**Results on MNIST and Contagio.** Figure 2(b) shows the results of BagFlip on MNIST and Contagio. When fixing $R$, the certified accuracy for the backdoor attack is much smaller than the certified accuracy for the trigger-less attack (Figure 1 and Figure 3 in Appendix E.2) because backdoor attacks are strictly stronger than trigger-less attacks. BagFlip cannot provide effective certificates for backdoor attacks on the more complex datasets CIFAR10 and EMBER, i.e., the certified radius is almost zero. **BagFlip can provide certificates against backdoor attacks on MNIST and Contagio, while BagFlip's certificates are not effective for CIFAR10 and EMBER.**

### 7.3 Computation Cost Analysis

We discuss the computation cost of BagFlip on the MNIST dataset and compare to other baselines.

**Training.** The cost of BagFlip during training is similar to all the baselines because BagFlip only adds noise in the training data. BagFlip and other baselines take about 16 hours on a single GPU to train $N = 1000$ classifiers on the MNIST dataset.

**Inference.** At inference time, BagFlip first evaluates the predictions of $N$ classifiers, and counts how many classifiers have the majority label ($N_1$) and how many have the runner-up label ($N_2$). Then, BagFlip uses a prepared lookup table to query the radius certified by $N_1$ and $N_2$. The inference time for each example contains the evaluation of $N$ classifiers and an O(1) table lookup. Hence, there is no difference between BagFlip and other baselines.

**Preparation.** BagFlip needs to prepare a table of size O($N^2$) to perform efficient lookup at inference time. The time complexity of preparing each table entry is presented in Sections 5 and 6. On the MNIST dataset, BagFlip with the relaxation proposed in Section 6 ($\delta = 10^{-4}$) needs 2 hours to prepare the lookup table on a single core. However, the precise BagFlip proposed in Section 5 needs 85 hours to prepare the lookup table. Bagging also uses a lookup table that can be built in 16 seconds on MNIST (Bagging only needs to do a binary search for each entry). FeatFlip needs approximately 8000 TB of memory to compute its table. Thus, FeatFlip is infeasible to run on the full MNIST dataset. FeatFlip is only evaluated on a subset of the MNIST-17 dataset containing only 100 training examples. RAB does not need to compute the lookup table because it has a closed-form solution for computing the certified radius.

**BagFlip has similar training and inference time compared to other baselines. The relaxation technique in Section 6 is useful to reduce the preparation time from 85 hours to 2 hours.** Even with the relaxation, BagFlip needs more preparation time than Bagging and RAB. We argue that the preparation of BagFlip is feasible because it only takes 12.5% of the time required by training.

## 8 Conclusion, Limitations, and Future Work

We presented BagFlip, a certified probabilistic approach that can effectively defend against both trigger-less and backdoor attacks. We foresee many future improvements to BagFlip. First, BagFlip treats both the data and the underlying machine learning models as closed boxes. Assuming a specific data distribution and training algorithm can further improve the computed certified radius. Second, BagFlip uniformly flips the features and the label, while it is desirable to adjust the noise levels for the label and important features for better normal accuracy according to the distribution of the data. Third, while probabilistic approaches need to retrain thousands of models after a fixed number of predictions, the deterministic approaches can reuse models for every prediction. Thus, it is interesting to develop a deterministic model-agnostic approach that can defend against both trigger-less and backdoor attacks.

# References

[1] Hojjat Aghakhani, Dongyu Meng, Yu-Xiang Wang, Christopher Kruegel, and Giovanni Vigna. Bullseye polytope: A scalable clean-label poisoning attack with improved transferability. In *IEEE European Symposium on Security and Privacy, EuroS&P 2021, Vienna, Austria, September 6-10, 2021*, pages 159–178. IEEE, 2021. doi: 10.1109/EuroSP51992.2021.00021. URL https://doi.org/10.1109/EuroSP51992.2021.00021.

[2] Hyrum S. Anderson and Phil Roth. EMBER: an open dataset for training static PE malware machine learning models. *CoRR*, abs/1804.04637, 2018. URL http://arxiv.org/abs/1804.04637.

[3] Battista Biggio, Blaine Nelson, and Pavel Laskov. Poisoning attacks against support vector machines. In *Proceedings of the 29th International Conference on Machine Learning, ICML 2012, Edinburgh, Scotland, UK, June 26 - July 1, 2012*. icml.cc / Omnipress, 2012. URL http://icml.cc/2012/papers/880.pdf.

[4] Ruoxin Chen, Jie Li, Chentao Wu, Bin Sheng, and Ping Li. A framework of randomized selection based certified defenses against data poisoning attacks. *CoRR*, abs/2009.08739, 2020. URL https://arxiv.org/abs/2009.08739.

[5] Xinyun Chen, Chang Liu, Bo Li, Kimberly Lu, and Dawn Song. Targeted backdoor attacks on deep learning systems using data poisoning. *CoRR*, abs/1712.05526, 2017. URL http://arxiv.org/abs/1712.05526.

[6] Jeremy M. Cohen, Elan Rosenfeld, and J. Zico Kolter. Certified adversarial robustness via randomized smoothing. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*, volume 97 of *Proceedings of Machine Learning Research*, pages 1310–1320. PMLR, 2019. URL http://proceedings.mlr.press/v97/cohen19c.html.

[7] Samuel Drews, Aws Albarghouthi, and Loris D'Antoni. Proving data-poisoning robustness in decision trees. In Alastair F. Donaldson and Emina Torlak, editors, *Proceedings of the 41st ACM SIGPLAN International Conference on Programming Language Design and Implementation, PLDI 2020, London, UK, June 15-20, 2020*, pages 1083–1097. ACM, 2020. doi: 10.1145/3385412.3385975. URL https://doi.org/10.1145/3385412.3385975.

[8] Krishnamurthy (Dj) Dvijotham, Jamie Hayes, Borja Balle, J. Zico Kolter, Chongli Qin, András György, Kai Xiao, Sven Gowal, and Pushmeet Kohli. A framework for robustness certification of smoothed classifiers using f-divergences. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net, 2020. URL https://openreview.net/forum?id=SJlKrkSFPH.

[9] Yansong Gao, Change Xu, Derui Wang, Shiping Chen, Damith Chinthana Ranasinghe, and Surya Nepal. STRIP: a defence against trojan attacks on deep neural networks. In David Balenson, editor, *Proceedings of the 35th Annual Computer Security Applications Conference, ACSAC 2019, San Juan, PR, USA, December 09-13, 2019*, pages 113–125. ACM, 2019. doi: 10.1145/3359789.3359790. URL https://doi.org/10.1145/3359789.3359790.

[10] Jonas Geiping, Liam Fowl, Gowthami Somepalli, Micah Goldblum, Michael Moeller, and Tom Goldstein. What doesn't kill you makes you robust(er): Adversarial training against poisons and backdoors. *CoRR*, abs/2102.13624, 2021. URL https://arxiv.org/abs/2102.13624.

[11] Jonas Geiping, Liam H. Fowl, W. Ronny Huang, Wojciech Czaja, Gavin Taylor, Michael Moeller, and Tom Goldstein. Witches' brew: Industrial scale data poisoning via gradient matching. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net, 2021. URL https://openreview.net/forum?id=01olnfLIbD.

[12] Tianyu Gu, Brendan Dolan-Gavitt, and Siddharth Garg. Badnets: Identifying vulnerabilities in the machine learning model supply chain. *CoRR*, abs/1708.06733, 2017. URL http://arxiv.org/abs/1708.06733.

[13] Jonathan Hayase, Weihao Kong, Raghav Somani, and Sewoong Oh. Defense against backdoor attacks via robust covariance estimation. In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, volume 139 of *Proceedings of Machine Learning Research*, pages 4129–4139. PMLR, 2021. URL http://proceedings.mlr.press/v139/hayase21a.html.

[14] Jinyuan Jia, Xiaoyu Cao, and Neil Zhenqiang Gong. Certified robustness of nearest neighbors against data poisoning attacks. *CoRR*, abs/2012.03765, 2020. URL https://arxiv.org/abs/2012.03765.

[15] Jinyuan Jia, Xiaoyu Cao, and Neil Zhenqiang Gong. Intrinsic certified robustness of bagging against data poisoning attacks. In *Thirty-Fifth AAAI Conference on Artificial Intelligence, AAAI 2021, Thirty-Third Conference on Innovative Applications of Artificial Intelligence, IAAI 2021, The Eleventh Symposium on Educational Advances in Artificial Intelligence, EAAI 2021, Virtual Event, February 2-9, 2021*, pages 7961–7969. AAAI Press, 2021. URL https://ojs.aaai.org/index.php/AAAI/article/view/16971.

[16] Pang Wei Koh, Jacob Steinhardt, and Percy Liang. Stronger data poisoning attacks break data sanitization defenses. *Mach. Learn.*, 111(1):1–47, 2022. doi: 10.1007/s10994-021-06119-y. URL https://doi.org/10.1007/s10994-021-06119-y.

[17] Aounon Kumar, Alexander Levine, Tom Goldstein, and Soheil Feizi. Curse of dimensionality on randomized smoothing for certifiable robustness. In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, volume 119 of *Proceedings of Machine Learning Research*, pages 5458–5467. PMLR, 2020. URL http://proceedings.mlr.press/v119/kumar20b.html.

[18] Guang-He Lee, Yang Yuan, Shiyu Chang, and Tommi S. Jaakkola. Tight certificates of adversarial robustness for randomly smoothed classifiers. In Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence d'Alché-Buc, Emily B. Fox, and Roman Garnett, editors, *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pages 4911–4922, 2019. URL https://proceedings.neurips.cc/paper/2019/hash/fa2e8c4385712f9a1d24c363a2cbe5b8-Abstract.html.

