

---

# Catastrophic Failures of Neural Active Learning on Heteroskedastic Distributions

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Models which can actively seek out the best quality training data hold the promise  
2 of more accurate, adaptable, and efficient machine learning. State-of-the-art tech-  
3 niques tend to prefer examples which are the most difficult to classify. While  
4 this works well on homogeneous datasets, we find that it can lead to catastrophic  
5 failures when performing active learning on multiple distributions which have  
6 different degrees of label noise (heteroskedasticity). Most active learning algo-  
7 rithms strongly prefer to draw from the distribution with more noise, even if its  
8 examples have no informative structure (such as solid color images). We find that  
9 active learning which encourages diversity and model uncertainty in the selected  
10 examples can significantly mitigate these failures. We hope these observations  
11 are immediately useful to practitioners and can lead to the construction of more  
12 realistic and challenging active learning benchmarks.

## 13 1 Introduction

14 In an active learning setup, a model has access to a pool of labeled data and a pool of unlabeled data.  
15 After training on the available labeled data, some selection rule is applied to identify a batch of  $k$   
16 unlabeled samples to be labeled and integrated into the training set before repeating the process. Under  
17 this paradigm, data are considered to be abundant but label acquisition is costly. The goal of an active  
18 learning algorithm is to identify unlabeled samples that, once labeled at used to fit model parameters,  
19 will elicit the most performant hypothesis possible given a fixed labeling budget.

20 In an effort to fulfill this objective, high-quality selection criteria generally favor a diverse set of  
21 examples where the model has a high degree of uncertainty. That is, we want to favor the selection of  
22 items that might lead to the most significant change from the model's current state, but we need to  
23 ensure we do not waste our labeling budget by selecting items that are similar to each other.

24 Preferring examples with high uncertainty often works well on homogeneous datasets, but can lead  
25 to catastrophic failure when training on a mixture of distinct distributions with different degrees of  
26 noise, as the active learning algorithm prefers the noisier distribution over the cleaner distribution.  
27 We refer to these as *Heteroskedastic Distributions*.

28 In this paper, we demonstrate that modern active learning algorithms designed for use with deep  
29 neural networks are significantly damaged when subjected to training data that's been corrupted by  
30 heteroskedastic noise. These approaches typically rely on notions of model improvement that are  
31 unable to disambiguate aleatoric from epistemic uncertainty, overselecting samples for which the  
32 model is unconfident but which are unlikely to best improve the current hypothesis. We show that this  
33 inefficiency can be partially reduced in two ways: by more heavily favor diversity and encouraging  
34 examples with high divergence between a conventionally trained model and an exponential moving  
35 average (EMA) of its iterates.

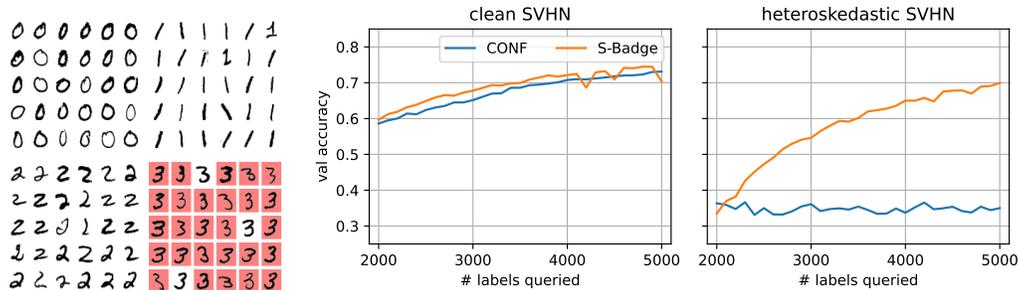


Figure 1: Overview of our setting and results. *Left*: We consider neural active learning settings with high heteroskedasticity; as an extreme example, the labels from one class are replaced with uniform noise (incorrect labels in red). Active learning algorithms have a strong tendency to prefer this random-label class over other classes. *Right*: Our proposed S-Badge algorithm, which encourages both model uncertainty and diversity in sampled points, performs similarly to least-confidence sampling on the original SVHN data, but greatly outperforms it on the heteroskedastic distributions.

36 Empirically, with this EMA divergence, which we measure on the penultimate layer representation,  
 37 seems to capture information about the model’s sensitivity to that point. Importantly, we observe  
 38 that this score is nearly zero on the noisy samples considered in this paper. When used to scale the  
 39 representations used in for the Badge algorithm, we find we can significantly reduce the frequency  
 40 with which noisy samples are selected, which we refer to as the S-Badge algorithm.

#### 41 1.1 Related work

42 **Neural active learning.** Active learning is an extremely well-researched area, with the richest  
 43 theory developed for the convex setting [1, 2]. More recently, however, there have been several  
 44 attempts to tractably generalize active learning to the deep regime. Such approaches can be thought  
 45 of as identifying batches of samples that cater more to either the model’s predictive uncertainty or to  
 46 the diversity of the selection.

47 In the former approach, a batch of points are selected in order of the model’s uncertainty about their  
 48 label. Many of these methods query samples that are nearest the decision boundary, an approach  
 49 that’s theoretically well understood in the linear regime when the batch size is 1 [3, 4, 5]. Some  
 50 deep learning-specific approaches have also been developed, including using the variance of dropout  
 51 samples to quantify uncertainty [6], and adversarial examples have been used to approximate the  
 52 distance between an unlabeled sample and the decision boundary. In the deep setting however, where  
 53 models are typically retrained from scratch after every round of selection, a larger batch size is usually  
 54 necessary for efficiency purposes.

55 For large batch sizes, algorithms that cater to diversity are usually more effective. In deep learning,  
 56 several methods take the representation obtained at the penultimate layer of the network, and aim to  
 57 identify a batch of samples that might summarize this space well [7, 8]. [9]. Other methods promote  
 58 diversity by minimizing an upper bound on some notion of model’s loss on unseen data [10, 11, 12,  
 59 13, 14]. This approach has also been taken to trade-off between diversity and uncertainty in deep  
 60 active learning [15, 16].

