

# Mitigating Adversarial Gray-Box Attacks Against Phishing Detectors

## Supplementary Material

### I. CONFIGURATION PARAMETERS

We provide in Table I the most influential parameter settings for each considered classifier of our evaluation.

### II. STANDARDIZED ATTACK TABLES

To facilitate understanding the effectiveness of our attacks, we provide the detailed results of the *attack* cases in Tables II to XIII. These tables have the same format as the those of the *no attack* case, i.e., Table IV and Table V of the main paper. In other words, Tables II to XIII show the F1-score, Accuracy, TPR and FPR of all the considered attacks (simple and complex) against the baseline and the

POC classifiers. We report the results for for GBA-1-GBA-3, and three of the seven variants of GBA- $\Delta$  (namely, when  $\Delta=10,40,70\%$ ).

By observing these tables we note that the FPR score is always the same between the *no-attack* and the *attack* scenarios. This is because our Gray Box attacks are evasion attacks, and aim to have a malicious (phishing) sample to be classified as benign. Hence, the false positive rate is not affected at all, because this metric would require a change in the “negative” (that is, benign) samples, which we do not touch in our evaluation. This reason motivates our focus on the Recall (or detection rate) metric for measuring the Impact of our attacks.

Table I: Most important parameters of each considered classifier. The DnW is an *EnsembleVotingClassifier* composed of DnW-1 and DnW-2.

Classifier	Type (as in <i>SciKit-learn</i> )	First Parameter	Second Parameter	Third Parameter
RF	RandomForestClassifier	<i>Estimators: 150</i>	<i>Criterion: Gini</i>	<i>maxFeatures: auto</i>
SVM	LinearSVC	<i>Loss: squaredHinge</i>	<i>Penalty: L2</i>	<i>C: 0.6476</i>
KNN	KNeighborsClassifier	<i>Neighbors: 1</i>	<i>Metric: minkowski</i>	<i>Algorithm: ballTree</i>
SGD	SGDClassifier	<i>Loss: squaredHinge</i>	<i>Penalty: ElasticNet</i>	<i>LearningRate: adaptive</i>
DT	DecisionTreeClassifier	<i>Splitter: best</i>	<i>Criterion: Gini</i>	<i>MaxDepth: none</i>
LR	LogisticRegression	<i>Penalty: L2</i>	<i>TOL: 0.0001</i>	<i>C: 1.0</i>
NB	GaussianNB	<i>Priors: None</i>	<i>Smoothing: 9e-07</i>	<i>N/A</i>
MLP	MLPClassifier	<i>Layers: 100</i>	<i>Activation: ReLU</i>	<i>Solver: adam</i>
AB	AdaBoostClassifier	<i>Estimators: 50</i>	<i>Algorithm: SAMME.R</i>	<i>LearningRate: 1.0</i>
ET	ExtraTreesClassifier	<i>Estimators: 340</i>	<i>Criterion: Entropy</i>	<i>MaxDepth: 40</i>
GB	GradientBoostingClassifier	<i>Estimators: 200</i>	<i>Criterion: MSE</i>	<i>Loss: deviance</i>
DnW-1	MLPClassifier	<i>Layers: 64, 16, 16, 4, 4</i>	<i>Activation: ReLU</i>	<i>Solver: adam</i>
DnW-2	MLPClassifier	<i>Layers: 256</i>	<i>Activation: ReLU</i>	<i>Solver: adam</i>
Bag	BaggingClassifier	<i>Estimators: 230</i>	<i>MaxSamples: 1000</i>	<i>BootstrapFeatures: true</i>





Table XII: Complex Attack Case (GBA- $\Delta$ ,  $\Delta = 40\%$ ): results of the POC classifiers for each dataset.