[19] Alexander Levine and Soheil Feizi. Deep partition aggregation: Provable defenses against general poisoning attacks. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net, 2021. URL https://openreview.net/forum?id=YUGG2tFuPM.

[20] Kang Liu, Brendan Dolan-Gavitt, and Siddharth Garg. Fine-pruning: Defending against backdooring attacks on deep neural networks. In Michael Bailey, Thorsten Holz, Manolis Stamatogiannakis, and Sotiris Ioannidis, editors, *Research in Attacks, Intrusions, and Defenses - 21st International Symposium, RAID 2018, Heraklion, Crete, Greece, September 10-12, 2018, Proceedings*, volume 11050 of *Lecture Notes in Computer Science*, pages 273–294. Springer, 2018. doi: 10.1007/978-3-030-00470-5\_13. URL https://doi.org/10.1007/978-3-030-00470-5_13.

[21] Yingqi Liu, Shiqing Ma, Yousra Aafer, Wen-Chuan Lee, Juan Zhai, Weihang Wang, and Xiangyu Zhang. Trojaning attack on neural networks. In *25th Annual Network and Distributed System Security Symposium, NDSS 2018, San Diego, California, USA, February 18-21, 2018*. The Internet Society, 2018. URL http://wp.internetsociety.org/ndss/wp-content/uploads/sites/25/2018/02/ndss2018_03A-5_Liu_paper.pdf.

[22] Yuzhe Ma, Xiaojin Zhu, and Justin Hsu. Data poisoning against differentially-private learners: Attacks and defenses. In Sarit Kraus, editor, *Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence, IJCAI 2019, Macao, China, August 10-16, 2019*, pages 4732–4738. ijcai.org, 2019. doi: 10.24963/ijcai.2019/657. URL https://doi.org/10.24963/ijcai.2019/657.

[23] Anna P. Meyer, Aws Albarghouthi, and Loris D'Antoni. Certifying robustness to programmable data bias in decision trees. In Marc'Aurelio Ranzato, Alina Beygelzimer, Yann N.

Dauphin, Percy Liang, and Jennifer Wortman Vaughan, editors, *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, pages 26276–26288, 2021. URL `https://proceedings.neurips.cc/paper/2021/hash/dcf531edc9b229acfe0f4b87e1e278dd-Abstract.html`.

[24] Blaine Nelson, Marco Barreno, Fuching Jack Chi, Anthony D. Joseph, Benjamin I. P. Rubinstein, Udam Saini, Charles Sutton, J. Doug Tygar, and Kai Xia. Exploiting machine learning to subvert your spam filter. In Fabian Monrose, editor, *First USENIX Workshop on Large-Scale Exploits and Emergent Threats, LEET '08, San Francisco, CA, USA, April 15, 2008, Proceedings*. USENIX Association, 2008. URL `http://www.usenix.org/events/leet08/tech/full_papers/nelson/nelson.pdf`.

[25] Fanchao Qi, Yuan Yao, Sophia Xu, Zhiyuan Liu, and Maosong Sun. Turn the combination lock: Learnable textual backdoor attacks via word substitution. In Chengqing Zong, Fei Xia, Wenjie Li, and Roberto Navigli, editors, *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing, ACL/IJCNLP 2021, (Volume 1: Long Papers), Virtual Event, August 1-6, 2021*, pages 4873–4883. Association for Computational Linguistics, 2021. doi: 10.18653/v1/2021.acl-long.377. URL `https://doi.org/10.18653/v1/2021.acl-long.377`.

[26] Ximing Qiao, Yukun Yang, and Hai Li. Defending neural backdoors via generative distribution modeling. In Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence d'Alché-Buc, Emily B. Fox, and Roman Garnett, editors, *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pages 14004–14013, 2019. URL `https://proceedings.neurips.cc/paper/2019/hash/78211247db84d96acf4e00092a7fba80-Abstract.html`.

[27] Goutham Ramakrishnan and Aws Albarghouthi. Backdoors in neural models of source code. *CoRR*, abs/2006.06841, 2020. URL `https://arxiv.org/abs/2006.06841`.

[28] Elan Rosenfeld, Ezra Winston, Pradeep Ravikumar, and J. Zico Kolter. Certified robustness to label-flipping attacks via randomized smoothing. In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, volume 119 of *Proceedings of Machine Learning Research*, pages 8230–8241. PMLR, 2020. URL `http://proceedings.mlr.press/v119/rosenfeld20b.html`.

[29] Aniruddha Saha, Akshayvarun Subramanya, and Hamed Pirsiavash. Hidden trigger backdoor attacks. In *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020*, pages 11957–11965. AAAI Press, 2020. URL `https://ojs.aaai.org/index.php/AAAI/article/view/6871`.

[30] Giorgio Severi, Jim Meyer, Scott E. Coull, and Alina Oprea. Explanation-guided backdoor poisoning attacks against malware classifiers. In Michael Bailey and Rachel Greenstadt, editors, *30th USENIX Security Symposium, USENIX Security 2021, August 11-13, 2021*, pages 1487–1504. USENIX Association, 2021. URL `https://www.usenix.org/conference/usenixsecurity21/presentation/severi`.

[31] Ali Shafahi, W. Ronny Huang, Mahyar Najibi, Octavian Suciu, Christoph Studer, Tudor Dumitras, and Tom Goldstein. Poison frogs! targeted clean-label poisoning attacks on neural networks. In Samy Bengio, Hanna M. Wallach, Hugo Larochelle, Kristen Grauman, Nicolò Cesa-Bianchi, and Roman Garnett, editors, *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada*, pages 6106–6116, 2018. URL `https://proceedings.neurips.cc/paper/2018/hash/22722a343513ed45f14905eb07621686-Abstract.html`.

[32] Florian Tramèr, Nicholas Carlini, Wieland Brendel, and Aleksander Madry. On adaptive attacks to adversarial example defenses. In Hugo Larochelle, Marc'Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin, editors, *Advances in Neural Information Processing*

*Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020. URL `https://proceedings.neurips.cc/pap er/2020/hash/11f38f8ecd71867b42433548d1078e38-Abstract.html`.

[33] Brandon Tran, Jerry Li, and Aleksander Madry. Spectral signatures in backdoor attacks. In Samy Bengio, Hanna M. Wallach, Hugo Larochelle, Kristen Grauman, Nicolò Cesa-Bianchi, and Roman Garnett, editors, *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada*, pages 8011–8021, 2018. URL `https://proceedings.neurips. cc/paper/2018/hash/280cf18baf4311c92aa5a042336587d3-Abstract.html`.

[34] Alexander Turner, Dimitris Tsipras, and Aleksander Madry. Label-consistent backdoor attacks. *CoRR*, abs/1912.02771, 2019. URL `http://arxiv.org/abs/1912.02771`.

[35] Binghui Wang, Xiaoyu Cao, Jinyuan Jia, and Neil Zhenqiang Gong. On certifying robustness against backdoor attacks via randomized smoothing. *CoRR*, abs/2002.11750, 2020. URL `https://arxiv.org/abs/2002.11750`.

[36] Bolun Wang, Yuanshun Yao, Shawn Shan, Huiying Li, Bimal Viswanath, Haitao Zheng, and Ben Y. Zhao. Neural cleanse: Identifying and mitigating backdoor attacks in neural networks. In *2019 IEEE Symposium on Security and Privacy, SP 2019, San Francisco, CA, USA, May 19-23, 2019*, pages 707–723. IEEE, 2019. doi: 10.1109/SP.2019.00031. URL `https://doi.org/10.1109/SP.2019.00031`.

[37] Hongyi Wang, Kartik Sreenivasan, Shashank Rajput, Harit Vishwakarma, Saurabh Agarwal, Jy-yong Sohn, Kangwook Lee, and Dimitris S. Papailiopoulos. Attack of the tails: Yes, you really can backdoor federated learning. In Hugo Larochelle, Marc'Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin, editors, *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020. URL `https://proceedings.neurips. cc/paper/2020/hash/b8ffa41d4e492f0fad2f13e29e1762eb-Abstract.html`.

[38] Wenxiao Wang, Alexander Levine, and Soheil Feizi. Improved certified defenses against data poisoning with (deterministic) finite aggregation. *CoRR*, abs/2202.02628, 2022. URL `https://arxiv.org/abs/2202.02628`.

[39] Maurice Weber, Xiaojun Xu, Bojan Karlas, Ce Zhang, and Bo Li. RAB: provable robustness against backdoor attacks. *CoRR*, abs/2003.08904, 2020. URL `https://arxiv.org/abs/ 2003.08904`.

[40] Huang Xiao, Battista Biggio, Gavin Brown, Giorgio Fumera, Claudia Eckert, and Fabio Roli. Is feature selection secure against training data poisoning? In Francis R. Bach and David M. Blei, editors, *Proceedings of the 32nd International Conference on Machine Learning, ICML 2015, Lille, France, 6-11 July 2015*, volume 37 of *JMLR Workshop and Conference Proceedings*, pages 1689–1698. JMLR.org, 2015. URL `http://proceedings.mlr.press/v37/xiao15.html`.

[41] Yuanshun Yao, Huiying Li, Haitao Zheng, and Ben Y. Zhao. Latent backdoor attacks on deep neural networks. In Lorenzo Cavallaro, Johannes Kinder, XiaoFeng Wang, and Jonathan Katz, editors, *Proceedings of the 2019 ACM SIGSAC Conference on Computer and Communications Security, CCS 2019, London, UK, November 11-15, 2019*, pages 2041–2055. ACM, 2019. doi: 10.1145/3319535.3354209. URL `https://doi.org/10.1145/3319535.3354209`.

[42] Chen Zhu, W. Ronny Huang, Hengduo Li, Gavin Taylor, Christoph Studer, and Tom Goldstein. Transferable clean-label poisoning attacks on deep neural nets. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*, volume 97 of *Proceedings of Machine Learning Research*, pages 7614–7623. PMLR, 2019. URL `http://proceedings. mlr.press/v97/zhu19a.html`.

