
Adversarial Robustness for Tabular Data through Cost and Utility Awareness

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Many machine learning applications (credit scoring, fraud detection, etc.) use
2 data in the *tabular domains*. Adversarial examples can be especially damaging
3 for these applications. Yet, existing works on adversarial robustness mainly focus
4 on machine-learning models in the image and text domains. We argue that due to
5 the differences between tabular data and images or text, existing threat models are
6 inappropriate for tabular domains. These models do not capture that cost can be
7 more important than imperceptibility, nor that the adversary could ascribe different
8 value to the utility obtained from deploying different adversarial examples. We
9 show that due to these differences the attack and defence methods used for images
10 and text cannot be directly applied to the tabular setup. We address these issues
11 by proposing new cost and utility-aware threat models tailored to capabilities and
12 constraints of attackers targeting tabular domains. We show that our approach is
13 effective on two tabular datasets corresponding to applications for which adversarial
14 examples can have economic and social implications.

15 1 Introduction

16 Adversarial examples are inputs deliberately crafted by an adversary to cause a classification mistake.
17 They pose a threat in applications for which such mistakes can have a negative impact in deployed
18 models (e.g., a financial loss [1] or a security breach [2–4]). The literature on adversarial examples
19 largely focuses on image [5–10] and text domains [11–16]. Yet, many of the applications where
20 adversarial examples are most damaging or helpful are not images or text. Fraud and abuse detection
21 systems [17], risk-scoring systems [1], operate on *tabular data*: A cocktail of categorical, ordinal, and
22 numeric features. As opposed to images, each of these features has its own different semantics. The
23 properties of the image domain have shaped the way adversarial examples and adversarial robustness
24 are approached in the literature [8], and have greatly influenced adversarial robustness research in the
25 text domain. In this paper, we argue that, in tabular domains, adversarial examples are of a different
26 nature and adversarial robustness has a different meaning.

27 We argue that two high-level differences need to be addressed: (a) “imperceptibility”, the main
28 constraint in existing image and text adversarial examples, is ill-defined and can be irrelevant for
29 tabular data; and (b) existing adversarial examples assume all adversarial inputs have the same value
30 for the adversary, while in tabular domains different adversarial examples can bring different gain to
31 the adversary. Authors in the literature commonly formalize the concept of “an example deliberately
32 crafted to cause a misclassification” as a *natural example*, i.e., an example coming from the data
33 distribution, that is *imperceptibly* modified by an adversary so that the classifier’s decision changes.
34 Typically, they formalize imperceptibility as closeness according to a mathematical distance such as
35 L_p [18, 19]. In tabular data, however, imperceptibility is not necessarily relevant. Let us consider the

36 following fraud detection toy example: An adversary aims to create a fraudulent financial transaction
37 (e.g., using stolen credit card credentials). in an app such as PayPal. Assume a fraud detector takes as
38 input two features: (1) `transaction amount`, and (2) `device` from which the transaction was sent.

39 In this example, *imperceptibility is not well-defined*. Is a modification to the `amount` feature from
40 \$200 to \$201 imperceptible? What increase or decrease would we consider perceptible? The issue is
41 even more apparent with categorical data, for which standard distances such as L_2 , L_∞ cannot even
42 capture imperceptibility: Is a change of the `device` feature from Android to an iPhone imperceptible?

43 Even if imperceptibility was well-defined, *imperceptibility might not be relevant*. Should we only be
44 concerned about adversaries making “imperceptible” changes, e.g., modifying `amount` from \$200
45 to \$201? What about attack vectors in which the adversary evades detection while changing the
46 transaction by a “perceptible” amount –from \$200 to \$2,000?

47 We argue that in tabular data the primary constraint should be *adversarial cost*, rather than imper-
48 ceptibility. Instead of looking at how visually or semantically similar feature vectors are, the focus
49 should be on *how costly it is for an adversary to enact a modification*. Costs capture the effort
50 of the adversary, e.g., financial or computational. “How much money does the adversary have to
51 spend to evade the detector?” better captures the possibility that an adversary deploys an attack than
52 establishing a threshold on the L_p distance the adversary would tolerate.

53 **Different tabular adversarial examples are of different value to the adversary.** In the literature,
54 except for Zhang and Evans [20], most formalizations of adversarial robustness implicitly consider
55 that all adversarial examples are equal in their importance [6, 10, 21–23]. In tabular data domains,
56 however, different adversarial examples can bring very different *gain* to the adversary. In the fraud
57 detection example, if a fraudulent transaction with `transaction amount` of \$2,000 successfully evades
58 the detector, it could be significantly more profitable than a transaction with `amount` of \$200.

59 Using cost as the primary constraint for adversarial examples provides a natural way to incorporate
60 the variability in adversarial gain. The adversary is expected to care about the profit that they would
61 obtain from the attack, i.e., the difference between the cost associated with crafting an adversarial
62 example, and the gain from successfully using it. We call this profit the *utility* of the attack. We show
63 how utility can be incorporated into the design of attacks to ensure their economic profitability, and
64 into the design of defences to ensure protection against adversaries that focus on profit.

65 In this paper, we introduce a framework to build adversarial examples tailored to tabular data. Our
66 contributions are:

- 67 • We propose two *adversarial objectives* for tabular data that address the limitations of
68 traditional adversary examples: a *cost-bounded* objective that substitutes standard impercep-
69 tibility constraints with adversarial costs; and a novel *utility-bounded* objective in which the
70 adversary adjusts their expenditure on different adversarial examples proportionally to the
71 potential gains from deploying them.
- 72 • We perform an empirical evaluation of attacks and defences with respect to proposed
73 objectives in realistic conditions demonstrating their applicability to real-world security
74 scenarios.

75 2 Adversarial Objectives in Tabular Data

76 **Notations.** The input domain’s *feature space* \mathbb{X} is composed of m features: $\mathbb{X} \subseteq \mathbb{X}_1 \times \mathbb{X}_2 \times \dots \times \mathbb{X}_m$.
77 For an example $x \in \mathbb{X}$, we denote the value of its i -th feature as x_i . Features x_i can be categorical,
78 ordinal, and numeric. Each example is associated to a binary label $y \in \{0, 1\}$. We assume the
79 adversary’s *target* to be a binary classifier $f(x) \in \{0, 1\}$ that aims to predict the class y to which an
80 example x belongs. In terms of capabilities, we assume the adversary can only perform modifications
81 that are within the domain constraints. E.g in the fraud-detection example, the adversary can change
82 the `transaction amount`, but the value must be positive. For a given initial labelled example (x, y) ,
83 we denote the set of feasible adversarial examples that can be reached within the capabilities of the
84 adversary as $\mathcal{F}(x, y) \subseteq \mathbb{X}$.

85 **Cost-Bounded Objective.** In the standard way to obtain an adversarial example [10], the adversary
86 aims to construct an example that maximizes the loss $\ell(\cdot, \cdot)$ incurred by the target classifier, while

87 keeping the L_p -distance from the initial example bounded:

$$\max_{x' \in \mathcal{F}(x, y)} \ell(\eta(x'), y) \quad \text{s.t. } \|x' - x\|_p \leq \varepsilon \quad (1)$$

88 This objective implicitly assumes that the adversary wants to keep the adversarial example as similar
89 to the initial example as possible in terms of the examples' feature values.

90 Formally, we associate a cost to the modifications needed to generate any adversarial example
91 $x' \in \mathcal{F}(x, y)$ (from the original example (x, y)). We encode this cost as a function $c : \mathbb{X} \times \mathbb{X} \rightarrow \mathbb{R}^+$.
92 We assume the generation cost is zero if and only if no change is enacted: $c(x, x) = 0$.

93 We assume that the cost-bounded adversary has a budget ε . The adversary aims to find any example
94 that flips the classifier's decision *and* that is within the cost budget:

$$\max_{x' \in \mathcal{F}(x)} \mathbf{1}[f(x') \neq y] \quad \text{s.t. } c(x, x') \leq \varepsilon \quad (2)$$

95 Alternatively, the adversary can optimize a standard surrogate objective which ensures that the
96 optimization problem can be solved in practice:

$$\max_{x \in \mathcal{F}(x, y)} \ell(\eta(x), y) \quad \text{s.t. } c(x, x') \leq \varepsilon, \quad (3)$$

97 **Utility-Bounded Objective.** The cost-bounded adversarial objective solves the issue of impercep-
98 tibility not being a suitable constraint for tabular data. It does not, however, tackle the problem of
99 heterogeneity of examples: the adversary cannot assign different importance to different adversarial
100 examples.

101 We propose to solve it by introducing the *gain* of an attack. The gain, $g(x^*) : \mathbb{X} \rightarrow \mathbb{R}^+$, represents
102 the reward (e.g., the revenue) that the adversary receives if their attack using adversarial example x^*
103 is successful. For example, in fraud detection gain would be a transaction amount, i.e. how much
104 money a fraudster can steal.

105 We also introduce the concept of *utility* of a successful attack as the net benefit of launching the
106 attack. We define the utility $u_{x, y}(x^*)$ of an attack mounted with adversarial example x^* as simply
107 the gain minus the costs:

$$u_{x, y}(x^*) \triangleq g(x^*) - c(x, x^*), \quad (4)$$

108 where (x, y) is the initial example.

109 The adversary can learn whether an example x^* evades the classifier or not (i.e., whether $f(x^*) \neq y$).
110 Then, they can decide to deploy an attack with an adversarial example x^* only if the utility of the
111 attack exceeds a given *margin* $\tau \geq 0$. Otherwise, the adversary discards this adversarial example.
112 Formally, we can model this process by using a *utility constraint* instead of a cost constraint:

$$\max_{x \in \mathcal{F}(x, y)} \mathbf{1}[f(x) \neq y] \quad \text{s.t. } u_{x, y}(x) \geq \tau \quad (5)$$

113 *Related work on adversarial costs.* Our generic cost-bounded objective is not the only possible
114 approach to model attacks in tabular domains. For example, works on adversarial robustness
115 in the context of decision tree-based classifiers often use per-feature constraints as adversarial
116 constraints [24–26]. At the low level, these constraints are formalized either as bounds on L_∞
117 distance [25, 26], or using functions determining constraints for each specific feature value [24]. In
118 these approaches feature constraints are independent. Such independence simplifies the problem, e.g.,
119 the usage of L_∞ independent constraints enables to split a multidimensional optimization problem
120 into a combination of simple one-dimension tasks [25]; or to limit the set of points affected by the
121 split change [24].

122 *Related work on utility-oriented adversaries.* The literature on *strategic classification* also considers
123 utility-oriented objectives [27–29] for their agents. In this body of work, however, agents are not
124 considered adversaries, and the gain is typically limited to $\{+1, -1\}$ reflecting the classifier decision.
125 Our model supports arbitrary gain, which enables us to model broader interests of the adversary such
126 as revenue. Only the work by Sundaram et al. [30] supports gains different from $+1$ or -1 , but they
127 focus on PAC-learning guarantees in the case of linear classifiers, whereas our goal is to provide
128 practical attack and defence algorithms for a wider family of classifiers.