61 **Data poisoning and distributional robustness.** A related body of work seeks to obtain models  
 62 and training procedures which are robust against *worst-case* perturbations to the data distribution.  
 63 For recent treatments of this topic and further references, see [17, 18]. A few recent works have  
 64 considered data poisoning in the active learning setting [19, 20], with defenses focusing on modifying  
 65 the setting rather than the algorithm. Overall, though the settings are compatible, the aim of this work  
 66 is to directly address the empirical performance of deep active learning with low-quality labels, rather  
 67 than a more pessimistic min-max robustness formulation.

68 **Heteroskedasticity in deep learning.** The issues of class imbalance and heteroskedasticity are of  
 69 interest in the supervised learning setting [21, 22, 23], in which various methods have been proposed

70 to make training more robust to these distributions. Our work seeks to initiate the study of the  
71 orthogonal (but analogous) issue in sample *selection*.

## 72 **2 Heteroskedastic Benchmarks for Neural Active Learning**

73 We introduce three new benchmarks for active learning on heteroskedastic distributions. In all cases,  
74 we introduce an additional set of  $N$  examples with purely random labels, which the model is trained  
75 on with it shuffled into the original data. In all cases, whether an example is one of these special noisy  
76 samples is not given to the model, but it is always reasonably easily predictable from the example’s  
77 features. This distinguishes our benchmarks from IID label noise, which is not predictable based on  
78 the example’s features.

- 79 • *Noisy-Blank*: We introduce  $N$  examples which are solid black ( $x = 0$ ) with a random  
80 example  $y \sim U(1, 10)$ .
- 81 • *Noisy-Diverse*: We increase the difficulty by introducing  $K = 100$  different types of  
82 examples, where each type is a random solid color and has a label randomly drawn from  
83 three different choices (unique to that type).  $N$  of these examples are sampled. This  
84 benchmark is designed to make the heteroskedastic distribution more diverse while still  
85 keeping the noisy examples simple.
- 86 • *Noisy-Class*: In our hardest setting, we take the examples in the dataset with a particular  
87 class  $y = 1$  and assign these examples uniformly random labels  $y \sim U(1, 10)$ . We then  
88 randomly repeat these examples to give  $N$  examples. In this case, the randomly labeled  
89 examples are challenging but still possible to identify.

## 90 **3 Methods and Experiments**

91 Here we experiment with several active learning algorithms, noising strategies, and model architec-  
92 tures. In all experiments, we use the same experimental settings from [16], starting with 2000 points,  
93 and query samples in batches of 100 points, until we got to a total of 5000 labeled points. We avoid  
94 warm-starting and retrain from a fresh initialization after each round of selection [24]. We ran with  
95 both a small 1-layer MLP with 512 hidden units as well as a ResNet18. We also added the additional  
96 noisy-labeled examples in all cases to make 80% of the examples from the noisy distribution and  
97 20% from the original distribution.

### 98 **3.1 Baselines**

99 We consider four baseline algorithms commonly used in the literature. Confidence sampling [25] and  
100 Margin sampling [26] are uncertainty-based strategies: Confidence sampling selects the  $k$  unlabeled  
101 points for which the most likely label has the smallest amount of probability mass, and Margin  
102 sampling selects the  $k$  points for which the difference in probability mass in the two most likely  
103 labels is smallest. The Coreset algorithm is a diversity-based approach that aims to select a batch of  
104 representative points, as measured in penultimate layer space of the current state of the model [7].  
105 BADGE is a hybridized approach, meant to strike a balance between uncertainty and diversity.  
106 BADGE represents data in a hallucinated gradient space before performing diverse selection using  
107 k-means++ [16].

### 108 **3.2 S-Badge: Increasing Sampling where Representations Change Across Training 109 Iterations**

110 We conjecture that while examples with noisy labels will have high loss and high predictive uncer-  
111 tainty, the model’s predictions will converge quickly and undergo little change later in training. By  
112 encouraging the selection process in active learning to not select these examples, we hope to improve  
113 performance in the heteroskedastic setting.

114 In addition to the main (or online) model ( $F_\theta$ ), we introduce an EMA model ( $F_\beta$ ) into the example-  
115 selection pipeline. The EMA model has the same architecture as the online model but uses a different  
116 set of parameters  $\beta$ , which are exponentially moving averages of  $\theta$ . That is,

$$\beta = \alpha \cdot \beta + (1 - \alpha) \cdot \theta \quad (1)$$

Table 1: Classification accuracy on SVHN with 5000 actively queried examples, with different heteroskedastic distribution corruptions. In parenthesis, we report the percentage of the examples (over the course of training) which the active learning selected from the non-corrupted original examples (100% being best, and 0% being worst).

ResNet	Clean	Noisy-Blank	Noisy-Diverse	Noisy-Class
Random	78.6%	64.0% (20.7%)	50.8% (19.6%)	25.3% (20.5%)
Confidence	68.5%	57.9% (0.00%)	55.9% (42.4%)	23.4% (14.2%)
Margin	73.6%	72.9% (82.6%)	60.8% (44.5%)	30.6% (15.9%)
Badge	76.3%	75.2% (99.0%)	57.0% (29.2%)	32.9% (30.1%)
S-Badge	76.0%	74.7% (99.0%)	67.1% (50.0%)	30.1% (36.9%)
MLP				
Random	70.6%	49.5% (20.3%)	46.5% (20.2%)	41.5% (19.9%)
Confidence	73.1%	35.0% (0.67%)	36.2% (12.9%)	40.9% (14.5%)
Margin	75.0%	65.2% (83.8%)	60.5% (60.0%)	46.7% (19.4%)
Badge	74.7%	69.1% (99.0%)	48.7% (19.8%)	43.8% (24.1%)
S-Badge	70.3%	70.0% (99.0%)	52.7% (22.5%)	47.7% (28.2%)

117 where  $\alpha$  is set to a high value of 0.999.