## Checklist

1. For all authors...

   (a) Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope? [Yes]

   (b) Did you describe the limitations of your work? [Yes] See Section 8.

   (c) Did you discuss any potential negative societal impacts of your work? [N/A]

   (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]

2. If you are including theoretical results...

   (a) Did you state the full set of assumptions of all theoretical results? [Yes]

   (b) Did you include complete proofs of all theoretical results? [Yes] We provide them in supplementary materials.

3. If you ran experiments...

   (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [Yes]

   (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes]

   (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [No] As our results are aggregated from thousands of models with the same statistical confidence level, we do not report the error bars.

   (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes]

4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...

   (a) If your work uses existing assets, did you cite the creators? [Yes]

   (b) Did you mention the license of the assets? [Yes] MIT license

   (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]

   (d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [N/A] They are either under MIT license or allowed for public usage.

   (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [N/A] There is no personally identifiable information or offensive content.

5. If you used crowdsourcing or conducted research with human subjects...

   (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]

   (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]

   (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]

# A    Probability Mass Functions of Distributions for $\text{FL}_s$, $\text{F}_s$, and $\text{L}$

**Distribution $\mu$ for $\text{F}_s$.**    The distribution $\mu(D, \mathbf{x})$ describes the outcomes of $\dot{D}$ and $\dot{\mathbf{x}}$. Each outcome $\dot{D}$ and $\dot{\mathbf{x}}$ can be represented as a combination of 1) selected indices $w_1, \ldots, w_k$, 2) smoothed training examples $\mathbf{x}'_{w_1}, \ldots, \mathbf{x}'_{w_k}$, and 3) smoothed test input $\mathbf{x}'$. The probability mass function (PMF) of $\mu(D, \mathbf{x})$ is,

$$p_{\mu(D,\mathbf{x})}(\dot{D}, \dot{\mathbf{x}}) = \frac{\rho^{(k+1)d}}{n^k} \left(\frac{\gamma}{\rho}\right)^{\sum_{i=1}^{k} \|\mathbf{x}'_{w_i} - \mathbf{x}_{w_i}\|_0 + \|\mathbf{x}' - \mathbf{x}\|_0}, \tag{10}$$

where $d$ is the dimension of the input feature, $\gamma = \frac{1-\rho}{K}$, $K + 1$ is the number of categories, and $\|\mathbf{x}' - \mathbf{x}\|_0$ is the $l_0$-norm, which counts the number of different features between $\mathbf{x}$ and $\mathbf{x}'$. Intuitively, Eq 10 is the multiply of the two PMFs of the two combined approaches because the bagging and the noise addition processes are independent.

**Distribution $\mu'$ for $\text{FL}_s$.**    For $\text{FL}_s$, we modify $\mu$ to $\mu'$ such that also flips the label of the selected instances. Concretely, $\dot{D}$ becomes $\{(\mathbf{x}'_{w_1}, y'_{w_1}), \ldots, (\mathbf{x}'_{w_k}, y'_{w_k})\}$, where $y'$ is a possibly flipped label with $\gamma$ probability. And the PMF of $\mu'(D, \mathbf{x})$ is,

$$p_{\mu'(D,\mathbf{x})}(\dot{D}, \dot{\mathbf{x}}) = \frac{\rho^{(k+1)d+k}}{n^k} \left(\frac{\gamma}{\rho}\right)^{\sum_{i=1}^{k} \|\mathbf{x}'_{w_i} - \mathbf{x}_{w_i}\|_0 + \mathbb{1}_{y_{w_i} \neq y'_{w_i}} + \|\mathbf{x}' - \mathbf{x}\|_0}, \tag{11}$$

where $\mathbb{1}_{y_{w_i} \neq y'_{w_i}}$ denotes whether the label $y'_{w_i}$ is flipped.

**Distribution $\mu''$ for $\text{L}$.**    For $\text{L}$, we modify $\mu$ to $\mu''$ such that it only flips the label of the selected instances. Also, it is not necessary to generate the smoothed test input because $\text{L}$ cannot flip input features. Concretely, $\dot{D}$ becomes $\{(\mathbf{x}_{w_1}, y'_{w_1}), \ldots, (\mathbf{x}_{w_k}, y'_{w_k})\}$, where $y'$ is a possibly flipped label with $\gamma$ probability. And the PMF of $\mu''(D)$ is,

$$p_{\mu''(D)}(\dot{D}) = \frac{\rho^k}{n^k} \left(\frac{\gamma}{\rho}\right)^{\sum_{i=1}^{k} \mathbb{1}_{y_{w_i} \neq y'_{w_i}}} \tag{12}$$

# B    Neyman–Pearson Lemma for the Multi-class Case

In multi-class case, instead of Eq 4, we need to certify

$$\forall \widetilde{D} \in S_r^\pi(D), \; \widetilde{\mathbf{x}} \in \pi(\mathbf{x}, y)_1. \; \Pr_{o \sim \widetilde{\mu}}(A(o) = y^*) > \Pr_{o \sim \widetilde{\mu}}(A(o) = y'), \tag{13}$$

where $y'$ is the runner-up prediction and $p'$ is the probability of predicting $y'$. Formally,

$$y' = \operatorname*{argmax}_{y \neq y^*} \Pr_{o \sim \mu}(A(o) = y), \quad p' = \Pr_{o \sim \mu}(A(o) = y')$$

We use Neyman–Pearson Lemma to check whether Eq 13 holds by computing a lower bound $\text{lb}$ of the LHS and an upper bound $\text{ub}$ of the RHS by solving the following optimization problems,

$$\text{lb} \triangleq \min_{A^? \in \mathcal{A}} \Pr_{o \sim \widetilde{\mu}}(A^?(o) = y^*), \quad \text{ub} \triangleq \max_{A^? \in \mathcal{A}} \Pr_{o \sim \widetilde{\mu}}(A^?(o) = y'), \tag{14}$$

$$s.t. \; \Pr_{o \sim \mu}(A^?(o) = y^*) = p^*, \; \Pr_{o \sim \mu}(A^?(o) = y') = p'$$

$$\widetilde{D} \in S_r^\pi(D), \widetilde{\mathbf{x}} \in \pi(\mathbf{x}, y)_1,$$

**Theorem B.1** (Neyman–Pearson Lemma for $\text{FL}_s$, $\text{F}_s$, $\text{L}$ in the Multi-class Case). *Let $\widetilde{D}$ and $\widetilde{\mathbf{x}}$ be a maximally perturbed dataset and test input, i.e., $|\widetilde{D} \backslash D| = r$, $\|\widetilde{\mathbf{x}} - \mathbf{x}\|_0 = s$, and $\|\widetilde{\mathbf{x}}_i - \mathbf{x}_i\|_0 + \mathbb{1}_{\widetilde{y}_i \neq y_i} = s$, for each perturbed example $(\widetilde{\mathbf{x}}_i, \widetilde{y}_i)$ in $\widetilde{D}$. Let $i_{\text{lb}} \triangleq \operatorname{argmin}_{i \in [1,m]} \sum_{j=1}^{i} p_\mu(\mathcal{L}_j) \geq p^*$ and*

$i_{\mathrm{ub}} \triangleq \operatorname{argmax}_{i \in [1,m]} \sum_{j=i}^{m} p_\mu(\mathcal{L}_j) \geq p'$, *The algorithm $\bar{A}^?$ is the minimizer of Eq 14 and its behaviors among $\forall i \in \{1, \ldots, m\}. \forall o \in \mathcal{L}_i$ are specified as*

$$\Pr(A^?(o) = y^*) = \begin{cases} 1, & i < i_{\mathrm{lb}} \\ \frac{p^* - \sum_{j=1}^{i-1} p_\mu(\mathcal{L}_j)}{p_\mu(\mathcal{L}_i)}, & i = i_{\mathrm{lb}} \\ 0, & i > i_{\mathrm{lb}} \end{cases} . \Pr(A^?(o) = y') = \begin{cases} 0, & i < i_{\mathrm{ub}} \\ \frac{p' - \sum_{j=i+1}^{m} p_\mu(\mathcal{L}_j)}{p_\mu(\mathcal{L}_i)}, & i = i_{\mathrm{ub}} \\ 1, & i > i_{\mathrm{ub}} \end{cases}$$

*Then,*

$$\mathrm{lb} = \sum_{j=1}^{i_{\mathrm{lb}}-1} p_{\widetilde{\mu}}(\mathcal{L}_j) + \left( p^* - \sum_{j=1}^{i_{\mathrm{lb}}-1} p_\mu(\mathcal{L}_j) \right) / \eta_{i_{\mathrm{lb}}}, \mathrm{ub} = \sum_{j=i_{\mathrm{ub}}+1}^{m} p_{\widetilde{\mu}}(\mathcal{L}_j) + \left( p' - \sum_{j=i_{\mathrm{ub}}+1}^{m} p_\mu(\mathcal{L}_j) \right) / \eta_{i_{\mathrm{ub}}}$$

Theorem B.1 is a direct application of Neyman–Pearson Lemma. And Lemma 4 in Lee et al. [18] proves the maximally perturbed dataset and test input achieve the worst-case bound lb and ub.

