

---

# A Lagrangian Duality Approach to Active Learning

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 We consider the batch active learning problem, where only a subset of the training  
2 data is labeled, and the goal is to query a batch of unlabeled samples to be labeled so  
3 as to maximally improve model performance. We formulate the learning problem  
4 using constrained optimization, where each constraint bounds the performance of  
5 the model on labeled samples. Considering a primal-dual approach, we optimize  
6 the primal variables, corresponding to the model parameters, as well as the dual  
7 variables, corresponding to the constraints. As each dual variable indicates how  
8 significantly the perturbation of the respective constraint affects the optimal value  
9 of the objective function, we use it as a proxy of the informativeness of the  
10 corresponding training sample. Our approach, which we refer to as Active Learning  
11 via Lagrangian duality, or ALLY, leverages this fact to select a diverse set of  
12 unlabeled samples with the highest estimated dual variables as our query set.  
13 We demonstrate the benefits of our approach in a variety of classification and  
14 regression tasks and also discuss its limitations depending on the capacity of the  
15 model used. We also show that ALLY can be used in a generative mode to create  
16 novel maximally-informative samples.

## 17 1 Introduction

18 Machine learning, and particularly deep learning, has seen tremendous progress in recent years  
19 in areas such as computer vision and natural language processing. One of the key drivers of  
20 such progress is the availability of massive, high-quality datasets, which enables training models  
21 comprising millions, or even billions, of parameters [1, 2, 3]. Nevertheless, in some areas, such as  
22 healthcare, obtaining *labeled* training data is challenging and/or expensive [4, 5, 6]. This has given  
23 rise to a class of approaches, collectively referred to as *active learning*, whose goal is to minimize the  
24 labeling effort for training machine learning models.

25 Active learning methods aim to improve data efficiency by querying the labels of samples presumed  
26 informative, in a feedback-driven fashion. In recent years, the pool-based active learning setting,  
27 in which queries are drawn from a large, static pool of unlabeled samples has drawn significant  
28 attention [7, 8, 9]. This is due to the abundance of such unlabeled pools and the compatibility of the  
29 pool-based setting with the training of deep neural networks.

30 Most active learners rely on defining a notion of *informativeness* of a given sample, such as model  
31 uncertainty [10, 11], expected model change [12, 13] or expected error reduction [14, 15]. However,  
32 in the batch setting, where multiple samples are queried simultaneously, not contemplating the  
33 information overlap between the samples can lead to sub-optimal queries. Consequently, batch  
34 *diversity* needs to be taken into account, often at the expense of individual sample informativeness  
35 [16].

36 In this paper, we demonstrate how a *constrained learning* formulation of the problem enables the use  
37 of *Lagrangian duality* for detecting informative samples. In particular, we bound the loss incurred by

38 each sample, and use the dual variables associated to these constraints as a measure of informativeness.  
39 We show that dual variables are directly related to the variations of the average optimal loss over the  
40 entire data distribution, which motivates our approach.

41 Through an iterative primal-dual strategy, we optimize the model parameters as well as the dual  
42 variables. We then leverage the learned embedding space [17, 18] to train a *dual regression head*  
43 that estimates the dual variable associated to each unlabeled sample. Our proposed active learning  
44 approach, which we refer to as Active Learning via Lagrangian dualitY, or ALLY can then be used to  
45 select a diverse and informative set of samples.

46 We evaluate the performance of ALLY on a suite of classification and regression tasks, and show that  
47 it performs better than or similarly to state-of-the-art batch active learning methods on a variety of  
48 tasks. We further demonstrate how the trained backbone, alongside the dual regression head, enable  
49 the generation of novel samples that can be optimized to be maximally informative, shedding light on  
50 the interpretability of the proposed active learning framework.

## 51 **2 Related Work**

### 52 **2.1 Active Learning**

53 The literature on active learning is voluminous and a myriad of strategies for the pool-based setting  
54 have been proposed [7]. In what follows, we describe some of the approaches most connected to our  
55 work.

56 A simple way of measuring model uncertainty is by computing the entropy of the predicted class  
57 distribution. Designed for the sequential case, Entropy Sampling [10] selects the unlabeled sample  
58 with highest associated output entropy. Among the relevant methods is BADGE [15], which employs  
59 a lower bound on the norm of the gradients in the final layer of the network as a measure of  
60 informativeness. BADGE balances diversity and informativeness by using the  $k$ -MEANS++ seeding  
61 algorithm to select a batch with large Gram determinant in the gradient space. BAIT [19] builds on  
62 traditional, Information Matrix based methods to efficiently select a batch that optimizes a bound on  
63 the MLE error in a two stage manner. Other methods that propose notions of informativeness are  
64 BALD [20], which uses the mutual information between predictions and model parameters as an  
65 uncertainty measure, and Learning Loss [21], which trains a loss prediction module and queries the  
66 samples that hypothetically generate high errors (and thus large model updates.)

67 Some diversity-promoting approaches are compatible with many informativeness measures. A  
68 popular approach is to cluster the samples of the unlabeled set and then select informative points  
69 from each cluster [22, 23, 24]. In [25], Monte Carlo sampling is used to simulate sequences of length  
70  $b$  of the sequential algorithm, and then a *best-matching* combination of the sequences is used to build  
71 a batch. A simpler approach is to select the  $b$  most informative points after a stochastic perturbation  
72 of the informativeness scores [26].

73 Some methods do not enforce informativeness explicitly, but rather query a set of data points that  
74 is maximally representative of the entire unlabeled set. Coreset [8], for instance, formulates pool-  
75 based active learning as a core-set selection problem, and aims to identify a set of points that  
76 geometrically covers the entire representation space. To do this, Coreset selects the batch that  
77 when added to the labeled set, minimizes the maximum distance between labeled and unlabeled  
78 examples. Coreset is compatible with deep neural networks and can be used in both regression and  
79 classification settings. Similarly, DAL [27] emphasizes representativeness by framing active learning  
80 as a binary classification task and selecting queries that maximize the similarity between the labeled  
81 and unlabeled set.

### 82 **2.2 Constrained Learning**

83 The need to tailor the behavior of machine learning systems has led to the development of a constrained  
84 learning theory. These developments [28, 29] have shown that, from a PAC (Probably Approximately  
85 Correct) perspective, learning under requirements is as hard as classical learning and that it can be  
86 done in practice through primal-dual learners. This has led to numerous applications across several  
87 areas of ML such as federated learning [30], fairness [28], stability of neural networks [31, 32] and  
88 adversarial robustness [33].

89 **3 Problem Formulation**

90 **3.1 Batch Active Learning**

91 Let  $\mathcal{D}$  denote a probability distribution over data pairs  $(\mathbf{x}, y)$ , where  $\mathbf{x} \in \mathcal{X} \subseteq \mathbb{R}^D$  represents a  
 92 feature vector (e.g., the pixels of an image) and  $y \in \mathcal{Y} \subseteq \mathbb{R}$  represents a label or measurement. In  
 93 classification tasks,  $\mathcal{Y}$  is a subset of  $\mathbb{N}$ , whereas in regression,  $\mathcal{Y} = \mathbb{R}$ .

