

## A. Proof of Theorem 1

Here we prove Theorem 1 for information gain score.

*Proof.*  $H(y)$  and  $H(y|x^{(j)} < \eta)$  are defined as

$$H(y) = -\frac{|\mathcal{I}_0|}{|\mathcal{I}|} \log\left(\frac{|\mathcal{I}_0|}{|\mathcal{I}|}\right) - \frac{|\mathcal{I}_1|}{|\mathcal{I}|} \log\left(\frac{|\mathcal{I}_1|}{|\mathcal{I}|}\right),$$

and

$$\begin{aligned} H(y|x^{(j)} < \eta) = & -\frac{|\mathcal{I}_L|}{|\mathcal{I}|} \left[ \frac{|\mathcal{I}_L \cap \mathcal{I}_0|}{|\mathcal{I}_L|} \log\left(\frac{|\mathcal{I}_L \cap \mathcal{I}_0|}{|\mathcal{I}_L|}\right) + \frac{|\mathcal{I}_L \cap \mathcal{I}_1|}{|\mathcal{I}_L|} \log\left(\frac{|\mathcal{I}_L \cap \mathcal{I}_1|}{|\mathcal{I}_L|}\right) \right] \\ & -\frac{|\mathcal{I}_R|}{|\mathcal{I}|} \left[ \frac{|\mathcal{I}_R \cap \mathcal{I}_0|}{|\mathcal{I}_R|} \log\left(\frac{|\mathcal{I}_R \cap \mathcal{I}_0|}{|\mathcal{I}_R|}\right) + \frac{|\mathcal{I}_R \cap \mathcal{I}_1|}{|\mathcal{I}_R|} \log\left(\frac{|\mathcal{I}_R \cap \mathcal{I}_1|}{|\mathcal{I}_R|}\right) \right]. \end{aligned}$$

For simplicity, we denote  $N_0 := |\mathcal{I}_0|$ ,  $N_1 := |\mathcal{I}_1|$ ,  $n_0 := |\mathcal{I}_L \cap \mathcal{I}_0|$  and  $n_1 := |\mathcal{I}_L \cap \mathcal{I}_1|$ . The information gain of this split can be written as a function of  $n_0$  and  $n_1$ :

$$\begin{aligned} IG = & C_1 \left[ n_0 \log\left(\frac{n_0}{N_0(n_1 + n_0)}\right) + n_1 \log\left(\frac{n_1}{N_1(n_1 + n_0)}\right) \right. \\ & + (N_0 - n_0) \log\left(\frac{N_0 - n_0}{N_0(N_1 + N_0 - n_1 - n_0)}\right) \\ & \left. + (N_1 - n_1) \log\left(\frac{N_1 - n_1}{N_1(N_1 + N_0 - n_1 - n_0)}\right) \right] + C_2, \end{aligned} \quad (5)$$

where  $C_1 > 0$  and  $C_2$  are constants with respect to  $n_0$ . Taking  $n_0$  as a continuous variable, we have

$$\frac{\partial IG}{\partial n_0} = C_1 \cdot \log\left(1 + \frac{n_0 N_1 - N_0 n_1}{(N_0 - n_0)(n_1 + n_0)}\right) \quad (6)$$

When  $\frac{\partial IG}{\partial n_0} < 0$ , perturbing one example in  $\Delta\mathcal{I}_R$  with label 0 to  $\mathcal{I}_L$  will increase  $n_0$  and decrease the information gain. It is easy to see that  $\frac{\partial IG}{\partial n_0} < 0$  if and only if  $\frac{n_0}{N_0} < \frac{n_1}{N_1}$ . This indicates that when  $\frac{n_0}{N_0} < \frac{n_1}{N_1}$  and  $\frac{n_0+1}{N_0} \leq \frac{n_1}{N_1}$ , perturbing one example with label 0 to  $\mathcal{I}_L$  will always decrease the information gain.  $\square$

Similarly, if  $\frac{n_1}{N_1} < \frac{n_0}{N_0}$  and  $\frac{n_1+1}{N_1} \leq \frac{n_0}{N_0}$ , perturbing one example in  $\Delta\mathcal{I}_R$  with label 1 to  $\mathcal{I}_L$  will decrease the information gain. As mentioned in the main text, to decrease the information gain score in Algorithm 1, the adversary needs to perturb examples in  $\Delta\mathcal{I}$  such that  $\frac{n_0}{N_0}$  and  $\frac{n_1}{N_1}$  are close to each other. Algorithm 3 gives an  $O(|\Delta\mathcal{I}|)$  method to find  $\Delta n_0^*$  and  $\Delta n_1^*$ , the optimal number of points in  $\Delta\mathcal{I}$  with label 0 and 1 to be added to the left.

