

# Robustness of Sketched Linear Classifiers to Adversarial Attacks

Ananth Mahadevan\*  
ananth.mahadevan@helsinki.fi  
University of Helsinki  
Helsinki, Finland

Yanhao Wang  
yhwang@dase.ecnu.edu.cn  
East China Normal University  
Shanghai, China

Arpit Merchant\*  
arpit.merchant@helsinki.fi  
University of Helsinki  
Helsinki, Finland

Michael Mathioudakis  
michael.mathioudakis@helsinki.fi  
University of Helsinki  
Helsinki, Finland

## ABSTRACT

Linear classifiers are well-known to be vulnerable to adversarial attacks: they may predict incorrect labels for input data that are adversarially modified with small perturbations. However, this phenomenon has not been properly understood in the context of sketch-based linear classifiers, typically used in memory-constrained paradigms, which rely on random projections of the features for model compression. In this paper, we propose novel Fast-Gradient-Sign Method (FGSM) attacks for sketched classifiers in full, partial, and black-box information settings with regards to their internal parameters. We perform extensive experiments on the MNIST dataset to characterize their robustness as a function of perturbation budget. Our results suggest that, in the full-information setting, these classifiers are less accurate on unaltered input than their uncompressed counterparts but just as susceptible to adversarial attacks. But in more realistic partial and black-box information settings, sketching improves robustness while having lower memory footprint.

## CCS CONCEPTS

• Theory of computation → Sketching and sampling; Adversarial learning.

## KEYWORDS

Sketching; Robustness; Adversarial Machine Learning

### ACM Reference Format:

Ananth Mahadevan, Arpit Merchant, Yanhao Wang, and Michael Mathioudakis. 2022. Robustness of Sketched Linear Classifiers to Adversarial Attacks. In *Proceedings of the 31st ACM Int'l Conference on Information and Knowledge Management (CIKM '22)*, Oct. 17–21, 2022, Atlanta, GA, USA. ACM, New York, NY, USA, 5 pages. <https://doi.org/10.1145/3511808.3557687>

## 1 INTRODUCTION

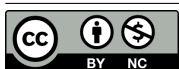
The resource-constrained learning paradigm aims to design classifiers that operate on high-dimensional feature spaces while having low memory footprint. It finds applications on small appliances like smart phones where memory-intensive language models for

speech recognition [13] or computer vision models for facial recognition [11] are required to be trained locally with limited RAM. In this paradigm, sketching [7] is a useful tool for compression employed by classifiers to (i) reduce memory footprint by hashing features to lower dimensions and (ii) maintain accurate weight estimates of an equivalent (i.e., trained on the same data) uncompressed classifier. However, in many data mining problems such as frequency estimation [5], heavy hitters [6], and  $k$ -means clustering [12], sketches have been shown to be brittle to adversarial inputs [1, 2, 10]. Similarly, linear classifiers are also well-known to be susceptible to adversarial data [9, 16, 18]. That is, small perturbations of data points can alter sketch estimates or the predicted labels of an uncompressed linear classifier. Nonetheless, the sensitivity of sketched classifiers to adversarial data has not been well studied. Therefore, our goal is to analyze the extent to which sketching adds robustness to linear classifiers against small perturbations, if at all, while offering its original recovery guarantees.

To this end, we focus on a state-of-the-art sketching-based linear classifier, i.e., the Weight-Median Sketched Classifier [17] (that we refer to as WM-SKETCH for short), as a compressed version of an online linear classifier (LEARNER). The WM-SKETCH internally defines an initial random count sketch [5] and then updates it using gradient descent with a regularized loss function. At inference time, WM-SKETCH predicts labels using either the count sketch directly or the weights in the original space of LEARNER recovered from the count sketch. As shown in [17], for appropriately chosen dimensions of the sketch matrix, the weights recovered from the sketched classifier closely approximate the weights of its uncompressed counterpart. This implies that adversarial perturbations designed for LEARNER may be adapted for WM-SKETCH.

In general, attack algorithms for uncompressed learners can be broadly classified into three categories namely, (i) gradient-based attacks such as Fast Gradient-Sign Method [9] (FGSM), (ii) score-based attacks such as local-search method [14], and (iii) decision-based attacks such as boundary method [3]. We refer interested readers to the surveys by Chakraborty et al. [4] and Pitropakis et al. [15] for further details. In this paper, we focus on characterizing the robustness of WM-SKETCH to FGSM-style adversarial attacks. This attack shifts the original feature vector by a small distance so as to position it on the opposite side of the decision boundary of the classifier being attacked. Generating adversarial examples in this manner requires a distance measure (e.g.,  $L_\infty$ -norm) that captures the size of the perturbation and complete knowledge of

\*Both authors contributed equally to this research.



