
Choosing Public Datasets for Private Machine Learning via Gradient Subspace Distance

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Differentially private stochastic gradient descent privatizes model training by
2 injecting noise into each iteration, where the noise magnitude increases with the
3 number of model parameters. Recent works suggest that we can reduce the noise by
4 leveraging public data for private machine learning, by projecting gradients onto a
5 subspace prescribed by the public data. However, given a choice of public datasets,
6 it is not clear which one may be most appropriate for the private task. We give an
7 algorithm for selecting a public dataset by measuring a low-dimensional subspace
8 distance between gradients of the public and private examples. The computational
9 and privacy cost overhead of our method is minimal. Empirical evaluation suggests
10 that trained model accuracy is monotone in this distance.

11 1 Introduction

12 Machine learning models have shown that they can memorize the information of their training data
13 [7]. Recent works have shown that attackers can recover many training samples from published
14 models through carefully designed attacks [3, 18]. This will cause critical privacy issues when the
15 models are trained on private data.

16 *Differential Privacy* (DP) [5] is a rigorous privacy criterion that provides theoretical guarantees
17 to the amount of information attackers can infer about any single training point. *Differentially*
18 *private stochastic gradient descent* (DPSGD) [1, 19, 2] is one of the most popular methods to
19 achieve differential privacy in deep learning. It makes two modifications to vanilla SGD: 1) clipping
20 per-sample gradients to ensure a bound on their ℓ_2 norms; 2) adding Gaussian noise to the gradient.

21 One downside of applying DP to machine learning is that we need to sacrifice the utility of machine
22 learning models to maintain privacy. Specifically, DPSGD adds random noise drawn from a spherical
23 Gaussian distribution, $\mathcal{N}(\mathbf{0}, \sigma^2 \mathbb{I}_p)$, where p is the model dimension, i.e., the number of model
24 parameters, and the variance σ^2 scales the noise. The magnitude of the noise introduced in each
25 step scales with the square root of the number of parameters p . For classic deep learning models for
26 Computer Vision tasks like ResNet, the added noise will be tens of times greater than the original
27 gradients, inevitably leading to worse utility.

28 Many works have proposed various methods to improve the utility of private machine learning [25,
29 8, 16, 24, 12]. One promising approach involves employing *public* data. Generally, there are two
30 ways of using public data in private training. One is transfer learning, where we pretrain the model
31 on a public dataset and then finetune the model on our target tasks (private data) [16, 1, 23, 15].
32 Another approach arises from the empirical observation that during the training process, the stochastic
33 gradients always stay in a lower-dimensional subspace of the high-dimensional gradient space p .
34 Based on this observation, some work suggests another approach to leverage public data: they use

35 the public data to find this lower-dimensional subspace and then project the noisy gradient onto this
 36 subspace [25, 8, 24, 12]. This generally improves utility over DPSGD without supplementary data.

37 However, this leaves an open question: *which public dataset should one select for a particular private*
 38 *task?* The ideal case may be if some part of the private dataset is public, as this avoids any distribution
 39 shift. But otherwise, we would like a way to quantify a public dataset’s fitness for use. Our main
 40 contribution is an algorithm for this purpose.

41 Our method needs a single batch of private and public examples from the dataset. Specifically,
 42 the algorithm performs three steps to derive the closeness between public data and private data: 1)
 43 compute the per-sample gradient of both public and private data, 2) find the gradient subspace of both
 44 private and public data by applying singular value decomposition (SVD), 3) compute the subspace
 45 distance d using Projection Metric [11]. Our algorithm gives a value that measures a type of distance
 46 between public and private data. Our empirical evaluation shows that the distance d derived from our
 47 algorithm follows the utility of the projection method monotonously, meaning that the distance d is a
 48 good indicator of public data’s utility.

49 2 Preliminaries

50 **Notation.** In this paper, we use p to denote the model dimension, i.e., the number of parameters in
 51 the model. k is the dimension of the lower-dimensional space we choose. m refers to the number of
 52 examples in a batch. We use superscript or subscript interchangeably to denote private or public data,
 53 like x_{priv} , V^{pub} .