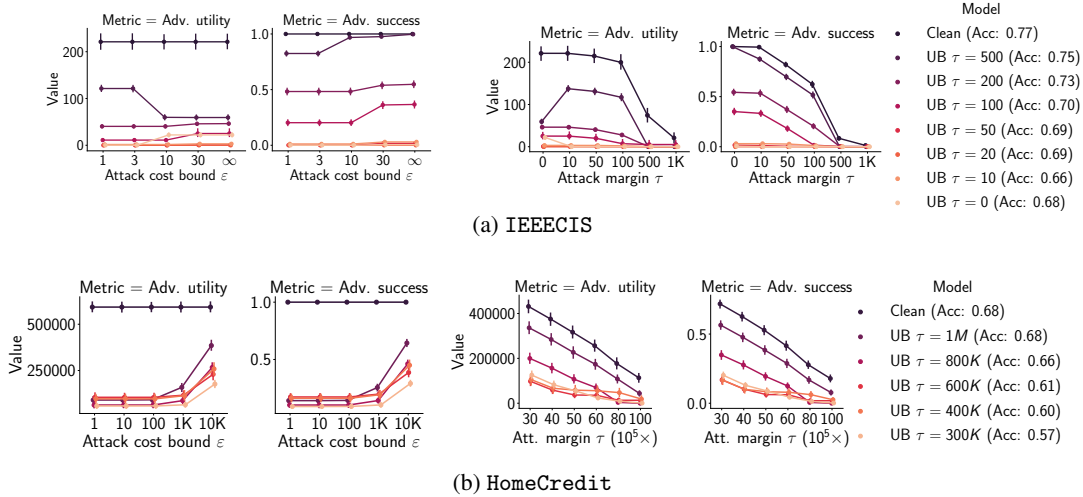


Figure 1: *Utility-Bounded adversarial training* for different adversarial utility margins τ . Evaluation against cost-bounded (left) and utility-bounded (right) adversaries. We show the adversary’s success and utility (y-axis) versus the adversary’s attack budget ϵ or desired margin τ (x-axis). On HomeCredit, the UB training decreases the performance of both UB and CB attacks, being robustness better against the former. Even when enabling a large profit margin ($\tau = 1M$) the attack success rate decreases by 40%, at the same time not affecting the accuracy.

129 3 Algorithms and Evaluation

130 In the full version of this work, we introduce algorithms for attacks and defences within the proposed
 131 adversarial models. We briefly summarize them next, with details provided in the Appendix.

132 **Attack Strategies** For the evaluation of our adversarial models, we implement attacks within both
 133 adversarial objectives using a greedy search algorithm. We describe the algorithm and its design
 134 choices in Appendix A. As a comparison baseline, we adapt the PGD algorithm [10], a common
 135 algorithm for generating adversarial examples, to our cost model, similarly to Ballet et al. [31]. In
 136 Appendix C.2 we show that the greedy algorithm is efficient and outperforms a PGD-based baseline.

137 **Defence With Adversarial Training** To train adversarially robust models, we relax the constraint
 138 sets of the original problems, simplifying them to weighted L_1 ball constraints. With such relaxed
 139 constraints, a PGD-based adversarial training [10] with projection onto the weighted L_1 ball becomes
 140 feasible. We detail this method in Appendix B. For the evaluation of the method, we use two
 141 datasets: HomeCredit [32] and IEEECIS [33], for which we estimate realistic cost and gain models
 142 (see Appendix E. In Fig. 1, we show the results of the evaluation for models trained against the
 143 utility-bounded adversary. These models show decent performance against cost-bounded, close to
 144 “classical”, adversaries. In Appendix C.3, we detail the experimental setup, and show the comparisons
 145 of training against both adversarial objectives, and a detailed study of accuracy-robustness tradeoffs.

146 4 Conclusions

147 In this paper, we revisited the problem of adversarial robustness when the target machine-learning
 148 model operates on tabular data. We introduced a new framework to design attacks and defences that
 149 account for the constraints existing in tabular adversarial scenarios: adversaries are limited by a budget
 150 to modify features, and adversaries may assign different utilities to different examples. Evaluating
 151 these attacks and defences on three real datasets we showed that our novel utility-based defence
 152 mechanism, not only generates models robust against utility-aware adversaries, but also against
 153 adversaries with a limited budget. On the contrary, performing adversarial training considering a
 154 cost-bounded adversary—as traditionally done in the literature—is a poor defence against adversaries
 155 focused on utility in some scenarios.

References

- 156
- 157 [1] Salah Ghamizi, Maxime Cordy, Martin Gubri, Mike Papadakis, Andrey Boytsov, Yves Le
158 Traon, and Anne Goujon. Search-based adversarial testing and improvement of constrained
159 credit scoring systems. In *ESEC/FSE*, 2020.
- 160 [2] Ambra Demontis, Marco Melis, Battista Biggio, Davide Maiorca, Daniel Arp, Konrad Rieck,
161 Igino Corona, Giorgio Giacinto, and Fabio Roli. Yes, machine learning can be more secure! A
162 case study on android malware detection. *CoRR*, 2017.
- 163 [3] Kathrin Grosse, Praveen Manoharan, Nicolas Papernot, Michael Backes, and Patrick D. Mc-
164 Daniel. On the (statistical) detection of adversarial examples. *CoRR*, 2017.
- 165 [4] Bojan Kolosnjaji, Ambra Demontis, Battista Biggio, Davide Maiorca, Giorgio Giacinto, Claudia
166 Eckert, and Fabio Roli. Adversarial malware binaries: Evading deep learning for malware
167 detection in executables. *CoRR*, 2018.
- 168 [5] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian J.
169 Goodfellow, and Rob Fergus. Intriguing properties of neural networks. *CoRR*, 2013.
- 170 [6] Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adver-
171 sarial examples. *CoRR*, 2014.
- 172 [7] Nicolas Papernot, Patrick D. McDaniel, Somesh Jha, Matt Fredrikson, Z. Berkay Celik, and
173 Ananthram Swami. The limitations of deep learning in adversarial settings. In *Euro S&P*, 2016.
- 174 [8] Seyed-Mohsen Moosavi-Dezfooli, Alhussein Fawzi, and Pascal Frossard. Deepfool: A simple
175 and accurate method to fool deep neural networks. In *CVPR*, 2016.
- 176 [9] Nicholas Carlini and David A. Wagner. Towards evaluating the robustness of neural networks.
177 In *S&P*, 2017.
- 178 [10] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu.
179 Towards deep learning models resistant to adversarial attacks. *CoRR*, 2017.
- 180 [11] Puyudi Yang, Jianbo Chen, Cho-Jui Hsieh, Jane-Ling Wang, and Michael I Jordan. Greedy
181 attack and Gumbel attack: Generating adversarial examples for discrete data. *JMLR*, 2020.
- 182 [12] Yutong Wang, Yufei Han, Hongyan Bao, Yun Shen, Fenglong Ma, Jin Li, and Xiangliang Zhang.
183 Attackability characterization of adversarial evasion attack on discrete data. In *KDD*, 2020.
- 184 [13] Boxin Wang, Hengzhi Pei, Han Liu, and Bo Li. AdvCodec: Towards a unified framework for
185 adversarial text generation. *CoRR*.
- 186 [14] Qi Lei, Lingfei Wu, Pin-Yu Chen, Alexandros G Dimakis, Inderjit S Dhillon, and Michael
187 Witbrock. Discrete adversarial attacks and submodular optimization with applications to text
188 classification. *CoRR*, 2018.
- 189 [15] Javid Ebrahimi, Anyi Rao, Daniel Lowd, and Dejing Dou. Hotflip: White-box adversarial
190 examples for text classification. In *ACL*, 2018.
- 191 [16] Bin Liang, Hongcheng Li, Miaoqiang Su, Pan Bian, Xirong Li, and Wenchang Shi. Deep text
192 classification can be fooled. In *IJCAI*, 2018.
- 193 [17] Michele Carminati, Luca Santini, Mario Polino, and Stefano Zanero. Evasion attacks against
194 banking fraud detection systems. In *RAID*, 2020.
- 195 [18] Mahmood Sharif, Lujo Bauer, and Michael K. Reiter. On the suitability of l_p -norms for creating
196 and preventing adversarial examples. *CoRR*, 2018.
- 197 [19] Wei Emma Zhang, Quan Z. Sheng, Ahoud Alhazmi, and Chenliang Li. Adversarial attacks on
198 deep learning models in natural language processing: A survey.
- 199 [20] Xiao Zhang and David Evans. Cost-sensitive robustness against adversarial examples. *CoRR*,
200 2018.

- 201 [21] Hongyang Zhang, Yaodong Yu, Jiantao Jiao, Eric P. Xing, Laurent El Ghaoui, and Michael I.
202 Jordan. Theoretically principled trade-off between robustness and accuracy. In *ICML*, 2019.
- 203 [22] Eric Wong, Frank Schmidt, Jan Hendrik Metzen, and J. Zico Kolter. Scaling provable adversarial
204 defenses. *CoRR*, 2018.
- 205 [23] Ali Shafahi, Mahyar Najibi, Amin Ghiasi, Zheng Xu, John P. Dickerson, Christoph Studer,
206 Larry S. Davis, Gavin Taylor, and Tom Goldstein. Adversarial training for free! In *NeurIPS*,
207 2019.
- 208 [24] Yizheng Chen, Shiqi Wang, Weifan Jiang, Asaf Cidon, and Suman Jana. Cost-aware robust tree
209 ensembles for security applications. In *USENIX*, 2021.
- 210 [25] Maksym Andriushchenko and Matthias Hein. Provably robust boosted decision stumps and
211 trees against adversarial attacks. 2019.
- 212 [26] Hongge Chen, Huan Zhang, Duane Boning, and Cho-Jui Hsieh. Robust decision trees against
213 adversarial examples. In *ICML*, 2019.
- 214 [27] Moritz Hardt, Nimrod Megiddo, Christos H. Papadimitriou, and Mary Wootters. Strategic
215 classification. In *ITCS*, 2016.
- 216 [28] Jinshuo Dong, Aaron Roth, Zachary Schutzman, Bo Waggoner, and Zhiwei Steven Wu. Strategic
217 classification from revealed preferences. In *EC*, 2018.
- 218 [29] Smitha Milli, John Miller, Anca D Dragan, and Moritz Hardt. The social cost of strategic
219 classification. In *FAT**, 2019.
- 220 [30] Ravi Sundaram, Anil Vullikanti, Haifeng Xu, and Fan Yao. Pac-learning for strategic classifica-
221 tion. In Marina Meila and Tong Zhang, editors, *ICML*, 2021.
- 222 [31] Vincent Ballet, Jonathan Aigrain, Thibault Laugel, Pascal Frossard, Marcin Detyniecki, et al.
223 Imperceptible adversarial attacks on tabular data. In *Robust AI in FS NeurIPS Workshop*, 2019.
- 224 [32] Kaggle. Home credit default risk, 2019. URL [https://www.kaggle.com/c/
225 home-credit-default-risk](https://www.kaggle.com/c/home-credit-default-risk).
- 226 [33] Kaggle. IEEE-CIS fraud detection, 2019. URL [https://www.kaggle.com/c/
227 ieee-fraud-detection](https://www.kaggle.com/c/ieee-fraud-detection).
- 228 [34] Bogdan Kulynych, Jamie Hayes, Nikita Samarin, and Carmela Troncoso. Evading classifiers in
229 discrete domains with provable optimality guarantees. *arXiv preprint arXiv:1810.10939*, 2018.
- 230 [35] Roni Stern, Rami Puzis, and Ariel Felner. Potential search: A bounded-cost search algorithm.
231 In *ICAPS*, 2011.
- 232 [36] Peter E. Hart, Nils J. Nilsson, and Bertram Raphael. A formal basis for the heuristic determina-
233 tion of minimum cost paths. *IEEE Trans. Sys. Sci. and Cybernetics*, 1968.
- 234 [37] Roni Stern, Ariel Felner, Jur van den Berg, Rami Puzis, Rajat Shah, and Ken Goldberg.
235 Potential-based bounded-cost search and anytime non-parametric A*. *Artificial Intelligence*,
236 2014.
- 237 [38] Richard E. Korf. Iterative-Deepening-A*: An optimal admissible tree search. In *Joint Confer-
238 ence on Artificial Intelligence*, 1985.
- 239 [39] Rina Dechter and Judea Pearl. Generalized best-first search strategies and the optimality of A*.
240 *J. ACM*, 1985.
- 241 [40] Samir Khuller, Anna Moss, and Joseph Naor. The budgeted maximum coverage problem. *Inf.
242 Process. Lett.*, 1999.
- 243 [41] Laurence A. Wolsey. Maximising real-valued submodular functions: Primal and dual heuristics
244 for location problems. *Math. Oper. Res.*, 1982.