118 We conjecture that the state difference between the EMA model and the online model can be a  
 119 helpful signal for querying unlabelled examples. To this end, we re-weight the gradients in the Badge  
 120 algorithm  $\frac{d\mathcal{L}}{dW}$  by the average value of this hidden state difference before running K-means++ seeding.  
 121 Where the hidden state difference is small, this reweighted gradient will be close to zero, and few of  
 122 these examples are likely to be selected.

### 123 3.3 Results

124 The results in Table 1 show consistent deterioration in test accuracy when selecting from heteroskedas-  
 125 tic distributions. We found that Badge and S-Badge, through their use of diversity in the selection  
 126 process, were nearly perfect in solving the *Noisy-Blank* task. However on *Noisy-Diverse*, we found  
 127 that S-Badge often significantly outperformed Badge. We also found that max-margin sampling was  
 128 a surprisingly effective baseline compared to least-confidence sampling.

### 129 3.4 Analysis

130 One of our more striking experimental results (Table 1) is that selecting examples with the lowest  
 131 prediction confidence can fail catastrophically on heteroskedastic distributions. In (Appendix A) we  
 132 provide a theory explaining why training only on high loss examples (which would have low confi-  
 133 dence given a well-calibrated model) can lead to poor performance on heteroskedastic distributions.

## 134 4 Conclusion

135 Neural Active Learning is an active area of research, with a plethora of new techniques competing  
 136 to achieve better results. Our work seeks to throw this research program a curve ball, by showing  
 137 that techniques which are competitive on homogeneous datasets with little label noise can fail  
 138 catastrophically when presented with diverse heteroskedastic distributions. Active learning techniques  
 139 which actively seek diversity fail less catastrophically, but still struggle when the noisy examples are  
 140 themselves somewhat diverse. Despite some techniques failing less than others, our results suggest  
 141 significant room for improvement and future research on this task.

## 142 References

- 143 [1] Burr Settles. Active learning literature survey. *University of Wisconsin, Madison*, 2010.
- 144 [2] Sanjoy Dasgupta. Two faces of active learning. *Theoretical computer science*, 2011.
- 145 [3] Simon Tong and Daphne Koller. Support vector machine active learning with applications to  
146 text classification. *Journal of Machine Learning Research*, 2001.
- 147 [4] Greg Schohn and David Cohn. Less is more: Active learning with support vector machines. In  
148 *International Conference on Machine Learning*, 2000.
- 149 [5] Gokhan Tur, Dilek Hakkani-Tür, and Robert E Schapire. Combining active and semi-supervised  
150 learning for spoken language understanding. *Speech Communication*, 2005.
- 151 [6] Yarin Gal, Riashat Islam, and Zoubin Ghahramani. Deep bayesian active learning with image  
152 data. In *International Conference on Machine Learning*, 2017.
- 153 [7] Ozan Sener and Silvio Savarese. Active learning for convolutional neural networks: A core-set  
154 approach. In *International Conference on Learning Representations*, 2018.
- 155 [8] Yonatan Geifman and Ran El-Yaniv. Deep active learning over the long tail. *arXiv:1711.00941*,  
156 2017.
- 157 [9] Daniel Gissin and Shai Shalev-Shwartz. Discriminative active learning. *arXiv:1907.06347*,  
158 2019.
- 159 [10] Yuhong Guo and Dale Schuurmans. Discriminative batch mode active learning. In *Neural  
160 Information Processing Systems*, 2008.
- 161 [11] Zheng Wang and Jieping Ye. Querying discriminative and representative samples for batch  
162 mode active learning. *Transactions on Knowledge Discovery from Data*, 2015.
- 163 [12] Yuxin Chen and Andreas Krause. Near-optimal batch mode active learning and adaptive  
164 submodular optimization. In *International Conference on Machine Learning*, 2013.
- 165 [13] Kai Wei, Rishabh Iyer, and Jeff Bilmes. Submodularity in data subset selection and active  
166 learning. In *International Conference on Machine Learning*, 2015.
- 167 [14] Andreas Kirsch, Joost van Amersfoort, and Yarin Gal. Batchbald: Efficient and diverse batch  
168 acquisition for deep bayesian active learning. In *Advances in Neural Information Processing  
169 Systems*, 2019.
- 170 [15] Jordan T Ash, Surbhi Goel, Akshay Krishnamurthy, and Sham Kakade. Gone fishing: Neural  
171 active learning with fisher embeddings. *arXiv preprint arXiv:2106.09675*, 2021.
- 172 [16] Jordan T Ash, Chicheng Zhang, Akshay Krishnamurthy, John Langford, and Alekh Agar-  
173 wal. Deep batch active learning by diverse, uncertain gradient lower bounds. *International  
174 Conference on Learning Representations*, 2020.
- 175 [17] Jacob Steinhardt, Pang Wei Koh, and Percy Liang. Certified defenses for data poisoning attacks.  
176 In *Proceedings of the 31st International Conference on Neural Information Processing Systems*,  
177 pages 3520–3532, 2017.
- 178 [18] Shiori Sagawa, Pang Wei Koh, Tatsunori B Hashimoto, and Percy Liang. Distributionally robust  
179 neural networks. In *International Conference on Learning Representations*, 2019.
- 180 [19] Jing Lin, Ryan Luley, and Kaiqi Xiong. Active learning under malicious mislabeling and  
181 poisoning attacks. *arXiv preprint arXiv:2101.00157*, 2021.
- 182 [20] Jose Rodrigo Sanchez Vicarte, Gang Wang, and Christopher W. Fletcher. Double-cross attacks:  
183 Subverting active learning systems. In *30th USENIX Security Symposium (USENIX Security  
184 21)*, pages 1593–1610. USENIX Association, August 2021.
- 185 [21] Kaidi Cao, Colin Wei, Adrien Gaidon, Nikos Arechiga, and Tengyu Ma. Learning imbalanced  
186 datasets with label-distribution-aware margin loss. *arXiv preprint arXiv:1906.07413*, 2019.

