Attributions #2 (PN:, PP:Avocado)

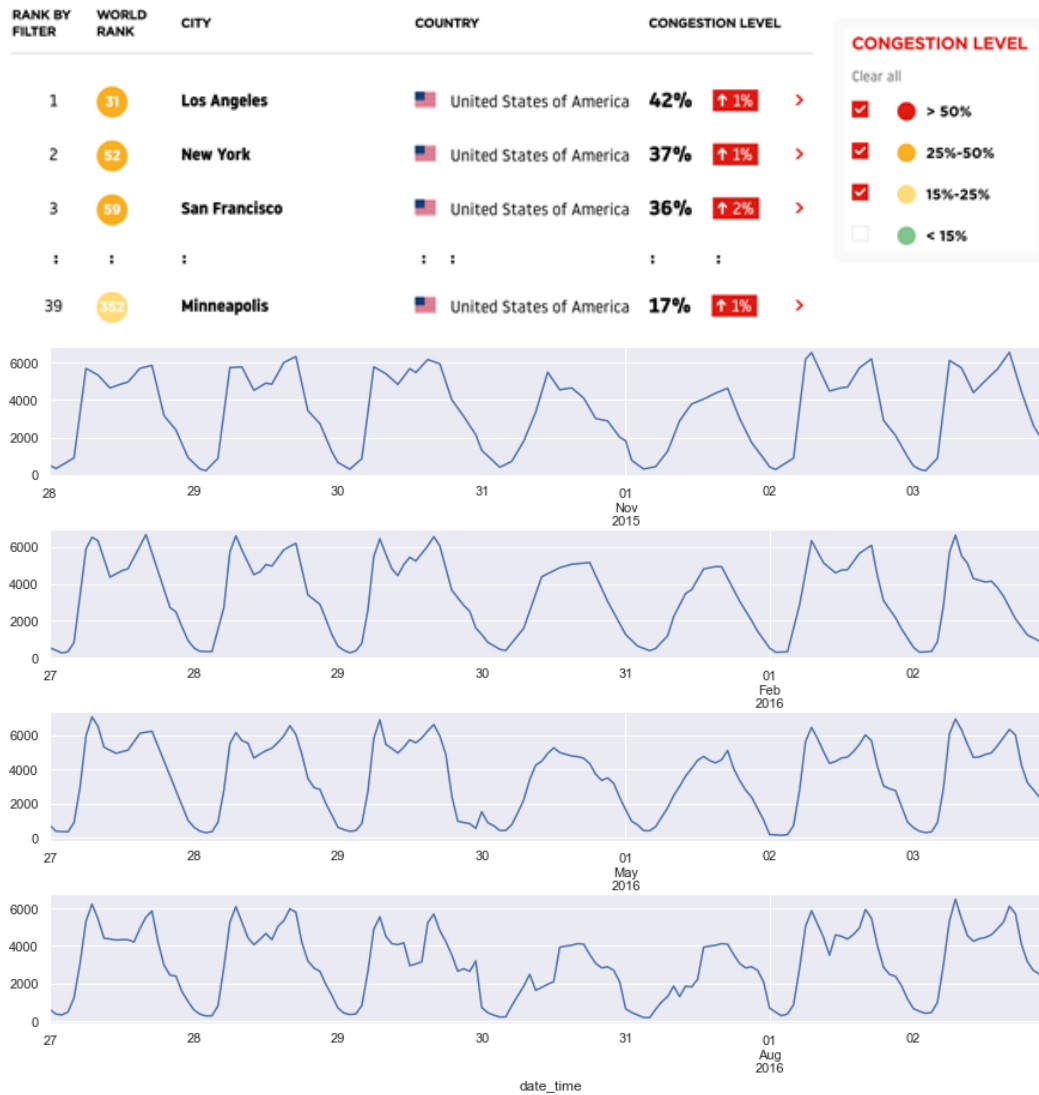


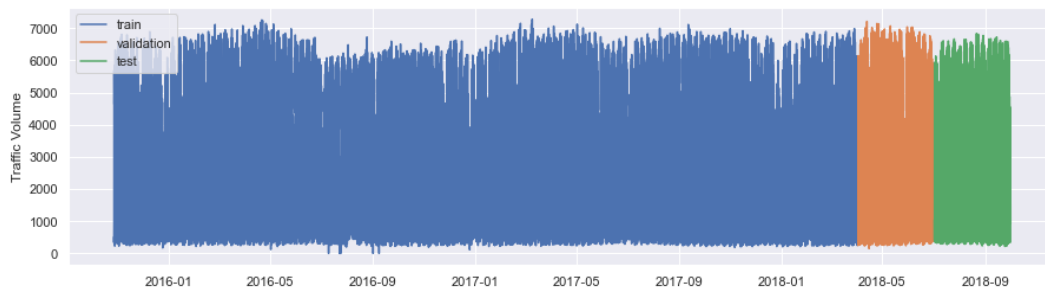
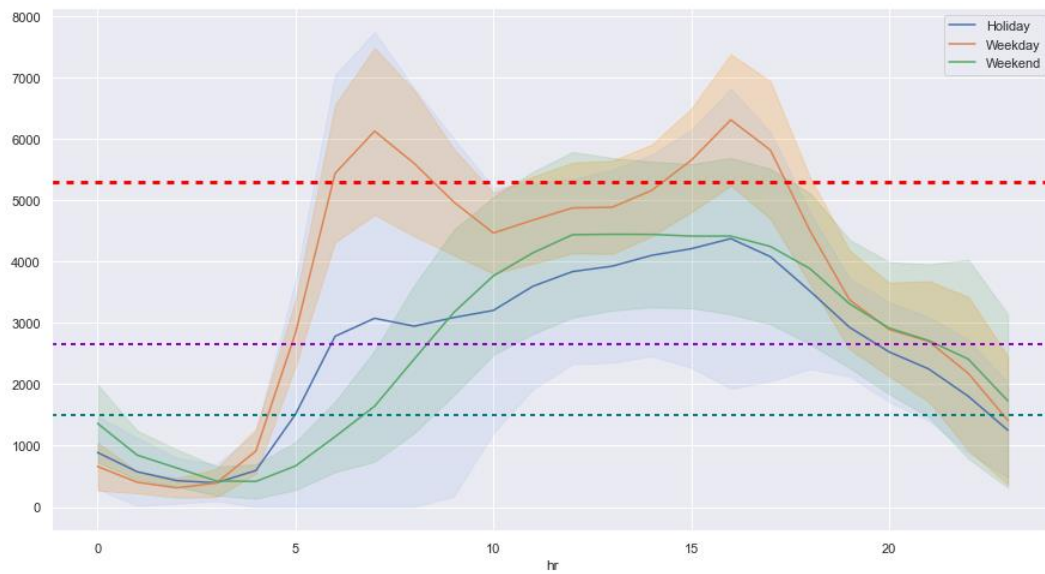
Perturbation-Based Attributions #3 (PN:Apple Golden, PP:)



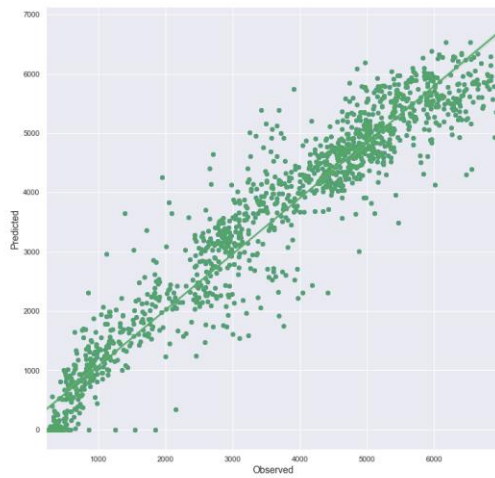


Chapter 9: Interpretation Methods for Multivariate Forecasting and Sensitivity Analysis



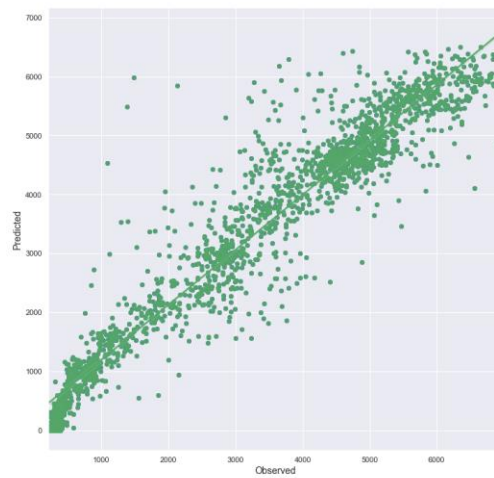


Traffic_Bidirectional_LSTM_672



RMSE_train: 537.1589
RMSE_test: 542.6649
r2: 0.9263

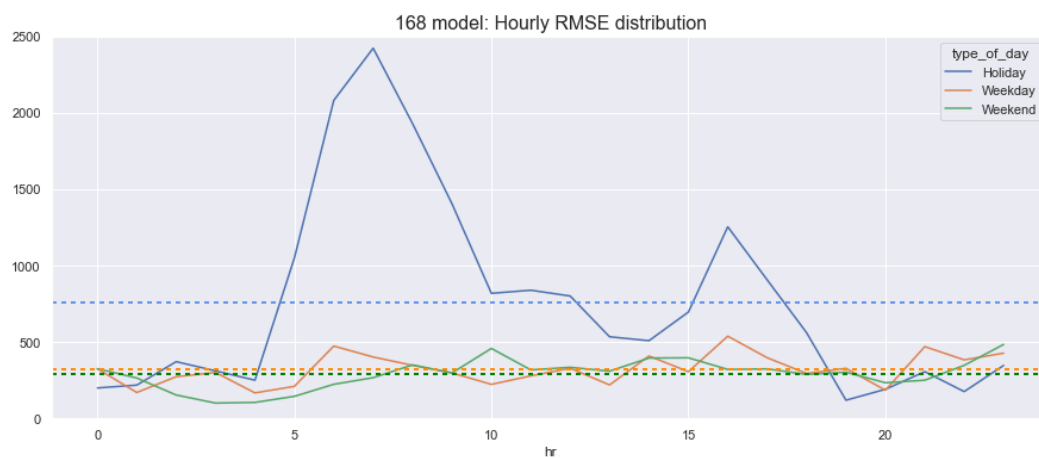
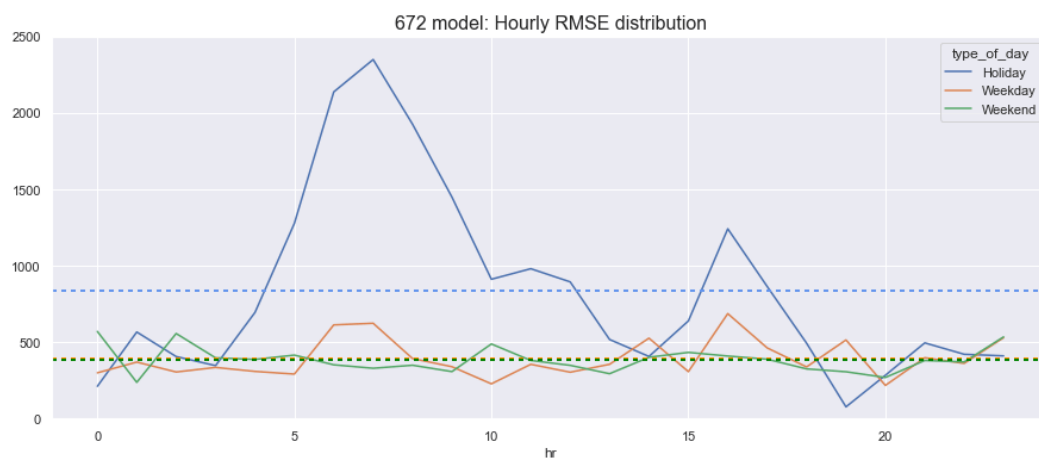
Traffic_LSTM_168



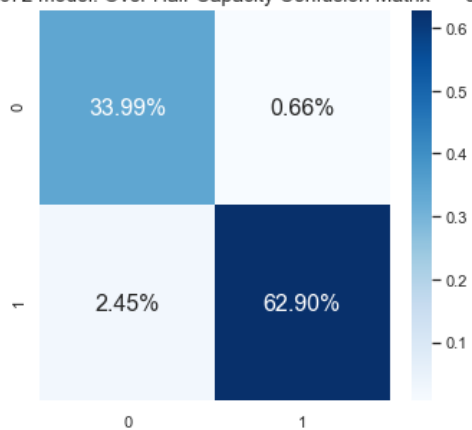
RMSE_train: 473.4274
RMSE_test: 561.8984
r2: 0.9187



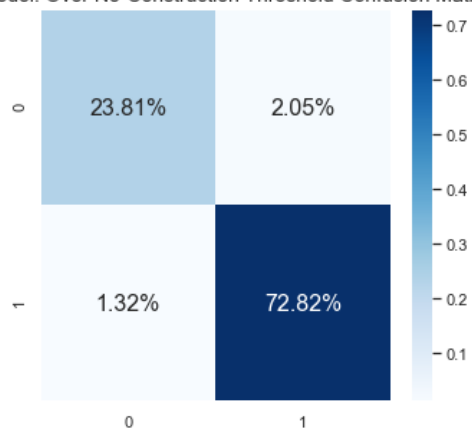




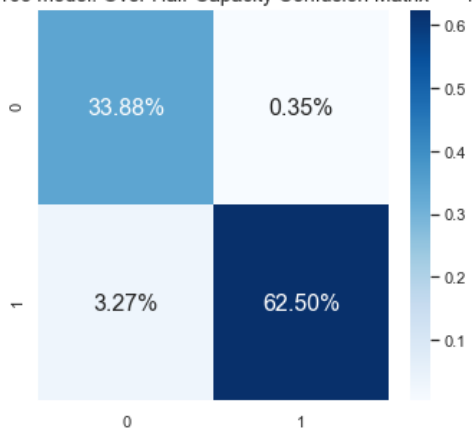
672 model: Over Half-Capacity Confusion Matrix



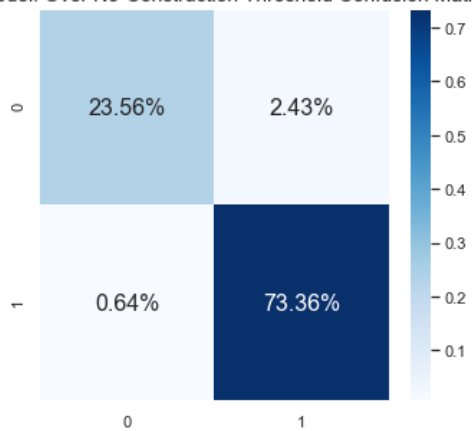
672 model: Over No-Construction Threshold Confusion Matrix



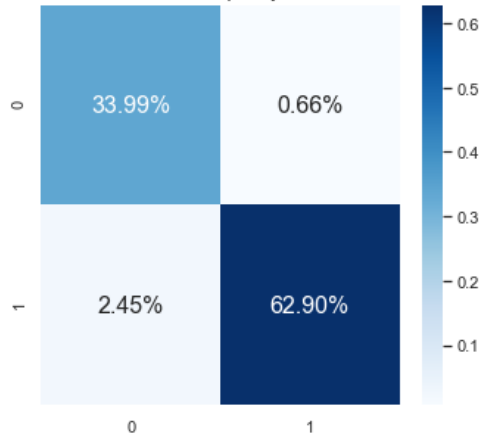
168 model: Over Half-Capacity Confusion Matrix



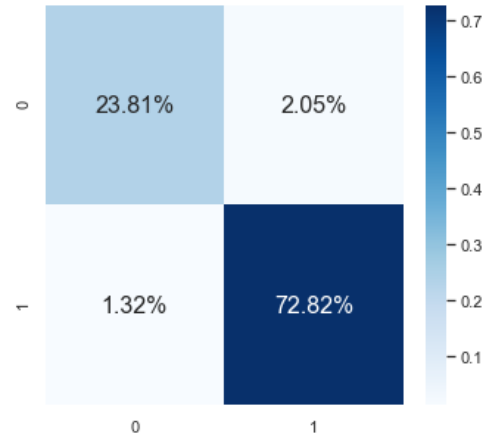
168 model: Over No-Construction Threshold Confusion Matrix