94 Initially, a set  $\mathcal{L} = \{(\mathbf{x}_i, y_i)\}_{i \in \mathcal{N}_{\mathcal{L}}}$  of data pairs, or labeled samples, is available, coming from a  
 95 probability distribution  $\mathcal{D}_{\mathcal{L}}$ . This set is used to learn a predictor  $f : \mathcal{X} \rightarrow \mathcal{Y}$  from a hypothesis  
 96 class  $\mathcal{F}$ . Then, a batch  $\mathcal{B}$  of samples, or *queries*, is selected from a pool of unlabeled samples  
 97  $\mathcal{U} = \{\mathbf{x}_i\}_{i \in \mathcal{N}_{\mathcal{U}}}$  and sent to an oracle for labeling. The goal is to select the batch that minimizes  
 98 the future expected loss. More precisely, we formulate the Batch Active Learning (BAL) problem as

$$\mathcal{B}^* = \arg \min_{\mathcal{B} \subseteq \mathcal{U} : |\mathcal{B}| \leq b} \min_{f \in \mathcal{F}} \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}_{\mathcal{L} \cup \mathcal{B}}} [\ell(f(\mathbf{x}), y)], \quad (\text{BAL})$$

99 where  $b$ , referred to as the *budget*, represents the maximum cardinality of  $\mathcal{B}$  and  $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$  is a  
 100 loss function (e.g., cross-entropy loss or mean-squared error).

101 This process is typically repeated multiple times. At each iteration, two main steps are performed:  
 102 (i) selecting  $\mathcal{B}_t$  and updating the sets:  $\mathcal{L}_t = \mathcal{L}_{t-1} \cup \mathcal{B}_t$  and  $\mathcal{U}_t = \mathcal{U}_{t-1} \setminus \mathcal{B}_t$ , and (ii) obtaining the  
 103 predictor  $f$  with the aggregate set of labeled samples  $\mathcal{L}_t$ . Steps (i) and (ii) correspond to the outer and  
 104 inner minimization problems in (BAL), respectively. In what follows, we focus on a single iteration,  
 105 and thus obviate the dependence on the iteration  $t$  to ease the notation.

106 **3.2 Constrained Statistical Learning**

107 Most active learning methods in the literature [8, 15, 20, 5, 7] formulate step (ii) above as an  
 108 *unconstrained* Statistical Risk Minimization (SRM) problem [34],

$$\min_{f \in \mathcal{F}} \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} [\ell(f(\mathbf{x}), y)]. \quad (\text{SRM})$$

109 Our approach, alternatively, uses a *Constrained* Statistical Learning (CSL) formulation,

$$P^* = \min_{f \in \mathcal{F}} \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} [\ell(f(\mathbf{x}), y)] \quad (\text{CSL-a})$$

$$\text{s.t. } \ell'(f(\mathbf{x}), y) \leq \epsilon(\mathbf{x}), \quad \mathcal{D}_{\mathbf{x}}\text{-a.e.} \quad (\text{CSL-b})$$

110 where  $\ell' : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$  is a secondary loss function,  $\epsilon : \mathcal{X} \rightarrow \mathbb{R}$  is a mapping from each data point to  
 111 a corresponding constraint upper bound, and  $\mathcal{D}_{\mathbf{x}}$  denotes the marginal distribution over  $\mathcal{X}$ . Note that  
 112 the objective function in (CSL-a) is the same as in (SRM), but the secondary loss is required to be  
 113 bounded  $\mathcal{D}_{\mathbf{x}}$ -almost everywhere.

114 Letting  $\lambda : \mathcal{X} \rightarrow \mathbb{R}^+$  denote the dual variable function, the Lagrangian associated to (CSL) can be  
 115 written as

$$\begin{aligned} L(f, \lambda) &= \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} [\ell(f(\mathbf{x}), y)] + \int_{\mathcal{X}, \mathcal{Y}} \lambda(\mathbf{x})(\ell'(f(\mathbf{x}), y) - \epsilon(\mathbf{x}))p(\mathbf{x}, y)d\mathbf{x}dy \\ &= \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \left[ \ell(f(\mathbf{x}), y) + \lambda(\mathbf{x})(\ell'(f(\mathbf{x}), y) - \epsilon(\mathbf{x})) \right], \end{aligned}$$

116 where it is implicitly assumed that the conditional distribution  $p(y|\mathbf{x})$  is a Dirac delta distribution,  
 117 i.e.,  $y$  is a deterministic function of  $\mathbf{x}$ . This leads to the dual problem,

$$D^* = \max_{\lambda \in \Lambda} \min_{f \in \mathcal{F}} L(f, \lambda(\mathbf{x})), \quad (\text{D-CSL})$$

118 where  $\Lambda := \{\lambda \mid \lambda(\mathbf{x}) \geq 0, \mathcal{D}_{\mathbf{x}}\text{-a.e.}\}$ . There are three main motivations for this infinite programming  
 119 formulation:

- 120 1. **Access to variations of  $P^*$ :** As we show in Theorem 3.2, this formulation gives us access  
 121 to  $\frac{\partial P^*}{\partial \epsilon(\mathbf{x})}$ , enabling the use of *dual variables* as an indicator of the *informativeness* of the  
 122 training samples.

- 123 2. **Resilience:** The most informative samples often lie in the tails of the distribution  $\mathfrak{D}$ . Those  
 124 samples appear less frequently in the dataset and thus, models can achieve low errors without  
 125 learning to classify/regress them correctly.
- 126 3. **Adaptive regularization:** For a fixed dual variable  $\lambda$ , the Lagrangian is simply a regular-  
 127 ized objective. Thus, the max-min formulation in D-CSL can be viewed as a regularized  
 128 minimization, where the regularization weight is updated during the training procedure  
 129 according to the degree of constraint satisfaction or violation.

130 The dual problem can be interpreted as finding the tightest lower bound on  $P^*$ . In the general case,  
 131  $D^* \leq P^*$ , which is known as weak duality. Nevertheless, under certain conditions,  $D^*$  attains  $P^*$   
 132 (strong duality) and we can derive a relation between the solution of (D-CSL) and the sensitivity of  
 133  $P^*$  with respect to  $\epsilon(\mathbf{x})$ . See Appendix A for more details.

134 In the following, we define the *Fréchet subdifferential* of a convex function, which allows us to justify  
 135 the use of dual variables as a measure of *informativeness* of a sample.

**Definition 3.1.** Let  $U, V$  be Banach spaces. The Fréchet subdifferential of a convex function  
 $P : U \rightarrow V$  at  $u \in U$  is defined as:

$$\partial P(u) = \{z \in U^* : P(v) - P(u) \geq \langle z, v - u \rangle \text{ for all } v \in U\},$$

136 where  $U^*$  denotes the topological dual space of  $U$ , and  $\langle z, v - u \rangle = \mathbb{E}_{\mathfrak{D}} [z(\mathbf{x})(v(\mathbf{x}) - u(\mathbf{x}))]$ .

137 Having the above definition, we state following theorem, which characterizes the variations of  $P^*$   
 138 (the optimum value of CSL) as a function of the constraint tightness  $\epsilon(\mathbf{x})$ .

139 **Theorem 3.2.** *If the problem (CSL) is strongly dual, then for any  $\mathbf{x} \in \mathcal{X}$ , we have*

$$-\lambda^*(\mathbf{x}) \in \partial P^*(\epsilon(\mathbf{x})),$$

140 where  $\partial P^*(\epsilon(\mathbf{x}))$  denotes the Fréchet subdifferential of  $P^*$  with respect to  $\epsilon(\mathbf{x})$ , and  $\lambda^*(\mathbf{x})$  is the  
 141 optimal dual variable associated to the constraint on  $\mathbf{x}$ .