- 187 [22] Jun Shu, Qi Xie, Lixuan Yi, Qian Zhao, Sanping Zhou, Zongben Xu, and Deyu Meng.  
188 Meta-weight-net: Learning an explicit mapping for sample weighting. *arXiv preprint*  
189 *arXiv:1902.07379*, 2019.
- 190 [23] Kaidi Cao, Yining Chen, Junwei Lu, Nikos Arechiga, Adrien Gaidon, and Tengyu Ma. Het-  
191 eroskedastic and imbalanced deep learning with adaptive regularization. *arXiv preprint*  
192 *arXiv:2006.15766*, 2020.
- 193 [24] Jordan T Ash and Ryan P Adams. On warm-starting neural network training. *Advances in*  
194 *Neural Information Processing Systems*, 2020.
- 195 [25] Dan Wang and Yi Shang. A new active labeling method for deep learning. In *International*  
196 *Joint Conference on Neural Networks*, 2014.
- 197 [26] Dan Roth and Kevin Small. Margin-based active learning for structured output spaces. In  
198 *European Conference on Machine Learning*, 2006.
- 199 [27] Kenji Kawaguchi and Haihao Lu. Ordered sgd: A new stochastic optimization framework for  
200 empirical risk minimization. In *International Conference on Artificial Intelligence and Statistics*,  
201 pages 669–679, 2020.
- 202 [28] Noah Golowich, Alexander Rakhlin, and Ohad Shamir. Size-independent sample complexity of  
203 neural networks. In *Conference On Learning Theory*, pages 297–299. PMLR, 2018.
- 204 [29] Peter L Bartlett and Shahar Mendelson. Rademacher and gaussian complexities: Risk bounds  
205 and structural results. *Journal of Machine Learning Research*, 3(Nov):463–482, 2002.
- 206 [30] Mehryar Mohri, Afshin Rostamizadeh, and Ameet Talwalkar. *Foundations of machine learning*.  
207 MIT press, 2012.
- 208 [31] Peter L Bartlett, Dylan J Foster, and Matus J Telgarsky. Spectrally-normalized margin bounds  
209 for neural networks. In *Advances in Neural Information Processing Systems*, pages 6240–6249,  
210 2017.
- 211 [32] Kenji Kawaguchi, Leslie Pack Kaelbling, and Yoshua Bengio. Generalization in deep learning.  
212 *In Mathematics of Deep Learning*, Cambridge University Press, to appear. Preprint available as:  
213 *MIT-CSAIL-TR-2018-014*, Massachusetts Institute of Technology, 2018.

## 214 A Generalization bound for biased query batches

### 215 A.1 Notation

216 Let  $\mathcal{D} = ((x_i, y_i))_{i=1}^n$  be a training dataset of  $n$  samples where  $x_i \in \mathcal{X} \subseteq \mathbb{R}^{d_x}$  is the input vector  
 217 and  $y_i \in \mathcal{Y} \subseteq \mathbb{R}^{d_y}$  is the target output vector for the  $i$ -th sample. A standard objective function is

$$L(\theta; \mathcal{D}) := \frac{1}{n} \sum_{i=1}^n L_i(\theta; \mathcal{D}), \quad (2)$$

218 where  $\theta \in \mathbb{R}^{d_\theta}$  is the parameter vector of the prediction model  $f(\cdot; \theta) : \mathbb{R}^{d_x} \rightarrow \mathbb{R}^{d_y}$ , and  $L_i(\theta; \mathcal{D}) :=$   
 219  $\ell(f(x_i; \theta), y_i)$  with the function  $\ell : \mathbb{R}^{d_y} \times \mathcal{Y} \rightarrow \mathbb{R}_{\geq 0}$  is the loss of the  $i$ -th sample.

220 Similarly to the notation of order statistics, we first introduce the notation of ordered indexes: given  
 221 a model parameter  $\theta$ , let  $L_{(1)}(\theta; \mathcal{D}) \geq L_{(2)}(\theta; \mathcal{D}) \geq \dots \geq L_{(n)}(\theta; \mathcal{D})$  be the decreasing values of  
 222 the individual losses  $L_1(\theta; \mathcal{D}), \dots, L_n(\theta; \mathcal{D})$ , where  $(j) \in \{1, \dots, n\}$  (for all  $j \in \{1, \dots, n\}$ ). That  
 223 is,  $\{(1), \dots, (n)\}$  as a perturbation of  $\{1, \dots, n\}$  defines the order of sample indexes by loss values.  
 224 Whenever we encounter ties on the values, we employ an arbitrary fixed tie-breaking rule in order to  
 225 ensure the uniqueness of such an order.

Denote  $r_i(\theta; \mathcal{D}) = \sum_{j=1}^n \mathbb{1}\{i = (j)\} \gamma_j$  where  $(j)$  depends on  $(\theta, \mathcal{D})$ . Given an arbitrary set  $\Theta \subseteq$   
 $\mathbb{R}^{d_\theta}$ , we define  $\mathfrak{R}_n(\Theta)$  as the (standard) Rademacher complexity of the set  $\{(x, y) \mapsto \ell(f(x; \theta), y) :$   
 $\theta \in \Theta\}$ :

$$\mathfrak{R}_n(\Theta) = \mathbb{E}_{\mathcal{D}, \xi} \left[ \sup_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n \xi_i \ell(f(\bar{x}_i; \theta), \bar{y}_i) \right],$$

