

Testing the Robustness of Learned Index Structures

Matthias Bachfischer
The University of Melbourne
Melbourne, Australia

bachfischer.matthias@googlemail.com

Renata Borovica-Gajic
The University of Melbourne
Melbourne, Australia

renata.borovica@unimelb.edu.au

Benjamin I. P. Rubinstein
The University of Melbourne
Melbourne, Australia

benjamin.rubinstein@unimelb.edu.au

ABSTRACT

While early empirical evidence has supported the case for learned index structures as having favourable average-case performance, little is known about their worst-case performance. By contrast, classical structures are known to achieve optimal worst-case behaviour. This work evaluates the robustness of learned index structures in the presence of adversarial workloads. To simulate adversarial workloads, we carry out a data poisoning attack on linear regression models that manipulates the cumulative distribution function (CDF) on which the learned index model is trained. The attack deteriorates the fit of the underlying ML model by injecting a set of poisoning keys into the training dataset, which leads to an increase in the prediction error of the model and thus deteriorates the overall performance of the learned index structure. We assess the performance of various regression methods and the learned index implementations ALEX and PGM-Index. We show that learned index structures can suffer from a significant performance deterioration of up to 20% when evaluated on poisoned vs. non-poisoned datasets.

KEYWORDS

Learned index, Database indexing, Adversarial Machine Learning

1 INTRODUCTION

In traditional database design, tree-based data structures such as B+ Trees and their variants have seen wide adoption due to their relative ease of implementation and optimal worst-case computational complexity guarantees. B+ Trees are general-purpose data structures that make no prior assumptions about the data distribution and do not take advantage of any patterns that might be specific to the data that the data structure stores. In practice however, real-world data often follows some underlying pattern that, if modeled appropriately, could significantly speed-up the data retrieval process. For instance, given a set of contiguous integer keys (*e.g.*, keys from 1 to 100 million), the key itself could be used as an offset for the index, thus reducing the time required to look-up any key in the dataset from $O(\log n)$ for a B+ Tree to $O(1)$. This increase in performance is what Kraska et al. hoped to achieve when they first introduced their work on Learned Index Structure (LIS) models [21]. Even though the concept of a LIS is still new, it has already led to a surge of inspiring results that leverage ideas from Machine Learning (ML), data structures, and database systems [7], [6], [31], [15], [1], [32], [5], [30], [23], [16], [11], [19], [8], [13], [27].

The core idea of a LIS is to model the functionality of a data structure as a prediction task, *i.e.*, given an input key, the objective of the LIS is to predict the key’s position in a key-sorted collection of key-value pairs. This approach allows the use of continuous functions to encode the data and leverage learning algorithms to approximate key lookup. Specifically, the approach proposed by Kraska et al. [21] is to approximate the Cumulative Distribution

Function (CDF) of the keys. Given a key k as an input, the CDF returns the probability that the chosen key takes a value less than or equal to k (*i.e.*, $P(X \leq \text{Key } k)$). Based on this observation, one can use the CDF to compute the number of keys less than the queried key k and infer the key’s memory location. In the context of a LIS, the CDF is used to provide a mapping from the key k to its position in a sorted array. The underlying learning task is one of supervised regression, or simply function interpolation [31].

While the LIS approach may be beneficial in a certain average-case sense, it also carries the risk of exploitation by a malicious adversary. Indeed adversarial analysis is an important tool for understanding worst-case performance and whether LIS approaches are worst-case competitive with classical structures.

2 RELATED WORK

2.1 Learned Index Structures

Since the first published work on learned index structures by Kraska et al. [21], the research community has explored a variety of ideas on how LIS could be used as a replacement for traditional index structures such as B+ Trees. Contrary to a B+ Tree, learned index structures rely on an underlying ML model to approximate the CDF. The most commonly used models are the Piecewise Linear Approximation (PLA) and Linear Spline (LS) models. The PLA model is a variant of the linear regression model that tries to approximate the CDF of the dataset by dividing it into variable-sized segments. The first and last keys of each segment are then used to construct a linear approximation of the data, resulting in a PLA model. As an alternative, LS models fit the data by approximating the CDF via linear spline points.