142 *Proof.* See Appendix B. □

143 For any  $\mathbf{x}_0 \in \mathcal{X}$ , let  $\delta_{\mathbf{x}_0}(\mathbf{x})$  be a bump function of radius  $\delta > 0$ , centered at  $\mathbf{x}_0$  (i.e., a continuous,  
 144 radially-decreasing function with support in  $\|\mathbf{x} - \mathbf{x}_0\| \leq \delta$ ). Theorem 3.2 implies that

$$P^*(\epsilon(\mathbf{x}) + t\delta_{\mathbf{x}_0}(\mathbf{x})) - P^*(\epsilon(\mathbf{x})) \geq \langle -\lambda^*(\mathbf{x}_0), t\delta_{\mathbf{x}_0}(\mathbf{x}) \rangle, \quad \forall t.$$

145 The problem (CSL) typically includes an infinite number of constraints. Theorem 3.2 implies that  
 146 the constraint whose perturbation has the most *potential impact* on the optimal value of (CSL) is the  
 147 constraint with the highest associated optimal dual variable. For instance, infinitesimally tightening  
 148 the constraint in a neighbourhood  $\mathbf{x}_0$  would restrict the feasible set, causing an increase of the optimal  
 149 value of (CSL) at a rate larger than  $\lambda^*(\mathbf{x}_0)$ . In that sense, the magnitude of the dual variables can  
 150 be used as a measure of informativeness. Similarly to non-support vectors in SVMs [35], samples  
 151 associated to inactive constraints (i.e.,  $\{\mathbf{x}_0 : \lambda^*(\mathbf{x}_0) = 0\}$ ), are considered uninformative.

## 152 4 Proposed Approach

153 In light of the results mentioned in Section 3.2 on the usefulness of dual variables in constrained  
 154 statistical learning as a measure of sample informativeness, we present our proposed method, ALLY,  
 155 in this section. We start by introducing a primal-dual procedure to empirically solve the constrained  
 156 learning problem, and we will then proceed to describe our active learning algorithm in detail.

### 157 4.1 Constrained Empirical Risk Minimization

158 The formulation in (D-CSL) poses two challenges: (i) the distribution  $\mathfrak{D}$  is usually unknown and  
 159 (ii) it is an infinite-dimensional problem since it optimizes over the functional spaces  $\mathcal{F}$  and  $\Lambda$ . We  
 160 handle the former by replacing expectations by their sample means over a set of labeled samples  
 161  $\mathcal{L} = \{(\mathbf{x}_i, y_i)\}_{i \in \mathcal{N}_{\mathcal{L}}}$ , as described in the classical Empirical Risk Minimization (ERM) theory

---

**Algorithm 1** Primal-dual constrained learning (PDCL)
 

---

```

1: Input: Labeled dataset  $\mathcal{L}$ , primal learning rate  $\eta_p$ , dual learning rate  $\eta_d$ , number of iterations  $T$ ,
   number of primal steps per iteration  $T_p$ , constraint vector  $\epsilon$ .
2: Initialize:  $\theta, \lambda \leftarrow \mathbf{0}$ .
3: for  $t = 1, \dots, T$  do
4:    $\theta \leftarrow \theta - \eta_p \nabla_{\theta} \hat{L}(\theta, \lambda)$  ( $\times T_p$ )           // Update primal variables ( $T_p$  SGD steps)
5:    $s_i \leftarrow \ell'(f_{\theta}(\mathbf{x}_i), y_i) - \epsilon_i, \forall i \in \mathcal{N}_{\mathcal{L}}$ .           // Evaluate constraint slacks
6:    $\lambda_i \leftarrow [\lambda_i + \eta_d s_i]_+, \forall i \in \mathcal{N}_{\mathcal{L}}$ .           // Update dual variables
7: end for
8: Return:  $\theta, \lambda$ .
  
```

---

162 [34, 36]. In order to resolve the latter, we introduce a *parameterization* of the hypothesis class  $\mathcal{F}$  as  
 163  $\mathcal{P} = \{f_{\theta} \mid \theta \in \Theta\}$ , while we create a separate dual variable  $\lambda_i$  for each training sample  $\mathbf{x}_i$ . These  
 164 modifications lead to the Constrained Empirical Risk Minimization (CERM) problem,

$$\hat{P}^* = \min_{\theta \in \Theta} \frac{1}{|\mathcal{N}_{\mathcal{L}}|} \sum_{i \in \mathcal{N}_{\mathcal{L}}} \ell(f_{\theta}(\mathbf{x}_i), y_i) \quad (\text{CERM-a})$$

$$\text{s.t. } \ell'(f_{\theta}(\mathbf{x}_i), y_i) \leq \epsilon_i, \forall i \in \mathcal{N}_{\mathcal{L}}. \quad (\text{CERM-b})$$

165 This, in turn, results in the corresponding empirical dual problem,

$$\hat{D}^* = \max_{\lambda \geq \mathbf{0}} \min_{\theta \in \Theta} \hat{L}(\theta, \lambda), \quad (\text{D-CERM})$$

166 where  $\lambda = \{\lambda_i\}_{i \in \mathcal{N}_{\mathcal{L}}}$ ,  $\lambda \geq \mathbf{0}$  represents element-wise non-negativity,  $\epsilon_i$  denotes the constraint upper  
 167 bound associated to the  $i^{\text{th}}$  point-wise constraint, and the *empirical* Lagrangian,  $\hat{L}(\theta, \lambda)$ , is defined as

$$\hat{L}(\theta, \lambda) = \frac{1}{|\mathcal{N}_{\mathcal{L}}|} \sum_{i \in \mathcal{N}_{\mathcal{L}}} \left[ \ell(f_{\theta}(\mathbf{x}_i), y_i) + \lambda_i [\ell'(f_{\theta}(\mathbf{x}_i), y_i) - \epsilon_i] \right].$$

168 The max-min problem (D-CERM) can be undertaken by alternating the minimization with respect to  
 169  $\theta$  and the maximization with respect to  $\lambda$  [37, 28, 38], which leads to the primal-dual constrained  
 170 learning procedure in Algorithm 1. Notice that  $\min_{\theta \in \Theta} \hat{L}(\theta, \lambda)$  is the minimum of a family of  
 171 affine functions on  $\lambda$ , and thus is concave. Consequently, the outer problem corresponds to the  
 172 maximization of a concave function and can be solved via gradient ascent. The inner minimization,  
 173 however, is generally non-convex, but there is empirical evidence that deep neural networks can attain  
 174 *good* local minima when trained with stochastic gradient descent [39, 40]. Some theoretical remarks  
 175 on Algorithm 1 can be found in Appendix C.

176 As shown in Algorithm 1, the dual variables accumulate the slacks (i.e., distances between the  
 177 per-sample secondary loss and constraint values) over the entire learning procedure. This allows  
 178 the dual variables to be used as a measure of informativeness, while at the same time affecting the  
 179 local optimum to which the algorithm converges. Quite interestingly, similar ideas on monitoring the  
 180 evolution of the loss for specific training samples in order to recognize impactful instances have been  
 181 used in several generalization analyses [41, 42].