## B. Gini Impurity Score

We also have a theorem for Gini impurity score similar to Theorem 1.

**Algorithm 3** Finding  $\Delta n_0^*$  and  $\Delta n_1^*$  to Minimize Information Gain or Gini Impurity

**Input:**  $N_0$  and  $N_1$ , number of instances with label 0 and 1.  $n_0^o$  and  $n_1^o$ , number of instances with label 0 and 1 that are certainly on the left.

**Input:**  $|\Delta\mathcal{I} \cap \mathcal{I}_0|$  and  $|\Delta\mathcal{I} \cap \mathcal{I}_1|$ , number of instances with label 0 and 1 that can be perturbed.

**Output:**  $\Delta n_0^*$ ,  $\Delta n_1^*$ , optimal number of points with label 0 and 1 in  $\Delta\mathcal{I}$  to be place on the left.

$\Delta n_0^* \leftarrow 0$ ,  $\Delta n_1^* \leftarrow 0$ ,  $\text{min\_diff} \leftarrow \left| \frac{n_0^o}{N_0} - \frac{n_1^o}{N_1} \right|$ ;

**for**  $\Delta n_0 \leftarrow 0$  **to**  $|\Delta\mathcal{I} \cap \mathcal{I}_0|$  **do**

$\text{ceil} \leftarrow \lceil \frac{N_1(n_0^o + \Delta n_0)}{N_0} \rceil - n_1^o$ ;

$\text{floor} \leftarrow \lfloor \frac{N_1(n_0^o + \Delta n_0)}{N_0} \rfloor - n_1^o$ ;

**for**  $\Delta n_1' \in \{\text{ceil}, \text{floor}\}$  **do**

$\Delta n_1 \leftarrow \max\{\min\{\Delta n_1', |\Delta\mathcal{I} \cap \mathcal{I}_1|\}, 0\}$ ;

**if**  $\text{min\_diff} > \left| \frac{\Delta n_0 + n_0^o}{N_0} - \frac{\Delta n_1 + n_1^o}{N_1} \right|$  **then**

$\Delta n_0^* \leftarrow \Delta n_0$ ,  $\Delta n_1^* \leftarrow \Delta n_1$ ,  $\text{min\_diff} \leftarrow$

$\left| \frac{\Delta n_0 + n_0^o}{N_0} - \frac{\Delta n_1 + n_1^o}{N_1} \right|$ ;

**end if**

**end for**

**end for**

Return  $\Delta n_0^*$  and  $\Delta n_1^*$ ;

**Theorem B.1.** If  $\frac{n_0}{N_0} < \frac{n_1}{N_1}$  and  $\frac{n_0+1}{N_0} \leq \frac{n_1}{N_1}$ , perturbing one example in  $\Delta\mathcal{I}_R$  with label 0 to  $\mathcal{I}_L$  will decrease the Gini impurity.

*Proof.* The Gini impurity score of a split with threshold  $\eta$  on feature  $j$  is

$$\begin{aligned} Gini = & \left(1 - \frac{|\mathcal{I}_0|^2}{|\mathcal{I}|^2} - \frac{|\mathcal{I}_1|^2}{|\mathcal{I}|^2}\right) \\ & - \frac{|\mathcal{I}_L|}{|\mathcal{I}|} \left(1 - \frac{|\mathcal{I}_0 \cap \mathcal{I}_L|^2}{|\mathcal{I}_L|^2} - \frac{|\mathcal{I}_1 \cap \mathcal{I}_L|^2}{|\mathcal{I}_L|^2}\right) \\ & - \frac{|\mathcal{I}_R|}{|\mathcal{I}|} \left(1 - \frac{|\mathcal{I}_0 \cap \mathcal{I}_R|^2}{|\mathcal{I}_R|^2} - \frac{|\mathcal{I}_1 \cap \mathcal{I}_R|^2}{|\mathcal{I}_R|^2}\right) \\ = & C_3 \left[ \frac{n_0^2 + n_1^2}{n_1 + n_0} + \frac{(N_0 - n_0)^2 + (N_1 - n_1)^2}{(N_0 + N_1 - n_0 - n_1)} \right] + C_4, \end{aligned} \quad (7)$$