Table 1: Overview of LIS implementations

Index	Model	Updates	Open-Source	Reference
RMI	Multiple	✗	✓	[21]
RadixSpline	LS	✗	✓	[19]
PGM-Index	PLA	✗	✓	[11]
ALEX	PLA	✓	✓	[7]
LIPP	PLA	✓	✓	[9]
COLIN	PLA	✓	✗	[35]

Table 1 overviews the most relevant learned index implementations published to date. A majority of indexes rely on a single type of model to construct the learned index: LS is used by RadixSpline [19] and PLA is used by Piecewise Geometric Model (PGM)-Index [11], Adaptive Learned indEX (ALEX) [7] and a variety of other competitors [9], [35]. The Recursive Model Index (RMI) [21] is the only LIS that supports a variety of model types, which gives the RMI a greater degree of flexibility, but also increases the complexity of tuning the RMI [24].

2.2 Adversarial Machine Learning

The term Adversarial Machine Learning describes the study of Machine Learning (ML) techniques against an adversarial opponent that aims to fool a model by supplying deceptive input. Adversarial ML has emerged in recent years as a new field, mostly driven by new advancements in computing capabilities [17], [2]. A main domain of focus has been computer vision.

Data poisoning attacks are widely studied within Adversarial ML and have been applied in a variety of contexts such as poisoning of neural network or recommender systems [14], [28]. To date, research on data poisoning has mostly been focused on classification and anomaly detection [3], [33], [4], while adversarial regression has largely remained underrepresented [22].

A gradient-based optimization framework for linear classifiers like Lasso or Ridge Regression was first introduced by Xiao et al. [34]. Jagielski et al. [18] have built upon this work and also proposed a defense mechanism against poisoning attacks called *TRIM*. Another novel attack algorithm on regression learning was proposed by Müller et al. [29]. It works by manipulating the training dataset in a way that causes maximum disturbance of the data points. In their experimental evaluation, the authors were able to observe that the Mean Squared Error (MSE) of the regressor increased by 150 percent after inserting only 2% of poisoned samples.

Poisoning attacks on LIS models differ significantly from previous attempts of poisoning linear regression models, because they require poisoning of the CDF. This task is challenging, because every insertion affects the values of all points in the dataset. The first researchers to study this new area of research were Kornaropoulos et al., who formulated two poisoning attacks on the hierarchical structure of the RMI model [20]. By leveraging the attacks described in [20], the researchers were able to increase the error of the poisoned RMI model by a factor of up to 300× compared to a non-poisoned model. We expand on this initial idea and perform a comprehensive empirical evaluation across a variety of models (e.g., SLR, LogTE, DLogTE, 2P, TheilSen, and LAD discussed in [9], as well as ALEX [7] and PGM [11]), open-sourcing our ready-to-use poisoning benchmark to the research community.

3 PRELIMINARIES

3.1 Terminology

In this work we focus on poisoning attacks against LIS models based on Piecewise Linear Approximation (PLA). To define the model, we let $D = \{x_i, y_i\}_{i=1}^n$ denote the data used by a learned index structure, with $x \in \mathbb{R}$ representing the input data vector and $y \in \mathbb{R}$ representing the output variable. In ordinary linear regression, predictions are made via a linear function $f(x, a, b) = a^T x + b$ with parameters $a \in \mathbb{R}^d$ and $b \in \mathbb{R}$ chosen to minimize average loss $\mathcal{L}(D, a, b) = \frac{1}{n} \sum_{i=1}^n (f(x_i, a, b) - y_i)^2$, also known as the Mean Squared Error.

In the context of our data poisoning attacks, a key is denoted by k and its key universe (range of potential keys) as \mathcal{K} , where $|\mathcal{K}| = m$. Similar to previous work on LIS models, it is assumed that keys are given as non-negative integer and that the total order of the keyset can always be derived. It is further assumed that each key is associated with a record and that the records are stored in an in-memory array that is sorted with respect to the key values.

3.2 Adversarial Model

3.2.1 Adversary's Goals. When executing a poisoning attack against a LIS model, the adversary's goal is to corrupt the learned index model during the training phase (i.e., index construction), so that its performance deteriorates during the test phase (i.e., prediction of the position of the given key). In this work, we focus on *poisoning availability attacks*, where the adversary's goal is to deteriorate the performance of a learning-based data structure. Specifically, the objective of the adversary is to generate a small number of *poisoning keys* that are used to augment the training dataset that consists of so-called *legitimate keys*. The assumption is that training a LIS model with both the *poisoning keys* and *legitimate keys* will result in a model whose performance is worse compared to a LIS model that is trained on only the *legitimate keys*.