226 where  $\bar{\mathcal{D}} = ((\bar{x}_i, \bar{y}_i))_{i=1}^n$ , and  $\xi_1, \dots, \xi_n$  are independent uniform random variables taking values  
 227 in  $\{-1, 1\}$  (i.e., Rademacher variables). Given a tuple  $(\ell, f, \Theta, \mathcal{X}, \mathcal{Y})$ , define  $M$  as the least upper  
 228 bound on the difference of individual loss values:  $|\ell(f(x; \theta), y) - \ell(f(x'; \theta), y')| \leq M$  for all  $\theta \in \Theta$   
 229 and all  $(x, y), (x', y') \in \mathcal{X} \times \mathcal{Y}$ . For example,  $M = 1$  if  $\ell$  is the 0-1 loss function.

$$\hat{\mathfrak{R}}_n(\Theta) = \mathbb{E}_{\xi} \left[ \sup_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n \xi_i \ell(f(x_i; \theta), y_i) \right],$$

### 230 A.2 Preliminaries

The previous paper [27] proves the following theoretical result. The stochastic optimization method  
 that uses a gradient estimator that is purposely biased toward those samples with the current top- $q$   
 losses (i.e., ordered SGD) implicitly minimizes a new objective function of

$$L_q(\theta; \mathcal{D}) = \frac{1}{q} \sum_{j=1}^n \gamma_j L_{(j)}(\theta; \mathcal{D}),$$

231 for any  $\mathcal{D}$  (including  $g(\mathcal{D})$ ), in the sense that such a gradient estimator is an unbiased estimator of a  
 232 (sub-) gradient of  $L_q(\theta; \mathcal{D})$ , instead of  $L(\theta; \mathcal{D})$ . Accordingly, the top- $q$ -biased stochastic optimization  
 233 method converges in terms of  $L_q$  instead of  $L$ .

Building up on this result, we consider generalization properties of the top- $q$ -biased stochastic  
 optimization with the presence of additional label noises in training data. We want to minimize the  
 expected loss,

$$\mathbb{E}_{(x, y) \sim \mathcal{P}} [\ell(f(x; \theta), y)]$$

by minimizing the training loss

$$L_q(\theta; g(\mathcal{D})),$$

where  $g(\mathcal{D}) = ((g_i^x(x_i), g_i^y(y_i)))_{i=1}^n$  is potentially corrupted by arbitrary noise and corruption effects  
 within arbitrary fixed functions  $g_i^x$  and  $g_i^y$  for  $i = 1 \dots, n$ , where  $(x_i, y_i) \sim \mathcal{P}$ . Thus, we want to  
 analyze the generalization gap:

$$\mathbb{E}_{(x, y) \sim \mathcal{P}} [\ell(f(x; \theta), y)] - L_q(\theta; g(\mathcal{D}))$$

234 The previous paper [27] showed the benefit of the top- $q$ -biased stochastic optimization method in  
 235 terms of generalization when  $g_i^x$  and  $g_i^y$  are identity functions and thus when the distributions are the  
 236 same for both expected loss and training loss. In contrast, in our setting, the distributions are different  
 237 for expected loss and training loss with potential noise corruptions through  $g_i^x$  and  $g_i^y$ .

### 238 A.3 Analysis

239 Theorem 1 presents a generalization bound for the top- $q$ -biased stochastic optimization:

240 **Theorem 1.** *Let  $\Theta$  be a fixed subset of  $\mathbb{R}^{d_\theta}$ . Then, for any  $\delta > 0$ , with probability at least  $1 - \delta$  over  
 241 an iid draw of  $n$  examples  $\mathcal{D} = ((x_i, y_i))_{i=1}^n$ , the following holds for all  $\theta \in \Theta$ :*

$$\mathbb{E}_{(x,y)}[\ell(f(x; \theta), y)] \leq L_q(\theta; g(\mathcal{D})) + 2\hat{\mathfrak{R}}_n(\Theta) + M \left(2 + \frac{s}{q}\right) \sqrt{\frac{\ln(2/\delta)}{2n}} - \mathcal{Q}_{n,q}(\Theta, g), \quad (3)$$

242 where  $\mathcal{Q}_{n,q}(\Theta, g) := \mathbb{E}_{\bar{\mathcal{D}}}[\inf_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n \binom{r_i(\theta; g(\bar{\mathcal{D}}))n}{q} \ell(f(g_i^x(\bar{x}_i); \theta), g_i^y(\bar{y}_i)) - \ell(f(\bar{x}_i; \theta), \bar{y}_i)]$ .

The expected error  $\mathbb{E}_{(x,y)}[\ell(f(x; \theta), y)]$  in the left-hand side of Equation (3) is a standard objective for generalization, whereas the right-hand side contains the data corruption function  $g$ . Here, we typically have  $\hat{\mathfrak{R}}_n(\Theta) = O(1/\sqrt{n})$  in terms of  $n$ . For example, consider the standard feedforward deep neural networks of the form  $f(x) = (\omega_T \circ \sigma_{T-1} \circ \omega_{T-1} \circ \sigma_{T-2} \cdots \sigma_1 \circ \omega_1)(x)$  where  $T$  is the number of layers,  $\omega_l(a) = W_l a$  with  $\|W_l\|_F \leq M_l$ , and  $\sigma_l$  is an element-wise nonlinear activation function that is 1-Lipschitz and positive homogeneous (e.g., ReLU). Then, if  $\|x\| \leq B$  for all  $x \in \mathcal{X}$ , using Theorem 1 of [28], we have that

$$\hat{\mathfrak{R}}_n(\Theta) \leq \frac{B(\sqrt{2 \log(2)T} + 1)(\prod_{l=1}^T M_l)}{\sqrt{n}}.$$

