

Ablation Study on Parameters and Algorithms for Calibrated Explanations

Tuwe Löfström

August 10, 2024

Abstract

This presents the information shown in [the ablation study](#) performed on 25 classification data sets. The main conclusion is that the choice of classifier influences the mean variance more than the size of the calibration set or the percentile sampling. Furthermore, runtime increases almost linearly with the number of percentiles sampled for numerical features, while calibration size has a much smaller effect. Factual explanations are faster than counterfactual explanations.

1 Ablation Analysis

The ablation analysis is focused on evaluating how the algorithm is affected by the calibration size and the number of percentiles sampled for numerical features. It is using a similar setup as the stability experiment, but with the following changes:

- The number of percentiles sampled for numerical features is varied between 1, 2, 3 (default), 4, and 9. The set of percentiles used are: [50], [33, 67], [25, 50, 75], [20, 40, 60, 80], [10, 20, 30, 40, 50, 60, 70, 80, 90]
- The calibration size is varied between 10%, 20% and 40% of the data not used for testing.
- Test size is fixed to 10% of the data.
- Only one repetition per percentile and calibration size is used.

Everything was run on 25 data sets. See the [Classification Experiment Ablation.py](#) for details on the experiment.

The tabulated results are the mean variance of the ablation measured per calibration size or percentile sampling. The variance is measured per instance and computed over the runs having the same calibration size/percentile sampling on the feature importance weight of the most influential feature, defined as the feature most often having highest absolute feature importance weight. The average variance is computed over the entire test set. The most influential feature is used since it is the feature that is most likely to be used in a decision but also the feature with the greatest expected variation (as a consequence of the weights having the highest absolute values).

1.1 Calibration Size

First out is a table with results per calibration size. Since different sampling sizes may result in different results for numerical features, the mean variance is only expected to be 0 for categorical-only datasets. The results are shown in [Table 1](#).

The most interesting observation from the results is that difference in mean variance is fairly low between the different calibration sizes. This indicates that the calibration size does not have a large impact on the feature importance weights. In fact, a smaller calibration set even tend to have a lower mean variance.

Dataset	Calibration Size			xGB			xGB			xGB			xGB			RF			RF			RF			RF		
	10%	20%	40%	CE	10%	20%	CE	10%	20%	40%	CCE	10%	20%	40%	CCE	CE	10%	20%	40%	CCE	10%	20%	40%	CCE	10%	20%	40%
colic	2.4e-05	1.8e-05	2.7e-05		2.4e-05	1.8e-05		6.9e-08	1.1e-06	7.4e-07		6.9e-08	1.1e-06	7.4e-07													
creditA	.00056	.00043	.0004		.00056	.00043	.0004	2.4e-05	4.1e-05	6e-05		2.4e-05	4.1e-05	6e-05													
diabetes	.00024	.00024	.00027		.00024	.00024	.00027	.00037	.00032	.0003		.00037	.00032	.0003													
german	2.2e-05	1.1e-05	2.2e-05		2.2e-05	1.1e-05	2.2e-05	5.1e-06	2.6e-06	1.8e-05		5.1e-06	2.6e-06	1.8e-05													
haberman	.0014	.00084	.00062		.0014	.00084	.00062	.00034	.00094	.00083		.00034	.00094	.00083													
heartC	7.6e-06	6.4e-05	4.7e-05		7.6e-06	6.4e-05	4.7e-05	2.7e-07	1.5e-05	1.4e-05		2.7e-07	1.5e-05	1.4e-05													
heartH	5.2e-05	3.2e-05	5.6e-05		5.2e-05	3.2e-05	5.6e-05	2.3e-06	1.3e-06	9e-07		2.3e-06	1.3e-06	9e-07													
hearts	.00015	.00026	.00025		.00015	.00026	.00025	3e-05	2.7e-05	5.5e-05		3e-05	2.7e-05	5.5e-05													
hepati	9.5e-05	.00017	.00014		9.5e-05	.00017	.00014	2e-05	4.7e-05	6.6e-05		2e-05	4.7e-05	6.6e-05													
iono	.00011	.00011	7.5e-05		.00011	.00011	7.5e-05	5.2e-05	5.2e-05	8.2e-05		5.2e-05	5.2e-05	8.2e-05													
je4042	.00016	.00022	.00026		.00016	.00022	.00026	3.5e-05	5e-05	4.5e-05		3.5e-05	5e-05	4.5e-05													
je4243	4.8e-05	5.9e-05	5.1e-05		4.8e-05	5.9e-05	5.1e-05	.00036	.00026	.00024		.00036	.00026	.00024													
kc1	.00025	.00017	.00016		.00025	.00017	.00016	.00037	.00022	.00017		.00037	.00022	.00017													
kc2	.00019	.0002	.00016		.00019	.0002	.00016	5.6e-06	3.4e-05	8.8e-05		5.6e-06	3.4e-05	8.8e-05													
kc3	5e-06	9e-06	7.2e-06		5e-06	9e-06	7.2e-06	7.3e-06	6.3e-06	4.5e-06		7.3e-06	6.3e-06	4.5e-06													
liver	.00014	.00013	.0001		.00014	.00013	.0001	.00061	.00044	.0005		.00061	.00044	.0005													
pc1req	.0	3.5e-35	2.9e-35		.0	3.5e-35	2.9e-35	6.3e-07	3.1e-07	2.1e-07		6.3e-07	3.1e-07	2.1e-07													
pc4	.00014	.00011	9.9e-05		.00014	.00011	9.9e-05	4.1e-05	4.7e-05	4.8e-05		4.1e-05	4.7e-05	4.8e-05													
sonar	1.3e-05	.00015	.0002		1.3e-05	.00015	.0002	4.8e-05	.00013	.00011		4.8e-05	.00013	.00011													
spect	1.4e-34	1.1e-34	7e-35		1.4e-34	1.1e-34	7e-35	.0	1.8e-35	2.3e-35		.0	1.8e-35	2.3e-35													
spectf	.00014	7e-05	8.5e-05		.00014	7e-05	8.5e-05	8.5e-05	5.9e-05	7.7e-05		8.5e-05	5.9e-05	7.7e-05													
transfusion	.00034	.00025	.00025		.00034	.00025	.00025	.00031	.00035	.00033		.00031	.00035	.00033													
ttt	2.8e-33	3.2e-33	3.3e-33		2.8e-33	3.2e-33	3.3e-33	1e-33	1.6e-33	1.7e-33		1e-33	1.6e-33	1.7e-33													
vote	1.3e-34	1.7e-34	1.2e-34		1.3e-34	1.7e-34	1.2e-34	1.3e-34	8.9e-35	1.1e-34		1.3e-34	8.9e-35	1.1e-34													
wbc	.00016	.00027	.00036		.00016	.00027	.00036	.00056	.00077	.00075		.00056	.00077	.00075													
Average	.00017	.00015	.00015		.00017	.00015	.00015	.00013	.00015	.00015		.00013	.00015	.00015													