- 245 [42] Eden Levy, Yael Mathov, Ziv Katzir, Asaf Shabtai, and Yuval Elovici. Not all datasets are born
246 equal: On heterogeneous data and adversarial examples. *CoRR*, 2020.
- 247 [43] Alex Kantchelian, J. D. Tygar, and Anthony D. Joseph. Evasion and hardening of tree ensemble
248 classifiers. In *ICML*, 2016.
- 249 [44] Francesco Cartella, Orlando Anunciacao, Yuki Funabiki, Daisuke Yamaguchi, Toru Akishita,
250 and Olivier Elshocht. Adversarial attacks for tabular data: Application to fraud detection and
251 imbalanced data. *SafeAI Workshop at AAAI*, 2021.
- 252 [45] Eric Wong, Leslie Rice, and J. Zico Kolter. Fast is better than free: Revisiting adversarial
253 training. In *ICLR*, 2020.
- 254 [46] Konstantinos Slavakis, Yannis Kopsinis, and Sergios Theodoridis. Adaptive algorithm for sparse
255 system identification using projections onto weighted ℓ_1 balls. In *ICASSP*, 2010.
- 256 [47] Guillaume Perez, Sebastian Ament, Carla Gomes, and Michel Barlaud. Efficient projection
257 algorithms onto the weighted ℓ_1 ball, 2020.
- 258 [48] Stefano Calzavara, Claudio Lucchese, Gabriele Tolomei, Seyum Assefa Abebe, and Salvatore
259 Orlando. Treant: training evasion-aware decision trees. *Data Min. Knowl. Discov.*, 2020.
- 260 [49] Daniël Vos and Sicco Verwer. Efficient training of robust decision trees against adversarial
261 examples. In Marina Meila and Tong Zhang, editors, *ICML*, 2021.
- 262 [50] Zafar Gilani, Ekaterina Kochmar, and Jon Crowcroft. Classification of twitter accounts into
263 automated agents and human users. In *ASONAM*, 2017.
- 264 [51] Experian. What is piggybacking credit, 2019. URL [https://www.experian.com/blogs/
265 ask-experian/what-is-piggybacking-credit/](https://www.experian.com/blogs/ask-experian/what-is-piggybacking-credit/).
- 266 [52] Sercan Ö Arik and Tomas Pfister. Tabnet: Attentive interpretable tabular learning. In *AAAI*,
267 2021.
- 268 [53] Ira Pohl. Heuristic search viewed as path finding in a graph. *Artif. Intell.*, 1970.
- 269 [54] Pratyush Maini, Eric Wong, and Zico Kolter. Adversarial robustness against the union of
270 multiple perturbation models. In *International Conference on Machine Learning*, 2020.

271 **A Finding adversarial examples in tabular domains**

272 In this section, we propose practical algorithms for finding adversarial examples suitable to achieve
273 the adversarial objectives we introduce in Section 2.

274 **A.1 Graphical Framework**

275 The optimization problems in Section 2 can seem daunting due to the large cardinality of $\mathcal{F}(x, y)$
276 when the feature space is large. To make the problems tractable, we transform them into graph-
277 search problems, following the approach by Kulynych et al. [34]. Consider a *state-space graph*
278 $\mathcal{G}(x) = (V, E)$. Each node corresponds to a feasible example in the feature space, $V = \mathcal{F}(x, y) \cup \{x\}$.
279 Edges between two nodes x and x' exist if and only if they differ in value of one feature: there exists
280 $i \in [n]$ such that $x_i \neq x'_i$, and $x_j = x'_j$ for all $j \neq i$. In other words, the immediate descendants of
281 a node in the graph consist of all feasible feature vectors that differ from the parent in exactly one
282 feature value.

283 Using this state-space graph abstraction, the objectives in Section 2 can be modelled as graph search
284 problems. Even though the graph size is exponential in the number of feature values, the search can
285 be efficient, because the search does not need the graph to be complete. Thus, we can construct the
286 graph on the fly.

287 Building the state-space graph is straightforward when features take discrete values. To encode
288 continuous features in the graph we discretize them by only considering changes to a continuous

289 feature i that lie within a finite subset of its domain \mathbb{X}_i — in particular, on a discrete grid. The search
 290 efficiency depends on the size of the grid. As the grid gets coarser, finding adversarial examples
 291 becomes easier. This efficiency comes at the cost of potentially missing adversarial examples that are
 292 not represented on the grid but could fulfil the adversarial constraints with less cost or higher utility.

293 A.2 Attacks as Graph Search

294 In the remainder of the paper we make the following assumptions about the adversarial model:

295 **Assumption 1** (Modular costs). *The adversary’s costs are modular: they decompose by features.*
 296 *Formally, changing the value of each feature i from x_i to x'_i has the associated cost $c_i(x_i, x'_i) > 0$,*
 297 *and the total cost of modifying x into x' is a sum of individual feature-modification costs:*

$$c(x, x') = \sum_i^n c_i(x_i, x'_i) \quad (6)$$

298 The state-space graph can encode modular costs by assigning weights to the graph edges. An edge
 299 between x and x' has an associated weight of $c_i(x_i, x'_i)$, where i is the index of the feature that differs
 300 between x and x' . For pairs of examples $x^{(0)}$ and $x^{(t)}$ that differ in more than one feature, the cost
 301 $c(x^{(0)}, x^{(t)})$ is the sum of the edge costs along the shortest path from $x^{(0)}$ to $x^{(t)}$.

302 **Assumption 2** (Constant gain). *For any initial example (x, y) , the adversary cannot change the gain:*

303

$$\forall x' \in \mathcal{F}(x, y) : g(x) = g(x') \quad (7)$$

304 This follows the approach in utility-oriented strategic classification (as detailed in Section 2). This
 305 assumption is not formally required for our attack algorithms, described next in this section, but we
 306 focus on this setting in our empirical evaluations in Appendix C.

307 **Strategies to find adversarial examples.** Under the constant per-instance gain, and modular-cost
 308 assumptions, the cost-bounded and utility-bounded adversaries look for any adversarial example
 309 that is within a (per-example) cost bound. These adversarial goals can be seen as instances of
 310 *bounded-cost search* [35].

311 We start with the *best-first search* (BFS) [36, 34], a flexible meta-algorithm that generalizes many
 312 common graph search algorithms. In its generic version (Algorithm 1) BFS keeps a bounded priority
 313 queue of *open nodes*. It iteratively pops the node v with the highest score value from the queue (best
 314 first), and adds its immediate descendants to the queue. This is repeated until the queue is empty. The
 315 algorithm returns the node with the highest score out of all popped nodes.

316 The BFS algorithm is parameterized by the *scoring function* $s : V \times V = \mathbb{X} \times \mathbb{X}$ and the size of
 317 the priority queue B . Different choices of the scoring function yield search algorithms suited for
 318 solving different graph-search problems, such as Potential Search for bounded-cost search [35, 37],
 319 and A* [38, 39] for finding the minimal-cost paths. When $B = \infty$, the algorithm might traverse the
 320 full graph and is capable of returning the optimal solution. As the size of B decreases, the optimality
 321 guarantees are lost. When $B = 1$ BFS becomes a *greedy* algorithm that myopically optimizes the
 322 scoring function. When $1 < B < \infty$ we get a *beam search* algorithm that keeps B best candidates at
 323 each iteration.

324 To achieve the adversarial objectives in Section 2, we propose to use a concrete instantiation of
 325 BFS, what we call the *Universal Greedy (UG)* algorithm. Inspired by heuristics for cost-bounded
 326 optimization of submodular functions [40, 41], we set the scoring function to balance the increase in
 327 the classifier’s score and the cost of the change:

$$s(v, t) = - \frac{\eta(t) - \eta(v)}{c(v, t)} \quad (8)$$

328 The minus sign appears because BFS expands the lowest scores first, and we need to maximize the
 329 score. We set the beam size to $B = 1$ (greedy), which enables us to find high-quality solutions to *both*
 330 cost-bounded and utility-bounded problems at reasonable computational costs (see Appendix C).

Algorithm 1 Best-First Search (BFS)

```
1: function BFSB,s,ε(x)
2:   open ← MINPRIORITYQUEUEB(x, 0)
3:   closed ← {}
4:   while open is not empty do
5:     v ← open.POP()
6:     if v ∉ closed then
7:       CLOSED ← CLOSED ∪ {v}
8:       if η(v) ≥ δ then return v
9:       S ← EXPAND(v)
10:      for t ∈ S do
11:        if t ∉ closed and c(x, t) ≤ ε then
12:          open.ADD(t, s(v, t))
```

331 A.3 Related Work on Attack Strategies

332 *Tabular domains.* Several works have proposed attacks on tabular data. Ballet et al. [31] propose to
333 apply existing continuous attacks to tabular datasets. The authors focus on crafting imperceptible
334 adversarial examples using standard methods from the image domain. They adapt these methods
335 such that less “important” features (low correlation with the target variable) can be perturbed to a
336 higher degree than other features. This corresponds to a special case within our framework, in which
337 the feature-modification costs depend on the feature importance with the difference that Ballet et al.
338 cannot guarantee that the proposed example will be feasible.

339 Levy et al. [42] suggest constructing a surrogate model capable of mimicking the target classifier. Part
340 of this surrogate model is a feature embedding function transforming tabular data to a homogeneous
341 continuous domain which aims to keep adversarial perturbations in the feasible set. Then, they apply
342 projected gradient descent to produce adversarial examples in the embedding space and map the
343 resulting examples to the tabular domains. As opposed to our methods, Levy et al. cannot provide
344 any guarantee that the produced adversarial example will lay in the feasible set.

345 Finally, Kantchelian et al. [43] propose a MILP-based and a coordinate-descent attack within different
346 L_p cost models against random-forest models.

347 Our attack differs from these three methods since they use L_p or similar bounds that do not capture
348 adversarial capabilities, whereas we use a cost bound that can capture realistic constraints.

349 In a concurrent work, Cartella et al. [44] propose to use a “custom” norm also based on feature
350 importance, similarly to Ballet et al. [31]. Cartella et al., however, use an adapted zero-order
351 optimization algorithm to find adversarial examples. Although their motivation is similar to ours, our
352 cost model is more general as we do not tie the costs to feature importance.

353 *Text domains.* Our universal greedy attack algorithm is similar to the methods for attacking classifiers
354 that operate on text [19, 11–16]. All these works, however, make use of adversarial constraints such
355 as restrictions on the number of modified words or sentences. These constraints do not apply to
356 tabular domains, as simply considering “number of changes” does not address the heterogeneity of
357 features. Our algorithms also differ from these approaches in that we incorporate complex adversarial
358 costs in the design of the algorithms. For example, the Greedy attack by Yang et al. [11], like us, uses
359 the target classifier’s confidence for choosing the best modifications to create adversarial examples
360 and allow to account for the number of modifications. Our framework not only considers the volume
361 of modifications but also their cost, better reflecting the adversary’s constraints.

362 B Defending from Adversarial Examples in Tabular Domains

363 The conventional approach to mitigate the risks of adversarial examples is adversarial training [6, 10].
364 In adversarial training, the training procedure includes adversarial examples along with natural ones.
365 In a standard approach by Madry et al. [10], for instance, these adversarial examples are constructed
366 by modifying natural examples x with perturbations constrained in a L_p -ball $d(x, x') < \varepsilon$, where d is
367 an L_p distance function.