where we use the same notation as in (5).  $C_3 > 0$  and  $C_4$  are constants with respect to  $n_0$ . Taking  $n_0$  as a continuous variable, we have

$$\frac{\partial Gini}{\partial n_0} = 2C_3 \frac{m_1 m_0 (n_0 m_1 + n_1 m_0 + 2n_1 m_1)}{(n_0 + n_1)^2 (m_0 + m_1)^2} \left( \frac{n_0}{m_0} - \frac{n_1}{m_1} \right), \quad (8)$$

where  $m_0 := N_0 - n_0$  and  $m_1 := N_1 - n_1$ . Then  $\frac{\partial Gini}{\partial n_0} < 0$  holds if  $\frac{n_0}{m_0} < \frac{n_1}{m_1}$ , which is equivalent to  $\frac{n_0}{N_0} < \frac{n_1}{N_1}$ .  $\square$

Since the conditions of Theorem 1 and Theorem B.1 are the same, Algorithm 1 and Algorithm 3 also work for tree-based models using Gini impurity score.

### C. Decision Boundaries of Robust and Natural Models

Figure 4 shows the decision boundaries and test accuracy of natural trees as well as robust trees with different  $\epsilon$  values on two dimensional synthetic datasets. All trees have depth 5 and we plot training examples in the figure. The results show that the decision boundaries of our robust decision trees are simpler than the decision boundaries in natural decision trees, agreeing with the regularization argument in the main text.

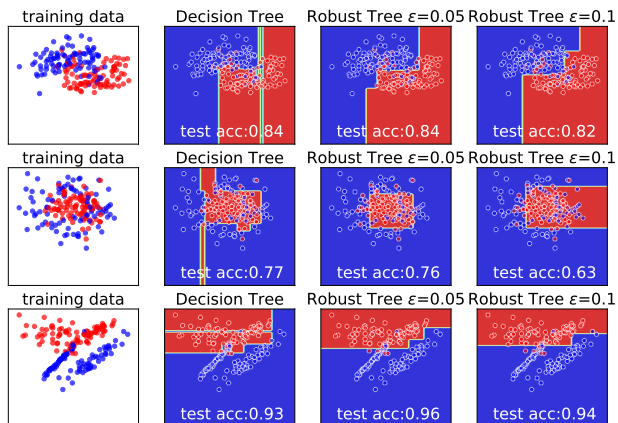


Figure 4. (Best viewed in color) The decision boundaries and test accuracy of natural decision trees and robust decision trees with depth 5 on synthetic datasets with two features.

### D. Omitted Results on $l_1$ and $l_2$ distortion

In Tables 4 and 5 we present the  $l_1$  and  $l_2$  distortions of vanilla (information gain based) decision trees and GBDT models obtained by Kantchelian’s  $l_1$  and  $l_2$  attacks. Again, only small or medium sized binary classification models can be evaluated by Kantchelian’s attack. From the results we can see that although our robust decision tree training algorithm is designed for  $l_\infty$  perturbations, it can also improve models  $l_1$  and  $l_2$  robustness significantly.

### E. Omitted Results on Models with Different Number of Trees

Figure 5 shows the  $l_\infty$  distortion and accuracy of Fashion-MNIST GBDT models with different number of trees. In Table 7 we present the test accuracy and  $l_\infty$  distortion of models with different number of trees obtained by Cheng’s  $l_\infty$  attack. For each dataset, models are generated during a single boosting run. We can see that the robustness of

robustly trained models consistently outperforms that of natural models with the same number of trees. Another interesting finding is that for MNIST and Fashion-MNIST datasets in Figures 3 (in the main text) and 5, models with more trees are generally more robust. This may not be true in other datasets; for example, results from Table 7 in the Appendix shows that on some other datasets, the natural GBDT models lose robustness when more trees are added.