## B.1 Computing the Certified Radius of BagFlip for Perturbations $\mathrm{FL}_s$ and $\mathrm{L}$

**Computing the Certified Radius for $\mathrm{FL}_s$**  In Eq 7, we need to consider the flipping of labels,

$$\underbrace{\left( \sum_{i=1}^{k} \|\mathbf{x}'_{w_i} - \mathbf{x}_{w_i}\|_0 + \mathbb{1}_{y'_{w_i} \neq y_{w_i}} + \|\mathbf{x}' - \mathbf{x}\|_0 \right)}_{\Delta} - \underbrace{\left( \sum_{i=1}^{k} \|\mathbf{x}'_{w_i} - \widetilde{\mathbf{x}}_{w_i}\|_0 + \mathbb{1}_{y'_{w_i} \neq \widetilde{y}_{w_i}} + \|\mathbf{x}' - \widetilde{\mathbf{x}}\|_0 \right)}_{\widetilde{\Delta}} = t$$

In Eq 8, the definition of $T(c, t)$ should be modified to,

$$\sum_{\substack{0 \leq u_1, \ldots, u_c \leq d+1 \\ 0 \leq u_0 \leq d}} \sum_{\substack{0 \leq v_1, \ldots, v_c \leq d+1 \\ 0 \leq v_0 \leq d \\ u_0 - v_0 + \ldots + u_c - v_c = t}} \prod_{i=1}^{c} L(u_i, v_i; s, d+1) b^{u_i} a^{d+1-u_i} L(u_0, v_0; s, d) b^{u_0} a^{d-u_0}$$

The algorithm associated with Theorem 5.4 should be modified as

$$T(0, t) = \sum_{u=\max(0,t)}^{\min(d,t+d)} L(u, u - t; s, d) b^u a^{d-u}, \forall -d \leq t \leq d,$$

$$T(c, t) = \sum_{t_1 = \max(-d, t-cd-c+1)}^{\min(d, t+cd+c-1)} T(c-1, t - t_1) G(t_1), \forall c > 0, -(c+1)d - c \leq t \leq (c+1)d + c,$$

where $G(t)$ is defined as

$$G(t) = \sum_{u=\max(0,t)}^{\min(d+1, t+d+1)} L(u, u - t; s, d+1) b^u a^{d+1-u}, \forall -d-1 \leq t \leq d+1$$

**Computing the Certified Radius for $\mathrm{L}$**  In Eq 7, we only need to consider the flipping of labels,

$$\underbrace{\left( \sum_{i=1}^{k} \mathbb{1}_{y'_{w_i} \neq y_{w_i}} \right)}_{\Delta} - \underbrace{\left( \sum_{i=1}^{k} \mathbb{1}_{y'_{w_i} \neq \widetilde{y}_{w_i}} \right)}_{\widetilde{\Delta}} = t$$

In the rest of the computation of $T(c, t)$, we set $d = 1$, i.e., consider the label as a one-dimension feature.

## B.2 Practical Perspectives

**Estimation of $p^*$ and $p'$.**  For each test example $\mathbf{x}$, we need to estimate $p^*$ and $p'$ of the smoothed algorithm $\bar{A}$ given the benign dataset $D$. We use Monte Carlo sampling to compute $p^*$ and $p'$. Specifically, for each test input $\mathbf{x}$, we train $N$ algorithms with the datasets $\dot{D}_1, \ldots, \dot{D}_N$ and evaluate

these algorithms on input $\dot{\mathbf{x}}_i$, where $(\dot{D}_i, \dot{\mathbf{x}}_i)$ is sampled from distribution $\mu(D, \mathbf{x})$. We count the predictions equal to label $y$ as $N_y = \sum_{i=1}^{N} \mathbb{1}_{A(\dot{D}_i, \dot{\mathbf{x}}_i)=y}$ and use the Clopper-Pearson interval to estimate $p^*$ and $p'$.

$$p^* = \mathrm{Beta}(\frac{\alpha}{|\mathcal{C}|}; N_{y^*}, N - N_{y^*} + 1), \quad p' = \mathrm{Beta}(\frac{1-\alpha}{|\mathcal{C}|}; N_{y'} + 1, N - N_{y'})$$

where $1 - \alpha$ is the confidence level, $|\mathcal{C}|$ is the number of different labels, and and $\mathrm{Beta}(\beta; \lambda, \theta)$ is the $\beta$th quantile of the Beta distribution of parameter $\lambda$ and $\theta$. We further tighten the estimation of $p'$ by $\min(p', 1 - p^*)$ because $p^* + p' \leq 1$ by definition.

However, it is computationally expensive to retrain $N$ algorithms for each test input. We can reuse the trained $N$ algorithms to estimate $p^*$ of $m$ test inputs with a simultaneous confidence level at least $1 - \alpha$ by using *Bonferroni correction*. Specifically, we evenly divide $\alpha$ to $\frac{\alpha}{m}$ when estimating for each test input. In the evaluation, we set $m$ to be the size of the test set.

**Computation of the certified radius.** Previous sections have introduced how to check whether Eq 4 holds for a specific $r$. We can use binary search to find the certified radius given the estimated $p^*$ and $p'$. Although checking Eq 4 has been reduced to polynomial time, it might be infeasible for in-time prediction in the real scenario. We propose to memoize the certified radius by enumerating all possibilities of pairs of $p^*$ and $p'$ beforehand so that checking Eq 4 can be done in $O(1)$. Notice that a pair of $p^*$ and $p'$ is determined by $N_{y^*}$ and $N_{y'}$. Recall that $N$ is the number of trained algorithms. $N_{y^*}$ and $N_{y'}$ can have $O(N^2)$ different pairs ($O(N)$ in the binary-classification case).

### B.3 A Relaxation of the Neyman–Pearson Lemma

We introduce a relaxation of the Neyman–Pearson lemma for the multi-class case.

**Theorem B.2** (Relaxation of the Lemma). *Define* lb *for the original subspaces* $\mathcal{L}_1, \ldots, \mathcal{L}_m$ *as in Theorem B.1. Define* $\mathrm{lb}_\delta$ *for* $\{\mathcal{L}_i\}_{i \in B}$ *as in Theorem B.1 by underapproximating* $p^*$ *as* $p^* - \sum_{i \notin B} p_\mu(\mathcal{L}_i)$. *Define* ub *for the original subspaces* $\mathcal{L}_1, \ldots, \mathcal{L}_m$ *as in Theorem B.1. Define* $\mathrm{ub}'$ *for underapproximated subspaces* $\{\mathcal{L}_i\}_{i \in B}$ *as in Theorem B.1 with* $p'$ *and let* $\mathrm{ub}_\delta = \mathrm{ub}' + \sum_{i \notin B} p_{\widetilde{\mu}}(\mathcal{L}_i)$. *Then, we have **Soundness:*** $\mathrm{lb}_\delta \leq \mathrm{lb}, \mathrm{ub}_\delta \geq \mathrm{ub}$ *and* $\delta$-***Tightness:*** *Let* $\delta \triangleq \sum_{i \notin B} p_{\widetilde{\mu}}(\mathcal{L}_i)$, *then* $\mathrm{lb}_\delta + \delta \geq \mathrm{lb}, \mathrm{ub}_\delta - \delta \geq \mathrm{ub}$.

## C Proofs

**Proof of Theorem 5.2.**

*Proof.* If we know an outcome $o \in \mathcal{L}_{c,t}$ has $\Delta$ distance from the clean data, then $o$ has $(k+1)d - \Delta$ features unchanged and $\Delta$ features changed when sampling from $\mu$. Thus, according to the definition of PMF in Eq 10, $p_\mu(o) = \frac{1}{n^k} \rho^{(k+1)d-\Delta} \gamma^\Delta$ for all $o \in \mathcal{L}_{c,t}$. Similarly, $p_{\widetilde{\mu}}(o) = \frac{1}{n^k} \rho^{(k+1)d-\widetilde{\Delta}} \gamma^{\widetilde{\Delta}}$ for all $o \in \mathcal{L}_{c,t}$. By the definition of likelihood ratios, we have $\forall o \in \mathcal{L}_{c,t}.\ \eta(o) = \eta_{c,t} = p_\mu(o)/p_{\widetilde{\mu}}(o) = \gamma^{\Delta-\widetilde{\Delta}} \rho^{\widetilde{\Delta}-\Delta} = \left(\frac{\gamma}{\rho}\right)^t$. $\qquad\square$

**Proof of Theorem 5.3.**

*Proof.* We first define a subset of $\mathcal{L}_{c,t}$ as $\mathcal{L}_{c,t,\Delta}$,

$$\mathcal{L}_{c,t,\Delta} = \{(\{(\mathbf{x}'_{w_i}, y_{w_i})\}_i, \mathbf{x}') \mid$$

$$\sum_{i=1}^{k} \mathbb{1}_{\mathbf{x}_{w_i} \neq \widetilde{\mathbf{x}}_{w_i}} = c,$$

$$\left(\sum_{i=1}^{k} \|\mathbf{x}'_{w_i} - \mathbf{x}_{w_i}\|_0 + \|\mathbf{x}' - \mathbf{x}\|_0\right) - \left(\sum_{i=1}^{k} \|\mathbf{x}'_{w_i} - \widetilde{\mathbf{x}}_{w_i}\|_0 + \|\mathbf{x}' - \widetilde{\mathbf{x}}\|_0\right) = t,$$

$$\left(\sum_{i=1}^{k} \|\mathbf{x}'_{w_i} - \mathbf{x}_{w_i}\|_0 + \|\mathbf{x}' - \mathbf{x}\|_0\right) = \Delta\}$$

18

Denote the size of $\mathcal{L}_{c,t,\Delta}$ as $|\mathcal{L}_{c,t,\Delta}|$, then $p_\mu(\mathcal{L}_{c,t})$ can be computed as

$$p_\mu(\mathcal{L}_{c,t}) = \sum_{0\leq\Delta\leq(c+1)d} p_\mu(\mathcal{L}_{c,t,\Delta}) = \sum_{0\leq\Delta\leq(c+1)d} \frac{1}{n^k}\gamma^\Delta\rho^{d-\Delta}|\mathcal{L}_{c,t,\Delta}| \tag{15}$$

Because every outcome in $\mathcal{L}_{c,t,\Delta}$ has the same probability mass and we only need to count the size of $\mathcal{L}_{c,t,\Delta}$.

$$|\mathcal{L}_{c,t,\Delta}| = \binom{k}{c}r^c(n-r)^{k-c} \sum_{\substack{0\leq\widetilde{\Delta}_0,\ldots,\widetilde{\Delta}_c\leq d \\ 0\leq\Delta_0,\ldots,\Delta_c\leq d \\ \Delta_0-\widetilde{\Delta}_0+\ldots+\Delta_c-\widetilde{\Delta}_c=t \\ \Delta_0+\ldots+\Delta_c=\Delta}} \prod_{i=0}^{c} L(\Delta_i,\widetilde{\Delta}_i; s,d), \tag{16}$$

where $L(\Delta,\widetilde{\Delta}; s,d)$ is defined as, similarly in the Lemma 5 of Lee et al. [18],

$$\sum_{i=\max(0,\widetilde{\Delta}-s)}^{\min(\Delta,d-s,\lfloor\frac{\Delta+\widetilde{\Delta}-s}{2}\rfloor)} (K-1)^j \binom{s}{j}\binom{s-j}{\Delta-i-j}K^i\binom{d-s}{i},$$

where $j = \Delta + \widetilde{\Delta} - 2i - s$.