243 In Theorem 1, we can see that a label noise corruption  $g$  can lead to the failure of the top- $q$ -biased  
 244 stochastic optimization via increasing the training loss  $L_q(\theta; g(\mathcal{D}))$  and decreasing the top- $q$ -biased  
 245 factor  $\mathcal{Q}_{n,q}(\Theta, g)$ . Here, if there is no corruption  $g$  (i.e., if  $g_i^x$  and  $g_i^y$  are identity functions), then  
 246 we have that  $\mathcal{Q}_{n,q}(\Theta, g) \geq 0$  because  $\mathcal{Q}_{n,q}(\Theta, g) = \mathbb{E}_{\bar{\mathcal{D}}}[\inf_{\theta \in \Theta} L_q(\theta; \bar{\mathcal{D}}) - L(\theta; \bar{\mathcal{D}})] \geq 0$  due to  
 247  $L_q(\theta; \bar{\mathcal{D}}) - L(\theta; \bar{\mathcal{D}}) \geq 0$  for any  $\theta$  and  $\bar{\mathcal{D}}$  when  $g_i^x$  and  $g_i^y$  are identity functions. Thus, the top- $q$ -biased  
 248 factor  $\mathcal{Q}_{n,q}(\Theta, g)$  can explain the improvement of the generalization of the top- $q$ -biased stochastic  
 249 optimization over the standard unbiased stochastic optimization. However, with the presence of  
 250 the corruption  $g$ ,  $\frac{r_i(\theta; g(\bar{\mathcal{D}}))n}{q} \ell(f(g_i^x(\bar{x}_i); \theta), g_i^y(\bar{y}_i))$  can be smaller than  $\ell(f(\bar{x}_i; \theta), \bar{y}_i)$  by fitting the  
 251 corrupted noise, resulting  $\mathcal{Q}_{n,q}(\Theta, g) < 0$ . This leads to a significant failure in the following sense:  
 252 the generalization gap  $(\mathbb{E}_{(x,y)}[\ell(f(x; \theta), y)] - L_q(\theta; g(\mathcal{D})))$  goes to zero as  $n$  approach infinity if  
 253  $\mathcal{Q}_{n,q}(\Theta, g) \geq 0$  with no data corruption, but the generalization gap no longer goes to zero as  $n$   
 254 approach infinity if  $\mathcal{Q}_{n,q}(\Theta, g) < 0$  with data corruption.

255 To see this, let us look at the asymptotic case when  $n \rightarrow \infty$ . Let  $\Theta$  be constrained such that  
 256  $\hat{\mathfrak{R}}_n(\Theta) \rightarrow 0$  as  $n \rightarrow \infty$ , which has been shown to be satisfied for various models and sets  $\Theta$ , including  
 257 the standard deep neural networks above [29, 30, 31, 32, 28]. The third term in the right-hand side  
 258 of Equation (3) disappear as  $n \rightarrow \infty$ . Thus, if there is no corruption (i.e., if  $g_i^x$  and  $g_i^y$  are identity  
 259 functions), it holds with high probability that  $\mathbb{E}_{(x,y)}[\ell(f(x; \theta), y)] \leq L_q(\theta; g(\mathcal{D})) - \mathcal{Q}_{n,q}(\Theta, g) \leq$   
 260  $L_q(\theta; g(\mathcal{D}))$ , where  $L_q(\theta; g(\mathcal{D}))$  is minimized by the top- $q$ -biased stochastic optimization. From this  
 261 viewpoint, the top- $q$ -biased stochastic optimization minimizes the expected error for generalization  
 262 when  $n \rightarrow \infty$ , if there is no corruption. However, if there is corruption,  $\mathbb{E}_{(x,y)}[\ell(f(x; \theta), y)] \leq$   
 263  $L_q(\theta; g(\mathcal{D})) - \mathcal{Q}_{n,q}(\Theta, g) \not\leq L_q(\theta; g(\mathcal{D}))$ , and hence  $\mathbb{E}_{(x,y)}[\ell(f(x; \theta), y)] - L_q(\theta; g(\mathcal{D})) \not\rightarrow 0$  even  
 264 in the asymptotic case.

### 265 A.4 Proof of Theorem 1

266 We first notice that the following proposition from [27] still holds with the corrupted data with the  
 267 same proof  $g(\mathcal{D})$ :

268 **Proposition 1.** *For any  $j \in \{1, \dots, n\}$ ,  $\gamma_j \leq \frac{s}{n}$ .*

269 We use this proposition in the following proof of Theorem 1 to bound the effect of replacing one  
 270 sample in a dataset.

*Proof of Theorem 1.* We find an upper bound on  $\sup_{\theta \in \Theta} \mathbb{E}_{(x,y)}[\ell(f(x; \theta), y)] - L_q(\theta; g(\mathcal{D}))$  based on McDiarmid's inequality. Define

$$\Phi(\mathcal{D}) = \sup_{\theta \in \Theta} \mathbb{E}_{(x,y)}[\ell(f(x; \theta), y)] - L_q(\theta; g(\mathcal{D})).$$

271 Our proof plan is to provide the upper bound on  $\Phi(\mathcal{D})$  by using McDiarmid's inequality. To apply  
 272 McDiarmid's inequality to  $\Phi(\mathcal{D})$ , we first show that  $\Phi(\mathcal{D})$  satisfies the remaining condition of  
 273 McDiarmid's inequality on the effect of changing one sample. Let  $\mathcal{D}$  and  $\mathcal{D}'$  be two datasets differing  
 274 by exactly one point of an arbitrary index  $i_0$ ; i.e.,  $\mathcal{D}_i = \mathcal{D}'_i$  for all  $i \neq i_0$  and  $\mathcal{D}_{i_0} \neq \mathcal{D}'_{i_0}$ . Since  $(j)$   
 275 depends on  $g(\mathcal{D})$ , we sometimes write  $(j; \mathcal{D}) = (j)$  to stress the dependence on  $\mathcal{D}$  under  $g$ . Then,  
 276 we provide an upper bound on  $\Phi(\mathcal{D}') - \Phi(\mathcal{D})$  as follows:

$$\begin{aligned} \Phi(\mathcal{D}') - \Phi(\mathcal{D}) &\leq \sup_{\theta \in \Theta} L_q(\theta; g(\mathcal{D})) - L_q(\theta; g(\mathcal{D}')) \\ &= \sup_{\theta \in \Theta} \frac{1}{q} \sum_{j=1}^n \gamma_j (L_{(j; \mathcal{D})}(\theta; g(\mathcal{D})) - L_{(j; \mathcal{D}')}(\theta; g(\mathcal{D}'))) \\ &\leq \sup_{\theta \in \Theta} \frac{1}{q} \sum_{j=1}^n |\gamma_j| |L_{(j; \mathcal{D})}(\theta; g(\mathcal{D})) - L_{(j; \mathcal{D}')}(\theta; g(\mathcal{D}'))| \\ &\leq \sup_{\theta \in \Theta} \frac{1}{q} \frac{s}{n} \sum_{j=1}^n |L_{(j; \mathcal{D})}(\theta; g(\mathcal{D})) - L_{(j; \mathcal{D}')}(\theta; g(\mathcal{D}'))| \end{aligned}$$