3.2.2 Adversary's Capabilities & Knowledge. Data poisoning attacks usually distinguish between two attack scenarios: *white-box* and *black-box* poisoning attacks.

In this work, we focus on *white-box* poisoning attacks where the attacker is assumed to have full access to the training data, i.e., the keyset K and the slope and intercept parameters a and b of the linear regression [25], [3], [20]. When performing the attack, the adversary is able to inject up to p maliciously-crafted poisoning keys into the training set prior to training the LIS model. The total number of data points in the training set is given by $N = n + p$, where n denotes the number of legitimate keys and p the number of poisoning keys in the training data. Similar to previous work on poisoning attacks, it is assumed that the adversary is able to control only a tiny fraction of the training set limited by the poisoning percentage α , where $\alpha = p/n$ [18], [34].

If we were to adapt this attack to the *black-box* scenario where the adversary has no direct access to the training data or training parameters of the model, the attacker would first need to infer the parameters of the model and subsequently use their estimates to perform the poisoning attack. Though more difficult to execute, *black-box* attacks allow better transferability of poisoning attacks against different training sets, as shown in [34] and [28].

3.2.3 Attack Evaluation Metric. In this research, the effect of poisoning the LIS is measured with respect to the mean lookup time in nanoseconds. To measure the performance impact of the poisoning attack, we calculate the ratio of the mean lookup time of a model that is trained on the legitimate data and the mean lookup time of a model trained on the poisoned data. For completeness, we also report the mean lookup time (in nanoseconds) across all poisoning thresholds for all the models considered.

4 TESTING THE ROBUSTNESS OF LEARNED INDEX STRUCTURES

To test the robustness of learned index structures, we execute a poisoning attack on linear regression models by attacking the CDF of the training data. The attack works by inserting a certain number of *poisoning keys* into the training dataset with the aim of increasing the approximation error of the regression and thus deteriorating the overall performance of the index.

To formulate the poisoning attack, we consider a LIS to consist of an index that is being constructed on a keyset K of size n , where

each key $k \in K$ has a rank r in the interval $[1, n]$. Here, r denotes the position of k in an ordered sequence of K . The objective of the LIS is to approximate the rank of the queried key by constructing a regression model on (k, r) , where the X-value is given by the key k , and the Y-value is denoted by the rank r . In other words, the function that the regression model approximates is the CDF of the input dataset.

Prior work on poisoning attacks on linear regression models was aimed at inserting maliciously-crafted poisoning keys that cause a “local change”, *i.e.*, inserting keys that do not affect the X- and Y-values of the legitimate points [18]. In the case of LIS models, the insertion of a single, maliciously-crafted key k_p will cause a shift in the rank of all keys larger than k_p . This change will in turn trigger a shift of the CDF, thus compounding the effect of the adversarial insertion.

To this date, the “compounding effect of adversarial insertion” has only been studied by the authors of [20], where they introduced a novel poisoning attack for linear regression on CDFs called *GreedyPoisoningRegressionCDF*. We follow their approach and describe the poisoning strategy of an adversary targeting linear regression models on CDFs in Definition 4.1. The parameter λ denotes the upper bound that limits the size of the poisoning keyset P and is chosen to be proportional to the size of the keyset. In the experiments described in Section 5.1, λ was set to a range of values between $\lambda = 0.01n$ and $\lambda = 0.20n$. For further details on the poisoning algorithm, the interested reader is referred to [20].

Definition 4.1 (Poisoning Linear Regression on CDF). Let K be the set of n integers that correspond to the keys and let P be the set of p integers that comprise the poisoning keys. The augmented set on which the linear regression model is trained is $\{(k'_1, r'_1), (k'_2, r'_2), \dots, (k'_n, r'_n)\}$, where $k'_i \in K \cup P$ and $r'_i \in [1, n + p]$. The goal of the adversary is to choose a set P of size at most λ so as to maximize the loss function of the augmented set $K \cup P$:

$$\arg \max_P \text{ s.t. } |P| \leq \lambda \left(\min_{a,b} \mathcal{L} \left(\{k'_i, r'_i\}_{i=1}^{n+p}, a, b \right) \right)$$