Definition 1 (Differential Privacy [5]). *A randomized algorithm \mathcal{A} is (ϵ, δ) -differential private if
 for any pair of datasets D, D' that differ in exactly one data point and for all subsets E of outputs, we
 have:*

$$\Pr[\mathcal{A}(D) \in E] \leq \Pr[\mathcal{A}(D') \in E] + \delta.$$

Definition 2 (Projection Metric [20, 6]). *The projection metric between two k -dimensional sub-
 spaces V_1, V_2 is defined as:*

$$d(V_1, V_2) = \left(\sum_{i=1}^k \sin^2 \theta_i \right)^{1/2} = \left(k - \sum_{i=1}^k \cos^2 \theta_i \right)^{1/2}$$

54 where θ_i ’s come from the principal angles between V_1 and V_2 .

55 Appendix A gives the formal definition of principle angles.

56 In this paper, we evaluate our method using GEP [24], the state-of-the-art private deep learning
 57 algorithm that leverages public data. We briefly describe their algorithm in Appendix A.

58 3 Methods

59 Now we define the problem formally. Suppose we have a task that consists of a private dataset X^{priv}
 60 and a private deep learning algorithm A that can leverage public data to improve model utility. We
 61 have a list of potential choice of public dataset $[X_1^{pub}, X_2^{pub}, \dots]$. We would like a metric that can
 62 prescribe which public dataset, when used with algorithm A on the private task X^{priv} , will have the
 63 best model utility.

64 At a high level, our method involves the following two steps: finding the gradient subspace of the
 65 data examples and computing the gradient subspace distance. The algorithm uses the same model A
 66 and a batch of unlabeled data examples from private and public datasets. Following standard DPSGD,
 67 the algorithm will first compute and store per-example gradients from each data example, that is
 68 $G_{priv}, G_{pub} \in \mathbb{R}^{m \times p}$. Then it computes the top- k singular vectors of both the private and public
 69 gradient matrix by performing singular value decomposition (SVD). Finally we use projection metric
 70 to derive the subspace distance d by taking the right singular vectors V_k^{pub}, V_k^{priv} from the previous
 71 step. The pseudo-code of our method is given in Algorithm 1.

72 Our algorithm is based on the empirical observation that the stochastic gradients stay in a lower-
 73 dimensional subspace during the training procedure of a deep learning model [10, 14]. We also

Algorithm 1 Gradient Subspace Distance

Input: Private examples x_{priv} , public examples x_{pub} , loss function \mathcal{L} , model weights w_0

Output: Distance between two image datasets d

- 1: $G_{priv} = \nabla \mathcal{L}(w_0, x_{priv})$ \triangleright Compute per-sample gradient matrix for private examples
 - 2: Compute top- k subspace of private gradient matrix V_k^{priv} :
 $U^{priv}, S^{priv}, V^{priv} \leftarrow \text{SVD}(G_{priv})$
 - 3: $G_{pub} = \nabla \mathcal{L}(w_0, x_{pub})$ \triangleright Compute per-sample gradient matrix for public examples
 - 4: Compute top- k subspace of private public matrix V_k^{pub} :
 $U^{pub}, S^{pub}, V^{pub} \leftarrow \text{SVD}(G_{pub})$
 - 5: $d = \text{ProjectionMetric}(V_k^{priv}, V_k^{pub})$
-

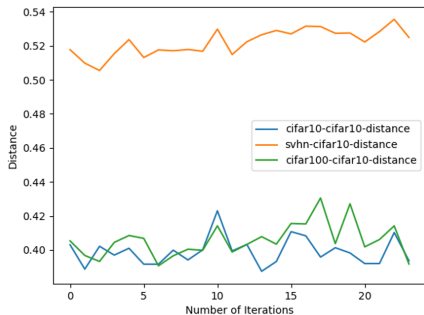
74 empirically evaluate this finding over different datasets and model settings. Details are given in
75 Appendix C.1. Such observation suggests that most of the information the gradient carries is
76 contented in much lower-dimensional space. Our method finds such subspace for private and public
77 data examples and then measures the distance between two subspaces.

78 We follow the conclusion in [11] and use projection metric [20, 6] to measure the subspace distance
79 between V_k^{pub} and V_k^{priv} . Intuitively, it considers all the principal angles by averaging them to show
80 intermediate characteristics between the two subspaces. It is suggested to be robust to the distribution
81 of data examples and enjoys great distance structure properties such as triangle inequality.