The binomial term in Eq 16 represents the different choices of selecting the $k$ indices with $c$ perturbed indices from a pool containing $r$ perturbed indices and $n-r$ clean indices. The rest of Eq 16 counts the number of different choices of flips. Specifically, we enumerate $\Delta$ and $\widetilde{\Delta}$ as $\sum_{i=0}^{c}\Delta_i$ and $\sum_{i=0}^{c}\widetilde{\Delta}_i$, where $i=0$ denotes the test input. And $L(\Delta_i,\widetilde{\Delta}_i; s,d)$ counts the number of different choices of flips for each example (see Lemma 5 of Lee et al. [18] for the derivation of $L(\Delta_i,\widetilde{\Delta}_i; s,d)$).

Then, plug Eq 16 into Eq 15, we have

$$p_\mu(\mathcal{L}_{c,t}) = \sum_{0\leq\Delta\leq(c+1)d}\frac{1}{n^k}\gamma^\Delta\rho^{d-\Delta}\binom{k}{c}r^c(n-r)^{k-c} \sum_{\substack{0\leq\widetilde{\Delta}_0,\ldots,\widetilde{\Delta}_c\leq d \\ 0\leq\Delta_0,\ldots,\Delta_c\leq d \\ \Delta_0-\widetilde{\Delta}_0+\ldots+\Delta_c-\widetilde{\Delta}_c=t \\ \Delta_0+\ldots+\Delta_c=\Delta}} \prod_{i=0}^{c} L(\Delta_i,\widetilde{\Delta}_i; s,d)$$

$$= \sum_{0\leq\Delta\leq(c+1)d}\frac{1}{n^k}\binom{k}{c}r^c(n-r)^{k-c} \sum_{\substack{0\leq\widetilde{\Delta}_0,\ldots,\widetilde{\Delta}_c\leq d \\ 0\leq\Delta_0,\ldots,\Delta_c\leq d \\ \Delta_0-\widetilde{\Delta}_0+\ldots+\Delta_c-\widetilde{\Delta}_c=t \\ \Delta_0+\ldots+\Delta_c=\Delta}} \prod_{i=0}^{c} L(\Delta_i,\widetilde{\Delta}_i; s,d)\gamma^{\Delta_i}\rho^{d-\Delta_i}$$

$$= \frac{1}{n^k}\binom{k}{c}r^c(n-r)^{k-c} \sum_{\substack{0\leq\widetilde{\Delta}_0,\ldots,\widetilde{\Delta}_c\leq d \\ 0\leq\Delta_0,\ldots,\Delta_c\leq d \\ \Delta_0-\widetilde{\Delta}_0+\ldots+\Delta_c-\widetilde{\Delta}_c=t}} \prod_{i=0}^{c} L(\Delta_i,\widetilde{\Delta}_i; s,d)\gamma^{\Delta_i}\rho^{d-\Delta_i}$$

$$= \text{Binom}(c; k, \frac{r}{n}) \sum_{\substack{0\leq\widetilde{\Delta}_0,\ldots,\widetilde{\Delta}_c\leq d \\ 0\leq\Delta_0,\ldots,\Delta_c\leq d \\ \Delta_0-\widetilde{\Delta}_0+\ldots+\Delta_c-\widetilde{\Delta}_c=t}} \prod_{i=0}^{c} L(\Delta_i,\widetilde{\Delta}_i; s,d)\gamma^{\Delta_i}\rho^{d-\Delta_i}$$

$\square$

**Proof of Theorem 5.4.**

*Proof.* From Eq 8 in Theorem 5.3, then for $c \geq 1$, we have

$$T(c,t) = \sum_{\substack{0 \leq \Delta_0, \ldots, \Delta_c \leq d}} \sum_{\substack{0 \leq \widetilde{\Delta}_0, \ldots, \widetilde{\Delta}_c \leq d \\ \Delta_0 - \widetilde{\Delta}_0 + \ldots + \Delta_c - \widetilde{\Delta}_c = t}} \prod_{i=0}^{c} L(\Delta_i, \widetilde{\Delta}_i; s, d) \gamma^{\Delta_i} \rho^{d-\Delta_i}$$

$$= \sum_{-d \leq t_1 \leq d} \sum_{0 \leq \Delta_c \leq d} \sum_{\substack{0 \leq \widetilde{\Delta}_c \leq d \\ \Delta_c - \widetilde{\Delta}_c = t_1}} L(\Delta_c, \widetilde{\Delta}_c; s, d) \gamma^{\Delta_c} \rho^{d-\Delta_c} \times$$

$$\sum_{\substack{0 \leq \Delta_0, \ldots, \Delta_{c-1} \leq d}} \sum_{\substack{0 \leq \widetilde{\Delta}_0, \ldots, \widetilde{\Delta}_{c-1} \leq d \\ \Delta_0 - \widetilde{\Delta}_0 + \ldots + \Delta_{c-1} - \widetilde{\Delta}_{c-1} = t - t_1}} \prod_{i=0}^{c-1} L(\Delta_i, \widetilde{\Delta}_i; s, d) \gamma^{\Delta_i} \rho^{d-\Delta_i}$$

$$= \sum_{-d \leq t_1 \leq d} T(0, t_1) T(c-1, t-t_1)$$

$\square$

**Proof of Theorem 6.1.**

Before giving the formal proof, we motivate the proof by the following knapsack problem, where each item can be divided arbitrarily. This allows a greedy algorithm to solve the problem, the same as the greedy process in Theorem 5.1.

**Example C.1.** *Suppose we have $m$ items with volume $p_\mu(\mathcal{L}_i)$ and cost $p_{\widetilde{\mu}}(\mathcal{L}_i)$. We have a knapsack with volume $p^*$. Determine the best strategy to fill the knapsack with the minimal cost* lb. *Note that each item can be divided arbitrarily.*

*The greedy algorithm sorts the item descendingly by "volume per cost" $p_\mu(\mathcal{L}_i)/p_{\widetilde{\mu}}(\mathcal{L}_i)$ (likelihood ratio) and select items until the knapsack is full. The last selected item $\mathcal{L}_{i_{\mathrm{lb}}}$ will be divided to fill the knapsack. Define the best solution as $S$ in this case and the minimal cost as* lb.

Now consider Theorem 6.1, which removes items $\{\mathcal{L}_i\}_{i \notin B}$ and reduces the volume of knapsack by the sum of the removed items' volume. Applying the greedy algorithm again, denote the best solution as $S'$ and the minimal cost in this case as $\mathrm{lb}_\delta$.

Soundness: The above process of removing items and reducing volume of knapsack is equivalent to just setting the cost of items in $\{\mathcal{L}_i\}_{i \notin B}$ as zero. Then, the new cost $\mathrm{lb}_\delta$ will be better than before (less than lb) because the cost of some items has been set to zero.

$\delta$-tightness: If we put removed items back to the reduced knapsack solution $S'$, this new solution is a valid selection in the original problem with cost $\mathrm{lb}_\delta + \sum_{i \notin B} p_{\widetilde{\mu}}(\mathcal{L}_i)$, and this cost cannot be less than the minimal cost lb, i.e., $\mathrm{lb}_\delta + \sum_{i \notin B} p_{\widetilde{\mu}}(\mathcal{L}_i) \geq \mathrm{lb}$.

*Proof.* We first consider a base case when $|B| = m - 1$, i.e., only one subspace is underapproximated. We denote the index of that subspace as $i'$. Consider $i_{\mathrm{lb}}$ computed in Theorem 5.1.

- If $i_{\mathrm{lb}} > i'$, then the likelihood ratio from $\{\mathcal{L}_i\}_{i \notin B}$ is used for computing lb in Theorem 5.1, meaning $\mathrm{lb}_\delta = \mathrm{lb} - p_{\widetilde{\mu}}(\mathcal{L}_{i'})$. This implies both soundness and $\delta$-tightness.

- If $i_{\mathrm{lb}} = i'$, then $\mathcal{L}_{i'}$ partially contributes to the computation of lb in Theorem 5.1. First, $\mathrm{lb}_\delta \leq \mathrm{lb}$ because the computation of $\mathrm{lb}_\delta$ can only sum up to $i_{\mathrm{lb}} - 1$ items (it cannot sum $i_{\mathrm{lb}}$th item), which implies $\mathrm{lb} \geq \sum_{i=1}^{i_{\mathrm{lb}}-1} p_{\widetilde{\mu}}(\mathcal{L}_i) \geq \mathrm{lb}_\delta$.