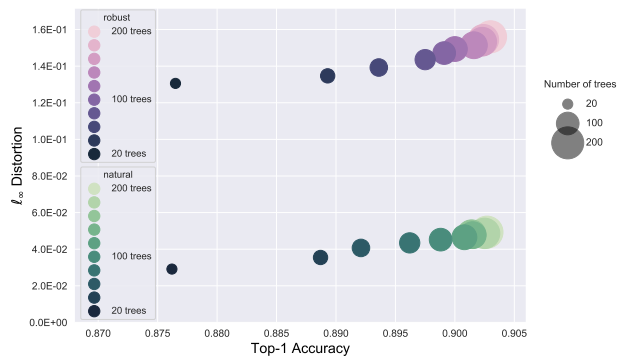


Figure 5. (Best viewed in color)  $l_\infty$  distortion vs. classification accuracy of GBDT models on Fashion-MNIST datasets with different numbers of trees (circle size). The adversarial examples are found by Cheng’s  $l_\infty$  attack. The robust training parameter  $\epsilon = 0.1$  for Fashion-MNIST. With robust training (purple) the distortion needed to fool a model increases dramatically with less than 1% accuracy loss.

### F. Reducing Depth Does Not Improve Robustness

One might hope that one can simply reduce the depth of trees to improve robustness since shallower trees provide stronger regularization effects. Unfortunately, this is not true. As demonstrated in Figure 6, the robustness of naturally trained GBDT models are much worse when compared to robust models, no matter how shallow they are or how many trees are in the ensemble. Also, when the number of trees in the ensemble model is limited, reducing tree depth will significantly lower the model accuracy.

### G. Random Forest Model Results

We test our robust training framework on random forest (RF) models and our results are in Table 6. In these experiments we build random forest models with 0.5 data sampling rate and 0.5 feature sampling rate. We test the robust and natural random forest model on three datasets and in each dataset, we tested 100 points using Cheng’s and Kantchelian’s  $l_\infty$  attacks. From the results we can see that our robust decision tree training framework can also significantly improve random forest model robustness.

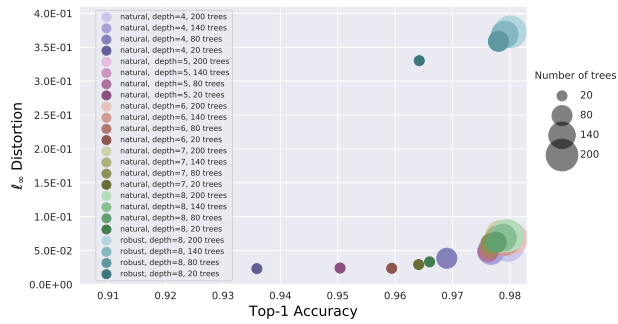


Figure 6. (Best viewed in color) Robustness vs. classification accuracy plot of GBDT models on MNIST dataset with different depth and different numbers of trees. The adversarial examples are found by Cheng’s  $\ell_\infty$  attack. The robust training parameter  $\epsilon = 0.3$ . Reducing the model depth cannot improve robustness effectively compared to our proposed robust training procedure.

## H. More MNIST and Fashion-MNIST Adversarial Examples

In Figure 7 we present more adversarial examples for MNIST and Fashion-MNIST datasets using GBDT models.

## Robust Decision Trees Against Adversarial Examples

Dataset	training	test	# of	# of	robust $\epsilon$	depth		test acc.		avg. $\ell_1$ dist. by Kantchelian's $\ell_1$ attack		avg. $\ell_2$ dist. by Kantchelian's $\ell_2$ attack		
	set size	set size	features	classes		robust	natural	robust	natural	robust	natural	robust	natural	
breast-cancer	546	137	10	2	0.3	5	5	.948	.942	<b>.534</b>		.270		.209
diabetes	614	154	8	2	0.2	5	5	.688	.747	<b>.204</b>		.075		.065
ionosphere	281	70	34	2	0.2	4	4	.986	.929	<b>.358</b>		.127		.106

Table 4. The test accuracy and robustness of information gain based single decision tree models. The robustness is evaluated by the average  $\ell_1$  and  $\ell_2$  distortions of adversarial examples found by Kantchelian's  $\ell_1$  and  $\ell_2$  attacks. Average  $\ell_\infty$  distortions of robust decision tree models found by the two attack methods are consistently larger than those of the naturally trained ones.