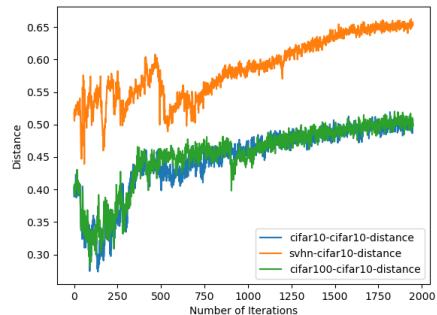
82 While one may have concern that such distance computation and comparison may have privacy
83 leakage, our method only needs a batch of private examples and use this batch once for distance
84 computation. There will be little privacy leakage during this process. Even if there are some extremely
85 private cases when we have to publish our choices of public data, we can spend some privacy budget
86 and apply some differential privacy mechanism such as exponential mechanism. The scoring function
87 would be the projection metric of \sqrt{k} -sensitivity.

88 4 Experiments

89 We evaluate our algorithm on three datasets widely used in Computer Vision: Fashion MNIST [22],
90 SVHN [17] and CIFAR-10 [13]. We also choose one medical image dataset: ChestX-ray14 [21], as
91 medical images are considered highly private-sensitive. A variety of datasets are chosen as public
92 data respectively. We use the state-of-the-art private deep learning algorithm that leverage public data,
93 GEP[24], for private training. The details of experiment settings are in Appendix B.



(a) Distance in 1 Epoch



(b) Distance over 100 Epoch

Figure 1: The trend of distance during the process of training a Resnet 20 model on CIFAR-10 using vanilla SGD. We follow a standard SGD training procedure and compute the distance between the current private batch and public examples at each iteration.

94 **Result.** We first empirically evaluate our distance measurement along the training process, as shown
 95 in Figure 1. Our empirical study shows that the relative distance between private and public datasets
 96 is uniform at most times over the training process. Based on this observation, our algorithm will only
 97 require a batch of private examples and one computation step that involves private gradient, meaning
 98 the privacy cost and computation overhead is minimal.

Table 1: GEP evaluation accuracy and corresponding distance in descending order. "-" means vanilla DP-SGD training.

Accuracy	Private Dataset	Public Dataset	Distance
58.63%	CIFAR-10	CIFAR-100	0.20
57.64%		CIFAR-10	0.24
56.75%		SVHN	0.28
52.16%		-	-
91.32%	SVHN	SVHN	0.25
89.29%		CIFAR-100	0.31
89.08%		MNIST-M	0.39
83.21%	FMNIST	-	-
85.25%		FMNIST	0.34
84.54%		FLOWER	0.43
83.91%		MNIST	0.50
79.77%	-	-	-

99 We compute the distance and evaluate using GEP for the chosen datasets. The evaluation results
 100 are given in Table 1. The empirical evaluation shows that the distance derived by our algorithm
 101 follows monotonously with the final trained model accuracy. Smaller distance implies that the private
 102 examples share more similarities with public examples, thus leading to better accuracy when we use
 103 those public data.

Table 2: GEP evaluation AUC and corresponding distance in descending order. "-" means vanilla DP-SGD training.

AUC	Private Dataset	Public Dataset	Distance
69.02%	ChestX-ray14	ChestX-ray14	0.15
66.62%		KagChest	0.36
64.90%		-	-
48.80%		CIFAR-100	0.55

104 For ChestX-ray14, we use the AUC metric because ChestX-ray14 is highly imbalanced where the
 105 "no finding" class takes up a large portion of the dataset. The evaluation results are given in Table
 106 2. For more complex tasks, a bad choice of public data, such as CIFAR-100 for ChestX-ray14, will
 107 result in worse utility than the DPSGD baseline. When practitioners want to leverage public data
 108 for private machine learning, it would be much more essential to use our algorithm to evaluate the
 109 quality of the public data before performing private training using algorithms like GEP.

110 5 Conclusion

111 While recent studies are focusing on leveraging public data for private machine learning, the quality
 112 of public data also matters and is still an open question. In this work, we propose a new algorithm
 113 that can help private deep learning practitioners to select public data at minimal time and privacy cost.
 114 The empirical evaluation suggests that our distance measurement is a good indicator of public data
 115 quality for private machine learning algorithms that leverage public examples.