This work is licensed under a Creative Commons Attribution-NonCommercial International 4.0 License.

the classifier’s internal parameters such as its weights, learning rate, etc. which might be a strong assumption in practice.

**Our Contributions.** For a target WM-SKETCH, we define three adversaries with FGSM-style attack protocols. Depending on the knowledge each adversary has about the target, the protocol involves choosing a surrogate classifier, crafting a perturbed example based on the surrogate, and attacking the target with it. First, in the white-box setting, the adversary knows that the target is a WM-SKETCH and has complete information about the target’s internal parameters, including its count sketch matrix. Here, the surrogate is the target itself. Second, in the grey-box setting, the adversary knows that the target is a WM-SKETCH and has complete information about all of WM-SKETCH’s internal parameters except for the count sketch matrix. Here, the adversary constructs a surrogate WM-SKETCH classifier with a different sketch matrix. And third, in the black-box setting, the adversary does not know that the target is a WM-SKETCH and hence has no access to the target’s internal parameters. Thus, it constructs a LEARNER as a surrogate. In each case, the robustness of the target WM-SKETCH is quantified by the accuracy on the adversarial input as a function of the perturbation budget. Our main contributions include:

- We analytically define FGSM-style attack protocols for WM-SKETCH in the context of the three adversaries discussed above.
- We conduct experiments on MNIST data to empirically quantify the robustness of WM-SKETCH against the three adversaries.
- In the white-box setting, we find that WM-SKETCH is just as susceptible to adversarial attacks as an uncompressed classifier.
- In the grey-box and black-box settings, however, sketching provides WM-SKETCH with improved robustness to adversarial attacks in addition to compression.

## 2 BACKGROUND: LINEAR CLASSIFIERS

**Uncompressed Learner.** Let  $[n]$  denote the set  $\{1, \dots, n\}$ . Denote LEARNER as an online linear classifier for binary classification. Let  $(\mathcal{X}, \mathcal{Y})$  be a stream of data points, where  $\mathbf{x} \in \mathcal{X} \subseteq \mathbb{R}^d$  is a feature vector and  $y \in \mathcal{Y} = \{-1, +1\}$  is the binary label for  $\mathbf{x}$ . A learning history of size  $T$  is a set of labeled data points  $\mathcal{H} = \{(\mathbf{x}^t, y^t)\}_{t \in [T]}$ . The parameters of LEARNER include (1) the hypothesis space  $\mathcal{W} \subseteq \mathbb{R}^d$ , where  $\|\mathcal{W}\|_2 \leq D$ , as well as the hypothesis (also referred to as weights)  $\mathbf{w}^t \in \mathcal{W}$  at any time  $t$  and (2) a time-decaying learning rate  $\eta^t > 0$ . Then, a regularized loss function at time  $t$  for labeled data point  $(\mathbf{x}^t, y^t)$  and hypothesis  $\mathbf{w} \in \mathcal{W}$  is defined as

$$L^t(\mathbf{x}^t, y^t, \mathbf{w}) = l(y^t \cdot \langle \mathbf{w}, \mathbf{x}^t \rangle) + \frac{\lambda}{2} \|\mathbf{w}\|_2^2, \quad (1)$$

where  $l(\cdot)$  is any convex, differentiable function and  $\lambda$  is a regularization parameter. LEARNER’s gradient descent update for labeled point  $(\mathbf{x}^t, y^t)$  is defined as

$$\mathbf{w}^{t+1} = \mathbf{w}^t - \eta^t \nabla L^t(\mathbf{x}^t, y^t, \mathbf{w}^t).$$

Denote  $\mathbf{w}_{\mathcal{H}}^* = \arg \min_{\mathbf{w} \in \mathcal{W}} \sum_{t=1}^{T_{\mathcal{H}}} L^t(\mathbf{x}^t, y^t, \mathbf{w})$  as LEARNER’s optimal hypothesis for history  $\mathcal{H}$ . At inference time, LEARNER’s prediction for  $\mathbf{x}^t$  is  $\hat{y}^t = \text{sign}(\langle \mathbf{w}^t, \mathbf{x}^t \rangle)$ . In this paper, we focus specifically on the logistic loss function  $l(a) = -\log(\sigma(a))$ , where  $\sigma(a) = 1/(1 + \exp(-a))$ . Note that our analyses can also be generalized to other convex, differentiable loss functions.