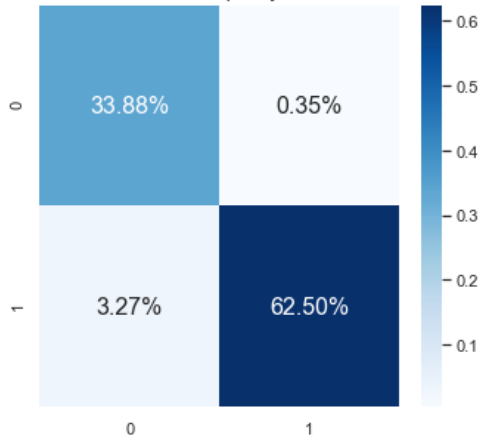
672 model: Over Half-Capacity Confusion Matrix



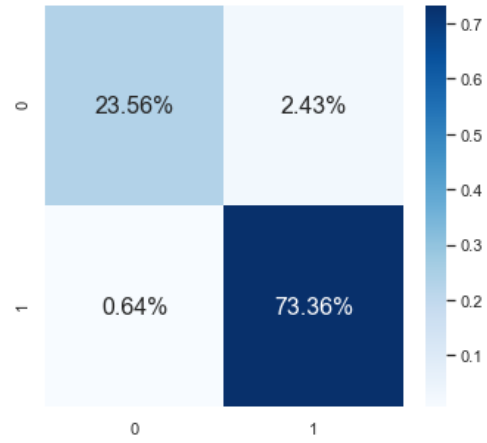
672 model: Over No-Construction Threshold Confusion Matrix

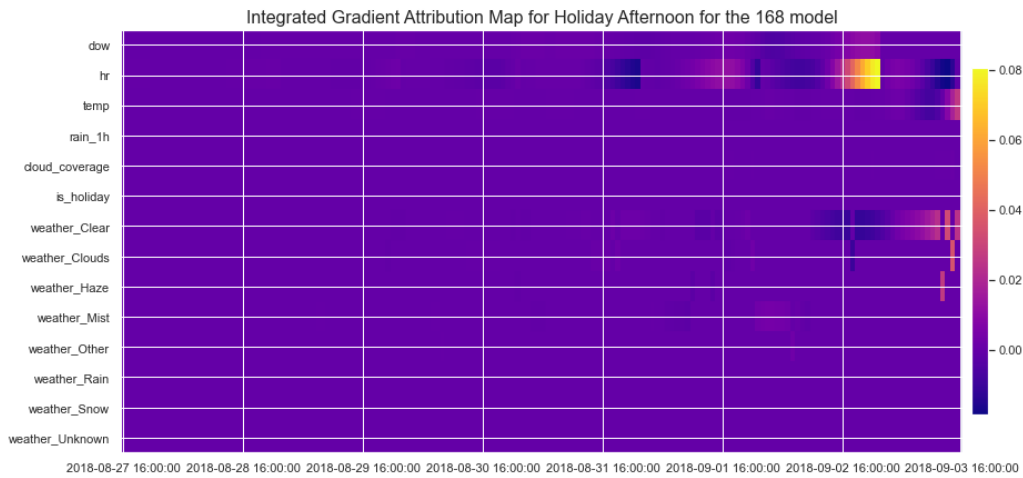
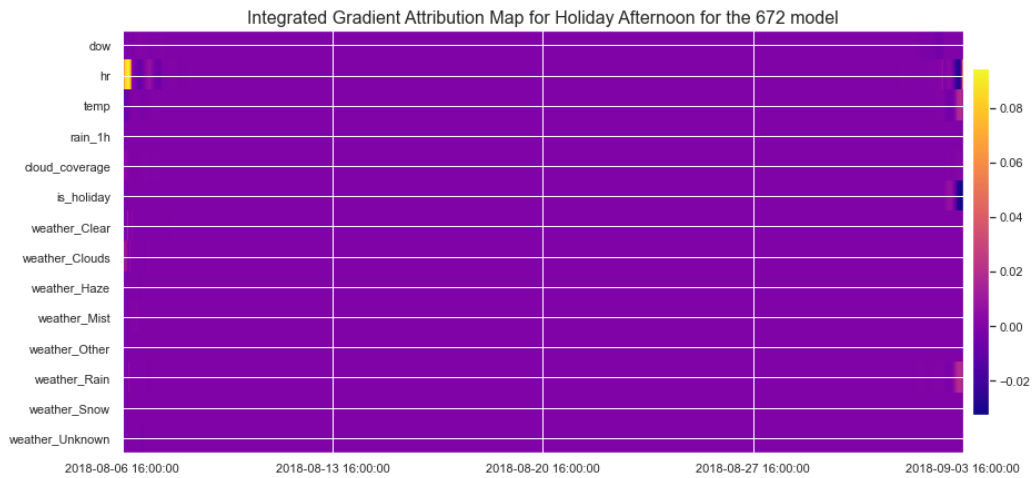


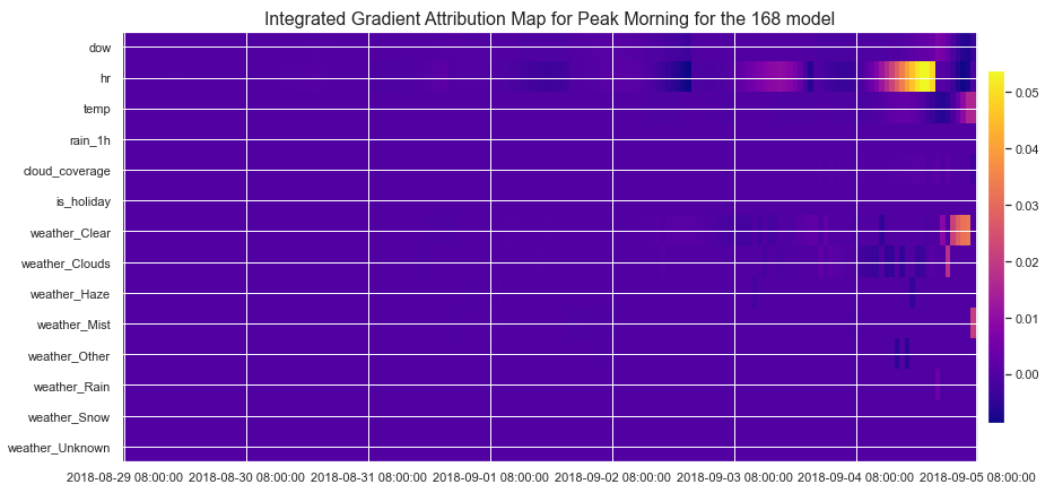
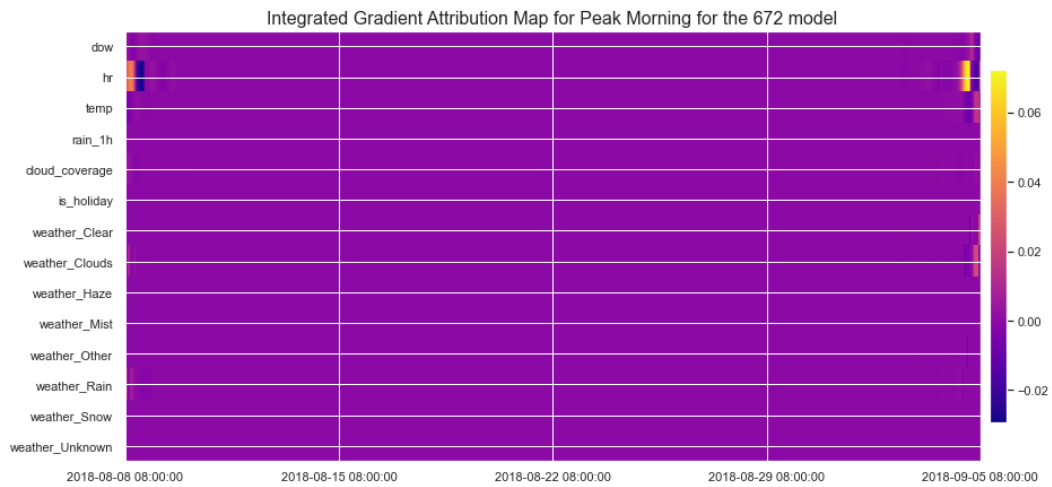
168 model: Over Half-Capacity Confusion Matrix

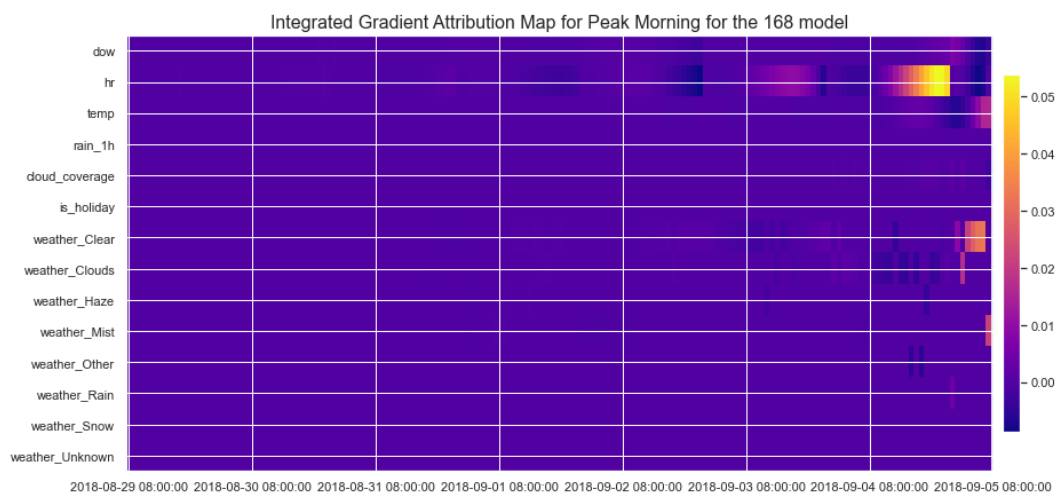
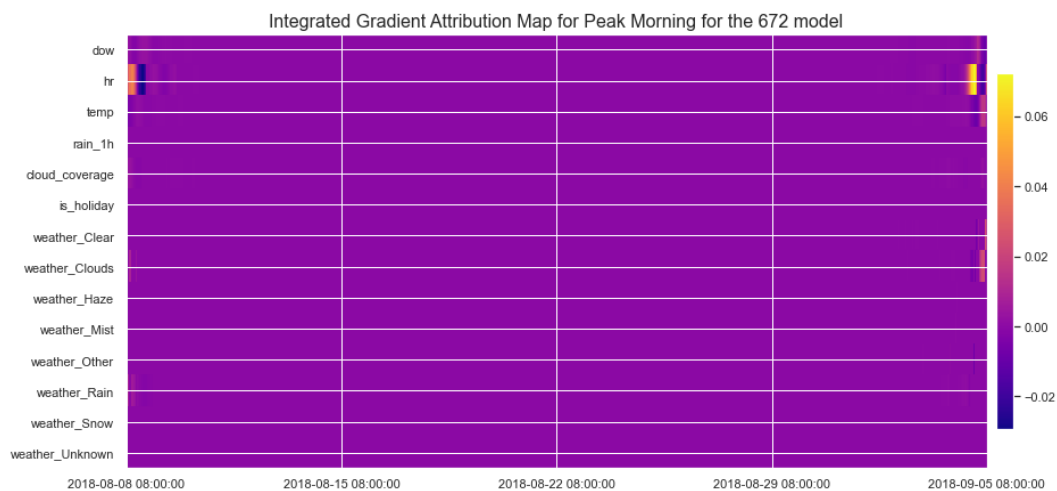


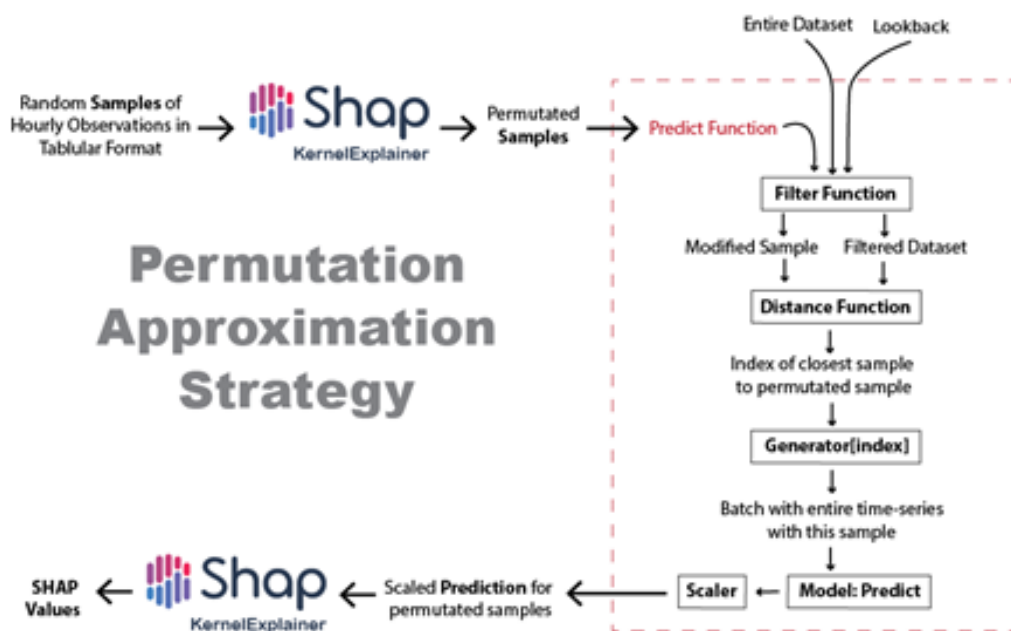
168 model: Over No-Construction Threshold Confusion Matrix