5 EXPERIMENTAL EVALUATION

5.1 Experimental Setup

To test the robustness of learned index structures, we set-up a flexible microbenchmark that allows us to quickly test the robustness of learned index structures against data poisoning attacks. The microbenchmark is based on the source code that was published by Eppert et al. [9]. We have extended the existing microbenchmark by implementing the CDF poisoning attack against different types of regression models and the learned index implementations ALEX and PGM-Index. The corresponding source code used in this work is available online.¹

The system architecture that we used for our experiments is shown in Figure 1. To simulate database workload, we have first generated a synthetic dataset consisting of 1000 keys. We subsequently executed the poisoning attack against the synthetic dataset while varying the poisoning threshold parameter from $p = 0.01$

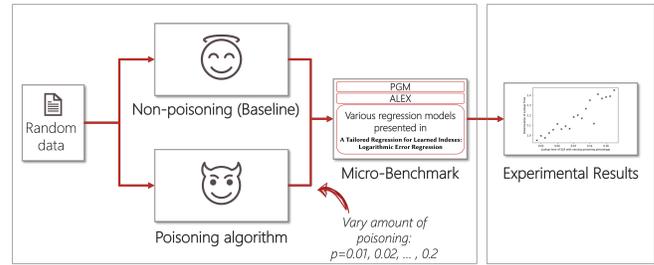


Figure 1: Benchmarking architecture used for experiments.

to $p = 0.20$. For each poisoning threshold, we obtained the corresponding set of poisoning keys and used the legitimate and poisoned keysets to measure the mean lookup time of the indexes on the legitimate (non-poisoned) dataset and the mean lookup time on the poisoned dataset.

To test the robustness of learned index structures, we focused on LIS models that approximate the CDF via PLA and whose source code is available as open-source. We have therefore decided to include ALEX, PGM-Index, and the regression models² discussed in [9] into our benchmark. Additional details on the indexes are provided in Table 2.

Method	Description	Parameters	Source
Regressions (SLR etc.)	[9]	-	[10]
ALEX	[7]	-	[26]
PGM-Index	[11]	max. error ϵ	[12]

Table 2: Overview of evaluated indexes.

5.2 Experimental Results

Figure 2 shows the results of the experiments. The performance deterioration of the LIS is calculated as the ratio between the mean lookup time in nanoseconds for the poisoned datasets and the mean lookup time for the legitimate (non-poisoned) dataset.

From the graphs, we can observe that simple linear regression (SLR) is particularly prone to the poisoning attack, as this regression model shows a steep increase in the mean lookup time when evaluated on the poisoned data. The performance of the competitors that optimize a different error function such as LogTE, DLogTE and 2P are more robust against adversarial attacks. For these regression models, the mean lookup time remains relatively stable even when the poisoning threshold is increased substantially.

Because SLR is the de-facto standard in learned index structures and used internally by the ALEX and the PGM-Index implementations, we would expect that these two models also exhibit a relatively high performance deterioration when evaluated on the poisoned dataset. Surprisingly, ALEX does not show any significant performance impact. This is most likely due to the usage of gapped arrays that allows the model to easily capture outliers in the data (up to a certain point). The performance of the PGM-Index deteriorates by a factor of up-to 1.3 \times .