277 where the first line follows the property of the supremum,  $\sup(a) - \sup(b) \leq \sup(a - b)$ , the second  
 278 line follows the definition of  $L_q$  where  $(j; \mathcal{D}) \neq (j; \mathcal{D}')$ , and the last line follows Proposition 1  
 279 ( $|\gamma_j| \leq \frac{s}{n}$ ).

280 We now bound the last term  $\sum_{j=1}^n |L_{(j; \mathcal{D})}(\theta; g(\mathcal{D})) - L_{(j; \mathcal{D}')}(\theta; g(\mathcal{D}'))|$ . This requires a careful  
 281 examination because  $|L_{(j; \mathcal{D})}(\theta; g(\mathcal{D})) - L_{(j; \mathcal{D}')}(\theta; g(\mathcal{D}'))| \neq 0$  for more than one index  $j$  (although  
 282  $\mathcal{D}$  and  $\mathcal{D}'$  differ only by exactly one point). This is because it is possible to have  $(j; \mathcal{D}) \neq (j; \mathcal{D}')$   
 283 for many indexes  $j$  where  $(j; \mathcal{D})$  in  $L_{(j; \mathcal{D})}(\theta; g(\mathcal{D}))$  and  $(j; \mathcal{D}')$  in  $L_{(j; \mathcal{D}')}(\theta; g(\mathcal{D}'))$ . To analyze  
 284 this effect, we now conduct case analysis. Define  $l(i; \mathcal{D})$  such that  $(j) = i$  where  $j = l(i; \mathcal{D})$ ; i.e.,  
 285  $L_i(\theta; g(\mathcal{D})) = L_{(l(i; \mathcal{D}))}(\theta; g(\mathcal{D}))$ .

286 Consider the case where  $l(i_0; \mathcal{D}') \geq l(i_0; \mathcal{D})$ . Let  $j_1 = l(i_0; \mathcal{D})$  and  $j_2 = l(i_0; \mathcal{D}')$ . Then,

$$\begin{aligned} &\sum_{j=1}^n |L_{(j)}(\theta; g(\mathcal{D})) - L_{(j)}(\theta; g(\mathcal{D}'))| \\ &= \sum_{j=j_1}^{j_2-1} |L_{(j)}(\theta; g(\mathcal{D})) - L_{(j)}(\theta; g(\mathcal{D}'))| + |L_{(j_2)}(\theta; g(\mathcal{D})) - L_{(j_2)}(\theta; g(\mathcal{D}'))| \\ &= \sum_{j=j_1}^{j_2-1} |L_{(j)}(\theta; g(\mathcal{D})) - L_{(j+1)}(\theta; g(\mathcal{D}))| + |L_{(j_2)}(\theta; g(\mathcal{D})) - L_{(j_2)}(\theta; g(\mathcal{D}'))| \\ &= \sum_{j=j_1}^{j_2-1} (L_{(j)}(\theta; g(\mathcal{D})) - L_{(j+1)}(\theta; g(\mathcal{D}))) + L_{(j_2)}(\theta; g(\mathcal{D})) - L_{(j_2)}(\theta; g(\mathcal{D}')) \\ &= L_{(j_1)}(\theta; g(\mathcal{D})) - L_{(j_2)}(\theta; g(\mathcal{D}')) \\ &\leq M, \end{aligned}$$

287 where the first line uses the fact that  $j_2 = l(i_0; \mathcal{D}') \geq l(i_0; \mathcal{D}) = j_1$  where  $i_0$  is the index of samples  
 288 differing in  $\mathcal{D}$  and  $\mathcal{D}'$ . The second line follows the equality  $(j; \mathcal{D}') = (j+1; \mathcal{D})$  from  $j_1$  to  $j_2 - 1$  in  
 289 this case. The third line follows the definition of the ordering of the indexes. The fourth line follows  
 290 the cancellations of the terms from the third line.

291 Consider the case where  $l(i_0; \mathcal{D}') < l(i_0; \mathcal{D})$ . Let  $j_1 = l(i_0; \mathcal{D}')$  and  $j_2 = l(i_0; \mathcal{D})$ . Then,

$$\sum_{j=1}^n |L_{(j)}(\theta; g(\mathcal{D})) - L_{(j)}(\theta; g(\mathcal{D}'))|$$

$$\begin{aligned}
&= |L_{(j_1)}(\theta; g(\mathcal{D})) - L_{(j_1)}(\theta; g(\mathcal{D}'))| + \sum_{j=j_1+1}^{j_2} |L_{(j)}(\theta; g(\mathcal{D})) - L_{(j)}(\theta; g(\mathcal{D}'))| \\
&= |L_{(j_1)}(\theta; g(\mathcal{D})) - L_{(j_1)}(\theta; g(\mathcal{D}'))| + \sum_{j=j_1+1}^{j_2} |L_{(j)}(\theta; g(\mathcal{D})) - L_{(j-1)}(\theta; g(\mathcal{D}))| \\
&= L_{(j_1)}(\theta; g(\mathcal{D})) - L_{(j_1)}(\theta; g(\mathcal{D}')) + \sum_{j=j_1+1}^{j_2} (L_{(j)}(\theta; g(\mathcal{D})) - L_{(j-1)}(\theta; g(\mathcal{D}))) \\
&= L_{(j_1)}(\theta; g(\mathcal{D}')) - L_{(j_2)}(\theta; g(\mathcal{D})) \\
&\leq M.
\end{aligned}$$

292 where the first line uses the fact that  $j_1 = l(i_0; \mathcal{D}') < l(i_0; \mathcal{D}) = j_2$  where  $i_0$  is the index of samples  
293 differing in  $\mathcal{D}$  and  $\mathcal{D}'$ . The second line follows the equality  $(j; \mathcal{D}') = (j-1; \mathcal{D})$  from  $j_1+1$  to  $j_2$  in  
294 this case. The third line follows the definition of the ordering of the indexes. The fourth line follows  
295 the cancellations of the terms from the third line.