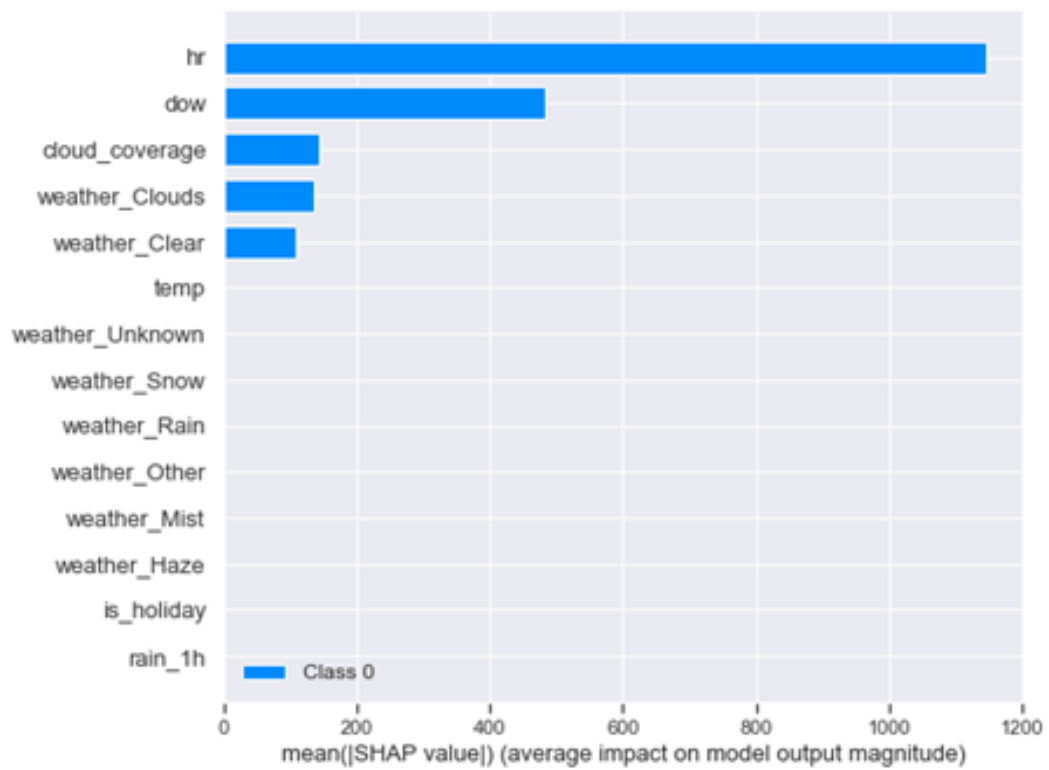




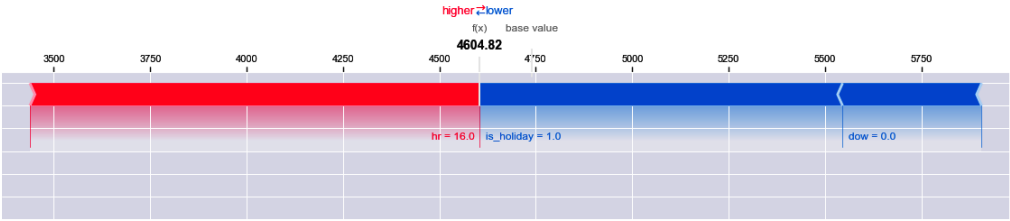




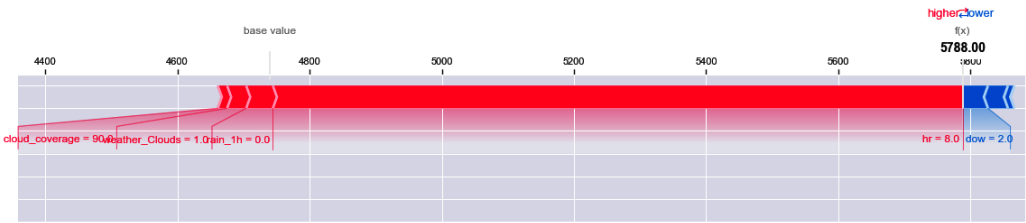




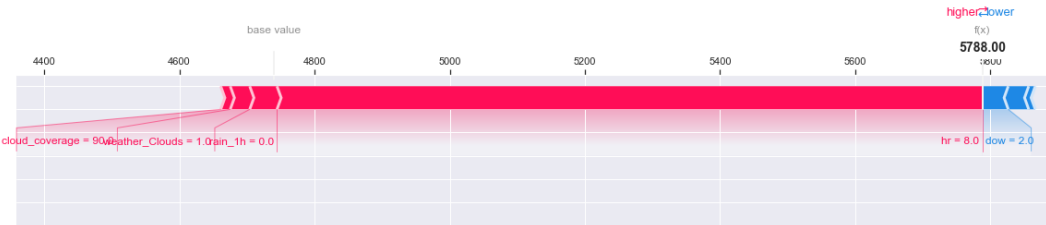
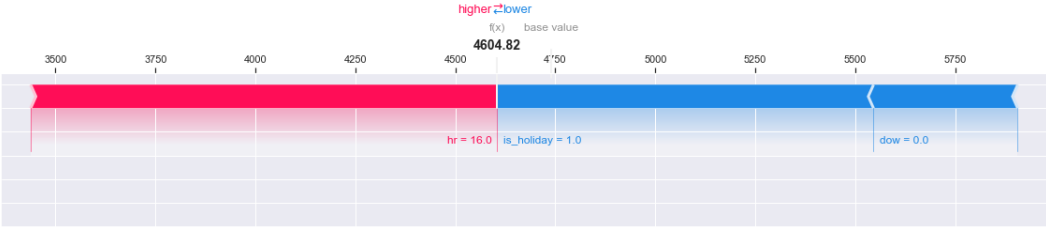
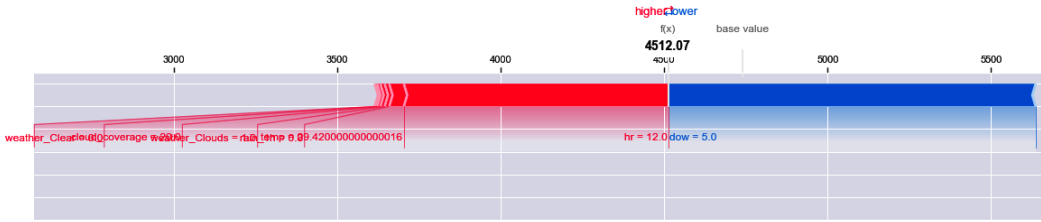
Holiday Afternoon



Peak Morning



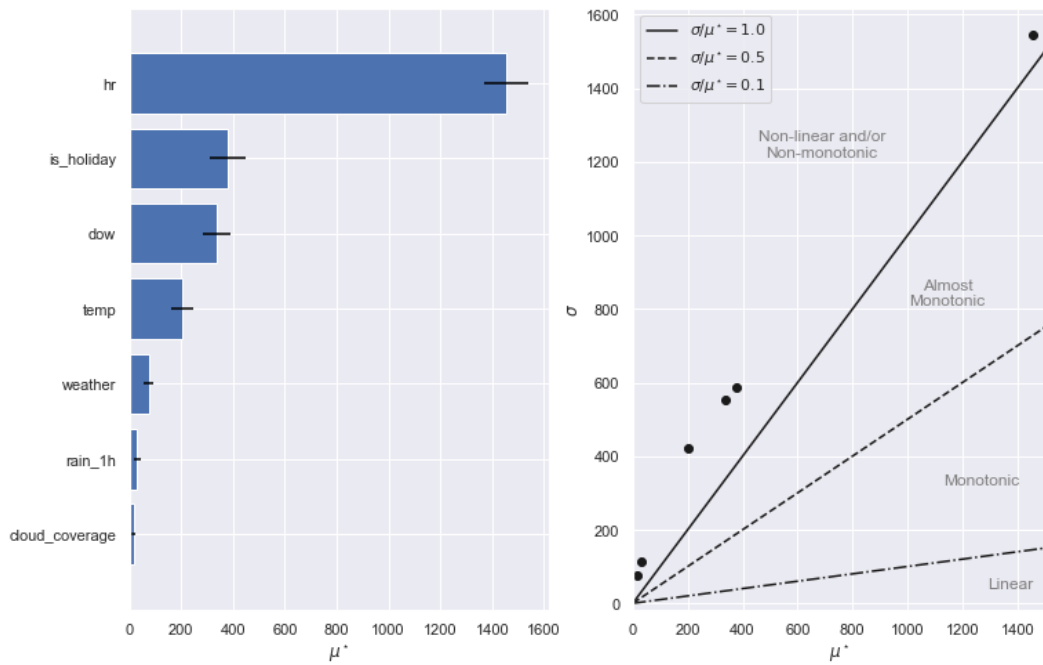
Hot Saturday



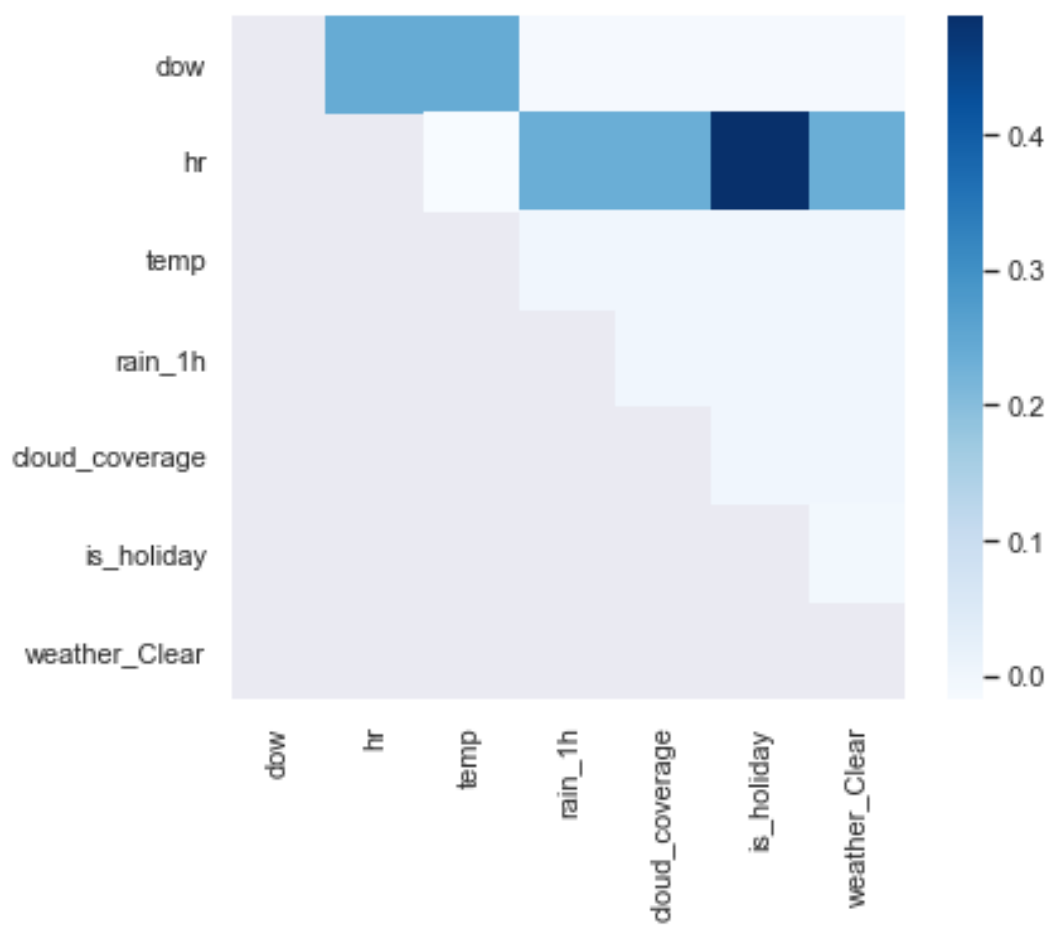


	count	mean	std	min	2.5%	50%	97.5%	max
dow	2232.000000	1.991935	1.415458	0.000000	0.000000	2.000000	4.000000	4.000000
hr	2232.000000	5.500000	7.933780	0.000000	0.000000	2.500000	23.000000	23.000000
temp	2232.000000	16.026438	5.380406	-2.570000	3.178750	16.935000	24.476750	30.458000
rain_1h	2232.000000	0.099628	0.603634	0.000000	0.000000	0.000000	1.451250	10.920000
cloud_coverage	2232.000000	29.178763	36.701417	0.000000	0.000000	1.000000	90.000000	100.000000
is_holiday	2232.000000	0.037634	0.190353	0.000000	0.000000	0.000000	1.000000	1.000000
weather_Clear	2232.000000	0.432348	0.495513	0.000000	0.000000	0.000000	1.000000	1.000000
weather_Clouds	2232.000000	0.207885	0.405885	0.000000	0.000000	0.000000	1.000000	1.000000
weather_Haze	2232.000000	0.010753	0.103159	0.000000	0.000000	0.000000	0.000000	1.000000
weather_Mist	2232.000000	0.104391	0.305835	0.000000	0.000000	0.000000	1.000000	1.000000
weather_Other	2232.000000	0.058244	0.234256	0.000000	0.000000	0.000000	1.000000	1.000000
weather_Rain	2232.000000	0.181452	0.385478	0.000000	0.000000	0.000000	1.000000	1.000000
weather_Snow	2232.000000	0.002240	0.047288	0.000000	0.000000	0.000000	0.000000	1.000000
weather_Unknown	2232.000000	0.002688	0.051789	0.000000	0.000000	0.000000	0.000000	1.000000

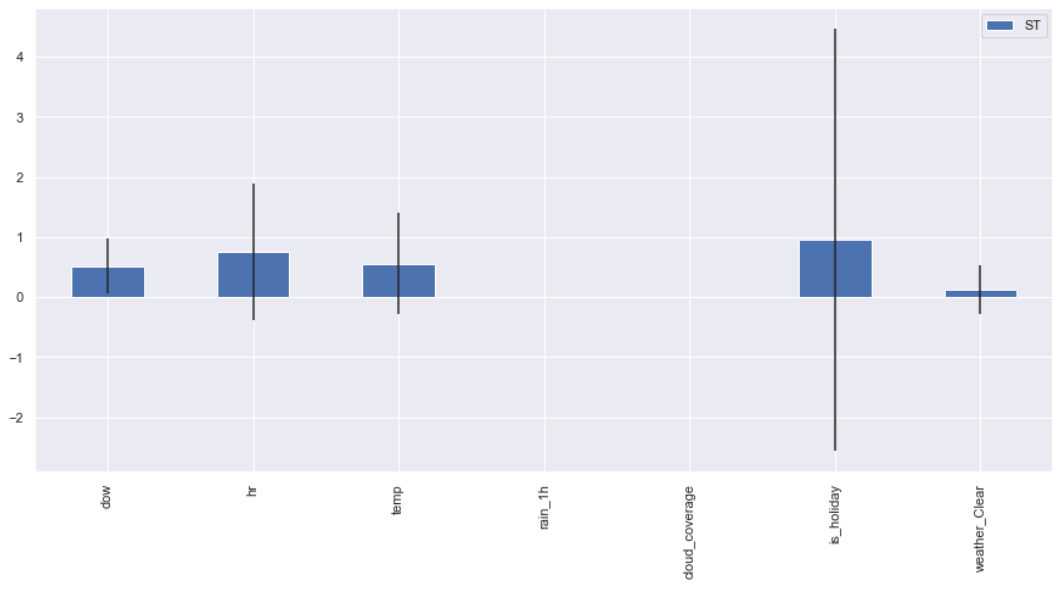
	features	μ	μ^*	σ
1	hr	-429.300110	1455.506958	1544.544312
5	is_holiday	-345.794861	379.520477	588.769897
0	dow	130.311508	336.568451	554.439819
2	temp	62.087799	202.984299	422.309845
6	weather	nan	75.732839	nan
3	rain_1h	-2.807377	30.730101	113.262093
4	cloud_coverage	9.897467	17.152805	74.319984

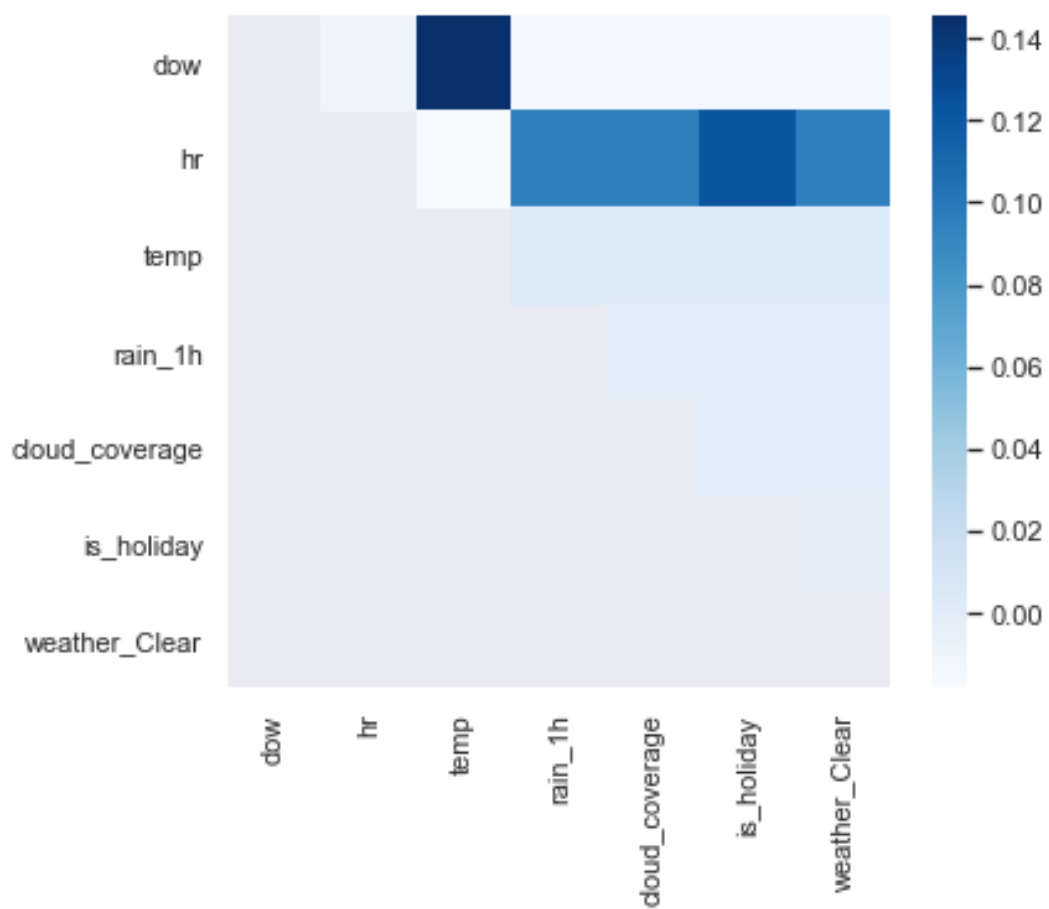


	features	1st	Total	Total Conf	Mean of Input
1	hr	0.009185	0.886824	0.912979	1.495931
2	temp	0.006123	0.506757	0.660847	14.059766
0	dow	0.009185	0.380068	0.366337	1.995599
5	is_holiday	0.003062	0.380068	0.479628	0.498047
6	weather_Clear	-0.003062	0.126689	0.314201	0.499023
3	rain_1h	0.000000	0.000000	0.000000	5.511458
4	cloud_coverage	0.000000	0.000000	0.000000	50.024740