## 182 4.2 ALLY: Active Learning via Lagrangian Duality

183 Our proposed active learning algorithm, ALLY, is presented in Algorithm 2 (for the case of  $b = 1$ ).  
 184 Given a set of labeled samples,  $\mathcal{L}$ , we first obtain the model parameters  $\theta^*$  and the dual variables  
 185 associated to samples in  $\mathcal{L}$  using the primal-dual constrained learning approach in Algorithm 1.  
 186 Taking a representation learning approach [17, 18, 43, 44], we then partition the model  $f_{\theta^*}$  to a  
 187 *backbone*  $f_{\phi^*} : \mathcal{X} \rightarrow \mathbb{R}^d$ , where  $d$  denotes the dimensionality of the *embedding space*, and a  
 188 *prediction head*  $f_{\psi^*} : \mathbb{R}^d \rightarrow \mathcal{Y}$ , such that  $f_{\theta^*} = f_{\phi^*} \circ f_{\psi^*}$  and  $\theta^* = \phi^* \cup \psi^*$ . In order to estimate  
 189 the informativeness of the samples in the unlabeled dataset  $\mathcal{U}$ , we train a *dual regression head*  
 190  $f_{\omega} : \mathbb{R}^d \rightarrow \mathbb{R}^+$  on the embeddings generated by  $f_{\phi^*}$  by minimizing the mean-squared error

$$L_{\lambda}(\omega) = \frac{1}{|\mathcal{N}_{\mathcal{L}}|} \sum_{i \in \mathcal{N}_{\mathcal{L}}} \|f_{\omega}(f_{\phi^*}(\mathbf{x}_i)) - \lambda_i^*\|^2, \quad (1)$$

---

**Algorithm 2** Active learning via Lagrangian duality (ALLY)

---

- 1: **Input:** Labeled set  $\mathcal{L}$ , unlabeled set  $\mathcal{U}$ , primal learning rate  $\eta_p$ , dual learning rate  $\eta_d$ , number of PDCL iterations  $T$ , number of primal steps per iteration  $T_p$ , constraint vector  $\epsilon$ .
  - 2:  $\theta^*, \lambda^* \leftarrow \text{PDCL}(\mathcal{L}, \eta_p, \eta_d, T, T_p, \epsilon)$ . // Run the Primal-Dual Algorithm
  - 3:  $\omega^* \leftarrow \arg \min_{\omega} \frac{1}{|\mathcal{N}_{\mathcal{L}}|} \sum_{i \in \mathcal{N}_{\mathcal{L}}} \|f_{\omega}(f_{\phi^*}(\mathbf{x}_i)) - \lambda_i^*\|^2$ . // Train the dual regression head
  - 4:  $j^* \leftarrow \arg \max_{j \in \mathcal{N}_{\mathcal{U}}} f_{\omega^*}(f_{\phi^*}(\mathbf{x}_j))$  // Find sample with highest dual variable
  - 5: **Return:**  $j^*$ .
- 

191 while the parameters  $\phi^*$ , hence the embeddings, are kept frozen. It should be noted that the idea of  
 192 mapping embeddings to dual variables is present in other machine learning settings [45]. Once the  
 193 dual regression head is trained, we evaluate it on the embeddings corresponding to the unlabeled  
 194 samples, and identify the sample with the highest predicted dual variable.

195 As explained in Section 2.1, in the batch setting, selecting the  $b$  samples with the highest associated  
 196 dual variables is not optimal, due to the potential information overlap of such samples [7, 16].  
 197 Our method is compatible with any batch diversity approach that takes informativeness scores and  
 198 unlabeled data points (or embeddings) as inputs.

### 199 4.3 Connection to BADGE [15]

200 The BADGE method [15] uses the gradient of the loss function with respect to the parameters of the  
 201 last layer -denoted by  $\theta_L$ - as a measure of informativeness, i.e.,

$$\frac{\partial \ell(f_{\theta}(\mathbf{x}), \hat{y}(\mathbf{x}))}{\partial \theta_L}, \tag{2}$$

202 where  $\hat{y}(\mathbf{x})$  is the *hypothetical* label of  $\mathbf{x}$ , defined as  $\hat{y}(\mathbf{x}) := \arg \max_{y \in \mathcal{Y}} [f_{\theta}(\mathbf{x})]_y$ . In contrast, as  
 203 discussed in Theorem 3.2, to evaluate the informativeness of a given sample, ALLY observes

$$\frac{\partial P^*}{\partial \epsilon(\mathbf{x})} = \frac{\partial P^*}{\partial \theta} \frac{\partial \theta}{\partial \epsilon(\mathbf{x})} = \frac{\partial \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} [\ell(f_{\theta^*}(\mathbf{x}), y)]}{\partial \theta} \frac{\partial \theta}{\partial \epsilon(\mathbf{x})}. \tag{3}$$

204 There are three main differences between these informativeness measures:

- 205 • ALLY uses the derivative of the average optimal loss *over the entire distribution*, whereas  
 206 BADGE only considers the point-wise derivative. (Nearly all strategies evaluate their scoring  
 207 function, or informativeness measure, on a single sample.)
- 208 • ALLY observes the gradient with respect to *all* model parameters, not only the ones in the  
 209 last layer.
- 210 • Aside from the derivative of the loss with respect to the model parameters, ALLY also  
 211 considers an additional term  $\frac{\partial \theta}{\partial \epsilon(\mathbf{x})}$ , which models how the results of the optimization (i.e.,  
 212 model parameters) change when the constraint function is perturbed.

## 213 5 Experimental Evaluation

### 214 5.1 Settings

215 We consider four image classification tasks and one biomedical, non-image regression task. In  
 216 the classification setting, we use standard datasets that commonly appear in the active learning  
 217 literature, namely STL-10 [46], CIFAR-10 [47], SVHN [48] and MNIST [49]. Lacking an established  
 218 benchmark regression dataset for active learning, we evaluate ALLY on the Parkinsons Telemonitoring  
 219 dataset (PTD) [50]. In this regression task, the goal is to predict UPDRS (Unified Parkinson’s Disease  
 220 Rating Scale) scores from dysphonia measurements such as variation in fundamental frequency. Since  
 221 measurements in this dataset are the result of a *costly* clinical trial that requires expert knowledge,  
 222 this task is a prime example in which active learning might be essential.

223 In all experiments, the initial labeled set  $\mathcal{L}_0$  consists of 100 randomly drawn samples, and the budget  
 224 is set to either  $b = 200$  or  $b = 1000$ . For STL-10, CIFAR-10, and SVHN, we use a ResNet-18

225 architecture [51] with an embedding size of 128. In the case of MNIST and PTD, which are simpler  
 226 tasks, we use a multi-layer perceptron (MLP) with two hidden layers, each with 256 neurons and  
 227 rectified linear unit (ReLU) activation, leading to an embedding size of 256. The dual regression head  
 228  $f_\omega$  is a MLP with 3 hidden layers of dimensions 64, 32 and 16, with ReLU activations and employs  
 229 batch normalization.

230 As done in [22, 23], to ensure diversity in the batch, we cluster the embeddings of the unlabeled  
 231 samples, i.e.,  $\{f_{\phi^*}(\mathbf{x}_j)\}_{j \in \mathcal{N}_U}$ , using the  $k$ -MEANS clustering algorithm [52], where  $k \leq b$  is a  
 232 hyperparameter. We then select the samples with the highest associated dual variables from each  
 233 cluster, while maintaining equity among the number of samples per cluster. As shown in Appendix F,  
 234 the performance of ALLY improves with an increased number of clusters  $k$ , since it leads to a more  
 235 diverse batch of selected samples. Therefore, in all our experiments, we set  $k = b$ .