  Next, we are going to prove $\mathrm{lb}_\delta + p_{\widetilde{\mu}}(\mathcal{L}_{i'}) \geq \mathrm{lb}$. Suppose the additional budget $p_\mu(\mathcal{L}_{i'})$ for lb selects additional subspaces (than underapproximated $\mathrm{lb}_\delta$) with likelihood ratio $\frac{p_\mu(\mathcal{L}_{i'-q})}{p_{\widetilde{\mu}}(\mathcal{L}_{i'-q})}, \ldots, \frac{p_\mu(\mathcal{L}_{i'-1})}{p_{\widetilde{\mu}}(\mathcal{L}_{i'-1})}, \frac{p_\mu(l_{i'})}{p_{\widetilde{\mu}}(l_{i'})}$ such that $p_\mu(\mathcal{L}_{i'}) = \sum_{j=1}^{q} p_\mu(\mathcal{L}_{i'-j}) + p_\mu(l_{i'})$, $\mathrm{lb} = \mathrm{lb}_\delta + \sum_{j=1}^{q} p_{\widetilde{\mu}}(\mathcal{L}_{i'-j}) + p_{\widetilde{\mu}}(l_{i'})$, and $l_{i'} \subseteq \mathcal{L}_{i'}$. Then we have $\mathrm{lb} \leq \mathrm{lb}_\delta + p_{\widetilde{\mu}}(\mathcal{L}_{i'})$ because $\frac{p_\mu(\mathcal{L}_{i'-q})}{p_{\widetilde{\mu}}(\mathcal{L}_{i'-q})} \geq \cdots \geq \frac{p_\mu(\mathcal{L}_{i'-1})}{p_{\widetilde{\mu}}(\mathcal{L}_{i'-1})} \geq \frac{p_\mu(l_{i'})}{p_{\widetilde{\mu}}(l_{i'})} = \frac{p_\mu(\mathcal{L}_{i'})}{p_{\widetilde{\mu}}(\mathcal{L}_{i'})}$ implies $\sum_{j=1}^{q} p_{\widetilde{\mu}}(\mathcal{L}_{i'-j}) + p_{\widetilde{\mu}}(l_{i'}) \leq$

$p_{\widetilde{\mu}}(\mathcal{L}_{i'})$. To see this implication, we have $p_{\widetilde{\mu}}(\mathcal{L}_{i'}) = p_\mu(\mathcal{L}_{i'})/\eta_{i'} = \sum_{j=1}^q p_\mu(\mathcal{L}_{i'-j})/\eta_{i'} + p_\mu(l_{i'})/\eta_{i'} \geq \sum_{j=1}^q p_\mu(\mathcal{L}_{i'-j})/\eta_{i'-j} + p_\mu(l_{i'})/\eta_{i'} = \sum_{j=1}^q p_{\widetilde{\mu}}(\mathcal{L}_{i'-j}) + p_{\widetilde{\mu}}(l_{i'})$.

- If $i_{\mathrm{lb}} < i'$, then $\mathrm{lb}_\delta \leq \mathrm{lb}$ because lb has more budget as $p^*$ than $p^* - p_\mu(\mathcal{L}_{i'})$. Suppose the additional budget $p_\mu(\mathcal{L}_{i'})$ for lb selects additional subspaces (than underapproximated $\mathrm{lb}_\delta$) with likelihood ratio $\frac{p_\mu(\mathcal{L}_{i_{\mathrm{lb}}-q})}{p_{\widetilde{\mu}}(\mathcal{L}_{i_{\mathrm{lb}}-q})}, \ldots, \frac{p_\mu(\mathcal{L}_{i_{\mathrm{lb}}})}{p_{\widetilde{\mu}}(\mathcal{L}_{i_{\mathrm{lb}}})}$ (we assume it selects the whole $\mathcal{L}_{i_{\mathrm{lb}}}$ for simplicity, and if it selects a subset of $\mathcal{L}_{i_{\mathrm{lb}}}$ can be proved similarly) such that $p_\mu(\mathcal{L}_{i'}) = \sum_{j=0}^q p_\mu(\mathcal{L}_{i_{\mathrm{lb}}-j})$ and $\mathrm{lb} = \mathrm{lb}_\delta + \sum_{j=0}^q p_{\widetilde{\mu}}(\mathcal{L}_{i_{\mathrm{lb}}-j})$. Then we have $\mathrm{lb} \leq \mathrm{lb}_\delta + p_{\widetilde{\mu}}(\mathcal{L}_{i'})$ because $\frac{p_\mu(\mathcal{L}_{i_{\mathrm{lb}}-j})}{p_{\widetilde{\mu}}(\mathcal{L}_{i_{\mathrm{lb}}-j})} \geq \ldots \geq \frac{p_\mu(\mathcal{L}_{i_{\mathrm{lb}}})}{p_{\widetilde{\mu}}(\mathcal{L}_{i_{\mathrm{lb}}})} \geq \frac{p_\mu(\mathcal{L}_{i'})}{p_{\widetilde{\mu}}(\mathcal{L}_{i'})}$ implies $\sum_{j=0}^q p_{\widetilde{\mu}}(\mathcal{L}_{i_{\mathrm{lb}}-j}) \leq p_{\widetilde{\mu}}(\mathcal{L}_{i'})$. The reason of implication can be proved in a similar way as above.

We then consider $|B| < m - 1$, i.e., more than one subspaces are underapproximated. We separate $\{\mathcal{L}_i\}_{i \notin B}$ into two parts, one $\{\mathcal{L}_{i'}\}$ contains any one of the subspaces, the other contains the rest $\{\mathcal{L}_i\}_{i \notin B \wedge i \neq i'}$.

Denote $\mathrm{lb}'_\delta$ for $\{\mathcal{L}_i\}_{i \in B \vee i=i'}$ as in Theorem 5.1 by underapproximating $p^*$ as $p^* - \sum_{i \notin B} p_\mu(\mathcal{L}_i) + p_\mu(\mathcal{L}_{i'})$. By inductive hypothesis, $\mathrm{lb} \geq \mathrm{lb}'_\delta$ and $\mathrm{lb} \leq \mathrm{lb}'_\delta + \sum_{i \notin B} p_\mu(\mathcal{L}_i) - p_\mu(\mathcal{L}_{i'})$. By the same process of the above proof when one subspace is underapproximated (comparing $\{\mathcal{L}_i\}_{i \in B \vee i=i'}$ with $\{\mathcal{L}_i\}_{i \in B}$), we have $\mathrm{lb}'_\delta \geq \mathrm{lb}_\delta$ and $\mathrm{lb}'_\delta \leq \mathrm{lb}_\delta + p_\mu(\mathcal{L}_{i'})$. Combining the results above, we have $\mathrm{lb} \geq \mathrm{lb}_\delta$ and $\mathrm{lb} \leq \mathrm{lb}_\delta + \sum_{i \notin B} p_\mu(\mathcal{L}_i)$.

$\square$

# D  A KL-divergence Bound on the Certified Radius

We can use KL divergence [8] to get a looser but computationally-cheaper bound on the certified radius. Here, we certify the trigger-less case for $\mathrm{F}_s$.

**Theorem D.1.** *Consider the binary classification case, Eq 4 holds if*

$$r < \frac{n \log(4p^*(1-p^*))}{2k \log(\frac{\gamma}{\rho})(\rho - \gamma)s}, \tag{17}$$

*where $n$ is the size of the training set, $r$ is the certified radius, $s$ is the number of the perturbed features, $k$ is the size of each bag, and $\rho, \gamma$ are the probabilities of a featuring remaining the same and being flipped.*

**Lemma D.1.** *Define $T$ as*

$$T = \sum_{u=0}^d \sum_{v=0}^d L(u,v;s,d)\gamma^u \rho^{d-u}(v-u) \tag{18}$$

*where $L(u,v;s,d)$ is the same quantity defined in Lee et al. [18]. Then, we have*

$$T = (\rho - \gamma)s$$

**Proof of the Theorem D.1.**

*Proof.* Denote $D' \sim \mu(D)$ and $D'' \sim \mu(\widetilde{D})$, from the theorem in Example 5 of [28], Eq 4 holds if

$$\mathrm{KL}(D''\|D') < -\frac{1}{2}\log(4p^*(1-p^*)) \tag{19}$$

Denote the PMF of selecting an index $w$ and flip $\mathbf{x}_w$ to $\mathbf{x}'_w$ by $\mu(D)$ as $p_\mu(\mathbf{x}'_w) = \frac{\rho^d}{n}\left(\frac{\gamma}{\rho}\right)^{\|\mathbf{x}'_w - \mathbf{x}_w\|_0}$. Similarly, the PMF of selecting an index $w$ and flip $\mathbf{x}_w$ to $\mathbf{x}'_w$ by $\mu(\widetilde{D})$ as $p_{\widetilde{\mu}}(\mathbf{x}'_w) = \frac{\rho^d}{n}\left(\frac{\gamma}{\rho}\right)^{\|\mathbf{x}'_w - \widetilde{\mathbf{x}}_w\|_0}$. We now calculate the KL divergence between the distribution generated from the

perturbed dataset $\widetilde{D}$ and the distribution generated from the original dataset $D$.

$$\mathrm{KL}(D''\|D')$$
$$=k\mathrm{KL}(D_1''\|D_1') \tag{20}$$
$$=k\sum_{w=1}^{|D|}\sum_{\mathbf{x}_w'\in[K]^d}p_{\widetilde{\mu}}(\mathbf{x}_w')\log\frac{p_{\widetilde{\mu}}(\mathbf{x}_w')}{p_\mu(\mathbf{x}_w')}$$
$$=k\sum_{\mathbf{x}_w\neq\widetilde{\mathbf{x}}_w}\sum_{\mathbf{x}_w'\in[K]^d}p_{\widetilde{\mu}}(\mathbf{x}_w')\log\frac{p_{\widetilde{\mu}}(\mathbf{x}_w')}{p_\mu(\mathbf{x}_w')} \tag{21}$$
$$\leq kr\sum_{\mathbf{x}_w'\in[K]^d}\frac{\rho^d}{n}\left(\frac{\gamma}{\rho}\right)^{\|\mathbf{x}_w'-\widetilde{\mathbf{x}}_w\|_0}\log(\frac{\gamma}{\rho})(\|\mathbf{x}_w'-\widetilde{\mathbf{x}}_w\|_0-\|\mathbf{x}_w'-\mathbf{x}_w\|_0)$$
$$=k\frac{r}{n}\rho^d\log(\frac{\gamma}{\rho})\sum_{\mathbf{x}_w'\in[K]^d}\left(\frac{\gamma}{\rho}\right)^{\|\mathbf{x}_w'-\widetilde{\mathbf{x}}_w\|_0}(\|\mathbf{x}_w'-\widetilde{\mathbf{x}}_w\|_0-\|\mathbf{x}_w'-\mathbf{x}_w\|_0) \tag{22}$$

where $\mathrm{KL}(D_1''\|D_1')$ is the KL divergence of the first selected instance. We have Eq 20 because each selected instance is independent. We have Eq 21 because the $p_\mu(\mathbf{x}_w')$ and $p_{\widetilde{\mu}}(\mathbf{x}_w')$ only differs when the $w$th instance is perturbed.