116 References

117 [1] Martin Abadi, Andy Chu, Ian Goodfellow, H. Brendan McMahan, Ilya Mironov, Kunal Talwar,
 118 and Li Zhang. Deep learning with differential privacy. In *Proceedings of the 2016 ACM SIGSAC*

- 119 *Conference on Computer and Communications Security, CCS '16*, page 308–318, New York,
120 NY, USA, 2016. Association for Computing Machinery.
- 121 [2] Raef Bassily, Adam D. Smith, and Abhradeep Thakurta. Private empirical risk minimization:
122 Efficient algorithms and tight error bounds. In *55th IEEE Annual Symposium on Foundations of*
123 *Computer Science, FOCS 2014, Philadelphia, PA, USA, October 18-21, 2014*, pages 464–473.
124 IEEE Computer Society, 2014.
- 125 [3] Nicholas Carlini, Florian Tramèr, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Kather-
126 ine Lee, Adam Roberts, Tom Brown, Dawn Song, Úlfar Erlingsson, Alina Oprea, and Colin
127 Raffel. Extracting training data from large language models. In *30th USENIX Security Sympo-*
128 *sium (USENIX Security 21)*, pages 2633–2650. USENIX Association, August 2021.
- 129 [4] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. ImageNet: A Large-Scale
130 Hierarchical Image Database. In *CVPR09*, 2009.
- 131 [5] Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith. Calibrating noise to sensi-
132 tivity in private data analysis. In Shai Halevi and Tal Rabin, editors, *Theory of Cryptography*,
133 pages 265–284, Berlin, Heidelberg, 2006. Springer Berlin Heidelberg.
- 134 [6] Alan Edelman, Tomás A. Arias, and Steven Thomas Smith. The geometry of algorithms with
135 orthogonality constraints. *SIAM J. Matrix Anal. Appl.*, 20(2):303–353, 1998.
- 136 [7] Matt Fredrikson, Somesh Jha, and Thomas Ristenpart. Model inversion attacks that exploit
137 confidence information and basic countermeasures. In *Proceedings of the 22nd ACM SIGSAC*
138 *Conference on Computer and Communications Security, CCS '15*, page 1322–1333, New York,
139 NY, USA, 2015. Association for Computing Machinery.
- 140 [8] Aditya Golatkar, Alessandro Achille, Yu-Xiang Wang, Aaron Roth, Michael Kearns, and
141 Stefano Soatto. Mixed differential privacy in computer vision. *CoRR*, abs/2203.11481, 2022.
- 142 [9] Gene H. Golub and Charles F. Van Loan. *Matrix Computations, Third Edition*. Johns Hopkins
143 University Press, 1996.
- 144 [10] Guy Gur-Ari, Daniel A. Roberts, and Ethan Dyer. Gradient descent happens in a tiny subspace.
145 *CoRR*, abs/1812.04754, 2018.
- 146 [11] Jihun Ham and Daniel D. Lee. Grassmann discriminant analysis: a unifying view on subspace-
147 based learning. In William W. Cohen, Andrew McCallum, and Sam T. Roweis, editors, *Machine*
148 *Learning, Proceedings of the Twenty-Fifth International Conference (ICML 2008), Helsinki,*
149 *Finland, June 5-9, 2008*, volume 307 of *ACM International Conference Proceeding Series*,
150 pages 376–383. ACM, 2008.
- 151 [12] Peter Kairouz, Mónica Ribero, Keith Rush, and Abhradeep Thakurta. (nearly) dimension
152 independent private ERM with adagrad rates via publicly estimated subspaces. In *Proceedings*
153 *of the 34th Annual Conference on Learning Theory, COLT '21*, pages 2717–2746, 2021.
- 154 [13] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images.
155 2009.
- 156 [14] Xinyan Li, Qilong Gu, Yingxue Zhou, Tiancong Chen, and Arindam Banerjee. Hessian based
157 analysis of SGD for deep nets: Dynamics and generalization. In Carlotta Demeniconi and
158 Nitesh V. Chawla, editors, *Proceedings of the 2020 SIAM International Conference on Data*
159 *Mining, SDM 2020, Cincinnati, Ohio, USA, May 7-9, 2020*, pages 190–198. SIAM, 2020.
- 160 [15] Xuechen Li, Florian Tramèr, Percy Liang, and Tatsunori Hashimoto. Large language models
161 can be strong differentially private learners. In *Proceedings of the 10th International Conference*
162 *on Learning Representations, ICLR '22*, 2022.
- 163 [16] Zelun Luo, Daniel J. Wu, Ehsan Adeli, and Li Fei-Fei. Scalable differential privacy with sparse
164 network finetuning. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR*
165 *2021, virtual, June 19-25, 2021*, pages 5059–5068. Computer Vision Foundation / IEEE, 2021.