¹<https://github.com/Bachfischer/LogarithmicErrorRegression>

²Regression models: SLR, LogTE, DLogTE, 2P, TheilSen, and LAD

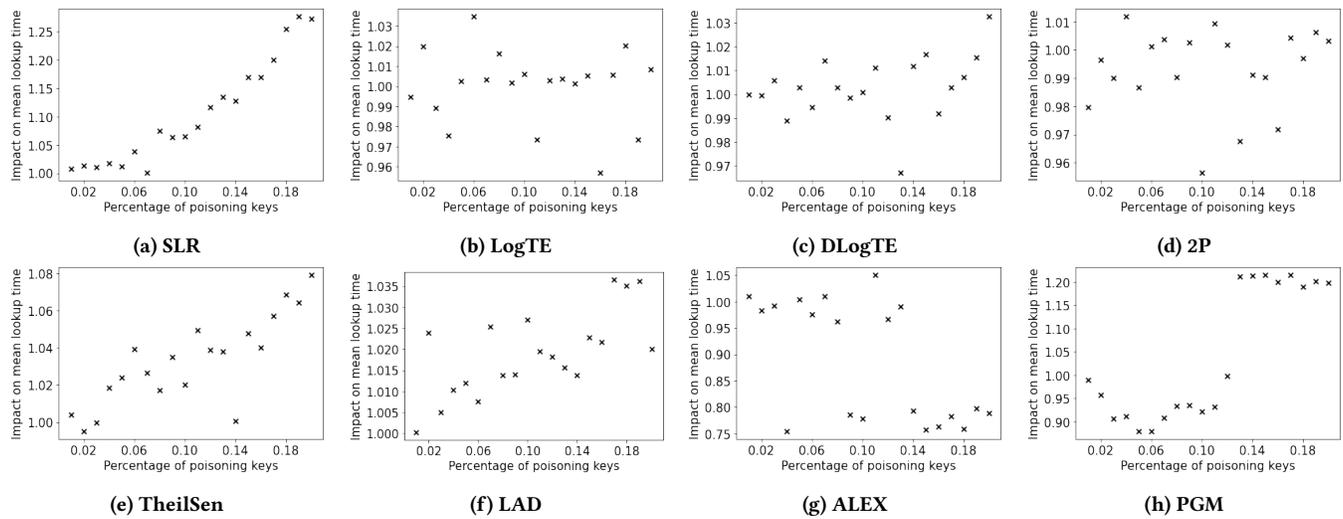


Figure 2: Performance deterioration of LIS models under poisoning attacks.

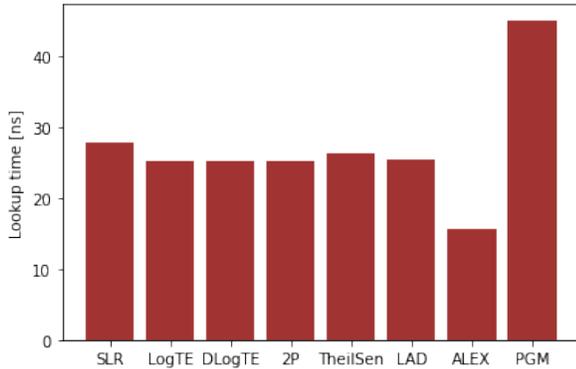


Figure 3: Mean lookup time avg. across all poisoning levels.

For further analysis, we have also calculated the overall mean lookup time for the evaluated learned indexes, averaged across all experiments. The results are shown in Figure 3. From Figure 3, we can observe that ALEX dominates all learned index structures. The performance of the regression models SLR, LogTE, DLogTE, 2P, TheilSen and LAD is also relatively similar, in a range between 30 - 40 nanoseconds, while the PGM-Index performs worst with a mean lookup time of > 50 nanoseconds.

6 CONCLUSIONS AND FUTURE OUTLOOK

In this research, we have tested the robustness of a variety of regression models as well as two learned index implementations ALEX and PGM-Index against adversarial workload. To simulate adversarial workload, we have executed a poisoning attack on a synthetic dataset consisting of 1000 keys. To evaluate the success of the poisoning attack, we measured the performance deterioration of the lookup time for various indexes. The results show that LIS models are prone to poisoning attacks and exhibit significantly worse performance after only a small subset of poisoning keys is introduced into the keysets.

The experiments described in this research focused exclusively on poisoning models that try to approximate the CDF via linear regression. Recent publications have introduced other models for constructing a LIS such as polynomial interpolation [31] and logarithmic error regression [9]. These novel model architectures would also provide an interesting target for poisoning attacks. While the poisoning attack used in this paper works by introducing poisoning keys prior to training the index models, future research may investigate how an adversary could leverage the update functionality of dynamic learned index models to insert and remove keys from a trained LIS model during runtime to deteriorate the fit of the LIS.