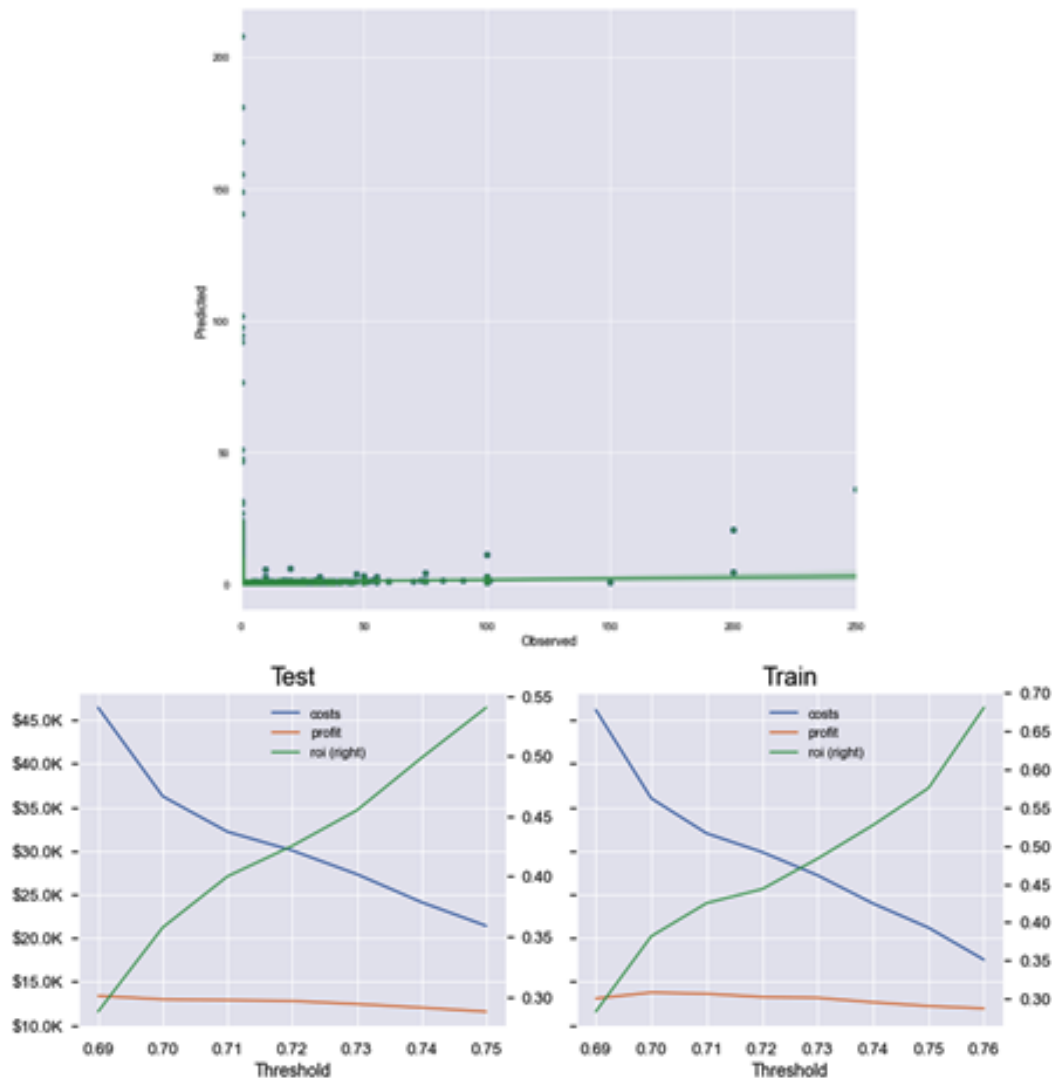


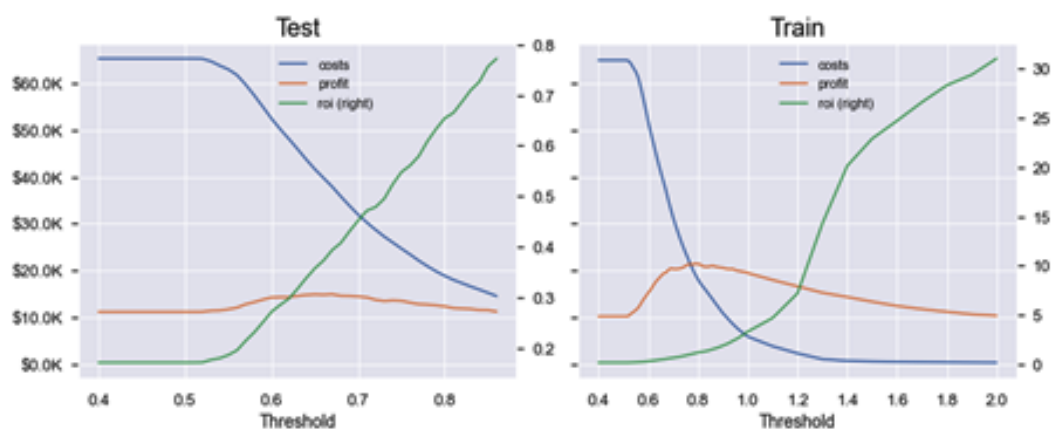
	features	1st	Total	Total Conf	Mean of Input
5	is_holiday	0.000852	0.953684	3.509326	0.498047
1	hr	0.010101	0.748595	1.132665	1.495931
2	temp	0.000677	0.552892	0.843215	14.059766
0	dow	0.009874	0.514826	0.452778	1.995599
6	weather_Clear	-0.002776	0.121222	0.404481	0.499023
4	cloud_coverage	-0.000000	0.000000	0.000000	50.024740
3	rain_1h	-0.000000	0.000000	0.000000	5.511458



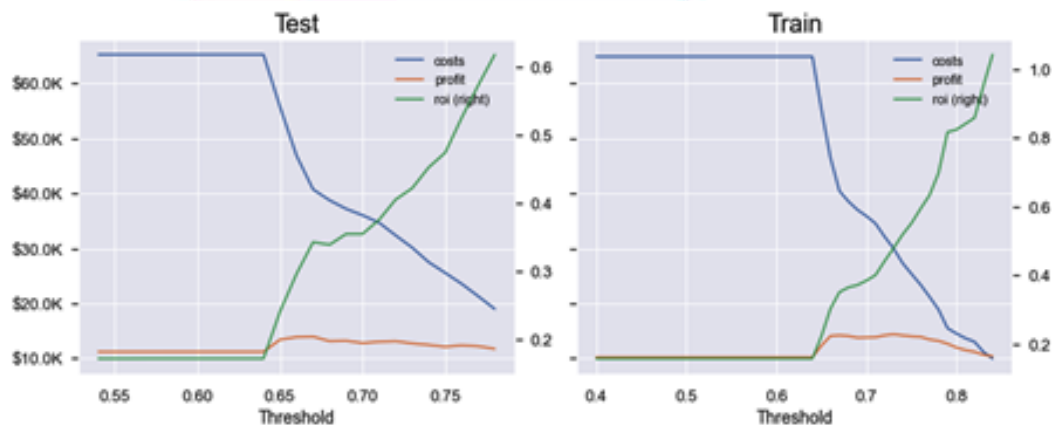


Chapter 10: Feature Selection and Engineering for Interpretability



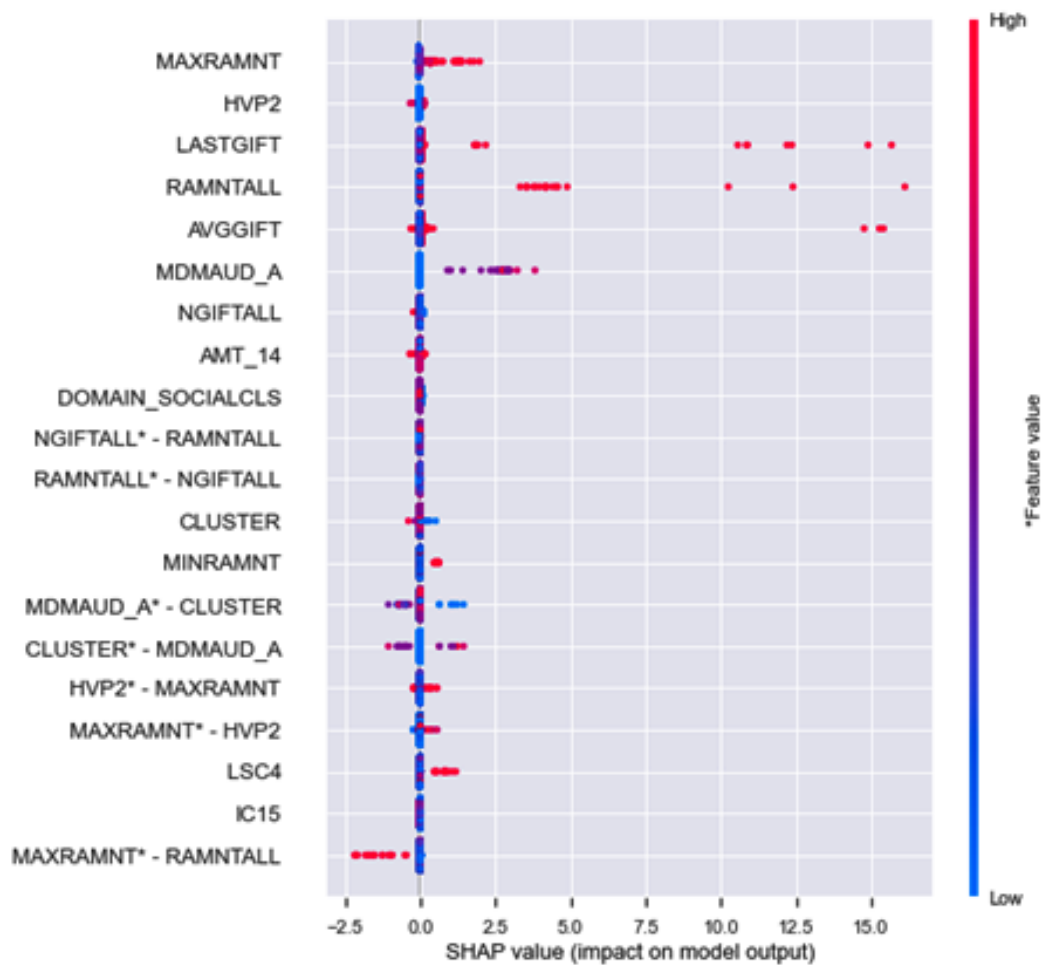


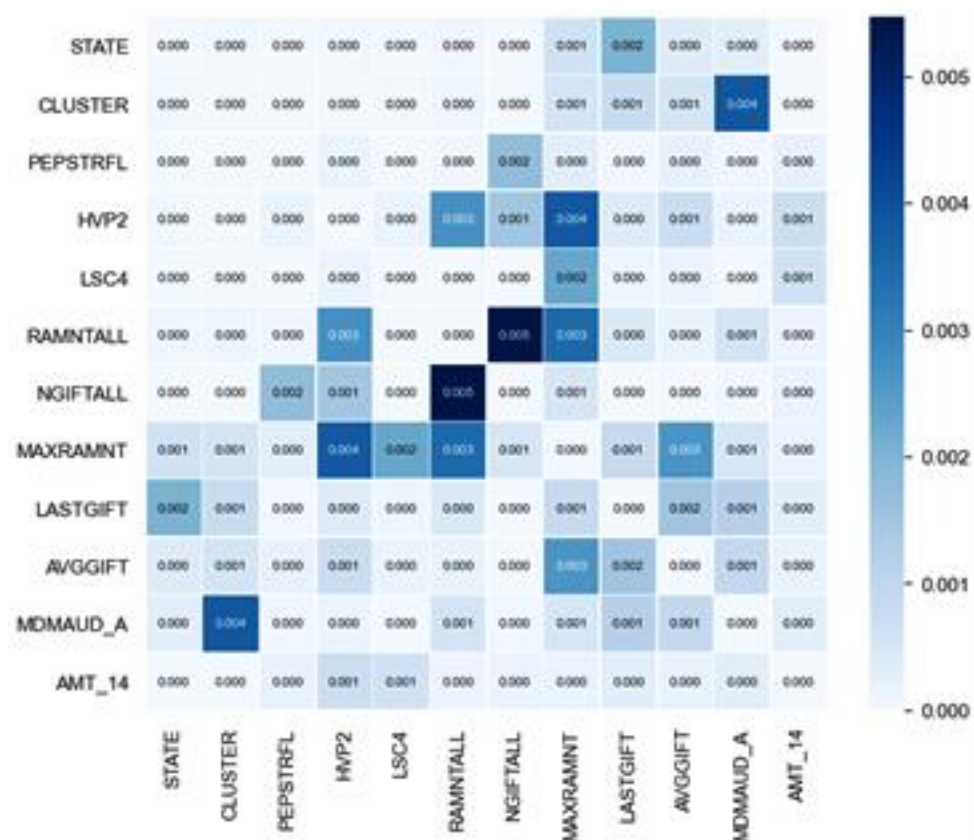
	depth	fs	rmse_train	rmse_test	max_profit_train	max_profit_test	max_roi	min_costs	speed	num_feat
rf_12_all	12	all	3.94	4.69	21521.98	14932.84	0.77	14532.28	2.89	415
rf_11_all	11	all	3.99	4.69	19904.00	15141.86	0.76	14928.04	2.73	398
rf_10_all	10	all	4.05	4.68	18603.92	14987.06	0.78	14396.28	2.43	383
rf_9_all	9	all	4.10	4.68	17453.14	14777.74	0.80	13997.12	2.19	346
rf_8_all	8	all	4.14	4.67	16439.72	14563.04	0.73	15308.84	1.94	315
rf_7_all	7	all	4.18	4.66	15435.32	14187.62	0.66	17164.56	1.71	277
rf_6_all	6	all	4.23	4.65	14651.12	13845.27	0.59	19305.20	1.41	240
rf_5_all	5	all	4.27	4.64	14242.32	13752.13	0.59	19199.12	1.22	201
rf_4_all	4	all	4.32	4.64	13715.90	13261.88	0.53	22392.40	1.00	160