**Sketched Learner.** Next, we consider the Weighted-Median Sketched Learner [17] (WM-SKETCH), which maintains a count sketch [5] over history  $\mathcal{H}$  and estimates LEARNER’s optimal hypothesis  $\mathbf{w}_{\mathcal{H}}^*$ . WM-SKETCH inherits the parameters  $\eta^t$ ,  $l(\cdot)$ , and  $\lambda$  from LEARNER and its other parameters include (1) the size  $k$  and depth  $s$  of count sketch, (2) the hypothesis  $\mathbf{z}^t \in \mathcal{Z} \subseteq \mathbb{R}^k$ , where  $\mathbf{z}^t$  is the count sketch at time  $t$  arranged as a vector, and (3) a scaled projection matrix  $R = \{-1/\sqrt{s}, +1/\sqrt{s}\}^{k \times d}$  created using random hash functions  $h_j : [d] \rightarrow [k/s]$  and  $\sigma_j : [d] \rightarrow [-1, +1]$  for each  $j \in [s]$ . In this case, the regularized loss function at time  $t$  for  $(\mathbf{x}^t, y^t)$  and hypothesis  $\mathbf{z} \in \mathcal{Z}$  is defined using feature projection as

$$\hat{L}^t(\mathbf{x}^t, y^t, \mathbf{z}) = l(y^t \cdot \langle \mathbf{z}, R\mathbf{x}^t \rangle) + \frac{\lambda}{2} \|\mathbf{z}\|_2^2. \quad (2)$$

WM-SKETCH’s gradient descent update for  $(\mathbf{x}^t, y^t)$  is given accordingly as  $\mathbf{z}^{t+1} = \mathbf{z}^t - \eta^t \nabla \hat{L}^t(\mathbf{x}^t, y^t, \mathbf{z}^t)$ . Furthermore, WM-SKETCH’s estimate  $\hat{\mathbf{w}}_{\text{WM}}^t$  of  $\mathbf{w}^t$  at time  $t$  is obtained as

$$\hat{\mathbf{w}}_i^t = \text{median} \left\{ \sqrt{s} \alpha \sigma_j(i) \mathbf{z}_{j(k/s)+h_j(i)}^t : j \in [s] \right\},$$

where  $\hat{\mathbf{w}}_i^t$  is the  $i$ -th weight of  $\hat{\mathbf{w}}_{\text{WM}}^t$  and  $\alpha = (1 - \eta^t \lambda)$  is a global scale parameter. Note that  $\mathbf{z}^t$ ’s index  $j(k/s) + h_j(i)$  represents the entry in COUNT-SKETCH’s  $j$ -th row and  $h_j(i)$ -th column. There are two ways for the prediction of  $\mathbf{x}^t$  using the sketched learner: (i) *weight recovery*, i.e.,  $\hat{y}_{\text{WR}}^t = \text{sign}(\langle \hat{\mathbf{w}}_{\text{WM}}^t, \mathbf{x}^t \rangle)$ , and (ii) *feature projection*, i.e.,  $\hat{y}_{\text{FP}}^t = \text{sign}(\langle \mathbf{z}^t, R\mathbf{x}^t \rangle)$ . WM-SKETCH’s optimal hypothesis for history  $\mathcal{H}$  of size  $T_{\mathcal{H}}$  is  $\mathbf{z}_{\mathcal{H}}^* = \arg \min_{\mathbf{z} \in \mathcal{Z}} \sum_{t=1}^{T_{\mathcal{H}}} \hat{L}^t(\mathbf{x}^t, y^t, \mathbf{z})$ . As shown in Theorem 1 of [17],  $\|\mathbf{z}_{\mathcal{H}}^* - R\mathbf{w}_{\mathcal{H}}^*\|_2^2$  is small. Further, with high probability over the choice of  $R$ , the optimal hypothesis can be recovered from  $\mathbf{z}_{\mathcal{H}}^*$  from COUNT-SKETCH within an error  $\epsilon$ , i.e.  $\mathbb{E}_{\mathcal{H}} \left[ \|\mathbf{w}_{\mathcal{H}}^* - \hat{\mathbf{w}}_{\text{WM}}^*\|_{\infty} \right] \leq \epsilon \|\mathbf{w}_{\mathcal{H}}^*\|_1$  (cf. Theorem 2 of [17]), where the expectation is over a random permutation of data points in the history  $\mathcal{H}$ . We use the feature projection  $\hat{y}_{\text{FP}}^t$  in our experiments<sup>1</sup>.