- 166 [17] Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Bo Wu, and Andrew Y Ng.
167 Reading digits in natural images with unsupervised feature learning. 2011.
- 168 [18] Reza Shokri, Marco Stronati, Congzheng Song, and Vitaly Shmatikov. Membership inference
169 attacks against machine learning models. In *2017 IEEE Symposium on Security and Privacy*
170 (*SP*), pages 3–18, 2017.
- 171 [19] Shuang Song, Kamalika Chaudhuri, and Anand D. Sarwate. Stochastic gradient descent with
172 differentially private updates. In *IEEE Global Conference on Signal and Information Processing,*
173 *GlobalSIP 2013, Austin, TX, USA, December 3-5, 2013*, pages 245–248. IEEE, 2013.
- 174 [20] Liwei Wang, Xiao Wang, and Jufu Feng. Subspace distance analysis with application to adaptive
175 bayesian algorithm for face recognition. *Pattern Recognit.*, 39(3):456–464, 2006.
- 176 [21] Xiaosong Wang, Yifan Peng, Le Lu, Zhiyong Lu, Mohammadhadi Bagheri, and Ronald Sum-
177 mers. Chestx-ray8: Hospital-scale chest x-ray database and benchmarks on weakly-supervised
178 classification and localization of common thorax diseases. In *2017 IEEE Conference on*
179 *Computer Vision and Pattern Recognition(CVPR)*, pages 3462–3471, 2017.
- 180 [22] Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for
181 benchmarking machine learning algorithms. *CoRR*, abs/1708.07747, 2017.
- 182 [23] Da Yu, Saurabh Naik, Arturs Backurs, Sivakanth Gopi, Huseyin A Inan, Gautam Kamath,
183 Janardhan Kulkarni, Yin Tat Lee, Andre Manoel, Lukas Wutschitz, Sergey Yekhanin, and
184 Huishuai Zhang. Differentially private fine-tuning of language models. In *Proceedings of the*
185 *10th International Conference on Learning Representations, ICLR ’22, 2022*.
- 186 [24] Da Yu, Huishuai Zhang, Wei Chen, and Tie-Yan Liu. Do not let privacy overbill utility:
187 Gradient embedding perturbation for private learning. In *9th International Conference on*
188 *Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net,
189 2021.
- 190 [25] Yingxue Zhou, Steven Wu, and Arindam Banerjee. Bypassing the ambient dimension: Pri-
191 vate {sgd} with gradient subspace identification. In *International Conference on Learning*
192 *Representations*, 2021.

193 **A Missing Preliminaries**

Definition 3 (Principal Angles [9]). Let V_1 and V_2 be two orthonormal matrices of $\mathbb{R}^{p \times k}$. The principal angles $0 \leq \theta_1 \leq \dots \leq \theta_k \leq \pi/2$ between two subspaces $\text{span}(V_1)$ and $\text{span}(V_2)$, are defined recursively by

$$\cos\theta_k = \max_{\mathbf{u}_k \in \text{span}(V_1)} \max_{\mathbf{v}_k \in \text{span}(V_2)} \mathbf{u}'_k \mathbf{v}_k, \text{ subject to}$$

$$\mathbf{u}'_k \mathbf{u}_k = 1, \mathbf{v}'_k \mathbf{v}_k = 1, \mathbf{u}'_k \mathbf{u}_i = 0, \mathbf{v}'_k \mathbf{v}_i = 0, (i = 1, \dots, k - 1)$$

194 The first principal angle θ_1 is the smallest angle between all pairs of unit vectors over two subspaces.
 195 The rest are similarly defined.

196 **Gradient Embedding Perturbation (GEP).** In this paper, we evaluate our method using GEP
 197 [24], the state-of-the-art private deep learning algorithm that leverages public data for private training.
 198 Here we briefly introduce their algorithm. GEP involves three steps: 1) it computes a set of the
 199 orthonormal basis for the lower-dimensional subspace; 2) GEP projects the private gradients to the
 200 subspace derived from step 1, thus dividing the private gradients into two parts: embedding gradients
 201 that contain most of the information carried by the gradient, and the remainder are called residual
 202 gradients; 3) GEP clips two parts of the gradients separately and perturbs them to achieve differential
 203 privacy.