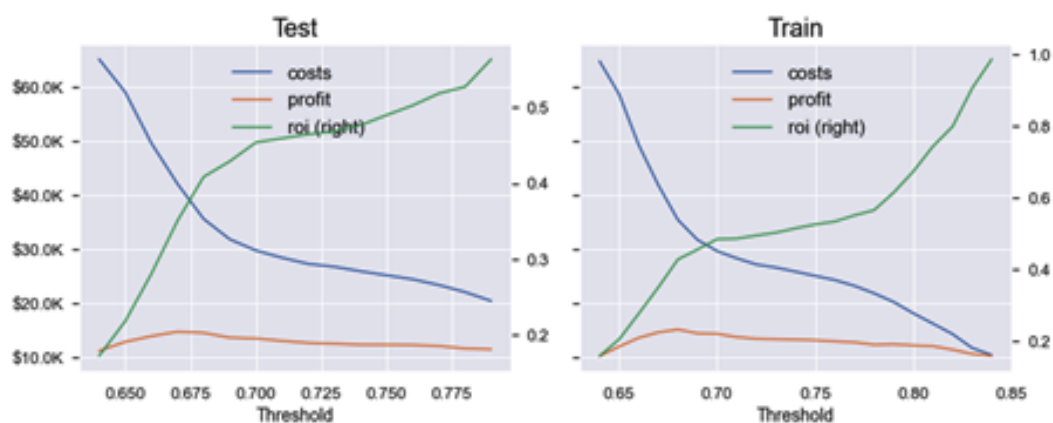
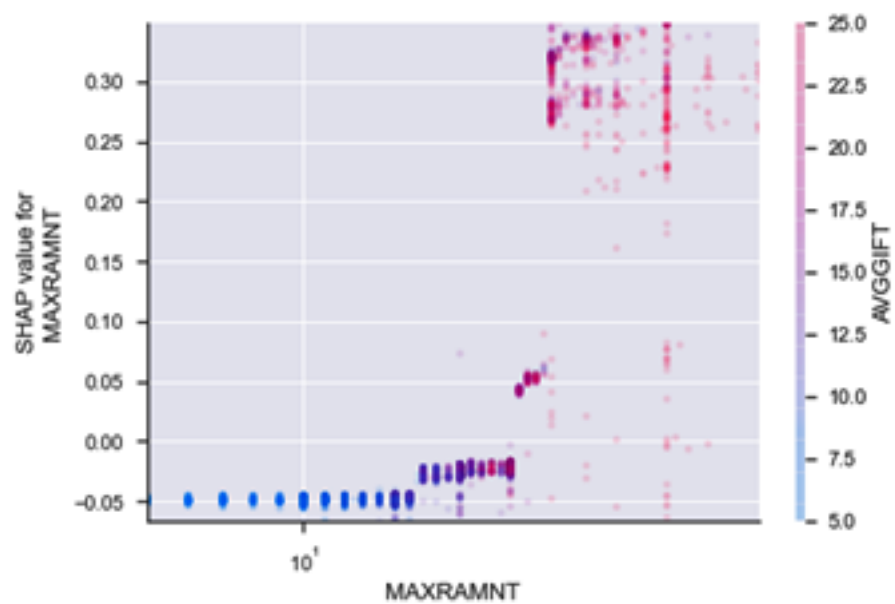


	depth	fs	rmse_train	rmse_test	max_profit_train	max_profit_test	max_roi	min_costs	speed	total_feat	num_feat
rf_11_all	11	all	3.99	4.69	19904.00	15141.86	0.76	14928.04	2.73	435	398
rf_10_all	10	all	4.05	4.68	18603.92	14987.06	0.78	14396.28	2.43	435	383
rf_12_all	12	all	3.94	4.69	21521.98	14932.84	0.77	14532.28	2.89	435	415
rf_11_f-corr	11	f-corr	3.98	4.67	19923.84	14894.94	0.77	14592.80	2.47	419	404
rf_9_all	9	all	4.10	4.68	17453.14	14777.74	0.80	13997.12	2.19	435	346
rf_8_all	8	all	4.14	4.67	16439.72	14563.04	0.73	15308.84	1.94	435	315
rf_5_f-mic	5	f-mic	4.31	4.57	14983.34	14481.39	0.62	18971.32	0.39	160	103
rf_7_all	7	all	4.18	4.66	15435.32	14187.62	0.66	17164.56	1.71	435	277
rf_6_all	6	all	4.23	4.65	14651.12	13845.27	0.59	19305.20	1.41	435	240
rf_5_all	5	all	4.27	4.64	14242.32	13752.13	0.59	19199.12	1.22	435	201
rf_4_all	4	all	4.32	4.64	13715.90	13261.88	0.53	22392.40	1.00	435	160
	depth	fs	rmse_train	rmse_test	max_profit_train	max_profit_test	max_roi	min_costs	speed	total_feat	num_feat
rf_11_all	11	all	3.99	4.69	19904.00	15141.86	0.76	14928.04	2.73	435	398
rf_10_all	10	all	4.05	4.68	18603.92	14987.06	0.78	14396.28	2.43	435	383
rf_12_all	12	all	3.94	4.69	21521.98	14932.84	0.77	14532.28	2.89	435	415
rf_11_f-corr	11	f-corr	3.98	4.67	19923.84	14894.94	0.77	14592.80	2.47	419	404
rf_9_all	9	all	4.10	4.68	17453.14	14777.74	0.80	13997.12	2.19	435	346
rf_5_e-llarsic	5	e-llarsic	4.28	4.45	15168.46	14768.37	0.56	20441.48	0.32	111	87
rf_8_all	8	all	4.14	4.67	16439.72	14563.04	0.73	15308.84	1.94	435	315
rf_5_f-mic	5	f-mic	4.31	4.57	14983.34	14481.39	0.62	18971.32	0.39	160	103
rf_6_h-rfe-lda	6	h-rfe-lda	4.25	4.48	15329.72	14351.74	0.71	15824.28	0.61	183	129
rf_6_e-logl2	6	e-logl2	4.28	4.60	15353.44	14199.90	0.67	16904.12	0.32	87	84
rf_7_all	7	all	4.18	4.66	15435.32	14187.62	0.66	17164.56	1.71	435	277
rf_6_all	6	all	4.23	4.65	14651.12	13845.27	0.59	19305.20	1.41	435	240
rf_5_all	5	all	4.27	4.64	14242.32	13752.13	0.59	19199.12	1.22	435	201
rf_4_e-llars	4	e-llars	4.36	4.45	14014.10	13633.19	0.52	22906.48	0.06	8	8
rf_6_h-rfe-rf	6	h-rfe-rf	4.40	4.78	13202.61	13347.15	0.41	28596.04	0.08	1	1
rf_4_all	4	all	4.32	4.64	13715.90	13261.88	0.53	22392.40	1.00	435	160
rf_3_e-lasso	3	e-lasso	4.46	4.49	14166.64	12930.30	0.51	22248.92	0.05	7	7

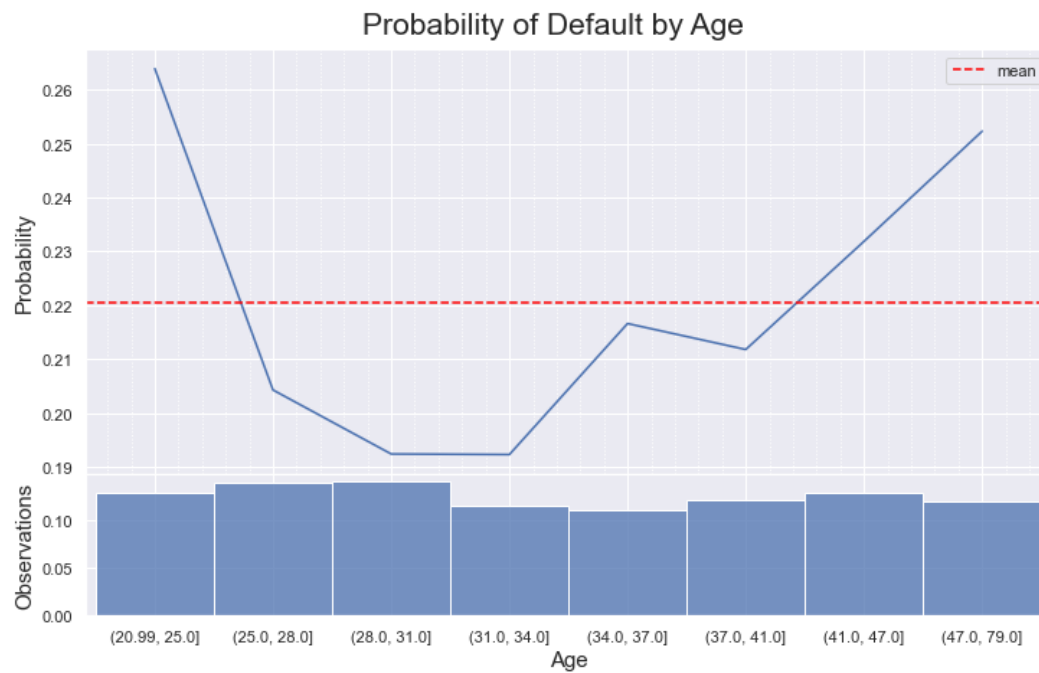
	depth	fs	rmse_train	rmse_test	max_profit_train	max_profit_test	max_roi	min_costs	speed	total_feat	num_feat
rf_5_e-llarsic	5	e-llarsic	4.28	4.45	15168.46	14768.37	0.56	20441.48	0.32	111	87
rf_5_f-mic	5	f-mic	4.31	4.57	14983.34	14481.39	0.62	18971.32	0.39	160	103
rf_6_h-rfe-lda	6	h-rfe-lda	4.25	4.48	15329.72	14351.74	0.71	15824.28	0.61	183	129
rf_6_a-shap	6	a-shap	4.23	4.52	15263.60	14282.20	0.61	18767.32	0.50	150	135
rf_5_a-ga-rf	5	a-ga-rf	4.39	4.45	14274.52	14220.53	0.69	13237.56	0.07	63	63
rf_6_e-logl2	6	e-logl2	4.28	4.60	15353.44	14199.90	0.67	16904.12	0.32	87	84
rf_6_all	6	all	4.23	4.65	14651.12	13845.27	0.59	19305.20	1.41	435	240
rf_5_w-sfs-lda	5	w-sfs-lda	4.43	4.63	14377.13	13801.55	0.51	22508.95	0.11	27	27
rf_5_all	5	all	4.27	4.64	14242.32	13752.13	0.59	19199.12	1.22	435	201
rf_4_e-llars	4	e-llars	4.36	4.45	14014.10	13633.19	0.52	22906.48	0.06	8	8
rf_6_a-pca	6	a-pca	4.30	4.46	14353.54	13351.57	0.47	23901.32	0.54	150	126
rf_6_h-rfe-rf	6	h-rfe-rf	4.40	4.78	13202.61	13347.15	0.41	28596.04	0.08	1	1
rf_4_all	4	all	4.32	4.64	13715.90	13261.88	0.53	22392.40	1.00	435	160
rf_6_w-sbs-et	6	w-sbs-et	4.34	4.53	14222.49	13119.17	0.71	14711.66	0.45	135	123
rf_3_e-lasso	3	e-lasso	4.46	4.49	14166.64	12930.30	0.51	22248.92	0.05	7	7