236 Regarding the role of the secondary loss, we opted for a generic formulation as the performed  
 237 sensitivity analysis in Theorem 3.2 holds for various choices of  $\ell'$ . However, in all our experiments,  
 238 we set  $\ell'(\cdot, \cdot) = \ell(\cdot, \cdot)$ . We believe that using unsupervised or self-supervised losses for  $\ell'$  is a  
 239 promising research direction, and we leave it for future work.

240 We compare our algorithm with Entropy sampling, BADGE, Coreset and BAIT. While Entropy  
 241 sampling focuses purely on uncertainty, BADGE and BAIT balance both diversity and informativeness  
 242 with different approaches. Coreset, also considered a state-of-the-art method, differs from the previous  
 243 methods in that it focuses on batch representativeness by framing active learning as a coreset selection  
 244 problem. It has been observed that, in some scenarios, several active learning methods fail to  
 245 consistently outperform Random Sampling [53, 54, 27]. We thus include it as one of the five  
 246 baselines. We adopt the PyTorch [55] implementation of the baselines from [56]. More details on the  
 247 experimental setting can be found in Appendix E.

248 **5.2 Interpreting the Informativeness Score of ALLY**

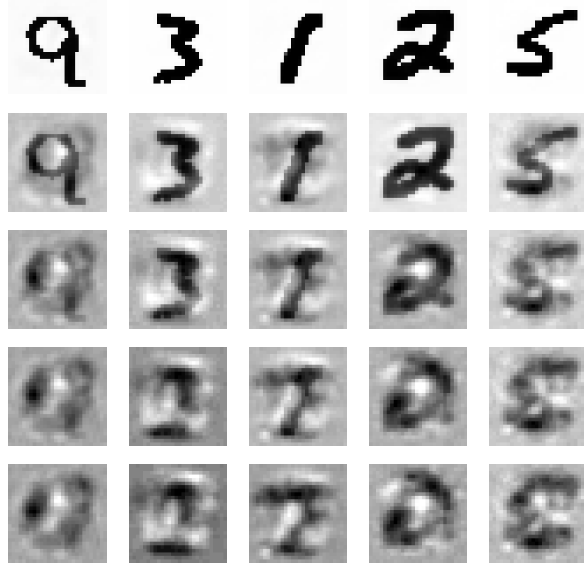


Figure 1: Sample generation by maximization of predicted dual variables. The top row shows the initial images from MNIST, to which the predictor associates a low dual variable. Rows 2-5 display the images resulting from subsequent iterations of gradient ascent. As the predicted dual variable increases, patterns corresponding to other classes appear, increasing the uncertainty on the true image label.

249 Our proposed framework allows us to leverage the trained backbone and dual regression head in  
 250 a *generative* manner to create novel samples that are most informative. Here, we focus on the  
 251 MNIST dataset, and use the trained model to generate synthetic images with maximal associated dual  
 252 variables. We begin by training a MLP on 10% of the MNIST dataset. Then, similarly to [57], we

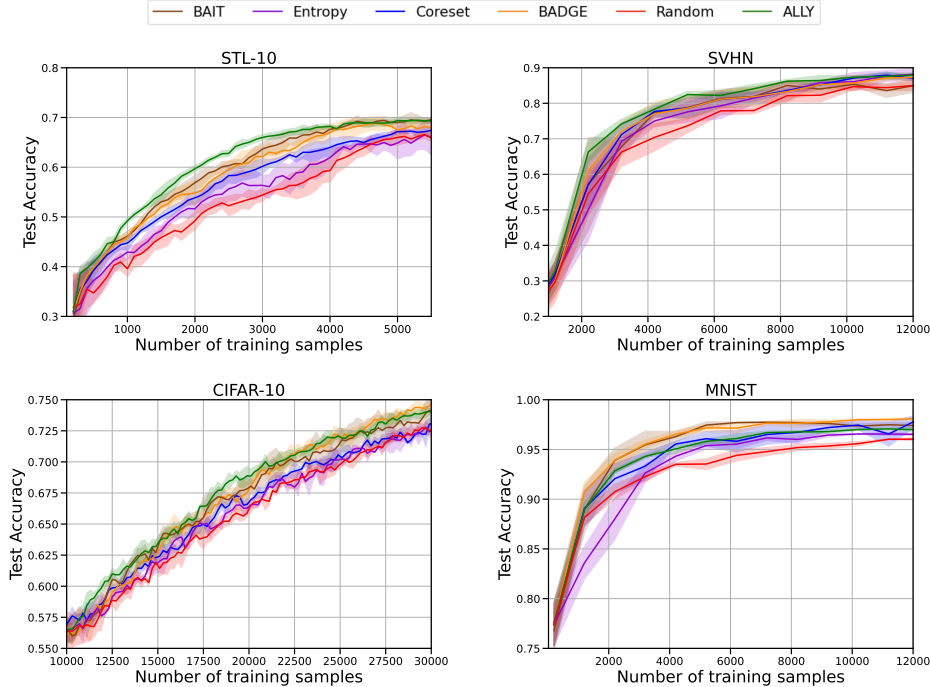


Figure 2: Accuracy in the test set as a function of the number of training samples in four classification settings and two budgets, 200 (left) and 1000 (right). Solid curves represent the mean across five different random seeds, while shaded regions correspond to the standard deviation.

253 perform gradient ascent on images that are initially considered uninformative, so as to maximize their  
 254 predicted dual variables. The progression of the resulting images is demonstrated in Figure 1.

255 As the predicted dual variable increases, patterns corresponding to other digits appear, increasing  
 256 the uncertainty on the true label of the image. For instance, the third column of Figure 1 shows  
 257 a handwritten ‘1’ that is progressively transformed into a blurred superposition of ‘7,’ ‘2’ and ‘1’.  
 258 Images in the last row of Figure 1 can be interpreted as lying in the tails of the distribution  $\mathcal{Q}$ , or  
 259 close to the decision boundary of the end-to-end model  $f_{\theta}$ .

### 260 5.3 Classification Results

261 The experiments on STL-10, CIFAR-10, SVHN and MNIST are all 10-class, image classification  
 262 tasks. We use the cross-entropy loss for both  $\ell(\cdot, \cdot)$  and  $\ell'(\cdot, \cdot)$  and set  $\epsilon(\mathbf{x}) = 0.2, \forall \mathbf{x}$ . As shown  
 263 in Figure 2, ALLY outperforms other baselines in STL-10, SVHN and CIFAR-10. In these three  
 264 datasets, the improvement in the number of samples needed by ALLY to achieve 97% of the final  
 265 accuracy, in comparison with the best baseline, is 9%, 8% and 2%, respectively. In MNIST, however,  
 266 BADGE and BAIT consistently outperform ALLY. This is partly due to the fact that the MLP, being  
 267 less expressive, yields embeddings of lower quality, hindering the prediction of dual variables.

### 268 5.4 Regression Results

269 We use mean-squared error for both  $\ell(\cdot, y)$  and  $\ell'(\cdot, y)$  and set  $\epsilon(\mathbf{x}) = 0.1, \forall \mathbf{x}$ . As seen in Figure  
 270 3, ALLY outperforms both Random and Coreset in this regression task, the gap being larger at the  
 271 beginning of the learning curve. Note that BADGE and Entropy are not applicable, since they are  
 272 limited to classification scenarios.<sup>1</sup>

<sup>1</sup>Although BAIT can be used for regression, the implementation code for the regression version of BAIT is not publicly available.