Dataset	training	test	# of	# of	# of	robust $\epsilon$	depth		test acc.		avg. $\ell_1$ dist. by Kantchelian's $\ell_1$ attack		dist. improv.	avg. $\ell_2$ dist. by Kantchelian's $\ell_2$ attack		dist. improv.		
	set size	set size	features	classes	trees		robust	natural	robust	natural	robust	natural		robust	natural			
breast-cancer	546	137	10	2	4	0.3	8	6	.978	.964	<b>.488</b>		.328	<b>1.49X</b>	<b>.431</b>		.251	<b>1.72X</b>
cod-rna	59,535	271,617	8	2	80	0.2	5	4	.880	.965	<b>.065</b>		.059	<b>1.10X</b>	<b>.062</b>		.047	<b>1.32X</b>
diabetes	614	154	8	2	20	0.2	5	5	.786	.773	<b>.150</b>		.081	<b>1.85X</b>	<b>.135</b>		.059	<b>2.29X</b>
ijcnn1	49,990	91,701	22	2	60	0.1	8	8	.959	.980	<b>.057</b>		.051	<b>1.12X</b>	<b>.048</b>		.042	<b>1.14X</b>
MNIST 2 vs. 6	11,876	1,990	784	2	1000	0.3	6	4	.997	.998	<b>1.843</b>		.721	<b>2.56X</b>	<b>.781</b>		.182	<b>4.29X</b>

Table 5. The test accuracy and robustness of GBDT models. Average  $\ell_1$  and  $\ell_2$  distortions of robust GBDT models are consistently larger than those of the naturally trained models. The robustness is evaluated by the average  $\ell_1$  and  $\ell_2$  distortions of adversarial examples found by Kantchelian's  $\ell_1$  and  $\ell_2$  attacks.

Dataset	training	test	# of	# of	# of	robust $\epsilon$	depth		test acc.		avg. $\ell_\infty$ dist. by Cheng's $\ell_\infty$ attack		dist. improv.	avg. $\ell_\infty$ dist. by Kantchelian's $\ell_\infty$ attack		dist. improv.		
	set size	set size	features	classes	trees		robust	natural	robust	natural	robust	natural		robust	natural			
breast-cancer	546	137	10	2	60	0.3	8	6	.993	.993	<b>.406</b>		.297	<b>1.37X</b>	<b>.396</b>		.244	<b>1.62X</b>
diabetes	614	154	8	2	60	0.2	5	5	.753	.760	<b>.185</b>		.093	<b>1.99X</b>	<b>.154</b>		.072	<b>2.14X</b>
MNIST 2 vs. 6	11,876	1,990	784	2	1000	0.3	6	4	.986	.983	<b>.445</b>		.180	<b>2.47X</b>	<b>.341</b>		.121	<b>2.82X</b>

Table 6. The test accuracy and robustness of random forest models. Average  $\ell_\infty$  distortion of our robust GBDT models are consistently larger than those of the naturally trained models. The robustness is evaluated by the average  $\ell_\infty$  distortion of adversarial examples found by Cheng's and Kantchelian's attacks.