204 **B Experiments Setting**

205 **Model Architecture.** For Fashion MNIST, we use a simple convolutional neural network with
 206 around 26000 parameters as in Table 3. For SVHN and CIFAR-10, we use ResNet20 which contains
 207 roughly 260,000 parameters. Batch normalization layers are replaced by group normalization layers
 208 for different private training, aligning with GEP settings. For ChestX-ray14, we use ResNet152
 209 which has been pretrained on ImageNet1k, a subset of the full ImageNet [4] dataset. We privately
 210 fine-tune its classification layer, which contains around 28,000 parameters. We use the same model
 211 architecture for subspace distance computation and GEP private training.

Table 3: Model architecture for Fashion MNIST.

Layer	Parameters
Conv2d	16 filters of 8x8, stride=2
Maxpooling2d	stride=2
Conv2d	32 filters 4x4, stride=2
Linear	32 units
Softmax	10 units

Table 4: Choices of public dataset for private dataset. The four datasets in the first row are private datasets. The datasets listed in the first columns are choices of public datasets. 'X' means we choose the two corresponding datasets as a pair of private/public dataset.

	CIFAR-10	SVHN	Fashion MNIST	ChestX-ray14
CIFAR-10	X			
CIFAR-100	X	X		X
SVHN	X	X		
MNIST_M		X		
Fashion MNIST			X	
Flower			X	
MNIST			X	
ChestX-ray14				X
KagChest				X

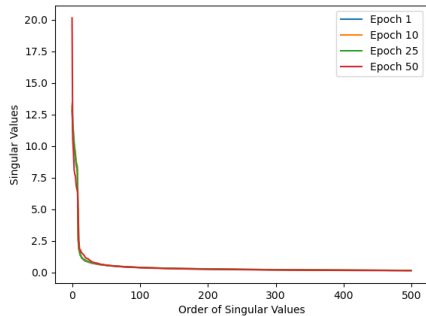
212 **Dataset Choice.** ChestX-ray14 consists of frontal view X-ray images with 14 different classes of
 213 lung disease. In our evaluation, there are 78,466 training examples and 20433 testing examples in
 214 ChestX-ray14. Our choices of public datasets for the four private datasets are described in Table 4.
 215 We sample 2000 examples from both private and public datasets for distance comparison using our
 216 algorithm. The same 2000 public examples are given to GEP for evaluation.

217 **Hyperparameter Setting.** For distance computation, we choose $k = 16$, that is, we only consider
 218 a 16-dimensional subspace. We follow the hyperparameter setting in the GEP paper for evaluation. In
 219 the GEP paper, they didn't evaluate GEP on the ChestX-ray14 dataset. In our evaluation, we choose
 220 $k = 100$ and clip norms are 3 and 1 for original and residual gradients, respectively. The learning
 221 rate for the SGD optimizer is set to 0.05. All other hyperparameters are set as default. We use $\epsilon = 2$
 222 and $\delta = 1e - 5$ for all the evaluations.

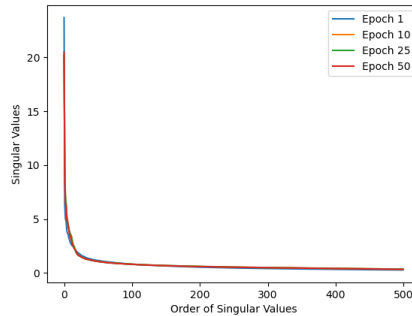
223 C More Experiments

224 C.1 Gradients are in a lower-dimensional subspace.

225 We evaluate the empirical observation that the stochastic gradients stay in a lower-dimensional
 226 subspace during the training procedure of a deep learning model [10, 14], as shown in Figure 2.
 227 Results show that only a tiny fraction of singular values are enormous. At the same time, the rest are
 228 close to 0, meaning that most of the gradients lie in a lower-dimensional subspace, corresponding to
 229 the top singular vectors.



(a) CIFAR-10



(b) ChestX-ray14

Figure 2: Top 500 singular values in the training procedure using vanilla SGD. Model architectures are in the Appendix B. Only a small fraction of singular values are extremely large while the rest are close to 0, meaning that most of the gradients lie in a lower-dimensional subspace, which corresponds to the top singular vectors.