Table 1: Calibration Size

Sample Size	xGB									RF								
	1	2	3	4	9	1	2	3	4	9	1	2	3	4	9			
Dataset	CE	CE	CE	CE	CE	CCE	CCE	CCE	CCE	CCE	CE	CE	CE	CE	CCE			
colic	.0026	.0027	.0027	.0027	.0027	.0026	.0027	.0027	.0027	.0027	.0043	.0043	.0043	.0043	.0043			
creditA	.007	.0066	.0064	.0062	.0061	.007	.0066	.0064	.0062	.0061	.003	.003	.003	.003	.0029			
diabetes	.0045	.0044	.0043	.0042	.0041	.0045	.0044	.0043	.0042	.0041	.0096	.0097	.0096	.0095	.0094			
german	.0011	.0011	.0011	.0011	.0011	.0011	.0011	.0011	.0011	.0011	.0027	.0027	.0026	.0026	.0026			
haberman	.012	.012	.012	.011	.011	.012	.012	.012	.011	.011	.0072	.0064	.0061	.0057	.0055			
heartC	.0043	.0043	.0043	.0043	.0043	.0043	.0043	.0043	.0043	.0043	.004	.004	.004	.004	.0041			
heartH	.0054	.0054	.0054	.0053	.0053	.0054	.0054	.0054	.0053	.0053	.0069	.0069	.0069	.0069	.0069			
hearts	.0096	.0087	.0084	.0083	.0082	.0096	.0087	.0084	.0083	.0082	.0054	.0053	.0052	.0052	.0052			
hepati	.0075	.0067	.0065	.0063	.0061	.0075	.0067	.0065	.0063	.0061	.001	.0011	.0012	.0011	.0011			
iono	.002	.002	.002	.002	.0021	.002	.002	.002	.002	.0021	.0041	.0041	.0041	.0041	.0039			
je4042	.0087	.0087	.0083	.0082	.008	.0087	.0087	.0083	.0082	.008	.0054	.0052	.0052	.0052	.0051			
je4243	.0027	.0027	.0026	.0026	.0026	.0027	.0027	.0026	.0026	.0026	.003	.0028	.0027	.0026	.0025			
ke1	.0064	.0064	.0063	.0062	.0061	.0064	.0064	.0063	.0062	.0061	.0048	.0041	.004	.0038	.0037			
ke2	.0045	.0043	.0042	.0042	.0042	.0045	.0043	.0042	.0042	.0042	.0039	.0038	.0038	.0037	.0037			
ke3	.00065	.00063	.00063	.00063	.00062	.00065	.00063	.00063	.00063	.00062	.0023	.0023	.0023	.0023	.0023			
liver	.0097	.0098	.0098	.0098	.0097	.0097	.0098	.0098	.0098	.0097	.0057	.0059	.0057	.0057	.0056			
pc1req	.0035	.0035	.0035	.0035	.0035	.0035	.0035	.0035	.0035	.0035	.004	.004	.004	.004	.004			
pc4	.00097	.00098	.001	.001	.001	.00097	.00098	.001	.001	.001	.00085	.00082	.00081	.00081	.00081			
sonar	.0095	.0087	.0084	.0084	.0083	.0095	.0087	.0084	.0084	.0083	.0048	.0048	.0047	.0048	.0048			
spect	.003	.003	.003	.003	.003	.003	.003	.003	.003	.003	.0021	.0021	.0021	.0021	.0021			
spectf	.0094	.0097	.0097	.0097	.0096	.0094	.0097	.0097	.0097	.0096	.0058	.0058	.0057	.0056	.0056			
transfusion	.0011	.001	.001	.00099	.00097	.0011	.001	.001	.00099	.00097	.0041	.0036	.0034	.0033	.0033			
ttt	.0037	.0037	.0037	.0037	.0037	.0037	.0037	.0037	.0037	.0037	.0017	.0017	.0017	.0017	.0017			
vote	.003	.003	.003	.003	.003	.003	.003	.003	.003	.003	.0023	.0023	.0023	.0023	.0023			
wbc	.0085	.0082	.0082	.0082	.0079	.0085	.0082	.0082	.0082	.0079	.002	.0018	.0016	.0016	.0015			
Average	.0053	.0051	.005	.005	.0049	.0053	.0051	.005	.005	.0049	.004	.0039	.0039	.0038	.0038			

