Unsupervised Learning under Latent Label Shift

Anonymous Author(s) Affiliation Address email

Abstract

1	What sorts of structure might enable a learner to discover classes from unla-
2	beled data? Traditional approaches rely on feature-space similarity and heroic
3	assumptions on the data. In this paper, we introduce unsupervised learning un-
4	der <i>Latent Label Shift</i> (LLS), where the label marginals $p_d(y)$ shift but the class
5	conditionals $p(\mathbf{x} y)$ do not. This work instantiates a new principle for identifying
6	classes: elements that shift together group together. For finite input spaces, we es-
7	tablish an isomorphism between LLS and topic modeling: inputs correspond to
8	words, domains to documents, and labels to topics. Addressing continuous data,
9	we prove that when each label's support contains a separable region, analogous to
10	an anchor word, oracle access to $p(d \mathbf{x})$ suffices to identify $p_d(y)$ and $p_d(y \mathbf{x})$ up
11	to permutation. Thus motivated, we introduce a practical algorithm that leverages
12	domain-discriminative models as follows: (i) push examples through domain dis-
13	criminator $p(d \mathbf{x})$; (ii) discretize the data by clustering examples in $p(d \mathbf{x})$ space;
14	(iii) perform non-negative matrix factorization on the discrete data; (iv) combine
15	the recovered $p(y d)$ with the discriminator outputs $p(d \mathbf{x})$ to compute $p_d(y x) \forall d$.
16	With semi-synthetic experiments, we show that our algorithm can leverage domain
17	information to improve state of the art unsupervised classification methods. We
18	reveal a failure mode of standard unsupervised classification methods when data-
19	space similarity does not indicate true groupings, and show empirically that our
20	method better handles this case. Our results establish a deep connection between
21	distribution shift and topic modeling, opening promising lines for future work.

22 1 Introduction

Discovering systems of categories from unlabeled data is a fundamental but ill-posed challenge in machine learning. Typical unsupervised learning methods group instances together based on featurespace similarity. Accordingly, given a collection of photographs of animals, a practitioner might hope that, in some appropriate feature space, images of animals of the same species should be somehow similar to each other. But why should we expect a clustering algorithm to recognize that dogs viewed in sunlight and dogs viewed at night belong to the same category? Why should we expect that butterflies and caterpillars should lie close together in feature space?

In this paper, we offer an alternative principle according to which we might identify a set of classes: 30 we exploit distribution shift across times and locations to reveal otherwise unrecognizable groupings 31 among examples. For example, if we noticed that whenever we found ourselves in a location where 32 butterflies are abundant, caterpillars were similarly abundant, and that whenever butterflies were 33 scarce, caterpillars had a similar drop in prevalence, we might conclude that the two were tied to the 34 same underlying concept, no matter how different they appear in feature space. In short, our principle 35 suggests that latent classes might be uncovered whenever instances that shift together group together. 36 Formalizing this intuition, we introduce the problem of unsupervised learning under Latent Label 37

Shift (LLS). Here, we assume access to a collection of domains $d \in \{1, \ldots, r\}$, where the mixture

Submitted to 36th Conference on Neural Information Processing Systems (NeurIPS 2022). Do not distribute.



Figure 1: Schematic of our DDFA algorithm. After training a domain discriminator, we (i) push all data through the discriminator; (ii) cluster the data based on discriminator outputs; (iii) solve the resulting discrete topic modeling problem and then combine $\hat{q}(d|x)$ and $\hat{q}(y,d)$ to estimate $\hat{p}_d(y|x)$.

proportions $p_d(y)$ vary across domains but the class conditional distribution p(x|y) is domain-39 invariant. Our goals are to recover the underlying classes up to permutation, and thus to identify 40 both the per-domain mixture proportions $p_d(y)$ and optimally adapted per-domain classifiers $p_d(y|x)$. 41 The essential feature of our setup is that only the true y's, as characterized by their class-conditional 42 distributions p(x|y), could account for the observed shifts in $p_d(x)$. We prove that under mild 43 assumptions, knowledge of this underlying structure is sufficient for inducing the full set of categories. 44 First, we focus on the *tabular setting*, demonstrating that when the input space is discrete and finite, 45 LLS is isomorphic to topic modeling [8]. Here, each distinct input x maps to a word each latent 46 *label y* maps to a *topic* and each domain d maps to a document. In this case, we can apply standard 47

⁴⁷ *identification results for topic modeling* [20, 4, 27, 32, 12] that rely only on the existence of anchor

words within each topic (for each label y_i there is at least one x in the support of y_i , that is not in

the support of any $y_j \neq y_i$). Here, standard methods based on Non-negative Matrix Factorization

(NMF) can recover each domain's underlying mixture proportion $p_d(y)$ and optimal predictor $p_d(y|x)$. [20, 32, 27]. However, the restriction to discrete inputs, while appropriate for topic modeling, proves

restrictive when our interests extend to high-dimensional continuous input spaces.