## 3 ADVERSARIAL ATTACKS

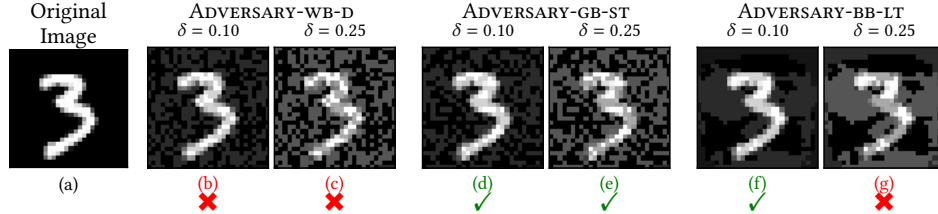
Our adversarial attacks are motivated by the question: *how does the performance of a given WM-SKETCH worsen when attacked by adversaries operating in the (i) white-box (full), (ii) grey-box (partial), and (iii) black-box (zero) information settings?* We formally specify an adversary in Section 3.1 and three attack protocols based on the aforementioned information settings in Section 3.2.

### 3.1 Adversary

For a target WM-SKETCH, an ADVERSARY’s parameters are:

- **Budget** ( $\delta$ ) denotes the maximum allowed perturbation of input point  $\mathbf{x}$  to obtain a perturbed point  $\tilde{\mathbf{x}}$ , i.e.,  $\|\tilde{\mathbf{x}} - \mathbf{x}\| \leq \delta$ .
- **Knowledge** denotes the available information about the target (also assume the access to the training history  $\mathcal{H}$ ), namely,
  - *White-Box* (WB): target is a WM-SKETCH, and access to all of the target’s parameters, including the sketching matrix  $R$  as well as  $s, k, \eta^t, \lambda$ , etc.
  - *Grey-Box* (GB): target is a WM-SKETCH, and access to all of the target’s parameters, except the sketching matrix  $R$ .

<sup>1</sup>The original implementation of [17] uses  $\hat{y}_{\text{WR}}^t$  for both training and prediction.



**Figure 1: Illustration of adversarial attacks.** The original image is as depicted in Figure (a). In Figures (b)-(c), (d)-(e), and (f)-(g), we show the perturbed images with budget  $\delta \in \{0.1, 0.25\}$  by ADVERSARY-WB-D, ADVERSARY-GB-ST, and ADVERSARY-BB-LT targeting a WM-SKETCH with  $k = 16$  and  $s = 2$ , respectively. Here, “✓” and “✗” indicate if the prediction is correct or not.

- *Black-Box (BB)*: access only to the target WM-SKETCH’s predicted label  $\hat{y}_{FP}$  for  $\mathbf{x}$ .

The ADVERSARY’s goal is to induce misclassification by finding a perturbed vector  $\tilde{\mathbf{x}}$  for a labeled data point  $(\mathbf{x}, y)$  given its *budget* and *knowledge*. This implies that the target WM-SKETCH predicts the correct label given  $\mathbf{x}$  but an incorrect label given  $\tilde{\mathbf{x}}$ . That is,  $\text{sign}(\langle \mathbf{z}^t, R\mathbf{x}^t \rangle) = y^t$ , but  $\text{sign}(\langle \mathbf{z}^t, R\tilde{\mathbf{x}}^t \rangle) \neq y^t$ .

### 3.2 Attack Protocols

The ADVERSARY’s attack protocol consists of three steps namely, (i) choosing a surrogate classifier, (ii) crafting a perturbed point  $\tilde{\mathbf{x}}$  for the surrogate, and (iii) attacking the target WM-SKETCH with  $\tilde{\mathbf{x}}$ .

**(i) Choosing a Surrogate.** Based on the ADVERSARY’s *Knowledge* (including the training history  $\mathcal{H}$ ), we have three surrogate options:

- *Direct (D)*: the target WM-SKETCH itself, in the white-box setting.
- *Sketch Transfer (ST)*: a WM-SKETCH trained with the same parameters as the target except a different  $R$ , in the grey-box setting.
- *Learner Transfer (LT)*: an independently trained LEARNER, in the black-box setting.