Table 2: Percentile Sampling

Calibration Size	CE				CCE		
	10%	20%	40%		10%	20%	40%
xGB	0.11	0.11	0.11		0.20	0.21	0.22
RF	0.11	0.11	0.11		0.20	0.21	0.22

Table 3: Computation Time - Calibration Size

Sample Size	CE						CCE				
	1	2	3	4	9		1	2	3	4	9
xGB	0.05	0.06	0.09	0.11	0.25		0.08	0.12	0.16	0.21	0.49
RF	0.04	0.06	0.08	0.10	0.25		0.07	0.12	0.16	0.21	0.48

Table 4: Computation Time - Sample Size

1.2 Percentile Sampling

Table 2 are the results per percentile sampling.

Even if there is some difference in the mean variance when varying the percentile sampling, the difference can mainly be attributed to the difference in underlying ML algorithm. There is a tendency that a larger set of percentiles tend to reduce the mean variance, which is expected. However, the tendency is not very strong.

2 Computing time

Now, let's look at the runtime taken to compute the explanations. The tabulated runtimes are the average time in seconds per instance. Table 3 show the summary over all data sets divided by calibration size and Table 4 show the summary over all data sets divided by percentile sampling (indicated by the number of samples drawn).

Detailed results per calibration size and percentile sampling is shown per algorithm in Tables 5 and 6.

The results regarding runtime are as can be expected and the observations are summarized below:

- The runtime increases almost linearly with the number of percentiles sampled for numerical features.
- The runtime increases with the calibration size, even if the difference in runtime is fairly small.
- CE is faster than CCE, as expected. The reason is that CCE will generally require additional calculations than CE, at least for numerical features.
- The runtime tend to increase with the number of features, even if it is not a linear increase. This is due to the fact that categorical features with many categories are more expensive to compute, at least as long as the sampling size is small. Consequently, the number of categorical features together with the number of categories per feature is more important than the total number of features, especially when the sampling size (only affecting numerical features) is small.
- With the old implementation (where one call was made to the learner's predict_proba for each perturbed instance), a great part of the difference in runtime could be attributed to the underlying model, indicating that they had a large difference in overhead. With the new implementation, the difference is negligible.

3 Conclusion

The individual algorithmic parameter that influence runtime most is the number of percentiles sampled for numerical features. As this tend to have a fairly small impact on the feature importance, it may be a reason to consider decreasing the number of percentiles sampled by default for numerical features (currently the default is 3: [25, 50, 75]). Using only the median ([50]) would on average reduce the runtime with almost half.

3.1 Final Note

The core algorithm was updated early August 2024. The results reported here are using the code committed in [chore: redirected __call__\(\) to explain\(\)](#). Average speedups per instance over all data sets range from 2-9 times faster depending on setup, with substantially higher speedups of up to almost 40 times faster for individual data sets and setups.

[illegible]

[illegible]