Then, to handle high-dimensional inputs, we propose Discriminate-Discretize-Factorize-Adjust 54 (DDFA), a general framework that proceeds in the following steps: (i) pool data from all domains to 55 produce a mixture distribution q(x, d); (ii) train a domain discriminative model f to predict q(d|x); 56 (iii) push all data through f, cluster examples in the pushforward distribution, and tabularize the 57 data based on cluster membership; (iv) solve the resulting discrete topic modeling problem (e.g., via 58 NMF), estimating q(y, d) up to permutation of the latent labels; (v) combine the predicted q(d|x)59 and q(y, d) to to estimate $p_d(y)$ and $p_d(y|x)$. In developing this approach, we draw inspiration from 60 recent works on distribution shift and learning from positive and unlabeled data that (i) leverage 61 black box predictors to perform dimensionality reduction [38, 23, 24]; and (ii) work with *anchor sets*, 62 separable subsets of continuous input spaces that belong to only one class's support [50, 39, 21, 6, 24]. 63 Our main theoretical result shows that domain discrimination provides a sufficient representation 64 for identifying all parameters of interest. Given oracle access to q(d|x) (which is identified without 65 66 labels), our procedure is asymptotically consistent. Our analysis reveals that the true q(d|x) maps all points in the same anchor set to a single point mass in the push-forward distribution. This motivates 67

⁶⁸ our practical approach of discretizing data by hunting for tight clusters in $\hat{q}(d|x)$ space.

In semi-synthetic experiments, we adapt existing image classification benchmarks to the LLS setting, 69 sampling without replacement to construct collections of label-shifted domains. We note that training 70 a domain discriminative classifier is a difficult task, and find that warm starting the initial layers of our 71 model with pretrained weights from unsupervised approaches can significantly boost performance. We 72 show that warm-started DDFA outperforms state-of-the-art (SOTA) unsupervised approaches when 73 domain marginals $p_d(y)$ are sufficiently sparse. In particular, we observe improvements of as much 74 as 30% accuracy over unsupervised SOTA on CIFAR-20. Further, on subsets of FieldGuide dataset, 75 where similarity between species and diversity within a species leads to failure of unsupervised 76 learning, we show that DDFA recovers the true distinctions. To be clear, these are not apples-to-77 apples comparisons: our methods are specifically tailored to the LLS setting. The takeaway is that the 78 structure of the LLS setting can be exploited to outperform the best unsupervised learning heuristics. 79

80 2 Related Work

Unsupervised Learning Standard unsupervised learning approaches for discovering labels often rely 81 82 on similarity in the original data space [40, 48]. While distances in feature space become meaningless for high-dimensional data, deep learning researchers have turned to similarity in a representations 83 space learned via self-supervised contrastive tasks [42, 19, 26, 11], or similarity in a feature space 84 learned end-to-end for a clustering task [9, 10, 45, 55]. Our problem setup closely resembles 85 independent component analysis (ICA), where one seeks to identify statistically independent signal 86 components from mixtures [33]. However, ICA's assumption of statistical independence among 87 88 the components does not obtain in our setup. In topic modeling [8, 4, 32, 12, 44], documents are 89 modeled as mixtures of topics, and topics as categorical distributions over a finite vocabulary. Topic models were pioneered by Latent Dirichlet Allocation (LDA) [8] and closely followed by papers 90 that relaxed assumptions on the distribution of topic mixing coefficients (pLSI) [31, 44]. The topic 91 modeling literature often draws on non-negative Matrix Factorization (NMF) methods [43, 51], which 92 decompose a given matrix into a product of two matrices with non-negative elements [18, 16, 25, 28]. 93 In both Topic Modeling and NMF, a fundamental problem has been to characterize the precise 94 conditions under which the system is uniquely identifiable [20, 4, 32, 12]. The anchor condition (also 95 referred to as separability) is known to be instrumental for identifying topic models [4, 12, 32, 20]. 96 In this work, we extend these ideas, leveraging separable subsets of each label's support (the anchor 97 98 sets) to produce anchor words in the discretized problem. Existing methods have attempted to extend latent variable modeling to continuous input domains by making assumptions about the functional 99 forms of the class-conditional densities, e.g., restricting to Gaussian mixtures [48, 47]. A second line 100 of approach involves finding an appropriate discretization of the continuous space [54]. 101