Formally, we study the following three adversaries: (a) ADVERSARY-WB-D (*White-Box, Direct*); (b) ADVERSARY-GB-ST (*Grey-Box, Sketch Transfer*); and (c) ADVERSARY-BB-LT (*Black-Box, Learner Transfer*).

**(ii) Crafting a Perturbed Point.** We focus on Fast-Gradient-Sign-Method (FGSM)-style attacks [9] for crafting the perturbed point for the surrogate. The key intuition behind this attack is to shift input  $\mathbf{x}$  by a maximum distance  $\delta$  such that it moves to the opposite side of the classifier’s decision boundary to induce misclassification.

The  $L_\infty$ -norm FGSM attack for a WM-SKETCH parameterized by  $R$  for input  $(\mathbf{x}, y)$  can be analytically derived as follows. Denote  $\tilde{\mathbf{x}} = \mathbf{x} + \mathbf{r}$  as the perturbed poin. Note, the regularization term  $\frac{\lambda}{2} \|\mathbf{z}\|_2^2$  is independent of  $\tilde{\mathbf{x}}$  (cf. Equation 2). We construct the attack for a general convex, differentiable loss function  $l(\cdot)$  and then instantiate specifically for logistic loss. Using first order approximation:

$$l(y \cdot \langle \mathbf{z}, R\tilde{\mathbf{x}} \rangle) = l(y \cdot \langle \mathbf{z}, R\mathbf{x} \rangle) + \nabla_{\tilde{\mathbf{x}}} l(y \cdot \langle \mathbf{z}, R\mathbf{x} \rangle)^\top \cdot \mathbf{r}$$

To compute  $\tilde{\mathbf{x}}$ , the smallest perturbation  $\mathbf{r}$  within allowed budget  $\delta$  which minimizes the loss is obtained as follows:

$$\begin{aligned} \arg \min_{\mathbf{r}} \quad & -\nabla_{\tilde{\mathbf{x}}} l(y \cdot \langle \mathbf{z}, R\mathbf{x} \rangle)^\top \cdot \mathbf{r} - l(y \cdot \langle \mathbf{z}, R\mathbf{x} \rangle) \\ \text{s.t.} \quad & \|\mathbf{r}\|_\infty \leq \delta \end{aligned} \quad (3)$$

Using Hölder’s inequality with  $\phi = -\nabla_{\tilde{\mathbf{x}}} l(y \cdot \langle \mathbf{z}, R\mathbf{x} \rangle)$ , we have

$$\phi^\top \mathbf{r} \geq -\|\mathbf{r}\|_\infty \|\phi\|_1 \geq -\delta \|\phi\|_1$$

This lower bound for  $\phi^\top \mathbf{r}$  is achieved when

$$\mathbf{r} = -\delta \cdot \text{sign}(\phi) = -\delta \cdot \text{sign}(-\nabla_{\tilde{\mathbf{x}}} l(y \cdot \langle \mathbf{z}, R\mathbf{x} \rangle))$$

In the case of the logistic loss function, we have

$$\nabla_{\tilde{\mathbf{x}}} l(y \cdot \langle \mathbf{z}, R\mathbf{x} \rangle) = -\sigma(-y \cdot \langle \mathbf{z}, R\mathbf{x} \rangle) \cdot y R^\top \mathbf{z}$$

**(iii) Attacking the Target.** Putting them all together, we have  $\mathbf{r} = -\delta \cdot \text{sign}(y R^\top \mathbf{z})$ . Then, ADVERSARY-WB-D’s perturbed point is:

$$\tilde{\mathbf{x}}_{WB-D} = \mathbf{x} - \delta \cdot \text{sign}(y R^\top \mathbf{z}) \quad (4)$$

Similarly, ADVERSARY-GB-ST perturbs an input point using a surrogate’s projection matrix  $\tilde{R}$  and COUNT-SKETCH  $\tilde{\mathbf{z}}$  instead of those of the target WM-SKETCH as follows:

$$\tilde{\mathbf{x}}_{GB-ST} = \mathbf{x} - \delta \cdot \text{sign}(y \tilde{R}^\top \tilde{\mathbf{z}}) \quad (5)$$