Classifier	LNU-Phish				DeltaPhish				Mendeley Phishing				UCI Phishing			
	FI	Acc	FPR	TPR	FI	Acc	FPR	TPR	FI	Acc	FPR	TPR	FI	Acc	FPR	TPR
RF	0.851	0.912	0.001	0.742	0.839	0.958	0.000	0.723	0.736	0.785	0.033	0.602	0.474	0.613	0.028	0.318
SVM	0.792	0.882	0.009	0.668	0.413	0.876	0.020	0.290	0.751	0.725	0.382	0.833	0.488	0.572	0.184	0.372
KNN	0.986	0.991	0.001	0.974	0.916	0.976	0.004	0.865	0.732	0.763	0.124	0.650	0.617	0.688	0.032	0.458
SGD	0.739	0.827	0.121	0.724	0.347	0.810	0.105	0.334	0.767	0.772	0.208	0.752	0.711	0.719	0.174	0.630
DT	0.921	0.950	0.004	0.860	0.739	0.938	0.000	0.586	0.583	0.676	0.104	0.454	0.597	0.678	0.026	0.435
LR	0.828	0.899	0.008	0.717	0.185	0.543	0.422	0.345	0.698	0.739	0.128	0.605	0.489	0.572	0.187	0.374
NB	0.841	0.894	0.071	0.826	0.484	0.825	0.125	0.544	0.752	0.683	0.599	0.968	0.673	0.683	0.209	0.594
MLP	0.824	0.899	0.001	0.702	0.580	0.904	0.014	0.440	0.641	0.689	0.180	0.557	0.282	0.524	0.045	0.170
AB	0.914	0.946	0.005	0.850	0.638	0.913	0.014	0.505	0.695	0.742	0.108	0.590	0.680	0.702	0.148	0.578
ET	0.977	0.985	0.000	0.955	0.777	0.945	0.000	0.636	0.704	0.766	0.029	0.559	0.958	0.955	0.028	0.940
GB	0.864	0.919	0.001	0.761	0.813	0.952	0.001	0.688	0.795	0.824	0.040	0.687	0.560	0.659	0.022	0.396
DnW	0.974	0.983	0.004	0.957	0.910	0.974	0.006	0.863	0.748	0.777	0.114	0.666	0.681	0.718	0.077	0.549
Bag	0.898	0.937	0.003	0.819	0.865	0.964	0.000	0.761	0.612	0.710	0.044	0.461	0.642	0.704	0.029	0.484

Table XIII: Complex Attack Case (GBA- $\Delta$ ,  $\Delta = 70\%$ ): results of the POC classifiers for each dataset.

Classifier	LNU-Phish				DeltaPhish				Mendeley Phishing				UCI Phishing			
	FI	Acc	FPR	TPR	FI	Acc	FPR	TPR	FI	Acc	FPR	TPR	FI	Acc	FPR	TPR
RF	0.623	0.815	0.001	0.454	0.646	0.921	0.000	0.477	0.494	0.655	0.033	0.339	0.930	0.927	0.028	0.890
SVM	0.536	0.782	0.009	0.372	0.450	0.881	0.020	0.323	0.446	0.509	0.382	0.398	0.756	0.753	0.184	0.700
KNN	0.981	0.988	0.001	0.965	0.843	0.958	0.004	0.746	0.374	0.569	0.124	0.259	0.647	0.706	0.032	0.491
SGD	0.502	0.722	0.121	0.414	0.442	0.828	0.105	0.451	0.779	0.782	0.208	0.772	0.840	0.827	0.174	0.828
DT	0.787	0.880	0.004	0.654	0.668	0.925	0.000	0.502	0.678	0.732	0.104	0.567	0.871	0.872	0.026	0.788
LR	0.717	0.848	0.008	0.568	0.239	0.560	0.422	0.457	0.434	0.594	0.128	0.313	0.858	0.843	0.187	0.867
NB	0.502	0.744	0.071	0.382	0.227	0.776	0.125	0.218	0.763	0.694	0.599	0.990	0.915	0.899	0.209	0.989
MLP	0.942	0.963	0.001	0.893	0.842	0.956	0.014	0.784	0.621	0.677	0.180	0.532	0.760	0.779	0.045	0.635
AB	0.764	0.869	0.005	0.624	0.438	0.883	0.014	0.302	0.706	0.749	0.108	0.605	0.817	0.810	0.148	0.775
ET	0.839	0.906	0.000	0.723	0.505	0.900	0.000	0.337	0.495	0.656	0.029	0.339	0.861	0.863	0.028	0.773
GB	0.793	0.884	0.001	0.658	0.703	0.931	0.001	0.545	0.599	0.703	0.040	0.444	0.957	0.953	0.022	0.933
DnW	0.852	0.912	0.004	0.748	0.772	0.942	0.006	0.649	0.510	0.635	0.114	0.382	0.534	0.629	0.077	0.387
Bag	0.972	0.982	0.003	0.952	0.807	0.951	0.000	0.676	0.492	0.650	0.044	0.340	0.735	0.765	0.029	0.595