Therefore, in both cases of  $l(i_0; \mathcal{D}') \geq l(i_0; \mathcal{D})$  and  $l(i_0; \mathcal{D}') < l(i_0; \mathcal{D})$ , we have that

$$\Phi(\mathcal{D}') - \Phi(\mathcal{D}) \leq \frac{s}{q} \frac{M}{n}.$$

Similarly,  $\Phi(\mathcal{D}) - \Phi(\mathcal{D}') \leq \frac{s}{q} \frac{M}{n}$ , and hence  $|\Phi(\mathcal{D}) - \Phi(\mathcal{D}')| \leq \frac{s}{q} \frac{M}{n}$ . Thus, by McDiarmid's  
inequality, for any  $\delta > 0$ , with probability at least  $1 - \delta$ ,

$$\Phi(\mathcal{D}) \leq \mathbb{E}_{\bar{\mathcal{D}}}[\Phi(\bar{\mathcal{D}})] + \frac{Ms}{q} \sqrt{\frac{\ln(1/\delta)}{2n}}.$$

296 Moreover, since

$$\sum_{i=1}^n r_i(\theta; g(\mathcal{D})) L_i(\theta; g(\mathcal{D})) = \sum_{j=1}^n \gamma_j \sum_{i=1}^n \mathbb{1}\{i = (j; \mathcal{D})\} L_i(\theta; g(\mathcal{D})) = \sum_{j=1}^n \gamma_j L_{(j)}(\theta; g(\mathcal{D})),$$

we have that

$$L_q(\theta; g(\mathcal{D})) = \frac{1}{q} \sum_{i=1}^n r_i(\theta; g(\mathcal{D})) L_i(\theta; g(\mathcal{D})).$$

297 Therefore,

$$\begin{aligned}
&\mathbb{E}_{\bar{\mathcal{D}}}[\Phi(\bar{\mathcal{D}})] \\
&= \mathbb{E}_{\bar{\mathcal{D}}} \left[ \sup_{\theta \in \Theta} \mathbb{E}_{(\bar{x}', \bar{y}')} [\ell(f(\bar{x}'; \theta), \bar{y}')] - L(\theta; \bar{\mathcal{D}}) + L(\theta; \bar{\mathcal{D}}) - L_q(\theta; g(\bar{\mathcal{D}})) \right] \\
&\leq \mathbb{E}_{\bar{\mathcal{D}}} \left[ \sup_{\theta \in \Theta} \mathbb{E}_{(\bar{x}', \bar{y}')} [\ell(f(\bar{x}'; \theta), \bar{y}')] - L(\theta; \bar{\mathcal{D}}) \right] - \mathcal{Q}_{n,q}(\Theta, g) \\
&\leq \mathbb{E}_{\bar{\mathcal{D}}, \bar{\mathcal{D}'}} \left[ \sup_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n (\ell(f(\bar{x}'_i; \theta), \bar{y}'_i) - \ell(f(\bar{x}_i; \theta), \bar{y}_i)) \right] - \mathcal{Q}_{n,q}(\Theta, g) \\
&\leq \mathbb{E}_{\xi, \bar{\mathcal{D}}, \bar{\mathcal{D}'}} \left[ \sup_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n \xi_i (\ell(f(\bar{x}'_i; \theta), \bar{y}'_i) - \ell(f(\bar{x}_i; \theta), \bar{y}_i)) \right] - \mathcal{Q}_{n,q}(\Theta, g) \\
&\leq 2\mathfrak{R}_n(\Theta) - \mathcal{Q}_{n,q}(\Theta, g).
\end{aligned}$$

298 where the third line and the last line follow the subadditivity of supremum, the forth line follows  
299 the Jensen's inequality and the convexity of the supremum, the fifth line follows that for each  
300  $\xi_i \in \{-1, +1\}$ , the distribution of each term  $\xi_i (\ell(f(\bar{x}'_i; \theta), \bar{y}'_i) - \ell(f(\bar{x}_i; \theta), \bar{y}_i))$  is the distribution of  
301  $(\ell(f(\bar{x}'_i; \theta), \bar{y}'_i) - \ell(f(\bar{x}_i; \theta), \bar{y}_i))$  since  $\bar{\mathcal{D}}$  and  $\bar{\mathcal{D}'}$  are drawn iid with the same distribution. Therefore,  
302 for any  $\delta > 0$ , with probability at least  $1 - \delta$ ,

$$\Phi(\mathcal{D}) \leq 2\mathfrak{R}_n(\Theta) - \mathcal{Q}_{n,q}(\Theta, g) + \frac{Ms}{q} \sqrt{\frac{\ln(1/\delta)}{2n}}.$$

Finally, since changing one data point in  $\mathcal{D}$  changes  $\hat{\mathfrak{R}}_n(\Theta)$  by at most  $M/m$ , McDiarmid's inequality implies that for any  $\delta > 0$ , with probability at least  $1 - \delta$ ,

$$\mathfrak{R}_n(\Theta) \leq \hat{\mathfrak{R}}_n(\Theta) + M\sqrt{\frac{\ln(1/\delta)}{2n}}.$$

303 By taking union bound, we obtain the statement of this theorem.

304

□