Robust Decision Trees Against Adversarial Examples

dataset (c)	train	test	feat.	# of trees	1		2		3		4		5		6		7		8		9		10	
					rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.
breast-cancer (2) $\epsilon = 0.3$ depth <sub>r</sub> = 8, depth <sub>n</sub> = 6	546	137	10	model	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.
				tst. acc.	.985	.942	.971	.964	.978	.956	.978	.964	.985	.964	.985	.964	.985	.971	.993	.971	.993	.971	1.00	.971
				$\ell_\infty$ dist.	.383	.215	.396	.229	.411	.216	.411	.215	.406	.226	.407	.229	.406	.248	.439	.234	.439	.238	.437	.241
covtype (7) $\epsilon = 0.2$ depth <sub>r</sub> = depth <sub>n</sub> = 8	400,000	181,000	54	model	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.
				tst. acc.	.775	.828	.809	.850	.832	.865	.847	.877	.858	.891	.867	.902	.875	.912	.882	.921	.889	.926	.894	.930
				$\ell_\infty$ dist.	.125	.066	.103	.064	.087	.062	.081	.061	.079	.060	.077	.059	.077	.058	.075	.056	.075	.056	.073	.055
cod-rna (2) $\epsilon = 0.2$ depth <sub>r</sub> = 5, depth <sub>n</sub> = 4	59,535	271,617	8	model	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.
				tst. acc.	.810	.947	.861	.959	.874	.963	.880	.965	.892	.966	.900	.967	.903	.967	.915	.967	.922	.967	.925	.968
				$\ell_\infty$ dist.	.077	.057	.066	.055	.063	.054	.062	.053	.059	.053	.057	.052	.056	.052	.056	.052	.056	.052	.058	.052
diabetes (2) $\epsilon = 0.2$ depth <sub>r</sub> = depth <sub>n</sub> = 5	614	154	8	model	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.
				tst. acc.	.760	.753	.760	.753	.766	.753	.773	.753	.773	.734	.779	.727	.779	.747	.779	.760	.779	.773	.786	.773
				$\ell_\infty$ dist.	.163	.066	.163	.065	.154	.071	.151	.071	.152	.073	.148	.072	.146	.067	.144	.062	.138	.062	.139	.060
Fashion-MNIST (10) $\epsilon = 0.1$ depth <sub>r</sub> = depth <sub>n</sub> = 8	60,000	10,000	784	model	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.
				tst. acc.	.877	.876	.889	.889	.894	.892	.898	.896	.899	.899	.900	.901	.902	.902	.902	.901	.902	.903	.903	.903
				$\ell_\infty$ dist.	.131	.029	.135	.035	.139	.041	.144	.043	.147	.045	.149	.047	.151	.048	.153	.048	.154	.049	.156	.049
HIGGS (2) $\epsilon = 0.05$ depth <sub>r</sub> = depth <sub>n</sub> = 8	10,500,000	500,000	28	model	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.
				tst. acc.	.676	.747	.688	.753	.700	.755	.702	.758	.705	.759	.709	.760	.711	.762	.712	.764	.716	.763	.718	.764
				$\ell_\infty$ dist.	.023	.013	.023	.014	.022	.014	.022	.014	.022	.014	.022	.014	.021	.015	.021	.015	.021	.015	.021	.015
ijcnn1 (2) $\epsilon = 0.1$ depth <sub>r</sub> = depth <sub>n</sub> = 8	49,990	91,701	22	model	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.
				tst. acc.	.933	.973	.942	.977	.947	.977	.952	.979	.958	.980	.959	.980	.962	.980	.964	.980	.967	.980	.968	.980
				$\ell_\infty$ dist.	.065	.048	.061	.047	.058	.048	.057	.047	.054	.048	.054	.047	.054	.047	.053	.047	.052	.047	.052	.047
MNIST (10) $\epsilon = 0.3$ depth <sub>r</sub> = depth <sub>n</sub> = 8	60,000	10,000	784	model	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.
				tst. acc.	.964	.966	.973	.975	.977	.977	.978	.978	.978	.978	.979	.979	.979	.979	.980	.979	.980	.979	.980	.980
				$\ell_\infty$ dist.	.330	.033	.343	.049	.352	.057	.359	.062	.363	.065	.367	.067	.369	.069	.370	.071	.371	.072	.373	.072
Sensorless (11) $\epsilon = 0.05$ depth <sub>r</sub> = depth <sub>n</sub> = 6	48,509	10,000	48	model	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.
				tst. acc.	.834	.977	.867	.983	.902	.987	.923	.991	.945	.992	.958	.994	.966	.996	.971	.996	.974	.997	.978	.997
				$\ell_\infty$ dist.	.037	.022	.036	.022	.035	.023	.035	.023	.035	.023	.035	.023	.035	.023	.035	.023	.035	.023	.035	.023
webspam (2) $\epsilon = 0.05$ depth <sub>r</sub> = depth <sub>n</sub> = 8	300,000	50,000	254	model	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.	rob.	nat.
				tst. acc.	.950	.976	.964	.983	.970	.986	.973	.989	.976	.990	.978	.990	.980	.991	.981	.991	.982	.992	.983	.992
				$\ell_\infty$ dist.	.049	.010	.048	.015	.049	.019	.049	.021	.049	.023	.049	.024	.049	.024	.049	.024	.048	.024	.049	.024

Table 7. The test accuracy and robustness of GBDT models. Here depth<sub>n</sub> is the depth of natural trees and depth<sub>r</sub> is the depth of robust trees. Robustness is evaluated by the average  $\ell_\infty$  distortion of adversarial examples found by Cheng’s attack (Cheng et al., 2019). The number in the parentheses after each dataset name is the number of classes. Models are generated during a single boosting run. We can see that the robustness of our robust models consistently outperforms that of natural models with the same number of trees.