And finally, ADVERSARY-BB-LT perturbs the input point  $\mathbf{x}$  based on a LEARNER with weights  $\mathbf{w}$  as follows:

$$\tilde{\mathbf{x}}_{BB-LT} = \mathbf{x} - \delta \cdot \text{sign}(y \mathbf{w}) \quad (6)$$

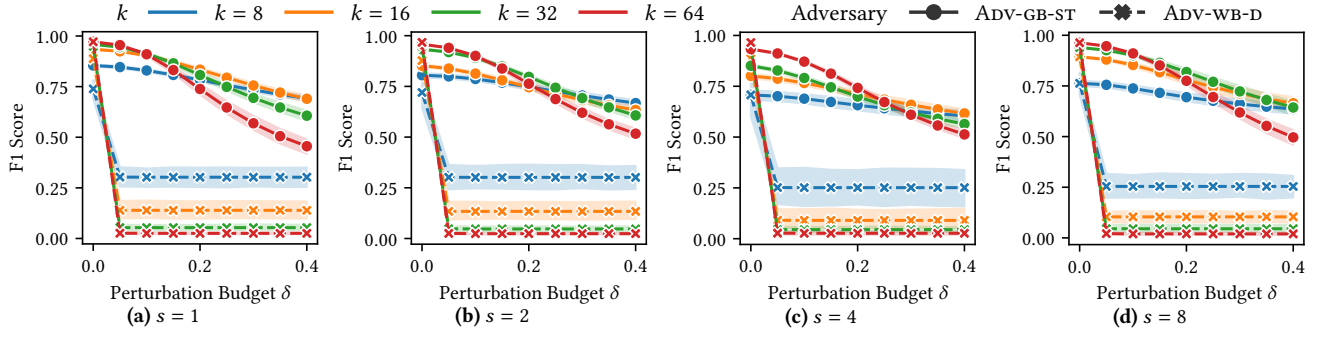
## 4 EXPERIMENTS

### 4.1 Experimental Setup

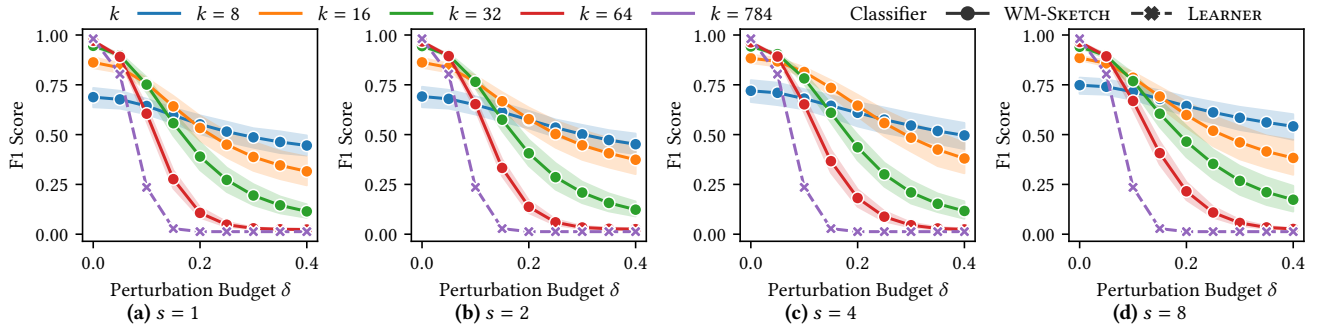
**Dataset.** We perform our experiments on the MNIST handwritten digits dataset [8] where the task is to identify a given image as the digit “3” or “7”. Each image is a 784-dimensional feature vector. We choose 12,396 and 2,038 images randomly for training and testing, respectively. For consistency, we randomly permute the training points to create a fixed history  $\mathcal{H}$  for all the classifiers.

**Classifiers.** We train an instance of LEARNER starting from an all-zero initial hypothesis with initial learning rate  $\eta^0 = 0.1$  and regularization parameter  $\lambda = 10^{-6}$ . We set  $l(\cdot)$  as the logistic loss function. We construct WM-SKETCH classifiers of four different sizes  $k \in \{8, 16, 32, 64\}$ . Smaller sizes imply higher compression ratios yet larger errors in estimates. For each size  $k$ , we use four sketch depths  $s \in \{1, 2, 4, 8\}$  and corresponding widths  $k/s$ . Furthermore, for each configuration of  $k$  and  $s$ , we create 20 unique  $R$  matrices with different random seeds. We report the results for WM-SKETCH that uses the feature projection for training and prediction.