368 The distance function and ε encode the *threat model* that adversarial training aims to defend against.
 369 The choice of the distance function depends on the characteristics of the input domain. In most
 370 previous works, the distance function aims at capturing imperceptibility within the given the bound
 371 ε . It is commonly assumed that if $d(x, x^*) < \varepsilon$, x^* is not substantially different from x , and the
 372 adversary would use x^* to attack. Otherwise, turning x to x^* results in a perceptible adversarial
 373 example that would be detected as malicious, and those examples are assumed to be outside of the
 374 threat model. As explained in Section 2, this approach does not apply to the tabular domains. In
 375 tabular domains, imperceptibility is not necessarily a relevant constraint. Instead, the adversary’s
 376 actions are constrained by feasibility and the cost of the transformation. Moreover, the tabular input
 377 domain is often a mix of discrete and continuous features, as opposed to continuous or quantized in,
 378 e.g., image domains, where adversarial examples are mostly studied.

379 Another difference between the image and tabular domains is the efficiency of generating adversarial
 380 examples. In images, adversarial examples used for training are generated using efficient methods
 381 such as Projected Gradient Descent (PGD) [10] or the Fast Gradient Sign Method (FGSM) [6, 45].
 382 These approaches produce adversarial examples fast, and enable the efficient implementation of
 383 adversarial training. Fast generation, however, is not possible for tabular domains. The algorithms
 384 to produce tabular adversarial introduced in Appendix A require thousands of inference operations
 385 over the target model. Under this condition, generating one example, which is all the adversary
 386 needs to perform an attack, may not be fast, but it is clearly feasible. Generating multiple adversarial
 387 examples per natural sample, however, in the dataset that the defender needs for adversarial training
 388 quickly becomes computationally infeasible, especially if the defender is computationally constrained.
 389 This computational cost constrains our capability of evaluation (Appendix C), for which we need
 390 to repeatedly run the defences. To make the generation of adversarial examples feasible during
 391 adversarial training, we introduce approximate versions of the attacks that rely on a relaxation of
 392 initial attack constraints.

393 B.1 Relaxing the Constraints

394 Following the setting of the standard Projected Gradient Descent (PGD) method [10], adversarial
 395 training for the cost-bounded adversary could be defined as follows:

$$\min_{\theta} \max_{x' \in \mathcal{F}(x, y)} \ell(\eta_{\theta}(x'), y) \quad \text{s.t. } c(x, x') \leq \varepsilon, \quad (9)$$

396 where η_{θ} is a parametric classifier and θ are its parameters.

397 To keep the computational requirements low, we relax the problem to optimize over a convex set,
 398 which enables us to adapt the PGD method. Let us define B_{ε} to be the constraint region of Eq. (9):

$$B_{\varepsilon}(x, y) \triangleq \{(x', y) \in \mathcal{F}(x, y) \mid c(x, x') \leq \varepsilon\}$$

399 We construct a relaxation of B_{ε} in two steps:

$$B_{\varepsilon} \xrightarrow{(1)} \bar{B}_{\varepsilon} \xrightarrow{(2)} \tilde{B}_{\varepsilon}$$

400 (1) *Continuous relaxation.* We map B_{ε} into a continuous space using an *encoding function*
 401 $\phi : \mathbb{X}^n \rightarrow \mathbb{R}^m$, and a *relaxed cost function* $\bar{c} : \mathbb{R}^m \times \mathbb{R}^m \rightarrow \mathbb{R}^+$. Continuous relaxation is
 402 defined as:

$$\bar{B}_{\varepsilon} \triangleq \{(\phi(x'), y) \mid \bar{c}(\phi(x), \phi(x')) \leq \varepsilon\}, \quad (10)$$

403 where $(x', y) \in \mathcal{F}(x, y)$. The pair (ϕ, \bar{c}) is designed to satisfy the following condition:

$$\forall (x', y) \in B_{\varepsilon}(x, y) : \bar{c}(\phi(x), \phi(x')) \leq c(x, x'), \quad (11)$$

404 ensuring that every example $(x', y) \in B_{\varepsilon}(x, y)$ is mapped to an element in the relaxed set,
 405 $\phi(x') \in B_{\varepsilon}(\phi(x), y)$. We denote the encoded value $\phi(x)$ as \bar{x} for convenience.

406 (2) *Convex cover.* To enable adversarial training using PGD, we need that elements of the
 407 relaxed set can be projected onto the constraint region. For this purpose, we cover \bar{B}_{ε}
 408 with a convex superset \tilde{B}_{ε} , e.g., a convex hull of \bar{B}_{ε} . The convex superset \tilde{B}_{ε} needs to be
 409 constructed such that there exists an efficient algorithm for *projection*. For a given (x, y) ,
 410 and a point $t \in \mathbb{R}^m$, we want to be able to efficiently solve $\min_{t' \in \tilde{B}_{\varepsilon}(x, y)} \|t - t'\|_2$.

Encoding and cost functions As we assume that the cost of modifications is modular (see Appendix A.2), we define the encoding (ϕ) and cost (\bar{c}) functions to be modular too:

$$\phi(x) = [\phi_1(x_1), \dots, \phi_n(x_n)]$$

$$\bar{c}(\phi(x), \phi(x')) = \sum_{i=1}^n \bar{c}_i(\phi_i(x_i), \phi_i(x'_i))$$

411 With this formulation, the problem of constructing suitable ϕ and \bar{c} functions is reduced to finding
 412 $\phi_i : \mathbb{X}_i \rightarrow \mathbb{R}^{m_i}$ and \bar{c}_i for each feature. If for all i both ϕ_i and \bar{c}_i fulfill (11), then the modular cost
 413 $\bar{c}(\bar{x}, \bar{x}')$ fulfills (11) as well.

414 In the following we introduce ϕ and \bar{c} functions for categorical and numeric features.

415 *Categorical features.* As encoding function $\phi(x_i)$ for categorical features we use standard one-hot
 416 encoding.

417 As the cost function for categorical features, we define \bar{c}_i :

$$\bar{c}_i(\bar{x}_i, \bar{x}'_i) = \min_{t \in \mathcal{F}(x, y)} c_i(x_i, t) \cdot \frac{1}{2} \|\bar{x}_i - \bar{x}'_i\|_1,$$

418 where $\mathcal{F}_i(x, y)$ is the set of feasible values of the feature i . For example, let x_i be a categorical feature
 419 with 4 possible values $X_i = \{a, b, c, d\}$, and let the minimal cost of change be 2. When $x_i = b$ and
 420 $x' = c$ ($\bar{x}_i = (0, 1, 0, 0)$, $\bar{x}'_i = (0, 0, 1, 0)$ after one-hot encoding). Then, $\bar{c}_i(\bar{x}_i, \bar{x}'_i) = \frac{1}{2} \cdot 2 \cdot 2 = 2$.

421 This cost function enables us to perform the two-step relaxation described before. First, it satisfies
 422 (11), and therefore the constraint region \bar{B}_ε includes all mapped examples of B_ε . Second, we can
 423 obtain the convex superset \tilde{B}_ε as a continuous L_1 ball around the mapped values $\bar{x} \in \bar{B}_\varepsilon$.

424 *Numeric features.* A *numeric feature* is a feature with values belonging to an ordered subset of \mathbb{R} (e.g.
 425 integer, real). In most cases, the identity function ($\phi(x_i) = x_i$) is sufficient for numerical features.
 426 However, more complex encoding functions could also be desirable. For example, when one needs to
 427 reduce numerical errors, which can be achieved by normalizing the feature values to $[-1, 1]$, or when
 428 the cost is non-linear.

429 In general, projecting onto arbitrary sets can be challenging. Specifically, the bounded region B_ε
 430 could be non-convex, e.g., hypothetically, if the cost is a pathological function such as the Dirichlet
 431 function. We therefore must limit the scope of possible adversarial cost functions that we can
 432 model during adversarial training to those that are compatible with efficient projection. For this, we
 433 introduce a cost model that covers a broad class of functions for which $c_i(x_i, x'_i)$ can be expressed as
 434 $K_i \cdot |\psi(x_i) - \psi(x'_i)|$, where K_i is a constant and $\psi(x)$ is an invertible function.

435 For instance, this model covers the following exponential cost model: $c(x, x') = K \cdot |e^x - e^{x'}|$. In
 436 this case, we can encode the features as $\phi(x_i) = \psi^{-1}(x_i) \triangleq \ln(x_i)$. This transformation enables us
 437 to account for certain non-linear cost functions c with respect to the input space using linear cost
 438 functions \bar{c} in the relaxed space \bar{B}_ε .

439 We define the cost function for numerical features as a piecewise-linear function, with different
 440 coefficients for increasing or decreasing the feature value:

$$\bar{c}_j(\bar{x}_j, \bar{x}'_j) = c_{j-}(x) \cdot [\bar{x}_j - \bar{x}'_j]^+ + c_{j+}(x) \cdot [\bar{x}'_j - \bar{x}_j]^+ \quad (12)$$

441 where $[t]^+$ returns t if $t > 0$, and 0 otherwise, and $c_{j-}(x)$ and $c_{j+}(x)$ encode the costs for decreasing
 442 and increasing the value of the feature j , respectively, and can vary from one initial example x to
 443 another.

444 Note that in this model the final cost of a modification could depend on the way in which this
 445 modification is achieved. A direct modification from x to x'' could have different cost than first
 446 modifying x to x' and then x' to x'' , i.e., $\bar{c}_i(\bar{x}, \bar{x}'') \neq \bar{c}_i(\bar{x}, \bar{x}') + \bar{c}_i(\bar{x}', \bar{x}'')$.

447 *Total cost.* Given the set of categorical feature indices, \mathcal{C} , and the set of numeric feature indices, \mathcal{I} ,
 448 the total cost function is:

$$\begin{aligned} \bar{c}(\bar{x}, \bar{x}') &= \sum_{i \in \mathcal{C}} \min_{t \in \mathcal{F}_i(x, y)} c_i(x_i, t) \cdot \frac{1}{2} \|\bar{x}_i - \bar{x}'_i\|_1 \\ &+ \sum_{j \in \mathcal{I}} c_{j-}(x) \cdot [\bar{x}_j - \bar{x}'_j]^+ + c_{j+}(x) \cdot [\bar{x}'_j - \bar{x}_j]^+ \end{aligned} \quad (13)$$

449 B.2 Adversarial Training with Projected Gradient Descent

450 Using the cost model introduced before, we redefine the training optimization problem in Eq. (9) to
 451 generate adversarial examples over a specific instantiation of the convex set \tilde{B} , as follows:

$$\min_{\theta} \max_{\bar{x}' \in \tilde{B}_{\varepsilon}(x, y)} \ell(\eta_{\theta}(\bar{x}'), y), \quad (14)$$

452 where we specify \tilde{B}_{ε} as:

$$\tilde{B}_{\varepsilon}(x, y) \triangleq \{\bar{x} + \delta \mid \delta \in \mathbb{R}^m \wedge \bar{c}(\bar{x}, \bar{x} + \delta) \leq \varepsilon\}. \quad (15)$$

453 Thus, we can rewrite Eq. (14):

$$\begin{aligned} \min_{\theta} \max_{\delta \in \mathbb{R}^m} \ell(\eta_{\theta}(\bar{x} + \delta), y) \\ \text{s.t. } \bar{c}(\bar{x}, \bar{x} + \delta) \leq \varepsilon \end{aligned} \quad (16)$$

454 This objective can be optimized using standard PGD-based adversarial training [10]. Due to the
 455 construction of our cost function in Eq. (13), we can use existing algorithms for projecting onto a
 456 weighted L_1 -ball [46, 47] with an appropriate choice of weights. As these approaches are standard,
 457 we omit them in the main body, and provide the details in Appendix D.

458 B.3 Adversarial Training against a Utility-Bounded Adversary

459 For the utility-bounded adversary we propose to use an objective similar to (16), applying individual
 460 constraints to different samples:

$$\begin{aligned} \min_{\theta} \max_{\delta \in \mathbb{R}^m} \ell(\eta_{\theta}(\bar{x} + \delta), y) \\ \text{s.t. } \bar{c}(\bar{x}, \bar{x} + \delta) \leq \varepsilon(x) \triangleq [g(\bar{x} + \delta)]_+ \end{aligned} \quad (17)$$

461 In this formulation, we use our assumption of invariant gain (see Appendix A.2), as $g(\bar{x} + \delta) = g(\bar{x})$.