REFERENCES

- [1] Hussam Abu-Libdeh, Deniz Altınbüken, Alex Beutel, Ed H Chi, Lyric Doshi, Tim Kraska, Andy Ly, and Christopher Olston. 2020. Learned Indexes for a Google-scale Disk-based Database. *arXiv preprint arXiv:2012.12501* (2020).
- [2] Marco Barreno, Blaine Nelson, Anthony D Joseph, and J Doug Tygar. 2010. The security of machine learning. *Machine Learning* 81, 2 (2010), 121–148.
- [3] Battista Biggio, Blaine Nelson, and Pavel Laskov. 2012. Poisoning attacks against support vector machines. *arXiv preprint arXiv:1206.6389* (2012).
- [4] Battista Biggio and Fabio Roli. 2018. Wild patterns: Ten years after the rise of adversarial machine learning. *Pattern Recognition* 84 (2018), 317–331.
- [5] Antonio Boffa, Paolo Ferragina, and Giorgio Vinciguerra. 2021. A “Learned” Approach to Quicken and Compress Rank/Select Dictionaries. In *2021 Proceedings of the Workshop on Algorithm Engineering and Experiments (ALENEX)*. SIAM, 46–59.
- [6] Yifan Dai, Yien Xu, Aishwarya Ganesan, Ramnathan Alagappan, Brian Kroth, Andrea Arpaci-Dusseau, and Remzi Arpaci-Dusseau. 2020. From wiskey to bourbon: A learned index for log-structured merge trees. In *14th USENIX Symposium on Operating Systems Design and Implementation (OSDI 20)*. 155–171.
- [7] Jialin Ding, Umar Farooq Minhas, Jia Yu, Chi Wang, Jaeyoung Do, Yinan Li, Hantian Zhang, Badrish Chandramouli, Johannes Gehrke, and Donald Kossmann. 2020. ALEX: an updatable adaptive learned index. In *Proceedings of the 2020 ACM SIGMOD International Conference on Management of Data*. 969–984.
- [8] Jialin Ding, Vikram Nathan, Mohammad Alizadeh, and Tim Kraska. 2020. Tsunami: A learned multi-dimensional index for correlated data and skewed workloads. *arXiv preprint arXiv:2006.13282* (2020).
- [9] Martin Eppert, Philipp Fent, and Thomas Neumann. 2021. A Tailored Regression for Learned Indexes: Logarithmic Error Regression. (2021).
- [10] Martin Eppert, Philipp Fent, and Thomas Neumann. 2021. *A Tailored Regression for Learned Indexes: Logarithmic Error Regression*. <https://github.com/umatin/LogarithmicErrorRegression>
- [11] Paolo Ferragina and Giorgio Vinciguerra. 2020. The PGM-index: a fully-dynamic compressed learned index with provable worst-case bounds. *Proceedings of the VLDB Endowment* 13, 10 (2020), 1162–1175.