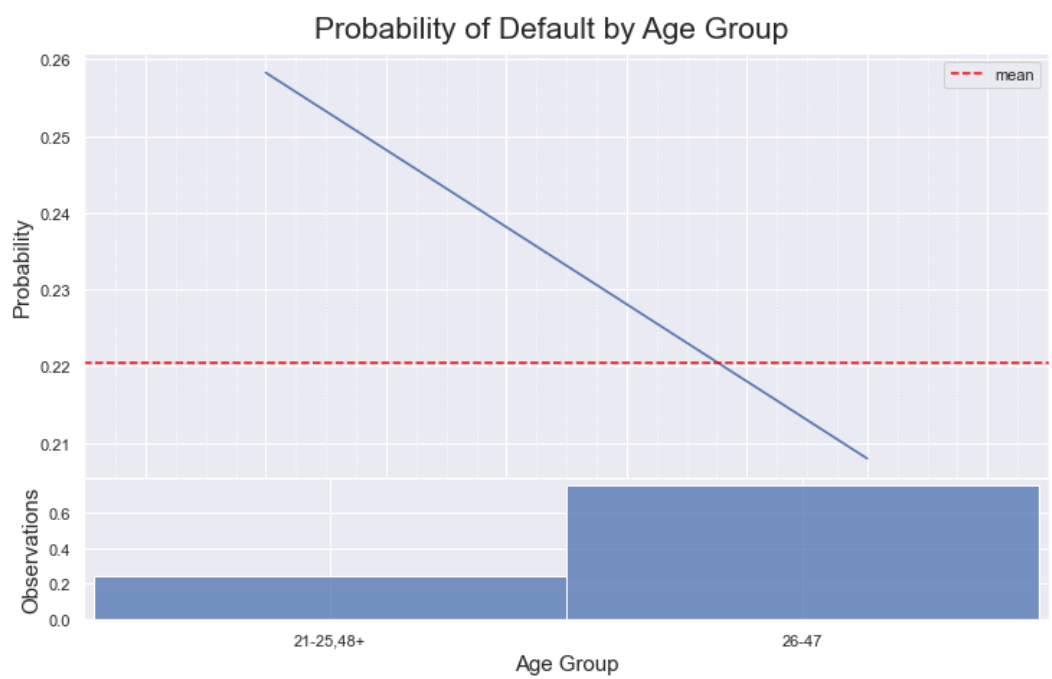


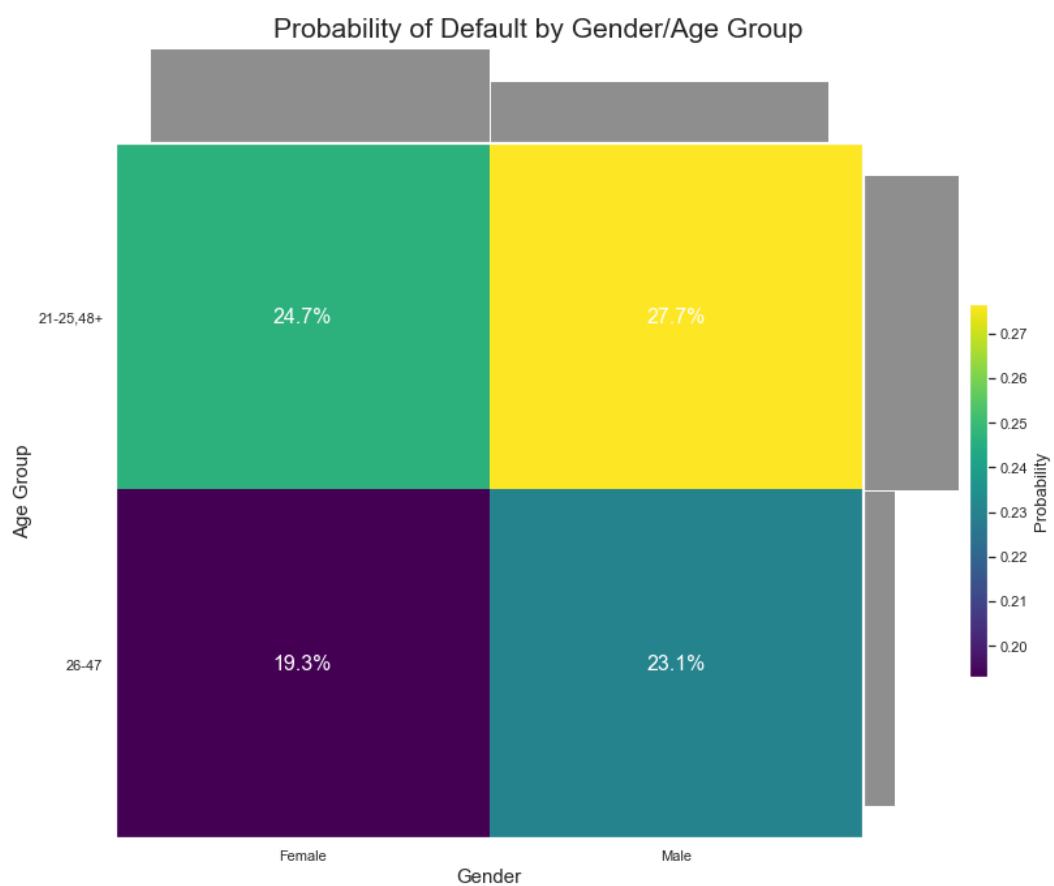


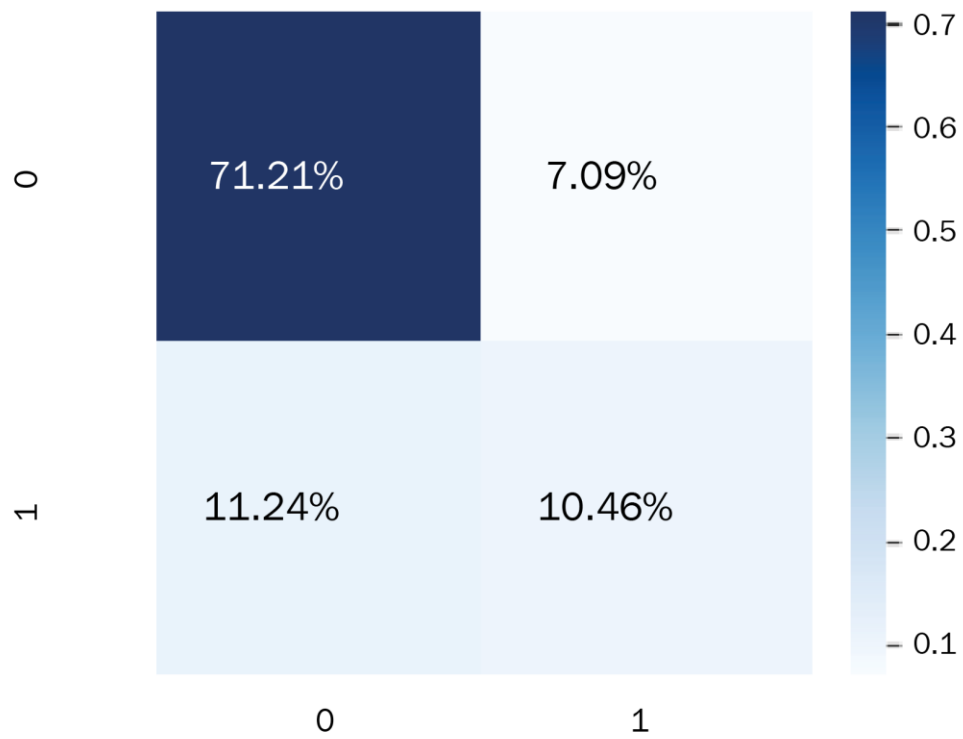


Chapter 11: Bias Mitigation and Causal Inference Methods

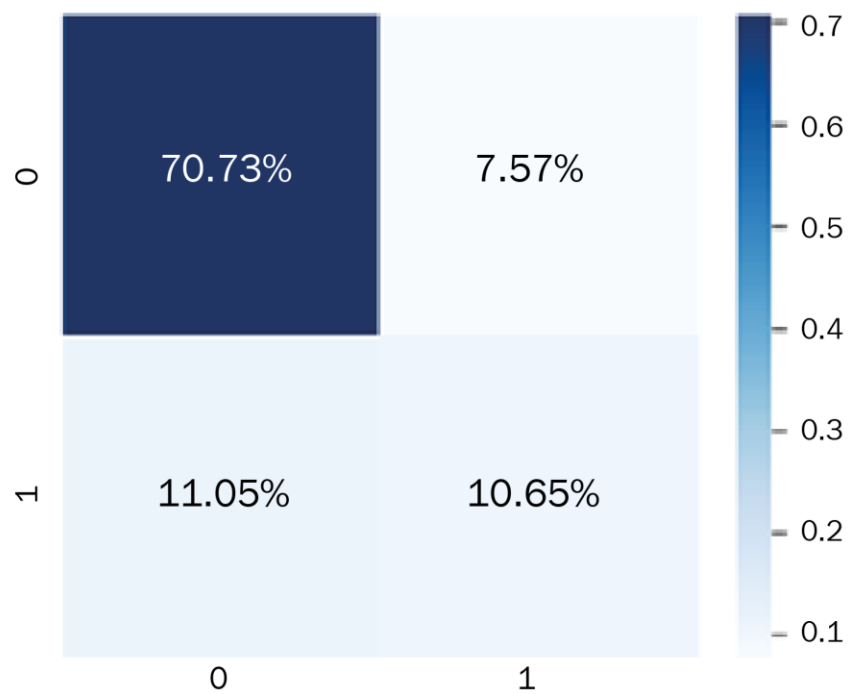




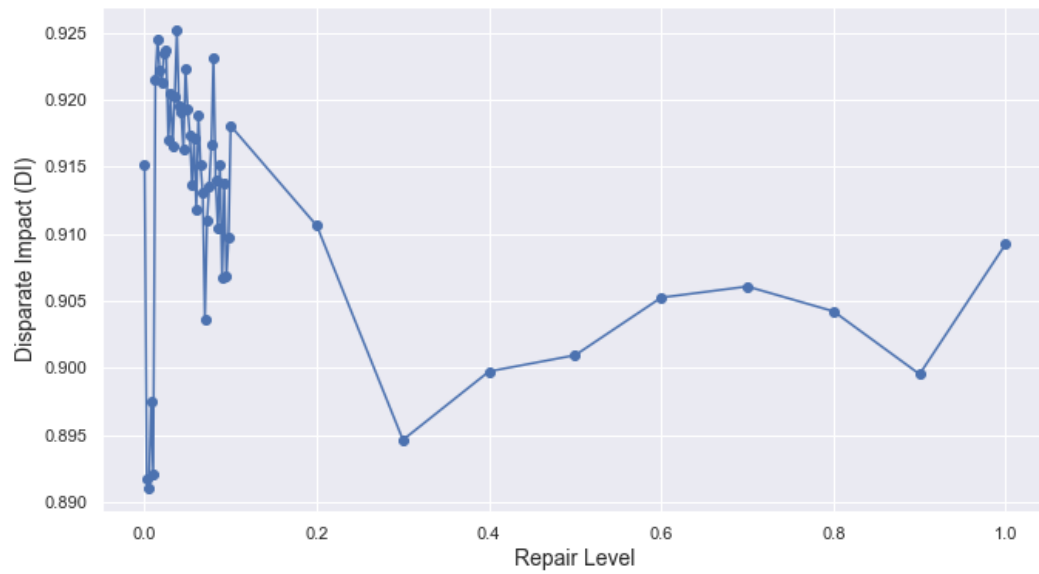




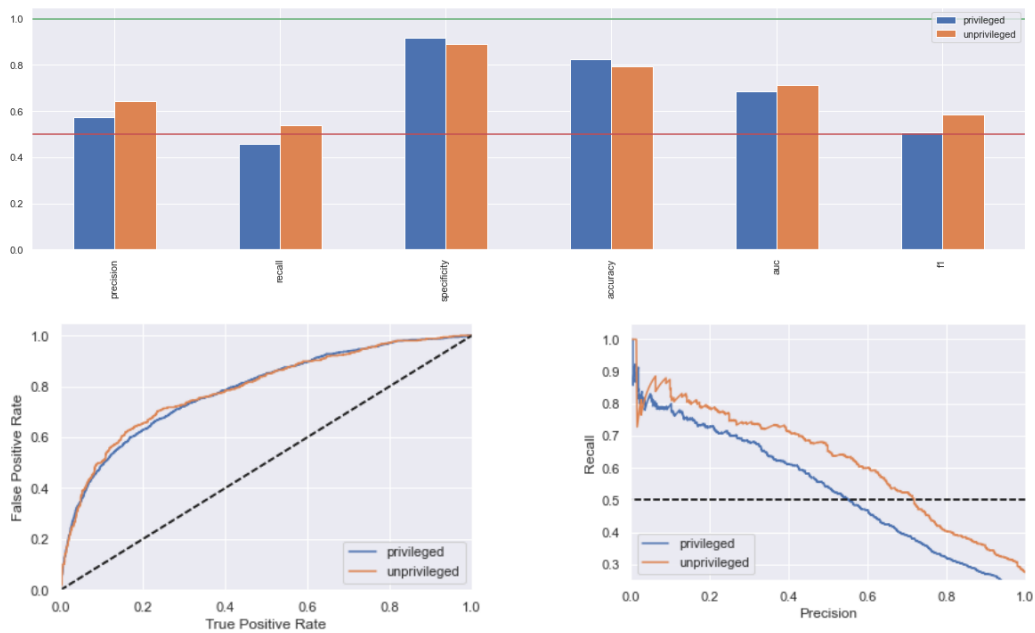
1. Accuracy_train: 0.8255 Accuracy_test: 0.8167
2. Precision_test: 0.5961 Recall_test: 0.4820
3. ROC-AUC_test: 0.7901 F1_test: 0.5330 MCC_test: 0.4243



1. Accuracy_train: 0.8240 Accuracy_test: 0.8138
2. Precision_test: 0.5847 Recall_test: 0.4908
3. ROC-AUC_test: 0.7886 F1_test: 0.5337 MCC_test: 0.4210

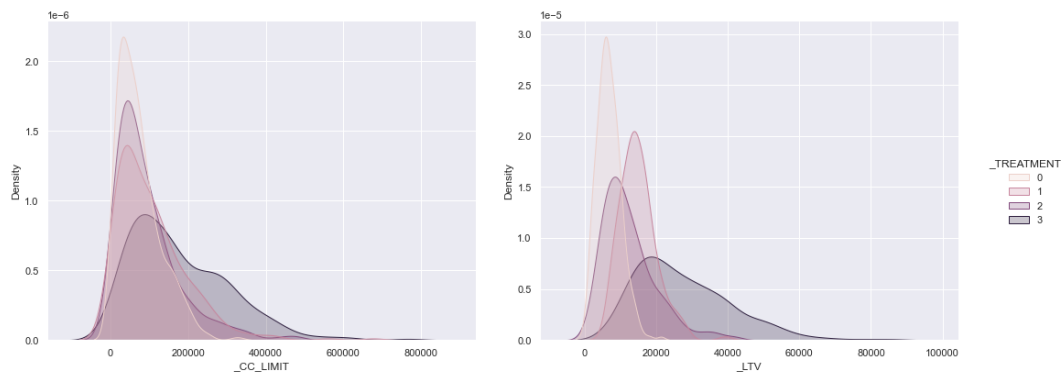
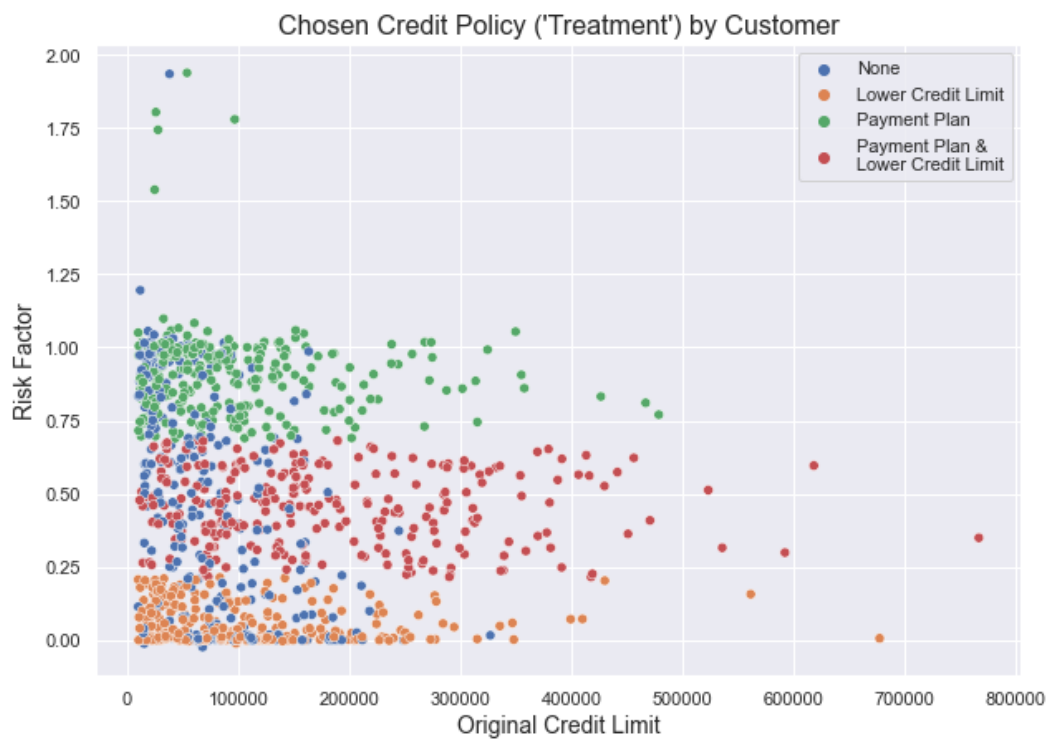


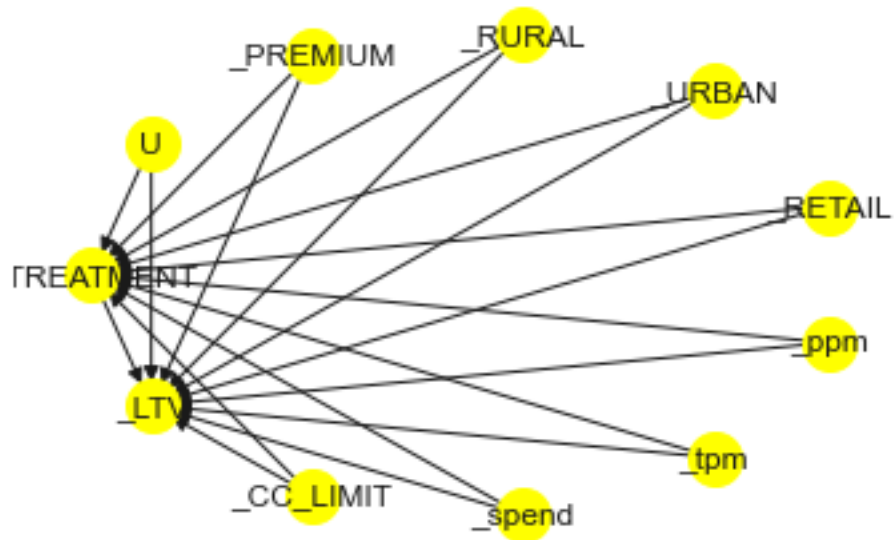
	accuracy_train	accuracy_test	f1_test	mcc_test	SPD	DI	AOD	EOD	DFBA
dt_2_gf	0.8214	0.8262	0.4812	0.4135	-0.0548	0.9388	-0.0430	-0.0216	0.2521
lgb_0_base	0.8255	0.8167	0.5330	0.4243	-0.0679	0.9193	-0.0550	-0.0265	0.2328
lgb_1_rw	0.8240	0.8138	0.5337	0.4210	-0.0371	0.9552	-0.0171	-0.0018	0.0349
lgb_1_dir	0.8237	0.8129	0.5301	0.4171	-0.0624	0.9252	-0.0493	-0.0214	0.2545
lgb_3_epp	0.8255	0.8101	0.5152	0.4025	-0.0260	0.9688	0.0022	-0.0021	0.0031
lgb_3_cpp	0.8255	0.2622	0.2129	-0.3055	-0.0711	0.7609	-0.0635	-0.1262	0.0432
log_2_pr	0.1912	0.1873	0.2844	-0.3363	0.0520	1.7627	0.0498	0.0235	0.3454



Credit Policy Experiment Outcomes







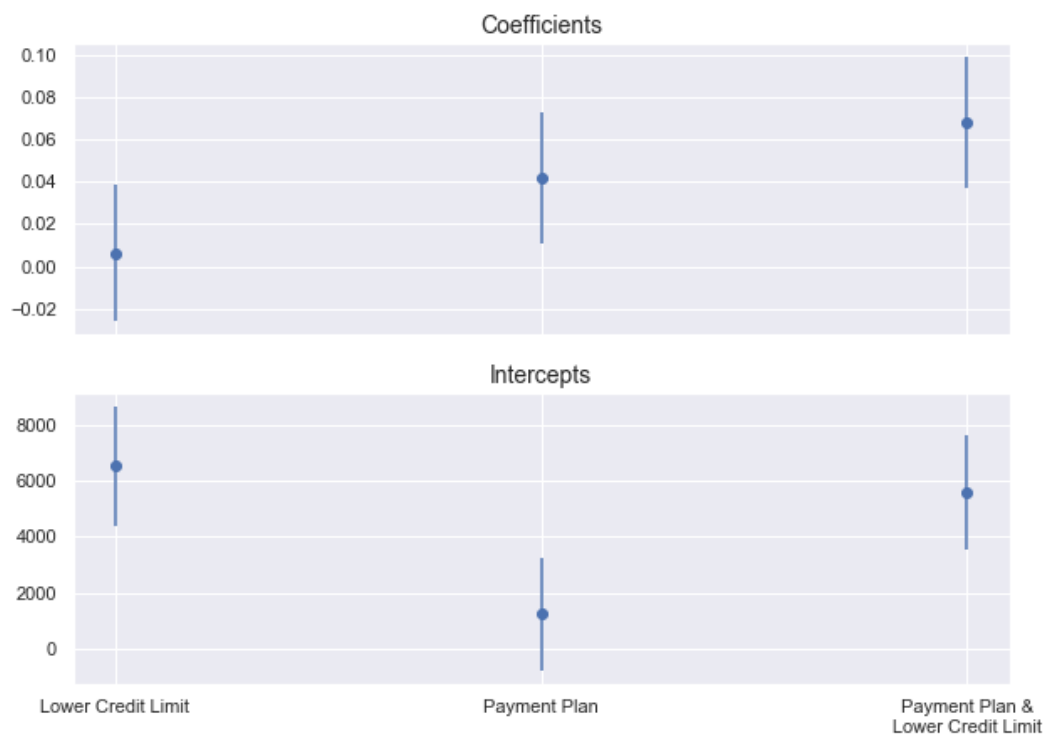
Treatment: Payment Plan &
Lower Credit Limit