Distribution Shift under the Label Shift Assumption The label shift assumption, where $p_d(y)$ 102 can vary but p(x|y) cannot, has been extensively studied in the domain adaptation literature [49, 52, 103 57, 38, 23] and also obtains in the problem of learning from positive and unlabeled data [22, 7, 24]. 104 For both problems, many classical approaches suffer from the curse of dimensionality, failing in 105 the settings where deep learning prevails. Our solution strategy draws inspiration from recent work 106 on label shift [38, 1, 5, 23] and PU learning [7, 39, 50, 24] that leverage black-box predictors to 107 produce sufficient low-dimensional representations for identifying target distributions of interest 108 (other works leverage black box predictors heuristically [34]). Key differences: While PU learning 109 requires identifying *one* new class for which we lack labeled examples provided that the positive 110 class contains an anchor set [24], LLS can identify an arbitrary number of classes (up to permutation) 111 from completely unlabeled data, provided a sufficient number of domains. 112

Domain Generalization The related problem of Domain Generalization (DG) also addresses learning with data drawn from multiple distributions and where the domain identifiers play a key role [41, 2]. However in DG, we are given *labeled* data from multiple domains, and our goal is to learn a classifier that can generalize to new domains. By contrast, in LLS, we work with unlabeled data only, leveraging the problem structure to identify the underlying labels.

118 **3** Unsupervised Learning under Latent Label Shift

Notation For a vector $v \in \mathbb{R}^p$, we use v_j to denote its j^{th} entry, and for an event E, we let $\mathbb{I}[E]$ denote the binary indicator of the event. By |A|, we denote the cardinality of set A. With [n], we denote the set $\{1, 2, \ldots, n\}$. We use $[A]_{i,j}$ to access the element at (i, j) in A. Let \mathcal{X} be the input space and $\mathcal{Y} = \{1, 2, \ldots, k\}$ the output space for multiclass classification. Throughout this paper, we use capital letters to denote random variables and small case letters to denote the corresponding values they take. For example, by X we denote the input random variable and by x, we denote a value that X may take.

We now formally introduce the problem of unsupervised learning under LLS. In LLS, we assume that we observe unlabeled data from r domains. Let $\mathcal{R} = \{1, 2, ..., r\}$ be the set of domains. By p_d , we denote the probability density (or mass) function for each domain $d \in \mathcal{R}$.

Definition 1 (Latent label shift). We observe data from r domains. While the label distribution among these domains can change, for all $d, d' \in \mathcal{R}$ and for all $(x, y) \in \mathcal{X} \times \mathcal{Y}$, we have $p_d(x|y) = p_{d'}(x|y)$.

Simply put, Definition 1 states that the conditional distribution $p_d(x|y)$ remains invariant across domains, i.e., they satisfy the label shift assumption. Thus, we can drop the subscript on this factor, denoting all $p_d(x|y)$ by p(x|y). Crucially, under LLS, $p_d(y)$ can vary across different domains.



Figure 2: Relationship under Q between observed D, observed X, and latent Y.

Under LLS, we observe unlabeled data with domain label $\{(x_1, d_1), (x_2, d_2), \dots, (x_n, d_n)\}$. Our goal breaks down into two tasks. Upto permutation of labels, we aim to (i) estimate the label marginal in each domain $p_d(y)$; and (ii) estimate the optimal per-domain predictor $p_d(y|x)$.