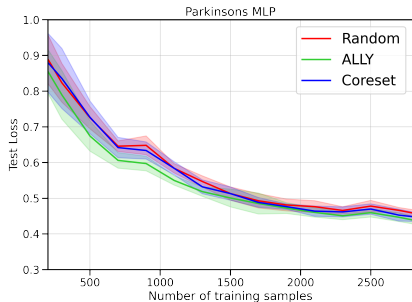


Figure 3: Mean-squared error in the test set as a function of the number of training samples in the Parkinson’s telemonitoring dataset.

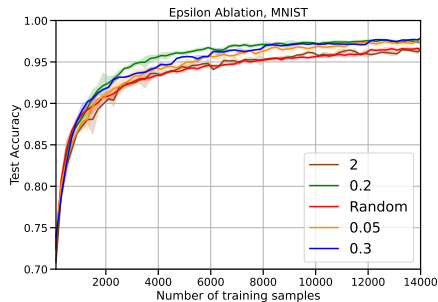


Figure 4: Ablation on the value of the constraint tightness  $\epsilon$  on the performance of ALLY in MNIST.

### 273 5.5 Ablation on the Constraint Tightness

274 Figure 4 illustrates the role  $\epsilon$  plays in the optimization procedure. For very large values of epsilon, the  
 275 constraint slacks become negative for all samples, and thus all dual variables become zero, making  
 276 them uninformative (analogous to an unconstrained problem) and the method performs similarly to  
 277 random sampling. Our experiments suggest that values in the range  $[1.1p_u, 1.3p_u]$ , where  $p_u$  is the  
 278 average loss observed when training the model without constraints, work well in practice.

## 279 6 Concluding Remarks

280 We presented ALLY, a principled batch active learning method based on Lagrangian duality. Our  
 281 method formulates the learning problem using constrained optimization and solves it in the Lagrangian  
 282 dual domain via a primal-dual approach. We then showed that the magnitude of the optimal dual  
 283 variables could be viewed as a measure of informativeness of the corresponding training sample, as it  
 284 indicates the sensibility of the optimum value of the objective function with respect to a perturbation  
 285 in the constraint.

286 Following the completion of the primal-dual learning phase, we leveraged the learned sample  
 287 representations, as well as their respective dual variables, to train a dual regression head. This  
 288 predictor is used to estimate the dual variables associated to unlabeled samples. We then promote  
 289 diversity in the batch by clustering the unlabeled sample embeddings selecting the samples with the  
 290 highest estimated dual variables from each cluster. We demonstrated that this principled method  
 291 outperforms state-of-the-art batch active learning algorithms in several classification and regression  
 292 experiments.

293 The image synthesis experiment shows that the trained model can shed light on the informativeness  
 294 measure induced by ALLY. In addition, it demonstrates that informative samples and outliers (such as  
 295 mislabeled samples) may be hard to distinguish. Recent empirical findings suggest that many active  
 296 learning algorithms consistently prefer to acquire samples that traditional models fail to learn [53].  
 297 Modifying ALLY in order to avoid sampling these so-called *collective outliers* (e.g., by setting an  
 298 upper bound on the dual variable associated to the queried samples) is a promising research direction  
 299 that we leave for future work. In our experiments, we have set the secondary loss to be identical  
 300 to the primary supervised loss (i.e., cross-entropy loss for classification, and mean-squared error  
 301 for regression). However, evaluating the performance of ALLY under alternative *unsupervised* or  
 302 *self-supervised* secondary losses is a promising future direction. Finally, it would be interesting to  
 303 evaluate the performance of ALLY under different diversity measures, comparing their computational  
 304 burden.

## 305 References

- 306 [1] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,  
307 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel  
308 Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler,  
309 Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott  
310 Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya  
311 Sutskever, and Dario Amodei. Language models are few-shot learners. In *Advances in Neural  
312 Information Processing Systems*, volume 33, pages 1877–1901, 2020.
- 313 [2] Dmitry Lepikhin, HyoukJoong Lee, Yuanzhong Xu, Dehao Chen, Orhan Firat, Yanping Huang,  
314 Maxim Krikun, Noam Shazeer, and Zhifeng Chen. GShard: Scaling giant models with condi-  
315 tional computation and automatic sharding. In *International Conference on Learning Represen-  
316 tations*, 2021. URL <https://openreview.net/forum?id=qrwe7XHTmYb>.
- 317 [3] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai,  
318 Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly,  
319 Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image  
320 recognition at scale. In *International Conference on Learning Representations*, 2021. URL  
321 <https://openreview.net/forum?id=YicbFdNTTy>.
- 322 [4] Ying Liu. Active learning with support vector machine applied to gene expression data for  
323 cancer classification. *Journal of chemical information and computer sciences*, 44 6:1936–41,  
324 2004.
- 325 [5] Steven C. H. Hoi, Rong Jin, Jianke Zhu, and Michael R. Lyu. Batch mode active learning  
326 and its application to medical image classification. In *Proceedings of the 23rd International  
327 Conference on Machine Learning*, ICML '06, page 417–424, New York, NY, USA, 2006.  
328 Association for Computing Machinery. ISBN 1595933832. doi: 10.1145/1143844.1143897.  
329 URL <https://doi.org/10.1145/1143844.1143897>.
- 330 [6] V. Nath, Dong Yang, Bennett A. Landman, Daguang Xu, and Holger R. Roth. Diminishing  
331 uncertainty within the training pool: Active learning for medical image segmentation. *IEEE  
332 Transactions on Medical Imaging*, 40:2534–2547, 2021.
- 333 [7] Burr Settles. Active learning literature survey. 2009.
- 334 [8] Ozan Sener and Silvio Savarese. Active learning for convolutional neural networks: A core-set  
335 approach, 2018.
- 336 [9] Gui Citovsky, Giulia DeSalvo, Claudio Gentile, Lazaros Karydas, Anand Rajagopalan, Afshin  
337 Rostamizadeh, and Sanjiv Kumar. Batch active learning at scale. *CoRR*, abs/2107.14263, 2021.  
338 URL <https://arxiv.org/abs/2107.14263>.
- 339 [10] David D. Lewis and William A. Gale. A sequential algorithm for training text classifiers. *CoRR*,  
340 abs/cmp-lg/9407020, 1994. URL <http://arxiv.org/abs/cmp-lg/9407020>.
- 341 [11] Ehsan Elhamifar, Guillermo Sapiro, Allen Yang, and S. Shankar Sasrty. A convex optimization  
342 framework for active learning. In *2013 IEEE International Conference on Computer Vision*,  
343 pages 209–216, 2013. doi: 10.1109/ICCV.2013.33.
- 344 [12] Burr Settles, Mark Craven, and Soumya Ray. Multiple-instance active learning. In J. Platt,  
345 D. Koller, Y. Singer, and S. Roweis, editors, *Advances in Neural Information Processing  
346 Systems*, volume 20. Curran Associates, Inc., 2008. URL [https://proceedings.neurips.  
347 cc/paper/2007/file/a1519de5b5d44b31a01de013b9b51a80-Paper.pdf](https://proceedings.neurips.cc/paper/2007/file/a1519de5b5d44b31a01de013b9b51a80-Paper.pdf).
- 348 [13] Wenbin Cai, Ya Zhang, and Jun Zhou. Maximizing expected model change for active learning  
349 in regression. In *2013 IEEE 13th International Conference on Data Mining*, pages 51–60, 2013.  
350 doi: 10.1109/ICDM.2013.104.
- 351 [14] Nicholas Roy and Andrew McCallum. Toward optimal active learning through sampling  
352 estimation of error reduction. In *Proceedings of the Eighteenth International Conference  
353 on Machine Learning*, ICML '01, page 441–448, San Francisco, CA, USA, 2001. Morgan  
354 Kaufmann Publishers Inc. ISBN 1558607781.
- 355 [15] Jordan T. Ash, Chicheng Zhang, Akshay Krishnamurthy, John Langford, and Alekh Agarwal.  
356 Deep batch active learning by diverse, uncertain gradient lower bounds. *CoRR*, abs/1906.03671,  
357 2019. URL <http://arxiv.org/abs/1906.03671>.