462 This objective aims to decrease the adversary’s utility by focusing the protection on samples with
 463 high gain. The main difference with respect to the cost-constrained objective in (16) is that here we
 464 use a different cost bound for different examples $\varepsilon(x)$. This formulation enables us to directly use the
 465 PGD-based adversarial training to defend against utility-bounded adversaries as well.

466 B.4 Related Work on Adversarial Training

467 To the best of our knowledge, there are no works on adversarial training for methods based on deep
 468 learning that tackle the tabular domains. We discuss existing methods and techniques with related
 469 goals.

470 *Adversarial robustness of decision trees.* Classifiers based on decision trees are prominently used in
 471 tabular domains. The adversarial robustness of such classifiers has been studied extensively [26, 25,
 472 48, 24, 49]. These works assume independent per-feature adversarial constraints, e.g., based on the
 473 L_{∞} metric. Our adversarial models, and thus our attacks and defences, are capable of capturing a
 474 broader class of adversarial cost functions that depend on feature modifications and better model the
 475 adversary’s constraints as we explain in Section 2.

476 C Experimental Evaluation

477 In this section, we show that our graph-based attacks can be used by adversaries to obtain profit, and
 478 that our proposed defences are effective at mitigating these attack’s harms.

479 **C.1 Experimental setup**

480 **C.1.1 Datasets**

481 We perform our experiments on three tabular datasets which represent real-world applications for
482 which adversarial examples can have social or economic implications:

- 483 • **TwitterBot** [50]. The dataset contains information about more than 3,400 Twitter accounts
484 either belonging to humans or bots. The task is to detect bot accounts. We assume that
485 the adversary is able to purchase bot accounts and interactions on darknet markets, thus
486 modifying the features that correspond to the account age, number of likes, and retweets.
- 487 • **IEEEECIS** [33]. The dataset contains information about around 600K financial transactions.
488 The task is to predict whether a transaction is fraudulent or benign. We model an adversary
489 that can modify three features for which we can outline the hypothetical method of possible
490 modification, and estimate its cost: payment-card type, email domain, and payment-device
491 type.
- 492 • **HomeCredit** [32]. The dataset contains financial information about 300K home-loan
493 applicants. The main task is to predict whether an applicant will repay the loan or default. We
494 use 33 features, selected based on the best solutions to the original Kaggle competition [32].
495 Of these, we assume that 28 can be modified by the adversary, e.g., the loan appointment
496 time. We also use a non-linear adversarial cost model for manipulating credit scores, inspired
497 by the practice of credit piggybacking [51].

498 **C.1.2 Models**

499 We evaluate our attacks against three types of ML models commonly applied to tabular data. First, an
500 L_2 -regularized *logistic regression (LR)* with a regularization parameter chosen using 5-fold cross-
501 validation. Second, *gradient-boosted decision trees (XGBoost)*. Third, *TabNet* architecture [52], a *deep*
502 *learning* model. TabNet is an attentive transformer neural network specifically designed for tabular
503 data. We optimize the number of steps as well as the capacity of TabNet’s fully connected layers
504 using grid search.

505 **C.1.3 Adversarial Features**

506 We assume that the feasible set consists of all positive values of numerical features and all possible
507 values of categorical features. For simplicity, we avoid features with mutual dependencies and treat
508 the adversarially modifiable features as independent. We detail the choice of the modifiable features
509 and their costs in Appendix E.2.

510 **C.1.4 Metrics**

511 To evaluate the effectiveness of the attacks and defences, we use three main metrics:

- 512 • **Adversary’s success rate:** The proportion of correctly classified examples from a test set
513 X_{test} for which adversarial examples successfully generated using the attack algorithm
514 $\mathcal{A}(x, y)$ evade the classifier:

$$\Pr_{(x,y) \sim X_{\text{test}}} [f(\mathcal{A}(x, y)) \neq y \wedge f(x) = y].$$

- 515 • **Adversarial cost:** Average cost of successful adversarial examples:

$$\mathbb{E}_{(x,y) \sim X_{\text{test}}} [c(x, \mathcal{A}(x, y)) \mid f(\mathcal{A}(x, y)) \neq y \wedge f(x) = y].$$

- 516 • **Adversarial utility:** Average utility (see Eq. (4)) of successful adversarial examples:

$$\mathbb{E}_{(x,y) \sim X_{\text{test}}} [u_{x,y}(\mathcal{A}(x, y)) \mid f(\mathcal{A}(x, y)) \neq y \wedge f(x) = y].$$

517 In all cases, we only consider correctly classified initial examples which enables us to distinguish
518 these security metrics from the target model’s accuracy. We introduce additional metrics in the
519 experiments when needed.

520 C.2 Attacks Evaluation

521 We evaluate the attack strategy proposed in Appendix A in terms of their effectiveness, and empirically
522 justify its design.

523 C.2.1 Design Choices of the Universal Greedy Algorithm

524 When designing attack algorithms in the BFS framework (see Algorithm 1) there are two main design
525 choices: the scoring function and the beam size. We explore different configurations and show that
526 our parameter choices for the Universal Greedy attack produce high-quality adversarial examples.

527 *Beam size.* We define the beam size of the Universal Greedy attack to be one. The other options that
528 we evaluate are 10 and 100. We evaluate by running three types of attacks: cost-bounded for three
529 cost bounds ε , and utility-bounded at the breakeven margin $\tau = 0$. The margin $\tau = 0$ is equivalent to
530 a cost-bounded attack with a variable cost bound equal to the gain of each initial example (denoted as
531 “Gain” in the tables).

532 We compute two metrics: Attack success, and the success-to-runtime ratio. This ratio represents how
533 much time is needed to achieve the same level of success rate using each choice of the beam size.
534 This metric is more informative for our evaluation than runtime, as runtime is just proportional to the
535 beam size.

536 For feasibility reasons, we use two datasets: `TwitterBot` and `IEEEECIS`. We aggregate the metrics
537 across the three models (LR, XGBT, TabNet), and report the average. The results on `TwitterBot`
538 are equivalent to the results on `IEEEECIS`, thus for conciseness we only report `IEEEECIS` results.

539 We find that the success rates are equal up to the percentage point for all choices of the beam size.
540 We show the detailed numeric results in Table 7 in the Appendix. As the smallest beam size of one is
541 the fastest to run, it demonstrates the best success/time ratio, therefore, is the best choice.

542 *Scoring function.* Recall from Eq. (8) that the scoring function is the cost-weighted increase in the
543 target classifier’s confidence:

$$s(v, t) = -\frac{\eta(t) - \eta(s)}{c(s, t)},$$

544 which aims to maximize the increase in classifier confidence at the lowest cost.

545 Suitable choices for the scoring function $s(v, t)$ could be:

- 546 • *A* algorithm* [38, 39, 34]: $s(v, t) = c(v, t) + \lambda \cdot h(t)$, where $h(t)$ is a heuristic function,
547 which estimates the remaining cost to a solution, and $\lambda > 0$ is a greediness parameter [53].
548 This scoring function balances the current known cost of a candidate and the estimated
549 remaining cost. We choose the model’s confidence for the positive class, $h(x) = \eta(x)$, as
550 heuristic function. Intuitively, this works as a heuristic, because the lower the confidence for
551 the positive class, the more likely we are close to a solution: an example classified as the
552 target class.
- 553 • *Potential Search (PS)* [35, 37]: $s(v, t) = h(t)/\varepsilon - c(v, t)$, which additionally takes into account
554 the cost bound ε , thus becoming more greedy (i.e., optimizing $s(v, t) = \lambda \cdot h(t)$ with
555 $\lambda \approx 1/\varepsilon$) when the cost of the current candidate leaves a lot of room within the ε budget. We
556 also choose $h(x) = \eta(x)$ as heuristic function.
- 557 • *Basic Greedy*: $s(v, t) = -\eta(t)/c(s, t)$, which aims to maximize the classifier’s confidence,
558 yet balance it with the incurred cost. Unlike Eq. (8), this scoring function does not care
559 about the relative increase of the confidence, only about its absolute value.

560 We evaluate the choice of the scoring function on the `TwitterBot` and `IEEEECIS` datasets, with the
561 beam size fixed to one. We run the cost-bounded and utility-bounded attacks in the same configuration
562 as before, and measure two metrics averaged over the models: Attack success, and attack success/time
563 ratio.

564 Table 1 shows the results. On `IEEEECIS`, the Universal Greedy outperforms the other choices in
565 terms of success rate and the success/time ratio. On the `TwitterBot` dataset, it outperforms the
566 other choices in the utility-bounded and unbounded attacks. For cost-bounded attacks, the Universal
567 Greedy offers very close performance to the best option, the Basic Greedy.

Cost bound → Scoring func. ↓	Adv. success, %				Cost bound → Scoring func. ↓	Adv. success, %			
	10	30	Gain	∞		1,000	10,000	Gain	∞
UG	45.32	56.57	56.22	68.20	UG	80.24	85.35	21.63	87.00
A*	42.37	55.62	55.34	53.47	A*	77.56	84.45	20.29	86.25
PS	45.32	55.14	56.18	N/A	PS	79.95	85.19	21.48	N/A
Basic Greedy	42.37	55.46	55.38	53.82	Basic Greedy	80.40	85.04	21.63	86.85

Cost bound → Scoring func. ↓	Success/time ratio				Cost bound → Scoring func. ↓	Success/time ratio			
	10	30	Gain	∞		1,000	10,000	Gain	∞
UG	3.78	4.80	2.53	2.06	UG	208.95	205.76	64.99	205.31
A*	3.29	3.83	1.89	1.15	A*	206.33	201.93	62.25	201.31
PS	3.78	4.01	2.26	N/A	PS	205.85	203.18	63.76	N/A
Basic Greedy	3.21	3.86	2.01	1.16	Basic Greedy	210.20	206.20	64.32	204.96

(a) IEEECIS

(b) Twitter bot

Table 1: *Effect of the scoring-function choice* for graph-based attacks on IEEECIS. In all settings, our Universal Greedy scoring function offers the best success rate and performance.

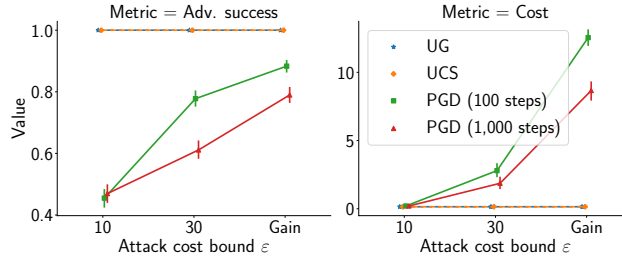


Figure 2: *Universal Greedy attack vs Baselines*. Left: Attack success rate (higher is better for the adversary). Right: Attack cost (lower is better for the adversary). For all cost bounds, our graph-based attack outperforms standard PGD and returns close to optimal-cost adversarial examples (obtained with Uniform-Cost Search, UCS).

568 C.2.2 Graph-based Attacks vs. Baselines

569 We compare the Universal Greedy (UG) algorithm against two baselines: previous work, and the
570 minimal-cost adversarial examples.

571 *Previous Work: PGD*. Since the introduced cost model differs from the existing approaches to
572 attacks on tabular data, we fundamentally cannot perform a fully apples-to-apples comparison against
573 existing attacks (see Appendix A). To compare against the high-level ideas from prior work, we
574 follow the spirit of the attack by Ballet et al. [31], which modifies the standard optimization problem
575 from Eq. (1) to use correlation-based weights. We adapt the standard L_1 -based PGD attack [10, 54]
576 to (1) support categorical features through discretization, and (2) use weighted L_1 norm following
577 our derivations in Appendix B.1. We provide a detailed description of this adaptation in Algorithm 4
578 in Appendix E.