The attacker can modify $\widetilde{D}$ to maximize Eq 22. And the Lemma 4 in Lee et al. [18] states that the maximal value is achieved when $\widetilde{\mathbf{x}}_w$ has exact $s$ features flipped to another value. Now suppose there are $s$ features flipped in $\widetilde{\mathbf{x}}_w$, we then need to compute Eq 22. If we denote $\|\mathbf{x}_w'-\widetilde{\mathbf{x}}_w\|_0$ as $u$ and $\|\mathbf{x}_w'-\mathbf{x}_w\|_0$ as $v$, and we count the size of the set $L(u,v;s,d)=\{\mathbf{x}_w'\in[K]^d\mid\|\mathbf{x}_w'-\widetilde{\mathbf{x}}_w\|_0=u,\|\mathbf{x}_w'-\mathbf{x}_w\|_0=v,\|\mathbf{x}_w-\widetilde{\mathbf{x}}_w\|_0=s\}$, then we can compute Eq 22 as

$$k\frac{r}{n}\rho^d\log(\frac{\rho}{\gamma})\sum_{u=0}^d\sum_{v=0}^d|L(u,v;s,d)|\left(\frac{\gamma}{\rho}\right)^u(v-u)$$

Let $T$ be defined as in Eq 18, we then have

$$\mathrm{KL}(D''\|D')=k\frac{r}{n}\log(\frac{\rho}{\gamma})T \tag{23}$$

Combine Eq 19, Eq 23, and Lemma D.1, we have

$$k\frac{r}{n}\log(\frac{\rho}{\gamma})(\rho-\gamma)s\leq\mathrm{KL}(D''\|D')<-\frac{1}{2}\log(4p^*(1-p^*))$$
$$r<\frac{n\log(4p^*(1-p^*))}{2k\log(\frac{\gamma}{\rho})(\rho-\gamma)s}$$

$\square$

**Proof of the Lemma D.1.**

*Proof.* Notice that the value of $T$ does not defend on the feature dimension $d$. Thus, we prove the lemma by induction on $d$. We further denote the value of $T$ under the feature dimension $d$ as $T_d$.

Let $L(u,v;s,d)$ be defined as in the Lemma 5 of Lee et al. [18],

$$\sum_{i=\max(0,v-s)}^{\min(u,d-s,\lfloor\frac{u+v-s}{2}\rfloor)}(K-1)^j\binom{s}{j}\binom{s-j}{u-i-j}K^i\binom{d-s}{i},$$

where $j=u+v-2i-s$.

**Base case.** Because $0\leq s\leq d$, when $d=0$, it is easy to see $T_0=s(\rho-\gamma)=0$.

**Induction case.** We first prove $T_{d+1}=s(\rho-\gamma)$ under a special case, where $s=d+1$, given the inductive hypothesis $T_d=s(\rho-\gamma)$ for $s=d$.

By definition

$$T_{d+1} = \sum_{u=0}^{d+1} \sum_{v=0}^{d+1} |L(u,v;s,d+1)| \gamma^u \rho^{d+1-u}(v-u), \tag{24}$$

By the definition of $L(u,v;s,d)$, we have the following equation when $0 \le s \le d$,

$$|L(u,v;d+1,d+1)| = (K-1)|L(u-1,v-1;d,d)| + |L(u-1,v;d,d)| + |L(u,v-1;d,d)| \tag{25}$$

and

$$\sum_{u=0}^{d} \sum_{v=0}^{d} |L(u,v;d,d)| \gamma^u \rho^{d-u} = 1 \tag{26}$$

Plug Eq 25 into Eq 24, we have the following equations when $s = d+1$,

$$T_{d+1} = \sum_{u=0}^{d+1} \sum_{v=0}^{d+1} (K-1)|L(u-1,v-1;d,d)| \gamma^u \rho^{d+1-u}(v-u) +$$

$$\sum_{u=0}^{d+1} \sum_{v=0}^{d+1} |L(u-1,v;d,d)| \gamma^u \rho^{d+1-u}(v-u) +$$

$$\sum_{u=0}^{d+1} \sum_{v=0}^{d+1} |L(u,v-1;d,d)| \gamma^u \rho^{d+1-u}(v-u)$$

$$= \gamma(K-1) \sum_{u=0}^{d} \sum_{v=0}^{d} |L(u,v;d,d)| \gamma^u \rho^{d-u}(v-u) + \tag{27}$$

$$\gamma \sum_{u=0}^{d} \sum_{v=0}^{d} |L(u,v;d,d)| \gamma^u \rho^{d-u}(v-u-1) + \tag{28}$$

$$\rho \sum_{u=0}^{d} \sum_{v=0}^{d} |L(u,v;d,d)| \gamma^u a^{d-u}(v-u+1) \tag{29}$$

$$= \sum_{u=0}^{d} \sum_{v=0}^{d} |L(u,v;d,d)| \gamma^u \rho^{d-u}[(v-u)+(\rho-\gamma)] \tag{30}$$

$$= d(\rho-\gamma) + (\rho-\gamma) \tag{31}$$

$$= s(\rho-\gamma)$$

We have Eq 25, Eq 27 and Eq 28 because $L(u,v;s,d) = 0$ when $u > d$, $v > d$, $u < 0$, or $v < 0$. We have Eq 30 because $\gamma(K-1) + \gamma + \rho = 1$ as $\gamma$ is defined as $\frac{1-\rho}{K}$. We have Eq 31 by plugging in Eq 26.

Next, we are going to show that $T_{d+1} = s(\rho-\gamma)$ for all $0 \le s \le d$, given the inductive hypothesis $T_d = s(\rho-\gamma)$ for all $0 \le s \le d$.

By the definition of $L(u,v;s,d)$, we have the following equations when $0 \le s \le d$,

$$|L(u,v;s,d+1)| = |L(u,v;s,d)| + K|L(u-1,v-1;s,d)| \tag{32}$$

23

Table 3: This paper compared to other approaches. RAB [39] can handle perturbation $\mathrm{F}_s^*$ that perturbs the input within a $l_2$-norm ball of radius $s$.

| Approach | Perturbation $\pi$ | Probability measure $\mu$ | Goal |
|---|---|---|---|
| Jia et al. [15], Chen et al. [4] | Any | Bagging | Trigger-less |
| Rosenfeld et al. [28] | L | Noise | Trigger-less |
| RAB [39] | $\mathrm{F}_s^*$ | Noise | Trigger-less, Backdoor |
| Wang et al. [35] | $\mathrm{FL}_s, \mathrm{F}_s, \mathrm{L}$ | Noise | Trigger-less, Backdoor |
| *This paper* | $\mathrm{FL}_s, \mathrm{F}_s, \mathrm{L}$ | Bagging+Noise | Trigger-less, Backdoor |

Plug Eq 32 into Eq 24, for all $0 \le s \le d$, we have

$$
\begin{aligned}
T_{d+1} &= \sum_{u=0}^{d+1} \sum_{v=0}^{d+1} |L(u,v;s,d)| \gamma^u \rho^{d+1-u}(v-u) + \\
&\quad \sum_{u=0}^{d+1} \sum_{v=0}^{d+1} K|L(u-1,v-1;s,d)| \gamma^u \rho^{d+1-u}(v-u) \\
&= \rho \sum_{u=0}^{d} \sum_{v=0}^{d} |L(u,v;s,d)| \gamma^u \rho^{d-u}(v-u) + \\
&\quad \gamma K \sum_{u=0}^{d} \sum_{v=0}^{d} |L(u,v;s,d)| \gamma^u \rho^{d-u}(v-u) \\
&= \rho s(\rho - \gamma) + \gamma K s(\rho - \gamma) \\
&= s(\rho - \gamma)
\end{aligned}
$$

$$(33)$$

$$(34)$$

We have Eq 33 and Eq 34 because $L(u,v;s,d) = 0$ when $u > d$, $v > d$, $u < 0$, or $v < 0$.

$\square$

# E   Experiments

## E.1   Dataset Details

MNIST is an image classification dataset containing 60,000 training and 10,000 test examples. Each example can be viewed as a vector containing 784 ($28 \times 28$) features.

CIFAR10 is an image classification dataset containing 50,000 training and 10,000 test examples. Each example can be viewed as a vector containing 3072 ($32 \times 32 \times 3$) features.

EMBER is a malware detection dataset containing 600,000 training and 200,000 test examples. Each example is a vector containing 2,351 features.

Contagio is a malware detection dataset, where each example is a vector containing 135 features. We partition the dataset into 6,000 training and 4,000 test examples.

MNIST-17 is a sub-dataset of MNIST, which contains 13,007 training and 2,163 test examples. MNIST-01 is a sub-dataset of MNIST, which contains 12,665 training and 2,115 test examples. CIFAR10-02 is a sub-dataset of CIFAR10, which contains 10,000 training and 2,000 test examples.