- 358 [16] K. Brinker. Incorporating diversity in active learning with support vector machines. In *Proceed-*  
359 *ings of the 20th International Conference on Machine Learning (ICML 2000)*, 2003.
- 360 [17] Yoshua Bengio, Aaron Courville, and Pascal Vincent. Representation learning: A review and  
361 new perspectives. *IEEE transactions on pattern analysis and machine intelligence*, 35(8):  
362 1798–1828, 2013.
- 363 [18] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive  
364 predictive coding. *arXiv preprint arXiv:1807.03748*, 2018.
- 365 [19] Jordan T. Ash, Surbhi Goel, Akshay Krishnamurthy, and Sham M. Kakade. Gone fishing:  
366 Neural active learning with fisher embeddings. *CoRR*, abs/2106.09675, 2021. URL <https://arxiv.org/abs/2106.09675>.
- 368 [20] Yarin Gal, Riashat Islam, and Zoubin Ghahramani. Deep bayesian active learning with image  
369 data. *CoRR*, abs/1703.02910, 2017. URL <http://arxiv.org/abs/1703.02910>.
- 370 [21] Donggeun Yoo and In So Kweon. Learning loss for active learning. *CoRR*, abs/1905.03677,  
371 2019. URL <http://arxiv.org/abs/1905.03677>.
- 372 [22] Fedor Zhdanov. Diverse mini-batch active learning. *CoRR*, abs/1901.05954, 2019. URL  
373 <http://arxiv.org/abs/1901.05954>.
- 374 [23] Zalán Bodó, Zsolt Minier, and Lehel Csató. Active learning with clustering. In Isabelle  
375 Guyon, Gavin Cawley, Gideon Dror, Vincent Lemaire, and Alexander Statnikov, editors, *Active*  
376 *Learning and Experimental Design workshop In conjunction with AISTATS 2010*, volume 16  
377 of *Proceedings of Machine Learning Research*, pages 127–139, Sardinia, Italy, 16 May 2011.  
378 PMLR. URL <https://proceedings.mlr.press/v16/bodo11a.html>.
- 379 [24] Xueying Zhan, Qingzhong Wang, Kuan-hao Huang, Haoyi Xiong, Dejing Dou, and Antoni B.  
380 Chan. A comparative survey of deep active learning, 2022. URL <https://arxiv.org/abs/2203.13450>.
- 382 [25] Javad Azimi, Alan Fern, Xiaoli Zhang-Fern, Glencora Borradaile, and Brent Heeringa. Batch  
383 active learning via coordinated matching, 2012. URL <https://arxiv.org/abs/1206.6458>.
- 384 [26] Andreas Kirsch, Sebastian Farquhar, and Yarin Gal. A simple baseline for batch active learning  
385 with stochastic acquisition functions. *CoRR*, abs/2106.12059, 2021. URL <https://arxiv.org/abs/2106.12059>.
- 387 [27] Daniel Gissin and Shai Shalev-Shwartz. Discriminative active learning. *CoRR*, abs/1907.06347,  
388 2019. URL <http://arxiv.org/abs/1907.06347>.
- 389 [28] Luiz F. O. Chamon and Alejandro Ribeiro. Probably approximately correct constrained learning,  
390 2021.
- 391 [29] Luiz F. O. Chamon, Santiago Paternain, Miguel Calvo-Fullana, and Alejandro Ribeiro. Con-  
392 strained learning with non-convex losses. *CoRR*, abs/2103.05134, 2021. URL <https://arxiv.org/abs/2103.05134>.
- 394 [30] Juan Cerviño, Luana Ruiz, and Alejandro Ribeiro. Training stable graph neural networks  
395 through constrained learning. In *ICASSP 2022 - 2022 IEEE International Conference on*  
396 *Acoustics, Speech and Signal Processing (ICASSP)*, pages 4223–4227, 2022. doi: 10.1109/  
397 ICASSP43922.2022.9746912.
- 398 [31] Zebang Shen, Juan Cervino, Hamed Hassani, and Alejandro Ribeiro. An agnostic approach to  
399 federated learning with class imbalance. In *International Conference on Learning Representa-*  
400 *tions*, 2022. URL <https://openreview.net/forum?id=Xo01bDt975>.
- 401 [32] Raghu Arghal, Eric Lei, and Shirin Saeedi Bidokhti. Robust graph neural networks via  
402 probabilistic lipschitz constraints. *CoRR*, abs/2112.07575, 2021. URL <https://arxiv.org/abs/2112.07575>.
- 404 [33] Alexander Robey, Luiz Chamon, George J. Pappas, Hamed Hassani, and Alejandro  
405 Ribeiro. Adversarial robustness with semi-infinite constrained learning. In M. Ran-  
406 zato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan, editors, *Ad-*  
407 *vances in Neural Information Processing Systems*, volume 34, pages 6198–6215. Curran  
408 Associates, Inc., 2021. URL <https://proceedings.neurips.cc/paper/2021/file/312ecfdfa8b239e076b114498ce21905-Paper.pdf>.
- 409