579 We run attacks using PGD with 100 and 1,000 steps, and compare it to UG (Appendix A) on the
580 TwitterBot and IEEECIS datasets. As PGD can only operate on differentiable models, in this
581 comparison we only evaluate the performance of the attacks against TabNet.

582 We run the cost-bounded attacks using two bounds ϵ , that are specific to each dataset (see Appendix E
583 for the exact attack parameters). As before, we also run a utility-bounded attack at the breakeven
584 margin $\tau = 0$. We measure the success rates of the attacks, as well as the average cost of the obtained
585 adversarial examples. For conciseness, we do not report the results on TwitterBot, as they find they
586 are equivalent to those on IEEECIS.

587 Fig. 2 shows that the UG attack consistently outperforms the PGD-based baseline both in terms of
588 the success rate, and the costs. Our attacks are superior even when the PGD-based baseline produces
589 feasible adversarial examples.

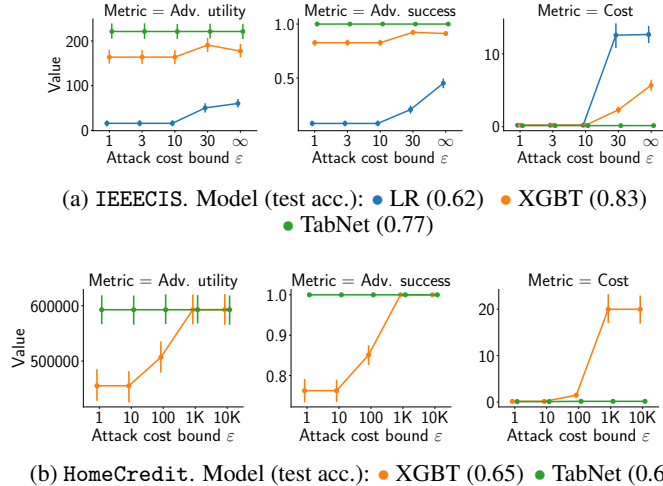


Figure 3: Results of cost-bounded graph-based attacks against three types of models. Left pane: Adversarial utility (higher is better for the adversary). Middle and right panes: See Fig. 2. On IEEECIS, the attack can achieve utility from approximately \$10 to \$125 per attack depending on the target model. On HomeCredit, the average utility ranges between \$400,000 and \$600,000.

590 *Minimal-Cost Adversarial Examples.* As UG is a greedy algorithm, we additionally evaluate how far
 591 are the obtained adversarial examples from the optimal ones in terms of cost. For this, we compare
 592 the results from UG to a standard Uniform-Cost Search (UCS) [34]. UCS is an instantiation of the
 593 BFS framework (see Appendix A) with unbounded beam size, and the scoring function equal to the
 594 cost: $s(v, t) = c(v, t)$. In our setting, UCS is guaranteed to return optimal solutions to the following
 595 optimization problem:

$$\min_{x' \in \mathcal{F}(x, y)} c(x, x') \quad \text{s.t. } f(x') \neq y$$

596 Fig. 2 shows that UG has almost no overhead over the minimal-cost adversarial examples on TabNet
 597 ($1.03\times$ overhead on average). In fact, the average and median cost overhead is $1.80\times$ and $1\times$ over
 598 all models, respectively. There exist some outlier examples, however, with over $100\times$ cost overhead.
 599 We provide more information on the distribution of cost overhead in Appendix E.

600 C.2.3 Performance against Undefended Models

601 Having shown that the attacks outperform the baseline, and the design choices are sound, we
 602 demonstrate that the attacks bring some *utility* to the adversary. In this section, we evaluate the
 603 attacks in a non-strategic setting: the models are not deliberately defended against the attacks.
 604 For conciseness, we only evaluate cost-bounded attacks, as the next section provides an extensive
 605 demonstration of utility-bounded attacks.

606 In *all* evaluated settings, the attacks have non-zero success rates and achieve non-zero adversarial
 607 utility. Fig. 3 show the results of cost-bounded attacks for IEEECIS and HomeCredit datasets. For
 608 utility-bounded attacks, we present the results in Fig. 7 in the Appendix due to the space constraints.
 609 We omit the results for LR on HomeCredit as it does not perform better than the random baseline.
 610 An average adversarial example obtained using the cost-bounded objective brings a profit of \$125 to
 611 the adversary when attacking the IEEECIS TabNet model, and close to 100% of examples in the test
 612 data can be modified into successful adversarial examples.

613 Although for all models we see non-zero success and utility, some models are less vulnerable than
 614 others—even without any protection. For example, the success rate of the adversary against LR on
 615 IEEECIS is much lower than against TabNet (at least 50p.p. lower). This model, however, is also
 616 comparatively inaccurate, with only 62% classification accuracy.

617 C.3 Evaluation of Our Defence Methods

618 We evaluate the defence mechanisms proposed in Appendix B in two scenarios. First, a scenario in
 619 which the adversary’s objective used by the defender for adversarial training—cost-bounded (CB) or

Table 2: Baseline performance

Accuracy	TwitterBot	IEEEICIS	HomeCredit
<i>Clean baseline</i>	0.775	0.755	0.680
<i>Robust baseline</i>	0.773	0.685	0.556
<i>Random baseline</i>	0.566	0.500	0.501

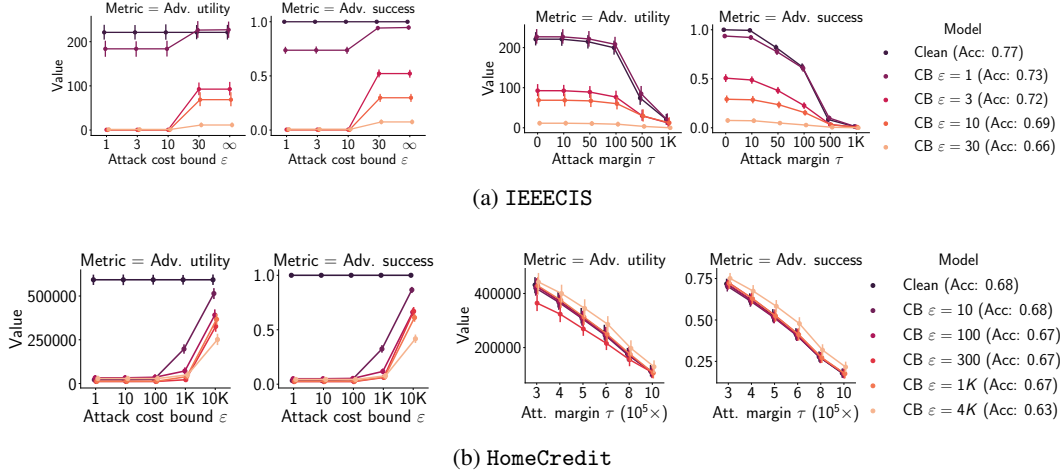


Figure 4: *Cost-bounded adversarial training* for different adversarial budgets ε . Evaluation against cost-bounded (left), and utility-bounded (right) attacks. We represent the adversary’s success and utility (y-axis) versus the adversary’s attack budget ε or desired utility margin τ (x-axis). CB attacks only have substantial success and profit when the adversary invests more than the budget assumed by the defender. UB attacks are thwarted for IEEECIS, but CB training is not significantly effective on HomeCredit, and for some models even enables higher adversary’s utility.

620 utility-bounded (UB)—matches the attack that will be deployed by the adversary. Second, a scenario
 621 in which the defender assumes the adversary’s objective incorrectly and uses a different attack than
 622 the adversary when performing adversarial training.

623 **Baselines.** We set two comparison baselines which provide boundaries for which a defence can be
 624 considered effective.

625 On the accuracy side, we consider the *clean baseline*: a model trained without any defence. It
 626 provides the best accuracy, but also the least robustness against attacks. Any defence that does not
 627 achieve at least the clean baseline’s *robustness* should not be considered, as the clean baseline would
 628 always provide better or equal accuracy, and hence a better robustness-accuracy trade-off.

629 On the robustness side, we consider the *robust baseline*: a model for which all features that can
 630 be changed by the adversary are masked with zeroes for training and testing. As this removes any
 631 adversarial input, this model is invulnerable to attacks within the assumed adversarial models. Any
 632 practical defence must outperform the robust baseline in terms of *accuracy*. Otherwise, the robust
 633 baseline would provide a better robustness-accuracy trade-off.

634 Table 2 shows the clean and robust baselines’ accuracy for the three datasets. On TwitterBot the
 635 robust baseline performs almost as well as the clean model. As there is no space for a better defence
 636 for TwitterBot, we only evaluate our defences for the IEEECIS and HomeCredit models.

637 We train our attacks and defences using the parameters listed in Table 3 in Appendix E.

638 C.3.1 Defender matches the adversary

639 We first evaluate the case in which the adversarial training used to generate the defence is perfectly
 640 tailored to the adversary’s objective.

641 **Cost-bounded Defence vs. Cost-bounded Attack.** We show in Fig. 4 the results when defender
 642 and adversary use CB objectives. For both IEEECIS and HomeCredit the CB trained defence

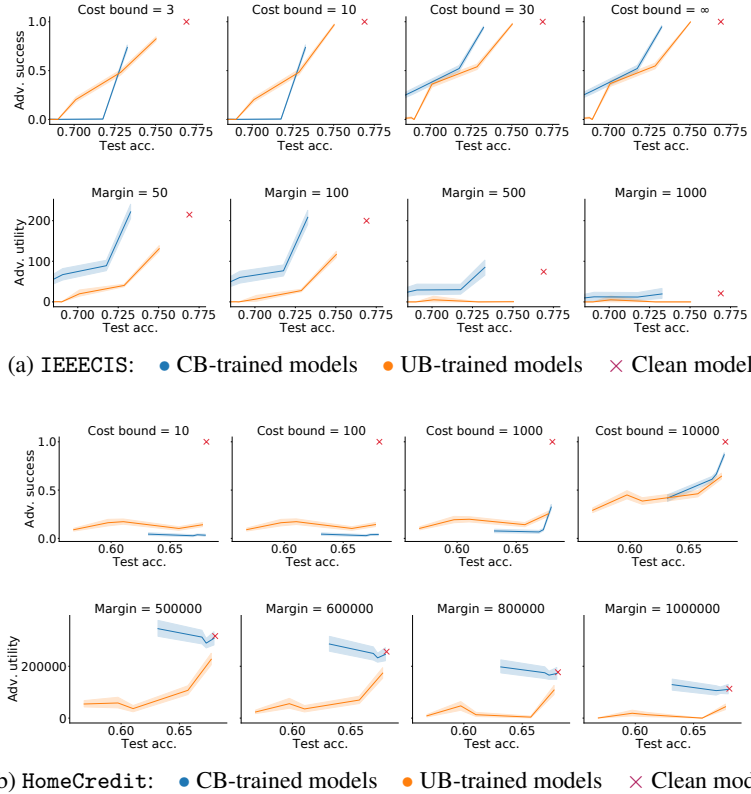


Figure 5: *Accuracy-robustness and utility-robustness trade-offs* for Cost-bounded and Utility-bounded adversarially trained models. The curves show accuracy (x-axis) and utility and success rate (x-axis) for the utility- and cost-bounded models presented in Fig. 4 and Fig. 1. When one curve is strictly below the other curve, it provides a better trade-off since it has better robustness for the same accuracy. Utility-bounded models consistently show better trade-offs for all utility-aware attacks. For CB attacks the situation is less consistent: for small cost-bounds CB defence outperforms utility-bounded one, while for the largest budgets utility bounded shows better results.