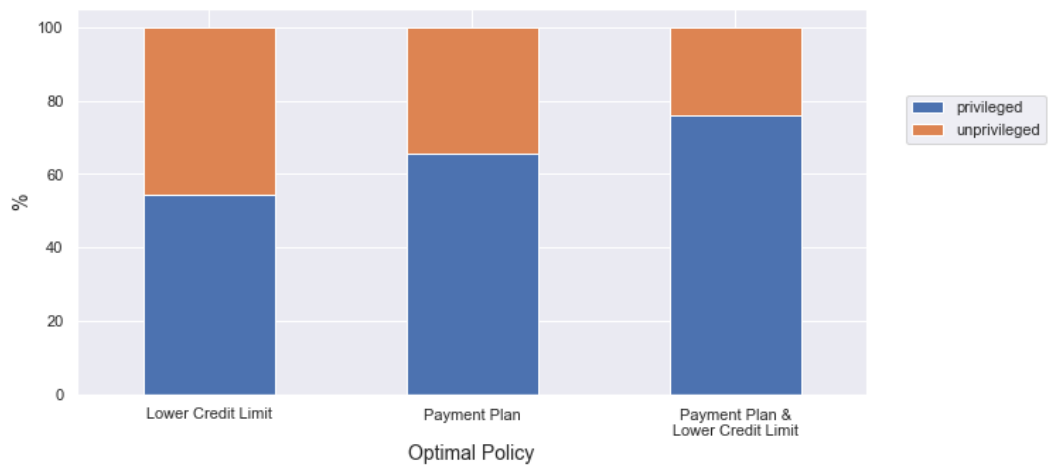
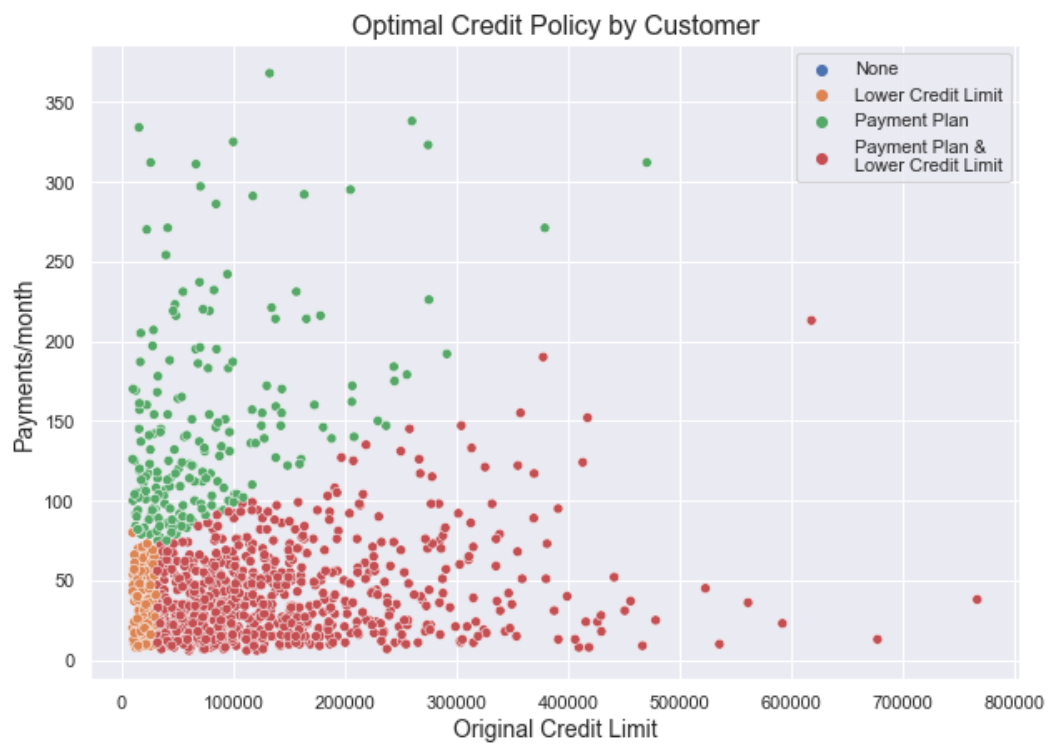
Coefficient Results

	point_estimate	stderr	zstat	pvalue	ci_lower	ci_upper
X0	0.095	0.03	3.171	0.002	0.046	0.145

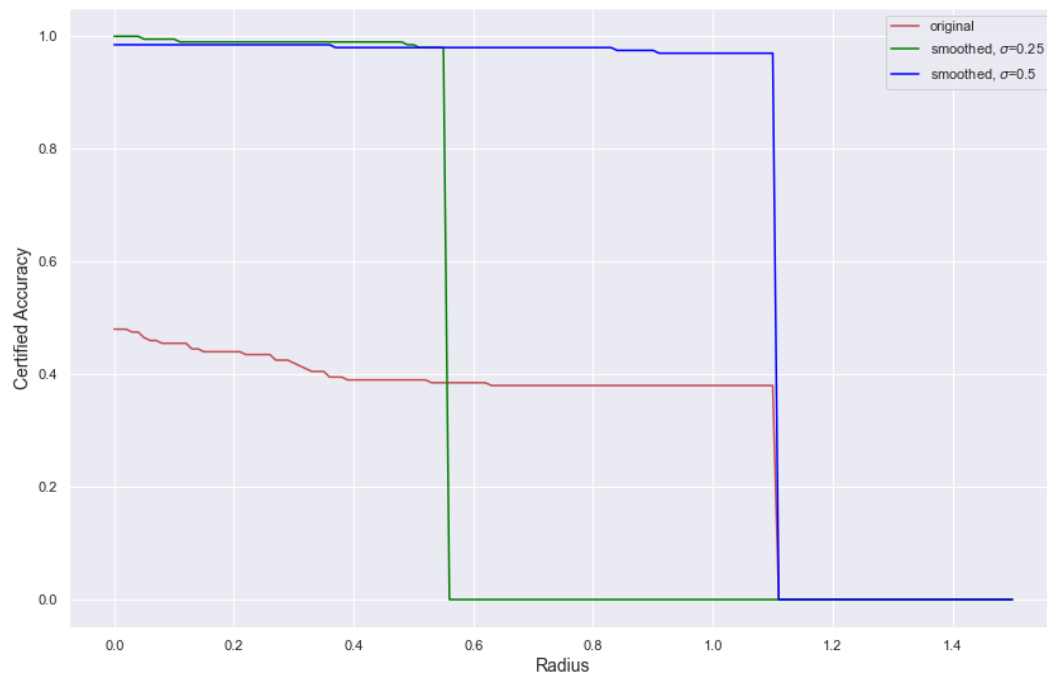
CATE Intercept Results

	point_estimate	stderr	zstat	pvalue	ci_lower	ci_upper
cate_intercept	4238.783	1907.747	2.222	0.026	1100.818	7376.748

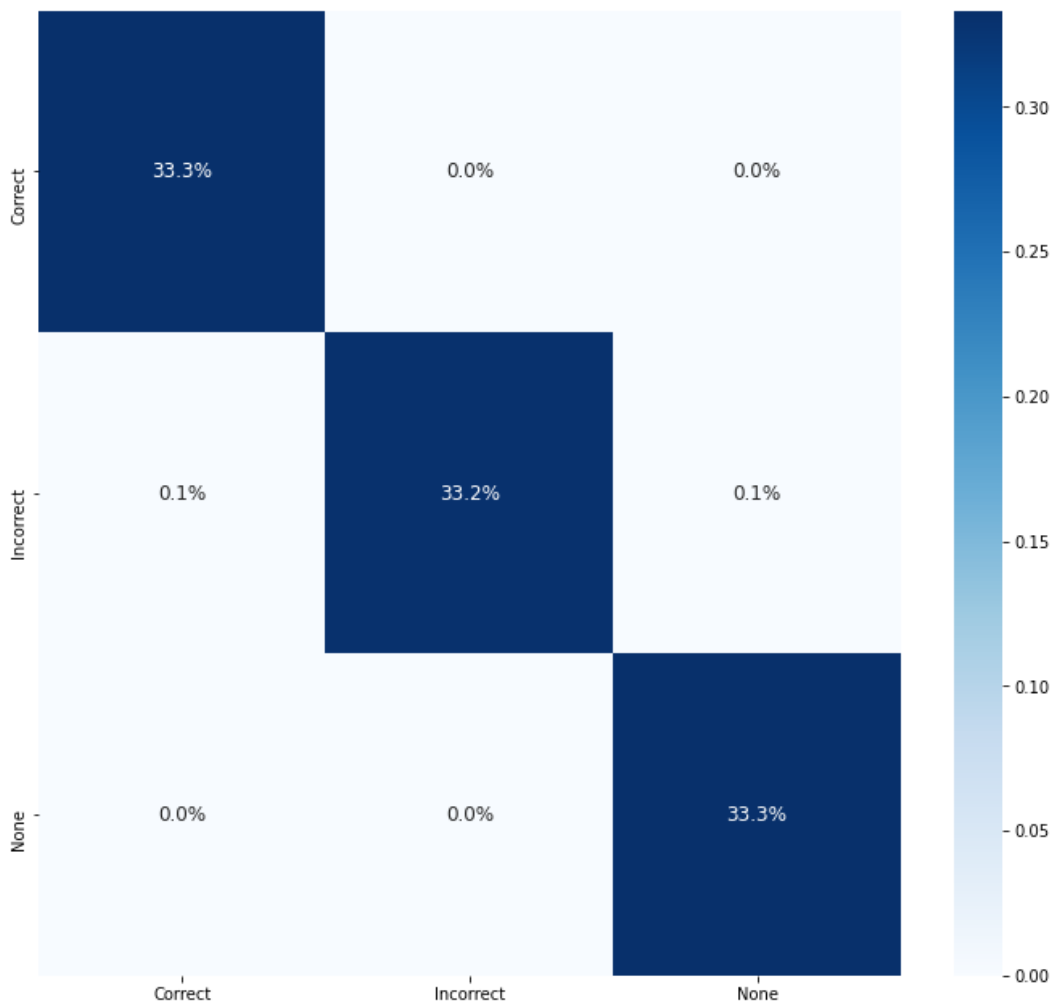




Chapter 13: Adversarial Robustness

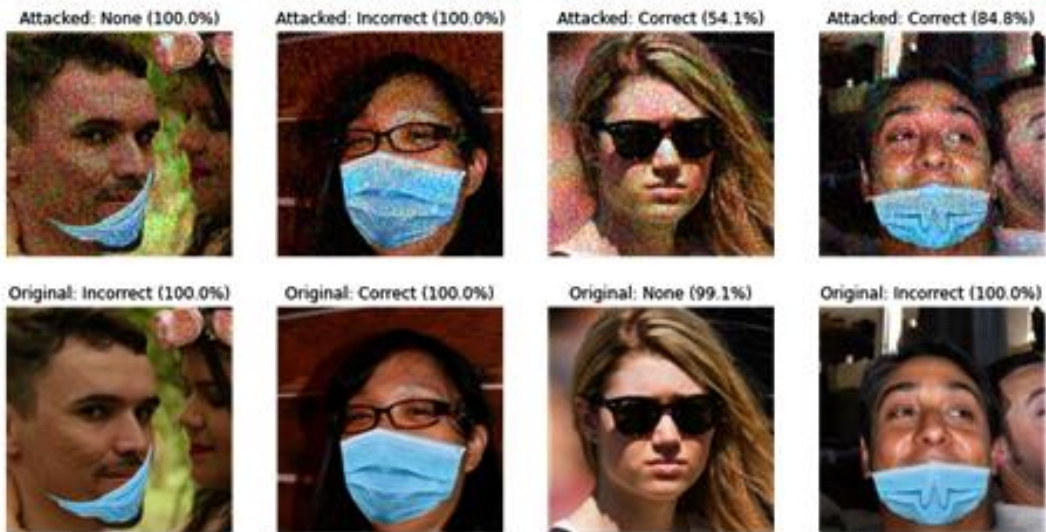




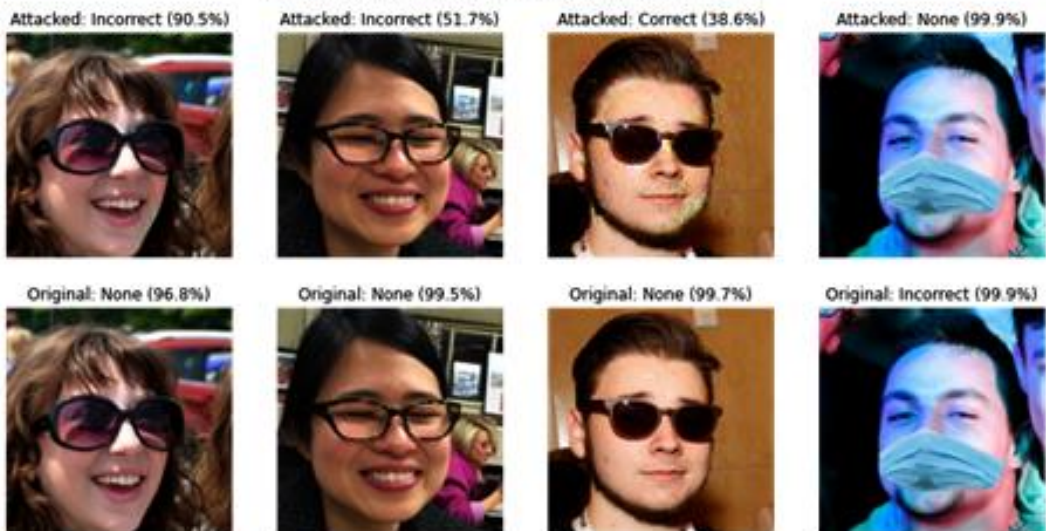


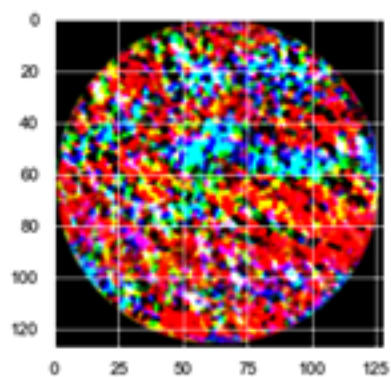
		Goal		
		Espionage	Sabotage	Fraud
Stage	Training	Inference (by poisoning)	Trojaning	
			Poisoning	
			Backdooring	
	Production	Inference	Reprogramming	
			Evasion	

FSGM Attack Average Perturbation: 0.092



C&W Inf Attack Average Perturbation: 0.003





AP Attack Average Perturbation: 0.080

Attacked: Incorrect (63.7%)



Attacked: Correct (52.3%)



Attacked: None (56.7%)



Attacked: Correct (67.1%)



Original: None (99.9%)



Original: None (100.0%)