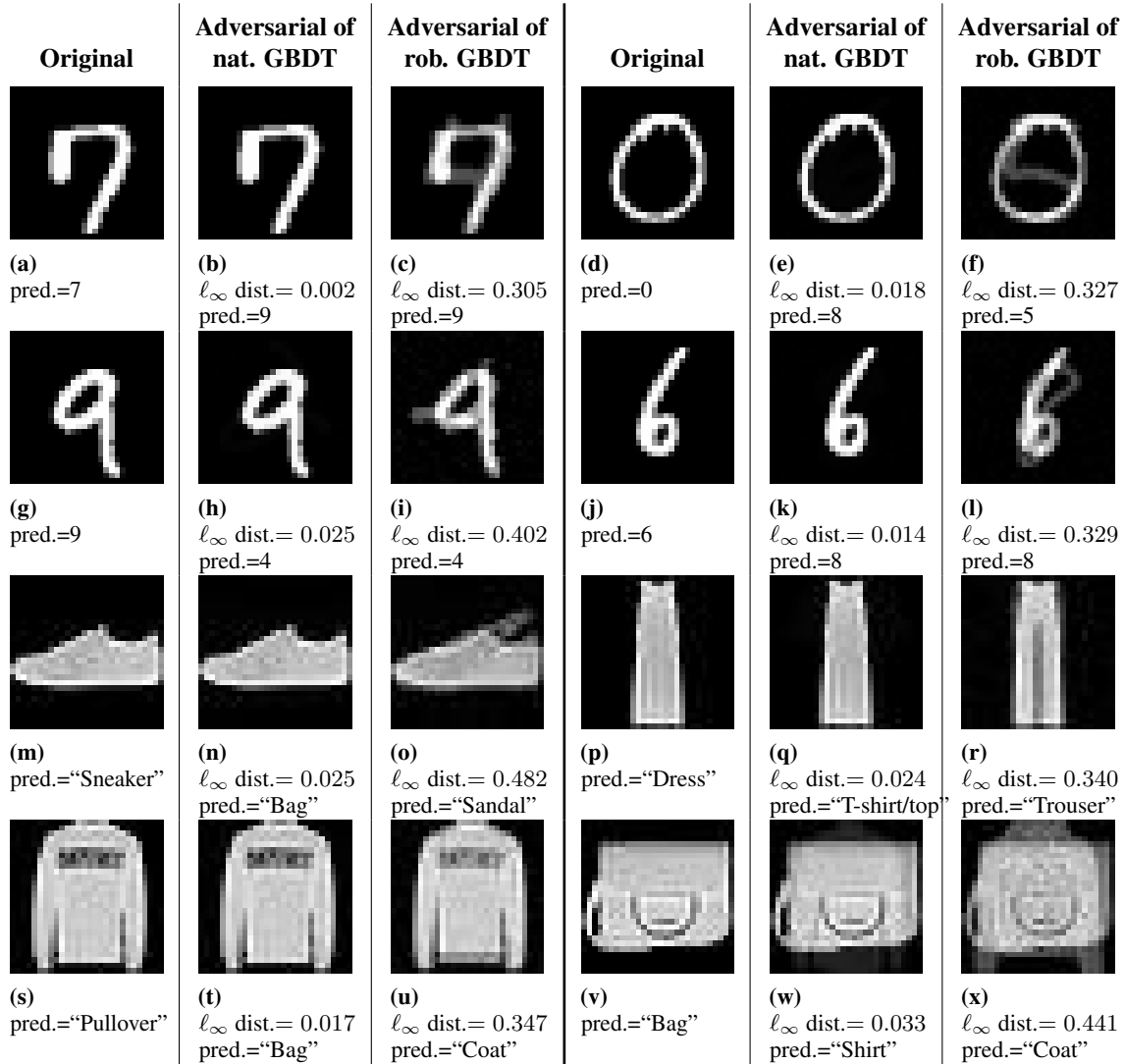


Figure 7. MNIST and Fashion-MNIST examples and their adversarial examples found using the untargeted Cheng’s  $\ell_\infty$  attack (Cheng et al., 2019) on 200-tree gradient boosted decision tree (GBDT) models trained using XGBoost with depth=8. For both MNIST and Fashion-MNIST robust models, we use  $\epsilon = 0.3$ .