643 is effective when the adversary uses CB attacks: the adversary only finds successful adversarial
 644 examples with positive utility if they invest more than the budget assumed by the defender. If the
 645 defender greatly underestimates the adversary’s budget of the adversary (e.g., training with $\varepsilon = 10$
 646 when the adversary’s budget is $\varepsilon = 1000$), the adversary obtains a high profit (close to 200K\$).
 647 Therefore, an effective defence requires an adequate estimation of the adversary’s capabilities.

648 **Utility-bounded Defence vs. Utility-bounded Attack.** Fig. 1 shows the results of our evaluation
 649 when the defender and the adversary use UB objectives. The defence is effective: it decreases both
 650 the success rate and the adversary’s utility on both datasets. On IEEECIS, the adversary can only
 651 succeed when their desired profit τ is smaller than the τ used to train the defence. On HomeCredit,
 652 we observe a similar behaviour, although when training for margins τ less than 500K model does not
 653 completely mitigate adversaries that wish to have larger profits. When the defender allows for large
 654 adversary’s profit margins (e.g., $\tau = 800K$ or $\tau = 1M$), the models become significantly robust
 655 with little accuracy loss.

656 C.3.2 Defender does not match the adversary

657 In the previous section, we show that if the defender correctly models the attacker’s objective, our
 658 defences offer good robustness. Next, we evaluate the performance when the defender’s model does
 659 not match the adversary’s objective. This is likely in realistic deployments, as the defender might not
 660 have any a priori knowledge of the adversary’s objective.

661 **Utility-bounded Defence vs. Cost-bounded Attack.** We show in Fig. 1 our evaluation results when
662 a CB adversary attacks a defence trained assuming UB objectives. For both datasets, the robustness
663 improves with respect to the clean baseline, even though robustness against CB adversaries is not
664 the defence goal. The improvement is more pronounced as the defender tightens the profit margin
665 (decrease in τ , being this effect much stronger on HomeCredit where even loose profit margins
666 provide significant robustness. The adversary can increase their success by increasing their budget ϵ .
667 Increasing the budget also increases the utility in HomeCredit. These experiments show that UB
668 training improves robustness *even when the adversary has a different objective*.

669 **Cost-bounded Defence vs. Utility-bounded Attack.** When we confront a UB adversary against
670 a CB defence, we observe a different behaviour (see Fig. 4). On IEEECIS CB adversarial training
671 increases the robustness of the model against utility-oriented adversaries—with greater effect as
672 the cost bound increases. However, when protecting against high adversary’s budgets ($\epsilon = 30$) the
673 impact on accuracy is too large and the robust baseline becomes preferable.

674 For HomeCredit the situation is worse. While performance is always above the robust baseline, we
675 observe little improvement with respect to the clean model. Even worse, for certain parameters the
676 utility of the adversary can even increase after the adversarial training (see the model trained with a
677 bound of $\epsilon = 4000$). We conclude that CB training might offer no guarantees if the adversary has a
678 different objective.

679 C.3.3 Robustness-Accuracy Trade-offs

680 In the previous sections, we evaluated the effectiveness of the defences depending on the adversary’s
681 and defender’s objectives. We now evaluate the trade-offs between defence effectiveness in reducing
682 the adversary’s success and the utility of the attacks, and the accuracy of the model.

683 As adversarial training penalizes changes in the model’s output caused by input feature perturbations,
684 it results in certain features having less influence on the output. These features cannot be used for
685 prediction to the same extent as features in the clean baseline, which leads to the degradation of the
686 model’s accuracy. On the positive side, these features can neither be used by the adversary—the
687 robust baseline being the extreme in which all features prone to manipulation are zeroed—reducing
688 the attack’s success and utility.

689 We show in Fig. 5 the trade-off between adversarial success and utility on the one side, and model
690 accuracy on the other side for IEEECIS (top) and HomeCredit (bottom) for all combinations of
691 the defender and adversarial objectives. For CB adversaries (top row for each dataset), it is not
692 clear which defence type is superior. Which defence provides better robustness for a given accuracy
693 depends on the adversary’s budget. On the contrary, for utility-bounded adversaries (bottom row
694 for each datasets), we consistently observe better robustness (less adversarial utility for the same
695 accuracy) for the utility-bounded defence compared to the cost-bounded. We conclude that in the
696 absence of knowledge of the adversary’s objective, utility-bounded defences are preferable. They
697 outperform CB adversarial training when the adversary is utility-oriented, and offer comparable
698 performance against CB attacks.

699 D Details on the Projection Algorithm and Adversarial Training

700 In this appendix, we describe our modifications to the traditional adversarial training pipeline.

701 D.1 Adversarial Training Procedure

702 Our training procedure is a version of the well-known adversarial training algorithm based on the
703 PGD method [10].

704 For every sample in a batch $(\phi(x^{(i)}), y^{(i)})_{i=1}^b$, we generate adversarial examples (lines 2–7) by finding
705 the perturbation $\delta^{(i)}$ (lines 4–7). The perturbations are normalized and multiplied by $\alpha = 2\epsilon/n$, to
706 improve the stability of the algorithm (line 6); and then projected to fulfill our relaxed problem in
707 Eq. (16) (line 7). We update the model weights (line 8), and return θ' .

Algorithm 2 Cost-bounded Adversarial Training Algorithm (single iteration)

Input: Model weights θ , batch of training examples $(\phi(x^{(i)}), y^{(i)})_{i=1}^b$, per-feature costs w_i , cost bound ε .

Output: Updated weights θ'

```

1:  $\alpha := 2\frac{\varepsilon}{n}$ 
2: for  $i$  in  $1..b$  do
3:    $\delta^{(i)} := 0$ 
4:   for  $t$  in  $1..n$  do
5:      $\nabla^{(i)} := \nabla_{\delta_i} \ell(f_\theta(\phi(x^{(i)}) + \delta^{(i)}), y_i)$ 
6:      $\delta^{(i)} := \delta^{(i)} + \alpha \frac{\nabla^{(i)}}{\|\nabla^{(i)}\|_1}$ 
7:      $\delta^{(i)} := P_{w, \varepsilon}(\phi(x^{(i)}) + \delta^{(i)})$ 
8:  $\theta' := \theta - \eta \nabla_\theta \frac{1}{b} \sum_{i=1}^b \ell(f_\theta(\phi(x^{(i)}) + \delta^{(i)}), y^{(i)})$ 

```

Return θ_{new}

708 D.2 Projection algorithm

709 We design an adapted projection algorithm to solve Eq. (14), presented in Algorithm 3. This algorithm
 710 is an extension of an existing sort-based weighted L_1 projection algorithm [46, 47]. It takes as input
 711 a sample \bar{x} and a perturbed sample \bar{x}' , and returns a valid perturbation vector δ such that $\bar{x} + \delta$ lies
 712 within the cost budget. With respect to the algorithm by Perez et al. [47], we introduce the capability
 713 to assign different weights based on a feature type and perturbation sign (line 2, in blue) to support
 714 our cost function in Eq. (13).

715 We now prove the correctness of this algorithm.

716 **Statement 1.** *Algorithm 3 is a valid projection algorithm onto the set \tilde{B}_ε , as defined in Eq. (15).*
 717 *Concretely, for a given \bar{x}, \bar{x}' , the algorithm returns δ^* such that:*

$$\delta^* = P_{\tilde{B}_\varepsilon(x, y)}(\bar{x}') \triangleq \arg \min_{\delta \in \mathbb{R}^m} \|\bar{x}' - (\bar{x} + \delta)\|_2$$

$$s.t. \bar{c}(\bar{x}, \bar{x} + \delta) \leq \varepsilon$$

718 *Proof of Statement 1.* First, if we keep either $c_{j+}(x)$ or $c_{j-}(x)$, the constraint becomes a weighted
 719 L_1 constraint, for which the complete proof is given by Perez et al. [47]. Then, we can recall the
 720 property that projection onto the weighted L_1 ball is equivalent to projection onto the simplex, if
 721 $\bar{c}(\bar{x}, \bar{x}') > \varepsilon$, and prove the similar property here.

Lemma 1. *For any \bar{x}, \bar{x}' , ε ,*

$$\delta^* = \arg \min_{\delta: \bar{c}(\bar{x}, \bar{x} + \delta) \leq \varepsilon} \|\bar{x}' - \bar{x} - \delta\|_2$$

$$\forall i \in [1..n] \implies \text{sign}(\delta_i) = \text{sign}(\bar{x}'_i - \bar{x}_i) \text{ or } 0$$

Proof. Proof by contradiction. Let us assume that the lemma does not hold and $\exists i : \text{sign}(\delta_i) =$
 $-\text{sign}(\bar{x}'_i - \bar{x}_i)$ and $\text{sign}(\delta_i) \neq 0$. Then, we can construct $\delta^* : \forall j \neq i, \delta_j^* = \delta_j$ and $\delta_i^* = -\delta_i$.

$$\|\bar{x}' - \bar{x} - \delta\|_2^2 = \|\bar{x}' - \bar{x} - \delta^*\|_2^2 - (\bar{x}'_i - \bar{x}_i - \delta_i)^2 + (\bar{x}'_i - \bar{x}_i - \delta_i^*)^2$$

Since $\text{sign}(\delta_i) = -\text{sign}(\bar{x}'_i - \bar{x}_i)$ and $\text{sign}(\delta_i) \neq 0$,

$$(\bar{x}'_i - \bar{x}_i - \delta_i)^2 > (\bar{x}'_i - \bar{x}_i - \delta_i^*)^2$$

Therefore,

$$\|\bar{x}' - \bar{x} - \delta^*\|_2^2 < \|\bar{x}' - \bar{x} - \delta\|_2^2$$

722 Which is a contradiction to the original statement. □

Algorithm 3 Cost-Ball Projection Algorithm

Input $\bar{x}, \bar{x}', c, \varepsilon, \mathcal{C}, \mathcal{I}$
Output $\delta^* = P_{\tilde{B}_\varepsilon(x,y)}(\bar{x}')$

- 1: $\delta = \bar{x}' - \bar{x}$
- 2: $w_i := \begin{cases} \min_{t \in \mathcal{F}_i(x,y)} c_i(x_i, t), & \text{if } i \in \mathcal{C} \\ c_{j-}(x), & \text{if } i \in \mathcal{I} \text{ and } \delta_i < 0 \\ c_{j+}(x), & \text{if } i \in \mathcal{I} \text{ and } \delta_i \geq 0 \end{cases}$
- 3: $z_i := \frac{\delta_i}{w_i}$
- 4: $\pi_z(\cdot) := \text{Permutation } \uparrow(z)$
- 5: $z_i := z_{\pi_z(i)}$
- 6: $J := \max \left\{ j : \frac{-\varepsilon + \sum_{i=j+1}^m w_{\pi_z(i)} \delta_{\pi_z(i)}}{\sum_{i=j+1}^m w_{\pi_z(i)}^2} > z_j \right\}$
- 7: $\lambda := \frac{-\varepsilon + \sum_{j=J+1}^m w_{\pi_z(j)} \delta_{\pi_z(j)}}{\sum_{j=J+1}^m w_{\pi_z(j)}^2}$
- 8: $\delta_i^* := \text{sign}(\delta_i) \max(\delta_i - w_i \lambda, 0)$

Return δ_i^*

The **highlighted** parts indicate the differences with respect to the sort-based weighted L_1 projection algorithm [47]. The function $\pi_z(i)$ denotes an outcome of permutation. Permutation $\uparrow(z)$ is a sort permutation in an ascending order.