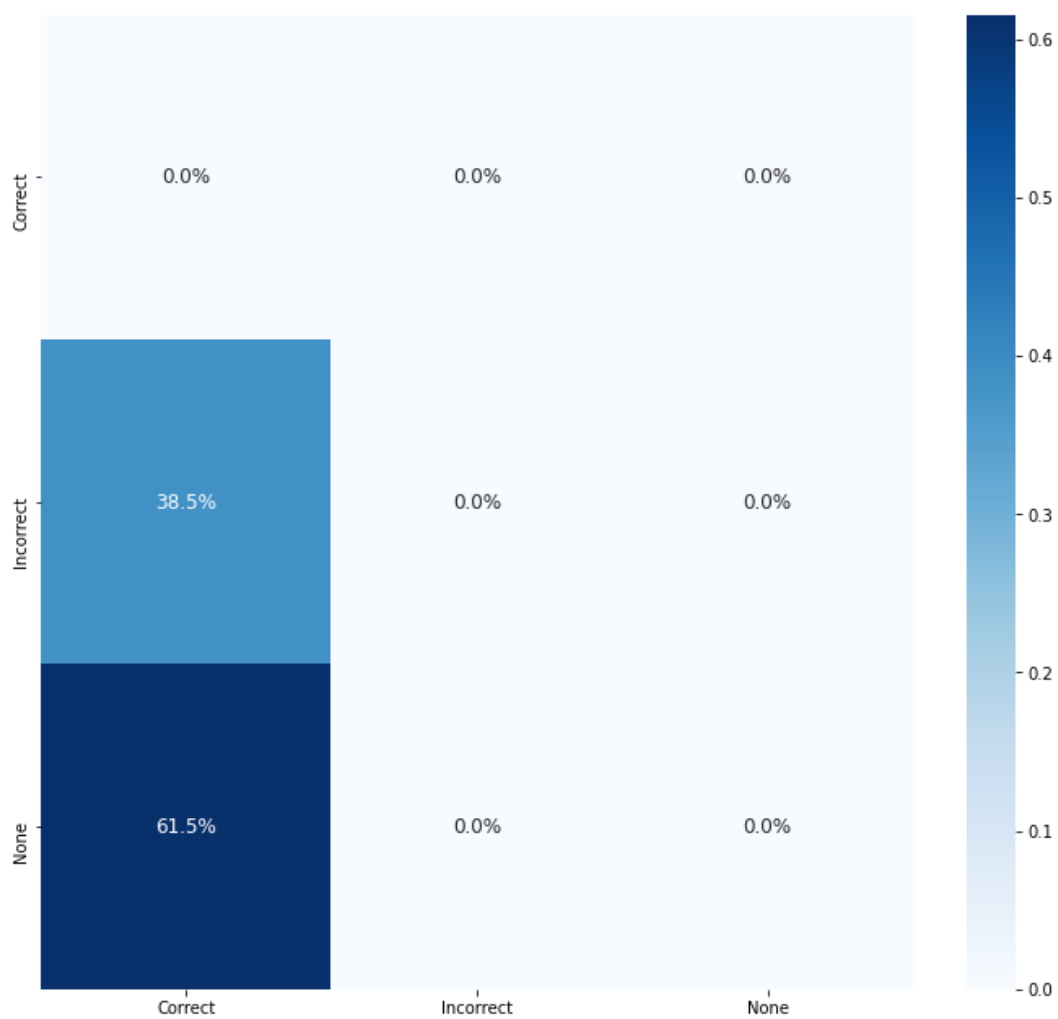


Original: Correct (100.0%)

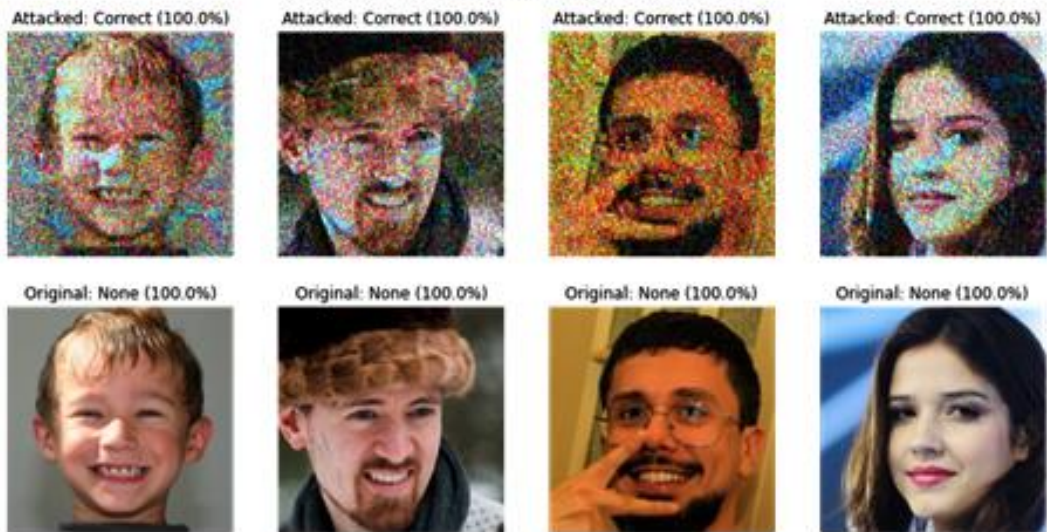


Original: Incorrect (100.0%)



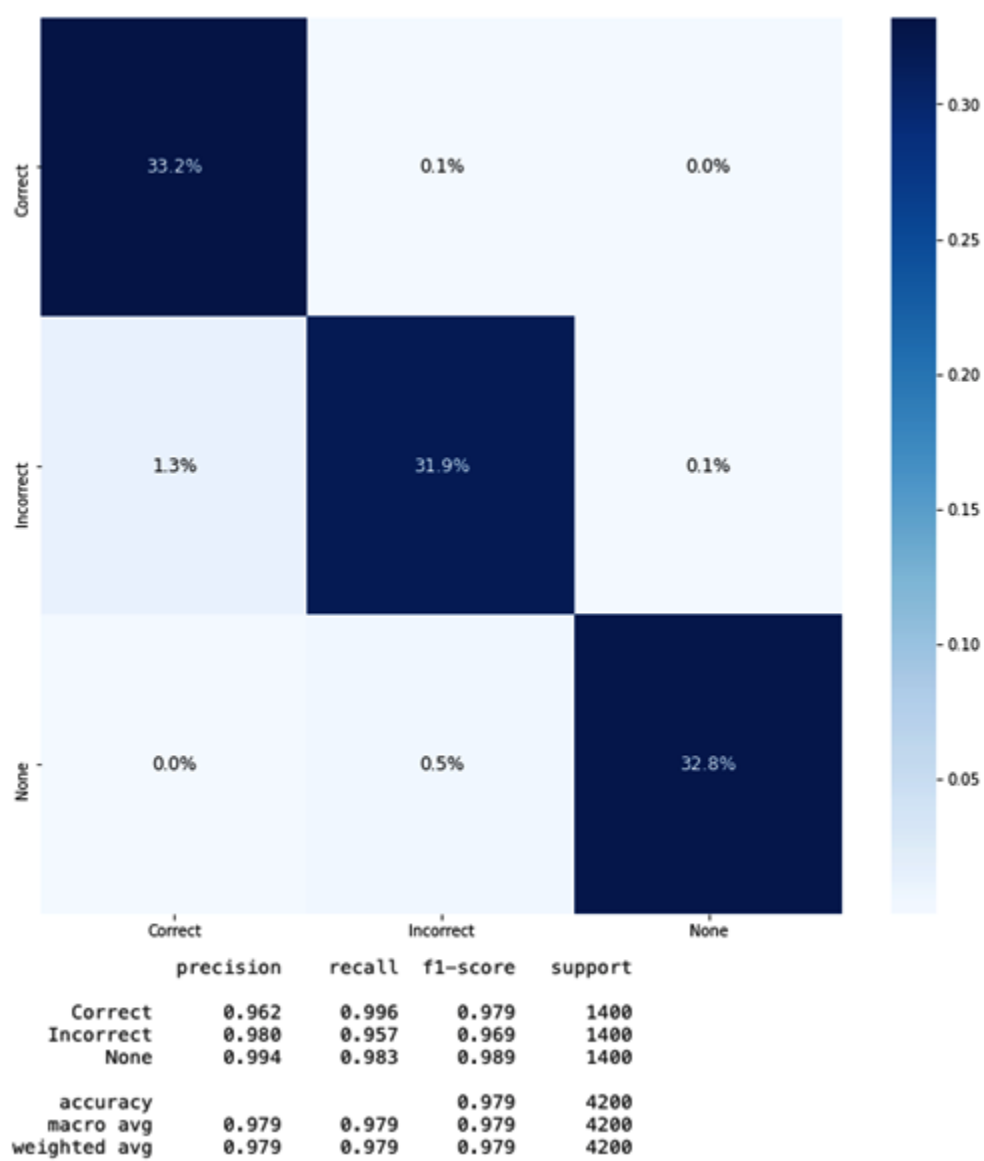


PGD Attack Average Perturbation: 0.147

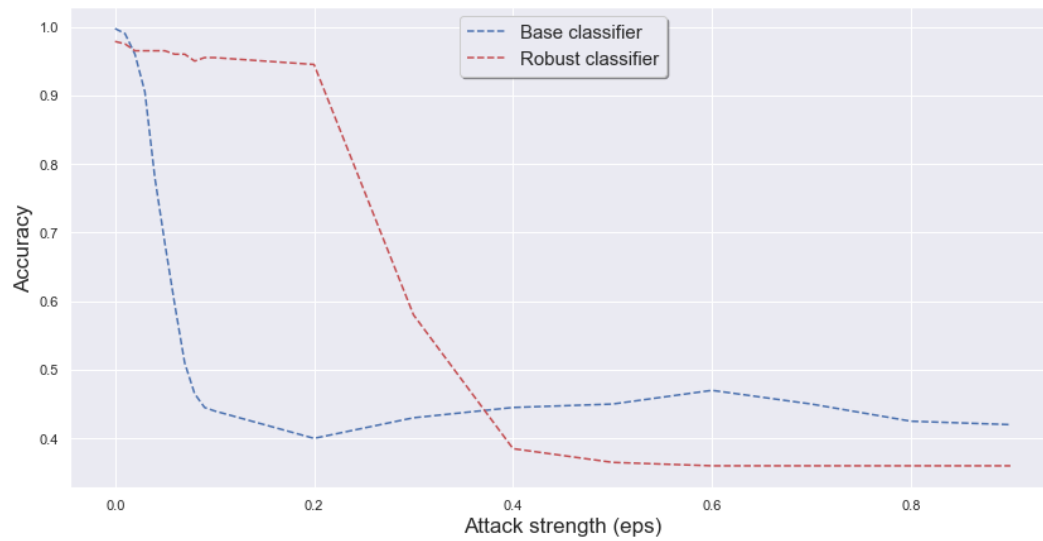


PGD Attack & Defended Average Perturbation: 0.069





FSGM Attack Average Perturbation: 0.035



Chapter 14: What's Next for Machine Learning Interpretability?



DIAGNOSTICS aim: EVALUATING MODEL / CHECKING ASSUMPTIONS / DETECTING PROBLEMS / UNDERSTANDING DATA

	Concerns	Interpretation Methods
FAIRNESS	Equity	• Class Balance 3 4 7 10 11 12
	Justice	• Comparing Metrics 7 11 12 (FPR, FNR)
	Diversity	• Comparing Plots 7 11 12 (Confusion Matrix, ROC Curve, PR Curve)
	Inclusion	• Group Fairness Metrics / Individual Fairness Metrics 11 (SPD, DI, AOD, EOD, DFBA, CDD) • Contour / Heat Probability Maps 12 • Sampling Bias Evaluations
ACCOUNTABILITY	Reliability	• Out-of-sample Evaluations 8
	Certainty	• Sensitivity Analysis 9 (Sobol, Morris, FAST)
	Security	• Causal Inference Methods 11 (DRL, DML, Forest Based, Meta-Learners)
	Safety	• Evasion Adversarial Robustness Evaluations 13 (FSGM, PGD, C&W, Adversarial Patches, Boundary, PDG, B&B, DeepFool..)
	Robustness	• Inference, Extraction & Poisoning Adversarial Robustness Evaluations
	Privacy	• Anomaly Detection / Metrics • Privacy Metrics
TRANSPARENCY	Interpretability	• Feature Importance Methods 1 2 3 4 5 8 9 10 12 (SHAP, Permutation, Model-specific)
	Explainability	• Dimensionality Reduction Methods 3 10 (PCA, t-SNE, VAE, DIP-VAE)
	Consistency	• Glass-box Models 3 (EBM, Skoped-Rules)
	Credibility	• Partial Dependence Plots & similar 4 5 7 9 11 12 (ICE, ALE, SHAP Dependence)
	Clarity	• White-box Surrogates 5 10 12 (Logistic Regression, Linear Regression, Rule Models, CART, KNN, ProfWeight)
		• Confirming with Statistical Tests & Correlations 5 10 12 (Spearman, Point-biserial, Cramér's V, Z-test) • Local Interpretation 6 7 9 (Decision Regions, ICE, Anchors, Counterfactuals, WIT, CEM, SHAP) • Deep Learning-specific 8 9 (IG, Saliency Maps, Grad-CAM, SmoothGrad, Semantic Segmentation) • Explainability Metrics

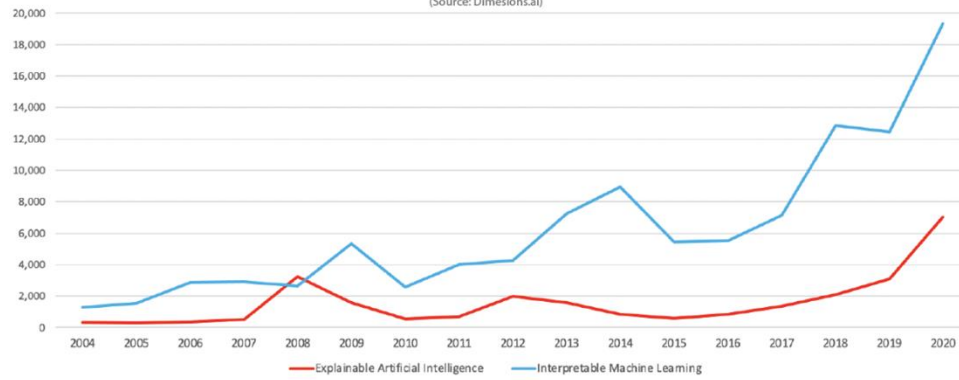
TREATMENT

aim:
FIXING & ANTICIPATING F.A.T PROBLEMS ("TUNING FOR INTERPRETABILITY")

	APPROACH	DATA	MODEL	PREDICTION
FAIRNESS	Mitigating Bias	Reweighting / DIR ¹¹ LFR / DIR / Unawareness ¹²	Cost-sensitive Learning ^{10 11 12} Prejudice Regul. / GerryFair ¹¹	Calibrating/Equalizing Odds ^{7 11} Reject Option Classification
	Placing Guardrails	Feature Engineering ^{10 12}	Monotonic Constraints ¹²	Prediction Abstention ^{11 13}
	Enhancing Reliability	Data Augmentation ^{8 11 13}	Adversarial Debiasing ¹¹	Fairness Model Certification
	Reducing Complexity	Feature Selection ¹⁰	Regularization ^{3 12}	
ACCOUNTABILITY	Enhancing Reliability	Drift Detection Data Augmentation ^{8 11 13}	Adversarial Training ¹⁵ Adv. Transformer Defenses Adversarial Robustness Certified Training & Inference ¹⁵	Adv. Postprocessing Defenses Adv. Detection Defenses Prediction Confidence Intervals
	Reducing Complexity	Feature Selection ¹⁰ Adv. Preprocessing Defenses ¹³	Regularization ^{3 12}	
	Mitigating Bias	Feature Engineering ^{10 12}	Monotonic Constraints ¹² (+ interaction/bi-variate constraints)	Calibrating/Equalizing Odds ^{7 11}
	Ensuring Privacy	Data Anonymization Differential Privacy	Federated Learning All Inference-attack Adversarial Defenses	Privacy-Preserving Inference
TRANSPARENCY	Reducing Complexity	Feature Selection ¹⁰	Regularization ^{3 12}	
	Enhancing Reliability	Feature Engineering ^{10 12}	Monotonic Constraints ¹²	Local Interpretation ^{6 7 8 9}

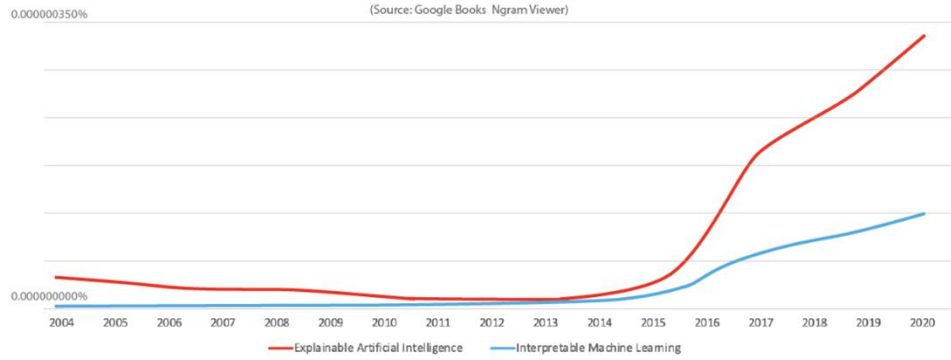
Academic Publications

(Source: Dimensions.ai)



Book Publications

(Source: Google Books Ngram Viewer)



Search Trends

(Source: Google Trends)