**Adversaries.** We vary  $\delta$  in the range  $\{0, 0.05, \dots, 0.4\}$ . For ADVERSARY-WB-D, we perturb the test set by  $\delta$  (cf. Equation 4). For ADVERSARY-GB-ST, we create one perturbed test set per  $\delta$  for each of the 20 WM-SKETCHES parameterized by their  $R$  (cf. Equation 5) to attack the remaining 19 WM-SKETCHES. And for ADVERSARY-BB-LT, we perturb the test set by  $\delta$  according to Equation 6 using LEARNER



**Figure 2: Accuracy of WM-SKETCH with sizes  $k \in \{8, 16, 32, 64\}$  for depths  $s \in \{1, 2, 4, 8\}$  as a function of budget  $\delta$  for the perturbed input constructed by ADVERSARY-WB-D (dashed lines) and ADVERSARY-GB-ST (solid lines) in full and partial information settings. WM-SKETCH is vulnerable against ADVERSARY-WB-D, but compression adds robustness against ADVERSARY-GB-ST.**



**Figure 3: Accuracy of WM-SKETCH classifiers with sizes  $k \in \{8, 16, 32, 64\}$  for depths  $s \in \{1, 2, 4, 8\}$  as well as LEARNER ( $k = 784$ ) as a function of budget  $\delta$  for the perturbed input constructed by ADVERSARY-BB-LT. Higher levels of compression provide increased robustness against adversarial attacks but lower accuracy on unperturbed input.**

(w) as the surrogate. In Figure 1, we present an illustrated instance of the original and ADVERSARY perturbed images.

**Implementation.** We adapt the original C++ implementation of WM-SKETCH by Tai et al. [17]. Our attack protocols are implemented in Python 3.9. Our code is available here. All experiments were run on a Linux server with 32 CPU cores and 50GB RAM.

## 4.2 Results and Analyses

The dashed lines in Figure 2 present the F1-scores of WM-SKETCHes when the input data is perturbed by ADVERSARY-WB-D (cf. Equation 4). The lines in different colors of the same sub-figure denote the results for WM-SKETCHes of different sizes (with the same depth  $s$ ). While the results for WM-SKETCHes of different depths are plotted in different sub-figures. WM-SKETCH achieves high accuracy on unperturbed data. However, even with small perturbation budget  $\delta \geq 0.05$ , its accuracy degrades significantly. This implies that WM-SKETCH is just as vulnerable to FGSM-style attacks as LEARNER.

The solid lines in Figure 2 depict the F1-scores of WM-SKETCHes given input data perturbed by ADVERSARY-GB-ST (cf. Equation 5). Lines with error bands represent the average and variance of the F1-scores on perturbed points obtained from 19 surrogate WM-SKETCHes of the same  $k$  yet different  $R$ 's. The largest WM-SKETCH ( $k = 64$ ) demonstrates the highest accuracy on unperturbed input but the least robustness compared to smaller sketches as  $\delta$  increases.

Figure 3 shows the robustness of WM-SKETCHes as a function of  $\delta$ -perturbed input from ADVERSARY-BB-LT (cf. Equation 6). The purple line confirms that LEARNER is misled by small perturbations while WM-SKETCHes are more robust. A key observation is that smaller compression levels (e.g.,  $k = 64$ ) lead to lower robustness but higher accuracy on unperturbed input, and vice versa.

## 5 CONCLUSION

In this paper, we initiate the study of the robustness of sketched classifiers. Specifically focusing on FGSM-style attacks that we design for WM-SKETCH [17], we find that (i) under full-information settings about WM-SKETCH's internal parameters, it is as brittle as uncompressed classifiers to small targeted perturbations, but (ii) in more realistic partial and black-box information settings, its random count sketch matrix adds robustness to transfer attacks at the expense of classification accuracy on unperturbed input. Our findings motivate further analysis of sketched classifiers to other attack paradigms and the design of countermeasures such as adversarial training, use of robust sketch alternatives, and ensemble methods to identify better tradeoffs between robustness and accuracy.

## ACKNOWLEDGMENTS

Michael Mathioudakis is supported by University of Helsinki and Academy of Finland Projects MLDB (322046) and HPC-HD (347747).