Algorithm 4 PGD-Based Attack

Input: P , initial example x , label y , costs w , cost bound ε .
Output: Adversarial example x^*

- 1: $\alpha := \frac{2\varepsilon}{n}$
- 2: $\delta := 0$
- 3: **for** j **in** $1..n$ **do**
- 4: $\nabla := \nabla_\delta \ell(f_\theta(\phi(x) + \delta), y)$
- 5: $\delta := \delta + \alpha \frac{\nabla}{\|\nabla\|_1}$
- 6: $\delta := P_{\mathcal{B}_{w,\varepsilon}}(\delta)$
 $x^* = P_{\mathcal{F}}(\delta)$

Return x^*

723 Based on this lemma we can see that, to find the projection of \bar{x}' , we can replace $\bar{c}(\bar{x}, \bar{x}')$ with the
724 following expression:

$$\bar{c}^*(\bar{x}, \bar{x}') = \sum_{i \in \mathcal{C}} \frac{1}{2} \|\bar{x}_i - \bar{x}'_i\|_1 \min_{t \in \mathcal{F}_i(x,y)} c_i(x_i, x'_i) + \sum_{j \in \mathcal{I}} c_{j*}(x) \cdot |\bar{x}'_j - \bar{x}_j|,$$

725 where c_{j*} is defined as follows:

$$c_{j*}(x) = \begin{cases} c_{j+}(x), & \text{if } \text{sign}(\bar{x}'_j - \bar{x}_j) \geq 0 \\ c_{j-}(x), & \text{if } \text{sign}(\bar{x}'_j - \bar{x}_j) < 0 \end{cases}$$

726 We can do so as both of these functions attain the same minimum value. \square

727 E Additional Experimental Details

728 In this appendix we provide the details of our evaluation aiming to improve the reproducibility of our
729 results.

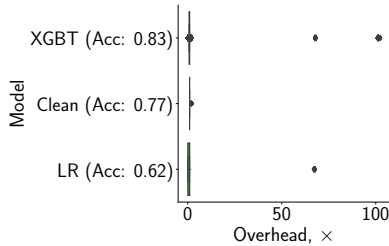


Figure 6: The distribution of cost overhead of adversarial examples obtained with UG over minimal-cost adversarial examples obtained with UCS on IEEECIS. Most UG adversarial examples have cost close to the minimal, although there exist outliers.

Table 3: IEEECIS and HomeCredit attack and defence parameters

Parameter	Value range
Adversarial Training (IEEEECIS)	
<i>Batch size</i>	2048
<i>Number of epochs</i>	400
<i>PGD iteration number</i>	20
<i>TabNet hyperparameters</i>	$N_D = 16, N_A = 16, N_{steps} = 4$
ϵ (for CB models)	[1, 3, 10, 30]
τ (for UB models)	[0, 10, 20, 50, 100, 200, 500]
Attacks (IEEEECIS)	
<i>Max. iterations</i>	100K
ϵ (for CB attacks)	[1, 3, 10, 30]
τ (for UB attacks)	[0, 10, 50, 500, 1000]
Adversarial Training (HomeCredit)	
<i>Batch size</i>	2048
<i>Num. of epochs</i>	100
<i>TabNet hyperparameters</i>	$N_D = 16, N_A = 16, N_{steps} = 4$
<i>Num. of PGD iterations</i>	20
ϵ (for CB models)	[1, 10, 100, 1000, 10000]
τ (for UB models)	[300K, 400K, 500K, 600K, 800K]
Attacks (HomeCredit)	
<i>Num. of iterations</i>	100
ϵ (for CB attacks)	[1, 10, 100, 1K, 10K]
τ (for UB attacks)	[10K, 300K, 400K, 500K, 600K, 800K]

730 E.1 Hyperparameter selection

731 We list our defence and attack parameters in Table 3. TabNet parameters are denoted according
 732 to the original paper [52]. We set the virtual batch size to 512. As training the clean baseline for
 733 HomeCredit was prone to overfitting, we reduced the training number of epochs to 100. Other
 734 hyperparameters were selected with a grid search.

735 E.2 Dataset Processing and Cost Models

736 E.2.1 TwitterBot

737 We use 19 numeric features from this dataset. We dropped 3 features, for which we cannot compute
 738 the effect of a transformation as we do not have access to the original tweets. We use the number of
 739 followers as the adversary’s gain. We assign costs of features based on estimated costs to purchase
 740 Twitter accounts of different characteristics on darknet markets.

741 **E.2.2** IEEECIS

742 We ascribe cost of changes, assuming that the adversary can change the device type and email address
 743 with a small cost. The device type can be changed with low effort using specific software on a mobile
 744 phone. Email domain can be changed with a registration of a new email address which typically
 745 cannot be automated. Although also low cost, it takes more time and effort than changing the device
 746 time. We reflect these assumptions ascribing the costs \$0.1 and \$0.2 to these changes. Changing
 747 the type of card requires obtaining a new card, which costs approximately \$20 in US-based darknet
 748 marketplaces in 2022, according to our research. We consider the transaction amount as a gain
 749 obtained by an adversary.

750 **E.2.3** HomeCredit

751 The main goal of the adversary in this task is receiving a credit approval, therefore, illustrative
 752 purposes, we set credit amount to be a gain of one sample. All features which can be used by an
 753 adversary are listed in Table 6 with the costs we ascribe to them. We assumed six groups of features
 754 and estimated the cost as follows:

- 755 • *Group 1*: features that an adversary can change with negligible effort such as email address,
 756 weekday or hour of the application. We ascribe \$0.1 cost to these transformations.
- 757 • *Group 2*: features associated to income. We use these as a numerical features to illustrate
 758 the flexibility of our method. We assume that to increase income by 1\$, the adversary needs
 759 to pay 1\$.
- 760 • *Group 3*: features associated to changing a phone number. Based on the US darknet
 761 marketplace prices we estimate that buying a phone number costs \$10.
- 762 • *Group 4*: features related documents which can be temporally changed in favor of an
 763 adversary. For example, a car can be transferred from one person to another for the
 764 application period and returned to the original owner after it. We ascribe a cost of \$100 to
 765 obtain these documents.
- 766 • *Group 5*: features that requires either document forging or permanent changes to a person’s
 767 status. For instance, buying a university diploma. These are expensive changes, and we
 768 estimate their cost in \$1 000.
- 769 • *Group 6*: features related to credit scores provided by unspecified external credit-scoring
 770 agencies. We estimate the cost of changes in this group with a manipulation model presented
 771 below.

772 **Credit-score manipulation.** In our feature set we include the features that contain credit scores from
 773 unspecified external credit-scoring agencies. One reported way of manipulating such credit scores is
 774 using credit piggybacking [51]. During piggybacking, a rating buyer finds a “donor” willing to share
 775 a credit for a certain fee. We introduce a model that captures costs of manipulating a credit score
 776 through piggybacking.

We assume that after one piggybacking manipulation the rating is averaged between “donor” and
 recipient, and that “donors” have the maximum rating (1.0). Then, the cost associated to increasing
 the rating from 0.5 to 0.75 is the same as that of increasing from 0.9 to 0.95. This cost cannot be
 represented by a linear function. Let the initial score value be x . The updated credit score after
 piggybacking is $x' = (x+1)/2$. If we repeat the operation n times, the score becomes:

$$x' = \frac{x + 2^n - 1}{2^n}$$

Thus, the number of required piggybacking operations can be computed from the desired final score
 x' as $n = \log_2 \frac{1-x}{1-x'}$, and the total cost is $c(x, x') = nC$, where C is the cost of one operation. We
 estimate to be \$10,000.

$$c(x, x') = C \log_2 \frac{1-x}{1-x'} = C(\log_2(1-x) - \log_2(1-x'))$$

777 To represent this cost function for adversarial training, we can use the encoding described in Ap-
 778 pendix B.1, setting $\phi(x) = \log_2(1-x)$. Then, the cost function becomes $\bar{c}(x, x') = C|\phi(x) - \phi(x')|$,
 779 which is suitable for our defence algorithm. It is worth mentioning that this cost is a lower bound of

Table 4: Costs of changing a feature in TwitterBot dataset

Feature	Estimated cost, \$
<i>likes_per_tweet</i>	0.025
<i>retweets_per_tweet</i>	0.025
<i>user_tweeted</i>	2
<i>user_replied</i>	2

Table 5: Costs of changing a feature in IEEECIS dataset

Feature	Estimated cost, \$
<i>DeviceType</i>	0.1
<i>P_emaildomain</i>	0.2
<i>card_type</i>	20

Table 6: Costs of changing a feature in HomeCredit

Feature	Estimated cost, \$
<i>NAME_CONTRACT_TYPE</i>	0.1
<i>NAME_TYPE_SUITE</i>	0.1
<i>FLAG_EMAIL</i>	0.1
<i>WEEKDAY_APPR_PROCESS_START</i>	0.1
<i>HOUR_APPR_PROCESS_START</i>	0.1
<i>AMT_INCOME_TOTAL</i>	1
<i>FLAG_EMP_PHONE</i>	10
<i>FLAG_WORK_PHONE</i>	10
<i>FLAG_CONT_MOBILE</i>	10
<i>FLAG_MOBIL</i>	10
<i>FLAG_OWN_CAR</i>	100
<i>FLAG_OWN_REALTY</i>	100
<i>REG_REGION_NOT_LIVE_REGION</i>	100
<i>REG_REGION_NOT_WORK_REGION</i>	100
<i>LIVE_REGION_NOT_WORK_REGION</i>	100
<i>REG_CITY_NOT_LIVE_CITY</i>	100
<i>REG_CITY_NOT_WORK_CITY</i>	100
<i>LIVE_CITY_NOT_WORK_CITY</i>	100
<i>NAME_INCOME_TYPE</i>	100
<i>CLUSTER_DAYS_EMPLOYED</i>	100
<i>NAME_HOUSING_TYPE</i>	100
<i>OCCUPATION_TYPE</i>	100
<i>ORGANIZATION_TYPE</i>	100
<i>NAME_EDUCATION_TYPE</i>	1000
<i>NAME_FAMILY_STATUS</i>	1000
<i>HAS_CHILDREN</i>	1000

780 the real cost, as the adversary can only do an integer number of operations. Nonetheless, it perfectly
781 fits our framework as Eq. (11) encompasses this cost model. This is not a fully realistic model, as we
782 cannot know how exactly credit score agencies compute the rating. However, it is reasonable, and
783 enables us to demonstrate how our framework’s support of non-linear costs.

Cost bound → Beam size ↓	Adv. success, %			
	10	30	Gain	∞
1	45.32	56.57	56.22	68.20
10	45.32	56.01	55.65	56.01
100	45.32	56.53	56.18	56.53

Cost bound → Beam size ↓	Success/time ratio			
	10	30	Gain	∞
1	3.78	4.80	2.53	2.06
10	2.14	2.25	1.31	1.15
100	0.66	0.65	0.65	0.66

Table 7: Effect of beam size B in the Universal Greedy algorithm on the IEEECIS dataset. The success rates are close for all choices of the beam size, thus the beam size of one offers the best performance in terms of runtime.

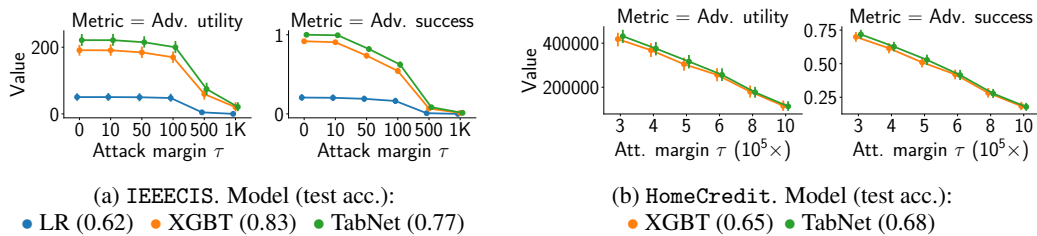


Figure 7: Results of utility-bounded graph-based attacks against three types of models. Left pane: Adversarial utility (higher is better for the adversary). Right pane: See Fig. 2. On IEEECIS, the attack can achieve utility from approximately up to approximately \$200 per attack against TabNet and XGBT. On HomeCredit, the average utility ranges between \$400,000 and \$200,000.