- [12] Paolo Ferragina and Giorgio Vinciguerra. 2021. *PGM-index: State-of-the-art learned data structure*. <https://github.com/gvinciguerra/PGM-index>
- [13] Alex Galakatos, Michael Markovitch, Carsten Binnig, Rodrigo Fonseca, and Tim Kraska. 2019. FITing-Tree: A Data-aware Index Structure. , 1189–1206 pages. <https://doi.org/10.1145/3299869.3319860>
- [14] Micah Goldblum, Dimitris Tsipras, Chulin Xie, Xinyun Chen, Avi Schwarzschild, Dawn Song, Aleksander Madry, Bo Li, and Tom Goldstein. 2020. Data Security for Machine Learning: Data Poisoning, Backdoor Attacks, and Defenses. *arXiv preprint arXiv:2012.10544* (2020).
- [15] Ali Hadian and Thomas Heinis. 2019. Interpolation-friendly B-trees: Bridging the Gap Between Algorithmic and Learned Indexes. (2019).
- [16] Darryl Ho, Jialin Ding, Sanchit Misra, Nesime Tatbul, Vikram Nathan, Vasimuddin Md, and Tim Kraska. 2019. LISA: towards learned DNA sequence search. *arXiv preprint arXiv:1910.04728* (2019).
- [17] Ling Huang, Anthony D Joseph, Blaine Nelson, Benjamin IP Rubinstein, and J Doug Tygar. 2011. Adversarial machine learning. In *Proceedings of the 4th ACM workshop on Security and artificial intelligence*. 43–58.
- [18] Matthew Jagielski, Alina Oprea, Battista Biggio, Chang Liu, Cristina Nita-Rotaru, and Bo Li. 2018. Manipulating machine learning: Poisoning attacks and countermeasures for regression learning. In *2018 IEEE Symposium on Security and Privacy (SP)*. IEEE, 19–35.
- [19] Andreas Kipf, Ryan Marcus, Alexander van Renen, Mihail Stoian, Alfons Kemper, Tim Kraska, and Thomas Neumann. 2020. RadixSpline: a single-pass learned index. In *Proceedings of the Third International Workshop on Exploiting Artificial Intelligence Techniques for Data Management*. 1–5.
- [20] Evgenios M Kornaropoulos, Silei Ren, and Roberto Tamassia. 2020. The Price of Tailoring the Index to Your Data: Poisoning Attacks on Learned Index Structures. *arXiv preprint arXiv:2008.00297* (2020).
- [21] Tim Kraska, Alex Beutel, Ed H Chi, Jeffrey Dean, and Neoklis Polyzotis. 2018. The case for learned index structures. In *Proceedings of the 2018 International Conference on Management of Data*. 489–504.
- [22] Chang Liu, Bo Li, Yevgeniy Vorobeychik, and Alina Oprea. 2017. Robust linear regression against training data poisoning. In *Proceedings of the 10th ACM Workshop on Artificial Intelligence and Security*. 91–102.
- [23] Stephen Macke, Alex Beutel, Tim Kraska, Maheswaran Sathiamoorthy, Derek Zhiyuan Cheng, and EH Chi. 2018. Lifting the curse of multidimensional data with learned existence indexes. In *Workshop on ML for Systems at NeurIPS*.
- [24] Ryan Marcus, Emily Zhang, and Tim Kraska. 2020. CDFShop: Exploring and Optimizing Learned Index Structures. In *Proceedings of the 2020 ACM SIGMOD International Conference on Management of Data*. 2789–2792.
- [25] Shike Mei and Xiaojin Zhu. 2015. Using machine teaching to identify optimal training-set attacks on machine learners. In *Twenty-Ninth AAAI Conference on Artificial Intelligence*.
- [26] Microsoft. 2021. *ALEX - A library for building an in-memory, Adaptive Learned index*. <https://github.com/microsoft/ALEX>
- [27] Michael Mitzenmacher. 2019. A model for learned bloom filters, and optimizing by sandwiching. *arXiv preprint arXiv:1901.00902* (2019).
- [28] Luis Muñoz-González, Battista Biggio, Ambra Demontis, Andrea Paudice, Vasin Wongrassamee, Emil C Lupu, and Fabio Roli. 2017. Towards poisoning of deep learning algorithms with back-gradient optimization. In *Proceedings of the 10th ACM Workshop on Artificial Intelligence and Security*. 27–38.
- [29] Nicolas Müller, Daniel Kowatsch, and Konstantin Böttinger. 2020. Data Poisoning Attacks on Regression Learning and Corresponding Defenses. In *2020 IEEE 25th Pacific Rim International Symposium on Dependable Computing (PRDC)*. IEEE, 80–89.
- [30] Vikram Nathan, Jialin Ding, Mohammad Alizadeh, and Tim Kraska. 2020. Learning multi-dimensional indexes. In *Proceedings of the 2020 ACM SIGMOD International Conference on Management of Data*. 985–1000.
- [31] Naufal Fikri Setiawan, Benjamin IP Rubinstein, and Renata Borovica-Gajic. 2020. Function interpolation for learned index structures. In *Australasian Database Conference*. Springer, 68–80.
- [32] Chuzhe Tang, Zhiyuan Dong, Minjie Wang, Zhaoguo Wang, and Haibo Chen. 2019. Learned indexes for dynamic workloads. *arXiv preprint arXiv:1902.00655* (2019).
- [33] Sanli Tang, Xiaolin Huang, Mingjian Chen, Chengjin Sun, and Jie Yang. 2019. Adversarial attack type i: Cheat classifiers by significant changes. *IEEE transactions on pattern analysis and machine intelligence* (2019).
- [34] Huang Xiao, Battista Biggio, Gavin Brown, Giorgio Fumera, Claudia Eckert, and Fabio Roli. 2015. Is feature selection secure against training data poisoning?. In *international conference on machine learning*. PMLR, 1689–1698.
- [35] Zhou Zhang, Pei-Quan Jin, Xiao-Liang Wang, Yan-Qi Lv, Shou-Hong Wan, and Xi-Ke Xie. 2021. COLIN: A Cache-Conscious Dynamic Learned Index with High Read/Write Performance. *Journal of Computer Science and Technology* 36, 4 (2021), 721–740.