## REFERENCES

- [1] Omri Ben-Eliezer, Rajesh Jayaram, David P. Woodruff, and Eylon Yogev. 2020. A Framework for Adversarially Robust Streaming Algorithms. In *Proceedings of the 39th ACM SIGMOD-SIGACT-SIGAI Symposium on Principles of Database Systems (PODS '20)*. ACM, New York, NY, USA, 63–80.
- [2] Vladimir Braverman, Avinatan Hassidim, Yossi Matias, Mariano Schain, Sandeep Silwal, and Samson Zhou. 2021. Adversarial Robustness of Streaming Algorithms through Importance Sampling. In *Advances in Neural Information Processing Systems*, Vol. 34. Curran Associates, Inc., Red Hook, NY, USA, 3544–3557.
- [3] Wieland Brendel, Jonas Rauber, and Matthias Bethge. 2018. Decision-Based Adversarial Attacks: Reliable Attacks Against Black-Box Machine Learning Models. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net, 12 pages.
- [4] Anirban Chakraborty, Manaar Alam, Vishal Dey, Anupam Chattopadhyay, and Debdeep Mukhopadhyay. 2018. Adversarial Attacks and Defences: A Survey. arXiv:1810.00069 [cs.LG]
- [5] Moses Charikar, Kevin C. Chen, and Martin Farach-Colton. 2004. Finding frequent items in data streams. *Theor. Comput. Sci.* 312, 1 (2004), 3–15.
- [6] Graham Cormode, Flip Korn, S. Muthukrishnan, and Divesh Srivastava. 2008. Finding Hierarchical Heavy Hitters in Streaming Data. *ACM Trans. Knowl. Discov. Data* 1, 4, Article 2 (2008), 48 pages.
- [7] Graham Cormode and S. Muthukrishnan. 2005. An improved data stream summary: the count-min sketch and its applications. *J. Algorithms* 55, 1 (2005), 58–75.
- [8] Li Deng. 2012. The MNIST Database of Handwritten Digit Images for Machine Learning Research. *IEEE Signal Process. Mag.* 29, 6 (2012), 141–142.
- [9] Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. 2015. Explaining and Harnessing Adversarial Examples. In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*. OpenReview.net, 11 pages.
- [10] Moritz Hardt and David P. Woodruff. 2013. How Robust Are Linear Sketches to Adaptive Inputs?. In *Proceedings of the Forty-Fifth Annual ACM Symposium on Theory of Computing (STOC '13)*. ACM, New York, NY, USA, 121–130.
- [11] Ashish Kapoor, Simon Baker, Sumit Basu, and Eric Horvitz. 2012. Memory constrained face recognition. In *2012 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2539–2546.
- [12] Stuart P. Lloyd. 1982. Least squares quantization in PCM. *IEEE Trans. Inf. Theory* 28, 2 (1982), 129–136.
- [13] Ian McGraw, Rohit Prabhavalkar, Raziell Alvarez, Montse Gonzalez Arenas, Kanishka Rao, David Rybach, Ouais Alsharif, Hasim Sak, Alexander Gruenstein, Françoise Beaufays, and Carolina Parada. 2016. Personalized speech recognition on mobile devices. In *2016 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 5955–5959.
- [14] Nina Narodytska and Shiva Prasad Kasiviswanathan. 2016. Simple Black-Box Adversarial Perturbations for Deep Networks. arXiv:2112.07030 [cs.LG]
- [15] Nikolaos Pitropakis, Emmanouil Panaousis, Thanassis Giannetsos, Eleftherios Anastasiadis, and George Loukas. 2019. A taxonomy and survey of attacks against machine learning. *Comput. Sci. Rev.* 34 (2019), 100199.
- [16] Jacob Steinhardt, Pang Wei Koh, and Percy Liang. 2017. Certified Defenses for Data Poisoning Attacks. In *Advances in Neural Information Processing Systems*, Vol. 30. Curran Associates, Inc., Red Hook, NY, USA, 3517–3529.
- [17] Kai Sheng Tai, Vatsal Sharan, Peter Bailis, and Gregory Valiant. 2018. Sketching Linear Classifiers over Data Streams. In *Proceedings of the 2018 International Conference on Management of Data (SIGMOD '18)*. ACM, New York, NY, USA, 757–772.
- [18] Peng Yang and Ping Li. 2021. Adversarial Kernel Sampling on Class-Imbalanced Data Streams. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management (CIKM '21)*. ACM, New York, NY, USA, 2352–2362.