- 410 [34] V.N. Vapnik. An overview of statistical learning theory. *IEEE Transactions on Neural Networks*,  
411 10(5):988–999, 1999. doi: 10.1109/72.788640.
- 412 [35] Corinna Cortes and Vladimir Vapnik. Support-vector networks. *Machine learning*, 20(3):  
413 273–297, 1995.
- 414 [36] Shai Shalev-Shwartz and Shai Ben-David. *Understanding Machine Learning: From Theory to*  
415 *Algorithms*. Cambridge University Press, USA, 2014. ISBN 1107057132.
- 416 [37] L. Hurwicz K. J. Arrow and H. Uzawa. Studies in linear and non-linear programming, by  
417 k. j. arrow, l. hurwicz and h. uzawa. stanford university press, 1958. 229 pages. *Canadian*  
418 *Mathematical Bulletin*, 3(3):196–198, 1960. doi: 10.1017/S0008439500025522.
- 419 [38] Ferdinando Fioretto, Pascal Van Hentenryck, Terrence W. K. Mak, Cuong Tran, Federico Baldo,  
420 and Michele Lombardi. Lagrangian duality for constrained deep learning. In *ECML/PKDD (5)*,  
421 pages 118–135, 2020. URL [https://doi.org/10.1007/978-3-030-67670-4\\_8](https://doi.org/10.1007/978-3-030-67670-4_8).
- 422 [39] Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding  
423 deep learning requires rethinking generalization. *CoRR*, abs/1611.03530, 2016. URL <http://arxiv.org/abs/1611.03530>.  
424
- 425 [40] Devansh Arpit, Stanisław Jastrzębski, Nicolas Ballas, David Krueger, Emmanuel Bengio,  
426 Maxinder S Kanwal, Tegan Maharaj, Asja Fischer, Aaron Courville, Yoshua Bengio, et al.  
427 A closer look at memorization in deep networks. In *International Conference on Machine*  
428 *Learning*, pages 233–242. PMLR, 2017.
- 429 [41] Mariya Toneva, Alessandro Sordani, Remi Tachet des Combes, Adam Trischler, Yoshua Bengio,  
430 and Geoffrey J. Gordon. An empirical study of example forgetting during deep neural network  
431 learning. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=BJ1xm30cKm>.  
432
- 433 [42] Angelos Katharopoulos and François Fleuret. Not all samples are created equal: Deep learning  
434 with importance sampling. *CoRR*, abs/1803.00942, 2018. URL <http://arxiv.org/abs/1803.00942>.  
435
- 436 [43] Yonglong Tian, Dilip Krishnan, and Phillip Isola. Contrastive multiview coding. In *Com-*  
437 *puter Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020,*  
438 *Proceedings, Part XI 16*, pages 776–794. Springer, 2020.
- 439 [44] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework  
440 for contrastive learning of visual representations. In *International conference on machine*  
441 *learning*, pages 1597–1607. PMLR, 2020.
- 442 [45] Harikrishna Narasimhan, Andrew Cotter, Yichen Zhou, Serena Wang, and Wenshuo Guo.  
443 Approximate heavily-constrained learning with lagrange multiplier models. In *Advances in*  
444 *Neural Information Processing Systems*, volume 33, pages 8693–8703, 2020.
- 445 [46] Adam Coates, Andrew Ng, and Honglak Lee. An analysis of single-layer networks in unsu-  
446 pervised feature learning. In Geoffrey Gordon, David Dunson, and Miroslav Dudík, editors,  
447 *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*,  
448 volume 15 of *Proceedings of Machine Learning Research*, pages 215–223. Fort Lauderdale, FL,  
449 USA, 11–13 Apr 2011. PMLR. URL <https://proceedings.mlr.press/v15/coates11a.html>.  
450
- 451 [47] Alex Krizhevsky. Learning multiple layers of features from tiny images. Technical report, 2009.
- 452 [48] Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Bo Wu, and Andrew Y. Ng.  
453 Reading digits in natural images with unsupervised feature learning. In *NIPS Workshop on Deep*  
454 *Learning and Unsupervised Feature Learning 2011*, 2011. URL [http://ufldl.stanford.edu/housenumbers/nips2011\\_housenumbers.pdf](http://ufldl.stanford.edu/housenumbers/nips2011_housenumbers.pdf).  
455
- 456 [49] Yann LeCun and Corinna Cortes. MNIST handwritten digit database. 2010. URL <http://yann.lecun.com/exdb/mnist/>.  
457
- 458 [50] Athanasios Tsanas, Max A. Little, Patrick E. McSharry, and Lorraine O. Ramig. Accurate tele-  
459 monitoring of parkinson’s disease progression by noninvasive speech tests. *IEEE Transactions*  
460 *on Biomedical Engineering*, 57(4):884–893, 2010. doi: 10.1109/TBME.2009.2036000.
- 461 [51] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image  
462 recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern*  
463 *Recognition (CVPR)*, June 2016.

- 464 [52] Stuart Lloyd. Least squares quantization in pcm. *IEEE transactions on information theory*, 28  
465 (2):129–137, 1982.
- 466 [53] Siddharth Karamcheti, Ranjay Krishna, Li Fei-Fei 0001, and Christopher D. Manning. Mind  
467 your outliers! investigating the negative impact of outliers on active learning for visual question  
468 answering. In Chengqing Zong, Fei Xia, Wenjie Li 0002, and Roberto Navigli, editors,  
469 *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and  
470 the 11th International Joint Conference on Natural Language Processing, ACL/IJCNLP 2021,  
471 (Volume 1: Long Papers), Virtual Event, August 1-6, 2021*, pages 7265–7281. Association for  
472 Computational Linguistics, 2021. ISBN 978-1-954085-52-7. URL [https://aclanthology.  
473 org/2021.acl-long.564](https://aclanthology.org/2021.acl-long.564).
- 474 [54] Yonatan Geifman and Ran El-Yaniv. Deep active learning over the long tail. *CoRR*,  
475 abs/1711.00941, 2017. URL <http://arxiv.org/abs/1711.00941>.
- 476 [55] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan,  
477 Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas  
478 Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy,  
479 Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-  
480 performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-  
481 Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*,  
482 pages 8024–8035. Curran Associates, Inc., 2019. URL [http://papers.nips.cc/paper/  
483 9015-pytorch-an-imperative-style-high-performance-deep-learning-library.  
484 pdf](http://papers.nips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf).
- 485 [56] Yao-Yuan Yang, Shao-Chuan Lee, Yu-An Chung, Tung-En Wu, Si-An Chen, and Hsuan-Tien  
486 Lin. libact: Pool-based active learning in python. Technical report, National Taiwan University,  
487 October 2017. URL <https://github.com/ntucllab/libact>. available as arXiv preprint  
488 <https://arxiv.org/abs/1710.00379>.
- 489 [57] Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adver-  
490 sarial examples. In Yoshua Bengio and Yann LeCun, editors, *3rd International Conference on  
491 Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track  
492 Proceedings*, 2015. URL <http://arxiv.org/abs/1412.6572>.

493 **Checklist**

- 494 1. For all authors...
- 495 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's  
496 contributions and scope? [Yes]
- 497 (b) Did you describe the limitations of your work? [Yes]
- 498 (c) Did you discuss any potential negative societal impacts of your work? [N/A]
- 499 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
500 them? [Yes]
- 501 2. If you are including theoretical results...
- 502 (a) Did you state the full set of assumptions of all theoretical results? [Yes]
- 503 (b) Did you include complete proofs of all theoretical results? [Yes]
- 504 3. If you ran experiments...
- 505 (a) Did you include the code, data, and instructions needed to reproduce the main experi-  
506 mental results (either in the supplemental material or as a URL)? [Yes]
- 507 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
508 were chosen)? [Yes]
- 509 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
510 ments multiple times)? [Yes]
- 511 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
512 of GPUs, internal cluster, or cloud provider)? [Yes]
- 513 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 514 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 515 (b) Did you mention the license of the assets? [N/A]
- 516 (c) Did you include any new assets either in the supplemental material or as a URL? [No]
- 517 (d) Did you discuss whether and how consent was obtained from people whose data you're  
518 using/curating? [N/A]
- 519 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
520 information or offensive content? [N/A]
- 521 5. If you used crowdsourcing or conducted research with human subjects...
- 522 (a) Did you include the full text of instructions given to participants and screenshots, if  
523 applicable? [N/A]
- 524 (b) Did you describe any potential participant risks, with links to Institutional Review  
525 Board (IRB) approvals, if applicable? [N/A]
- 526 (c) Did you include the estimated hourly wage paid to participants and the total amount  
527 spent on participant compensation? [N/A]