For CIFAR10, we categorize each feature into 4 categories. For the rest of the datasets, we binarize each feature. A special case is L, where we do not categorize features.

## E.2   Experiment Results Details

We train all models on a server running Ubuntu 18.04.5 LTS with two V100 32GB GPUs and Intel Xeon Gold 5115 CPUs running at 2.40GHz. For computing the certified radius, we run experiments across hundreds of machines in high throughput computing center.

### E.2.1 Defend Trigger-less Attacks

**Comparison to Bagging**  We show full results of comparison on $F_s$ in Figures 3 4 and 5. The results are similar as described in Section 7.

We additionally compare BagFlip with Bagging on $FL_1$ using MNIST-17 and $L$ using MNIST, and show the results in Figure 6. BagFlip still outperforms Bagging on $FL_1$ using MNIST-17. However, Bagging outperforms BagFlip on $L$ because when the attacker is only able to perturb the label, then $s = 1$ is equal to $s = \infty$ and flipping the labels hurts the accuracy.

**Comparison to LabelFlip**  We compare two configurations of BagFlip to LabelFlip using MNIST and show results of BagFlip in Figure 7. The results show that LabelFlip achieves less than $60\%$ normal accuracy, while BagFlip-0.95 (BagFlip-0.9) achieves 89.2% (88.6%) normal accuracy, respectively. BagFlip achieves higher certified accuracy than LabelFlip across all $R$. In particular, the certified accuracy of LabelFlip drops to less than $20\%$ when $R = 0.83$, while BagFlip-0.95 (BagFlip-0.9) still achieves 38.9% (36.2%) certified accuracy, respectively. **BagFlip has higher normal accuracy and certified accuracy than LabelFlip.**

### E.2.2 Defend Backdoor Attacks

We set $k = 100$ for MNIST-17 when comparing to FeatFlip. We set $k = 50, 200$ for MNIST-01 and CIFAR10-02 respectively when comparing to RAB. And we set $k = 100, 1000, 3000, 30$ for MNIST, CIFAR10, EMBER, and Contagio respectively when evaluating BagFlip on $F_1$. We use BadNets to modify 10% of examples in the training set.

**Comparison to RAB**  We show the comparison with full configurations of RAB-$\sigma$ in Figures 8(a) and 8(b), where $\sigma = 0.5, 1, 2$ are different Gaussian noise levels. Note that RAB's curves are not visible because the certified radius is almost zero anywhere.

**Results on EMBER and CIFAR10**  BagFlip cannot compute effective certificates for CIFAR10, i.e., the certified accuracy is zero even at $R = 0$, thus we do not show the figure for CIFAR10. Figure 9 shows the results of BagFlip on EMBER. BagFlip cannot compute effective certificates for EMBER, neither. We leave the improvement of BagFlip as a future work.



Figure 9: BagFlip-0.95 on EMBER against backdoor attack with $F_1$. Dashed lines show normal accuracy.

(a) $F_1$ on MNIST

(b) $F_\infty$ on MNIST

(c) $F_1$ on EMBER

(d) $F_\infty$ on EMBER
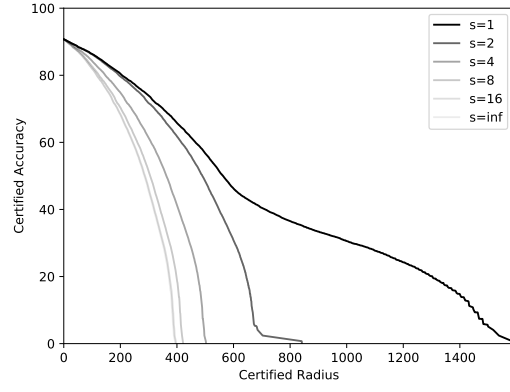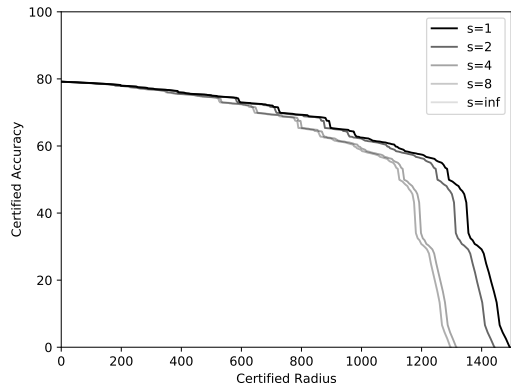
(e) $F_1$ on Contagio

(f) $F_\infty$ on Contagio

Figure 3: Compared to Bagging [15]. The horizontal dashed lines show the normal accuracy. The solid lines show the certified accuracy at different $R$. BagFlip-$a$ shows the result of the noise level $a$. The blue line shows the uncategorized version of bagging.
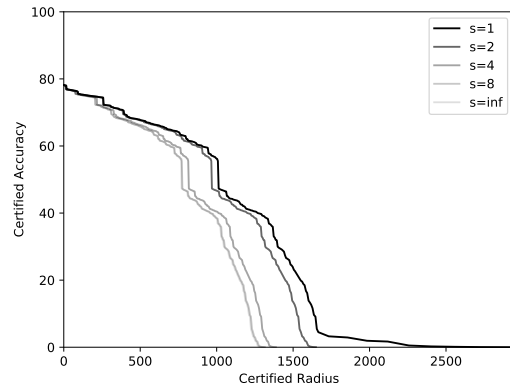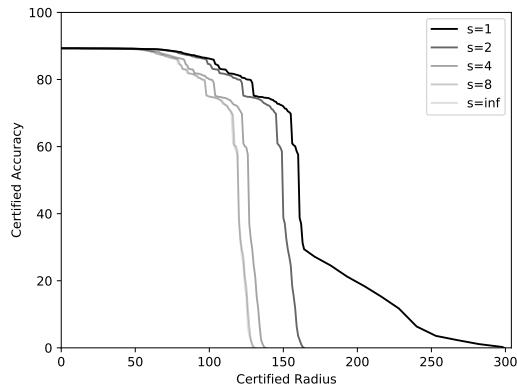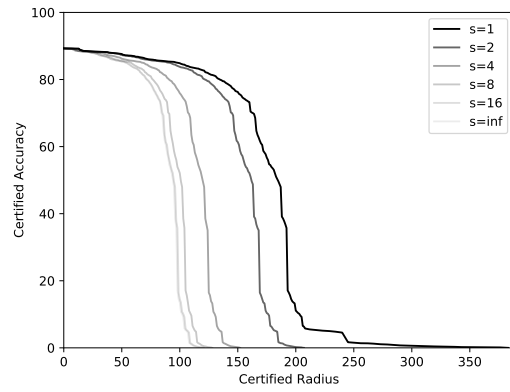
26

(a) BagFlip-0.9 on MNIST

(b) BagFlip-0.8 on MNIST

(c) BagFlip-0.95 on EMBER

(d) BagFlip-0.9 on EMBER

(e) BagFlip-0.9 on Contagio

(f) BagFlip-0.8 on Contagio

Figure 4: Results of BagFlip on different $s$ and different datasets.

27

(a) $F_1$ on CIFAR10

(b) $F_1$ on CIFAR10

(c) BagFlip-0.95 on CIFAR10

(d) BagFlip-0.9 on CIFAR10

Figure 5: Results of BagFlip on CIFAR10.



(a) Compared to Bagging on $FL_1$ on MNIST-17

(b) Compared to Bagging on $L$ using MNIST

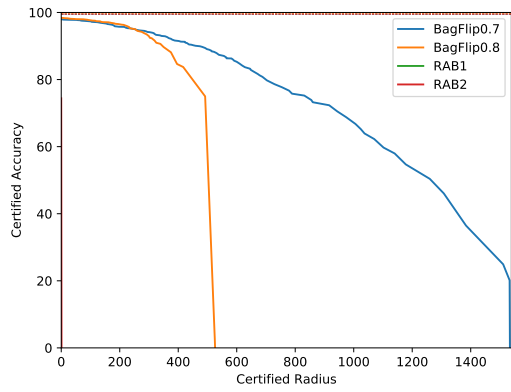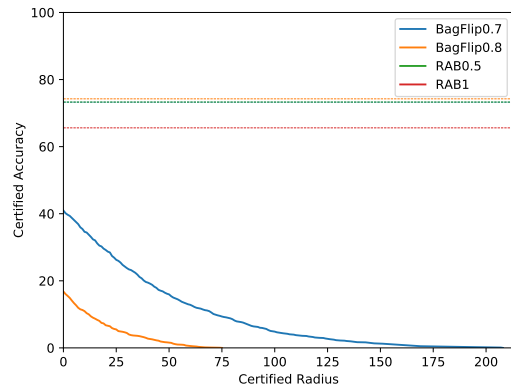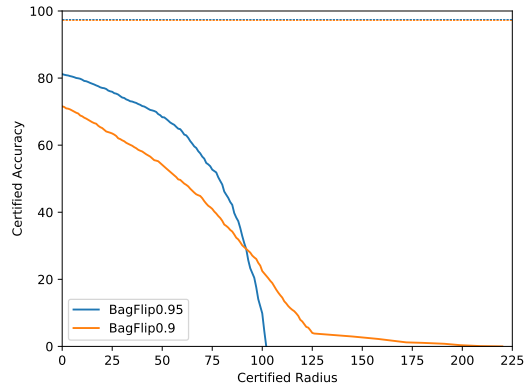Figure 6: Compared to Bagging on $FL_1$ on MNIST-17 and $L$ using MNIST.

Figure 7: Results of BagFlip on L.



(a) Compared to RAB on MNIST-01



(b) Compared to RAB on CIFAR10-02



(c) Compared to FeatFlip on MNIST-17

Figure 8: Compared to FeatFlip and RAB.