Mixing distribution Q A key step in our algorithm will be to train a domain discriminative model. 136 Towards this end we define Q, a distribution over $\mathcal{X} \times \mathcal{Y} \times \mathcal{R}$, constructed by taking a uniform mixture 137 over all domains. By q, we denote the probability density (or mass) function of Q. Define Q such 138 that $q(x, y|D = d) = p_d(x, y)$, i.e., when we condition on D = d we recover the joint distribution 139 over $\mathcal{X} \times \mathcal{Y}$ specific to that domain d. For all $d \in \mathcal{R}$, we define $\gamma_d = q(d)$, i.e., the prevalence of each 140 domain in our distribution Q. Notice that q(x, y) is a mixture over the distributions $\{p_d(x, y)\}_{d \in \mathcal{R}}$, 141 with $\{\gamma_d\}_{d\in\mathcal{R}}$ as the corresponding mixture coefficients. Under LLS (Definition 1), X does not 142 depend on D when conditioned on Y (Fig. 2). 143

Additional notation for the discrete case To begin, we setup notation for discrete input spaces with $|\mathcal{X}| = m$. Without loss of generality, we assume that $\mathcal{X} = \{1, 2, ..., m\}$. The label shift assumption allows us to formulate the label marginal estimation problem in matrix form. Let $\mathbf{Q}_{X|D}$ be an $m \times r$ matrix such that $[\mathbf{Q}_{X|D}]_{i,d} = p_d(X = i)$, i.e., the *d*-th column of $\mathbf{Q}_{X|D}$ is $p_d(x)$. Let $\mathbf{Q}_{X|Y}$ be an $m \times k$ matrix such that $[\mathbf{Q}_{X|Y}]_{i,j} = p(X = i|Y = j)$, the *j*-th column is a distribution over *X* given Y = j. Similarly, define $\mathbf{Q}_{Y|D}$ as a $k \times r$ matrix whose *d*-th column is the domain marginal $p_d(y)$. Now with Definition 1, we have $p_d(x) = \sum_y p_d(x, y) = \sum_y p_d(x|y)p_d(y) = \sum_y p(x|y)p_d(y)$. Since this is true $\forall d \in \mathcal{R}$, we can express this in a matrix form as $\mathbf{Q}_{X|D} = \mathbf{Q}_{X|Y}\mathbf{Q}_{Y|D}$.

Additional assumptions Before we present identifiability results for the LLS problem, we introduce four additional assumptions required throughout the paper:

- A.1 There are at least as many domains as classes, i.e., $|\mathcal{R}| \ge |\mathcal{Y}|$.
- A.2 The matrix formed by label marginals (as columns) across different domains is full-rank, i.e., rank $(\mathbf{Q}_{Y|D}) = k$.
- A.3 Equal representation of domains, i.e., for all $d \in \mathcal{R}, \gamma_d = 1/r$.

158 A.4 Fix $\epsilon > 0$. For all $y \in \mathcal{Y}$, there exists a subdomain $A_y \subseteq \mathcal{X}$, such that $q(A_y) \ge \epsilon$ with 159 $q(A_y|y) > 0$ and for all $y' \in \mathcal{Y} \setminus \{y\}$, $q(A_y|y') = 0$. We refer to this assumption as ϵ -anchor 160 sub-domain condition.

We now comment on the assumptions. A.1–A.2 are benign, these assumptions just imply that the 161 matrix $\mathbf{Q}_{Y|D}$ is full row rank. Without loss of generality, A.3 can be assumed when dealing with 162 data from a collection of domains. When this condition is not satisfied, one could just re-sample data 163 points uniformly at random from each domain d. Intuitively, A.4 states that for each label $y \in \mathcal{Y}$, we 164 have some subset of inputs that only belong to that class y. To avoid vanishing probability of this 165 subset, we ensure at least ϵ probability mass in our mixing distribution Q. The anchor word condition 166 is related to the positive sub-domain in PU learning, which requires that there exists a subset of \mathcal{X} in 167 which all examples only belong to the positive class [50, 39, 21, 6]. 168

169 4 Theoretical Analysis

In this section, we establish identifiability of LLS problem. We begin by considering the case where the input space is discrete and formalize the isomorphism to topic modeling. Then we establish the identifiability of the system in this discrete setting by appealing to existing results in topic modeling [32]. Finally, extending results from discrete case, we provide novel analysis to establish our identifiability result for the continuous setting.

Isomorphism to topic modeling Recall that for the discrete input setting, we have the matrix formulation: $\mathbf{Q}_{X|D} = \mathbf{Q}_{X|Y}\mathbf{Q}_{Y|D}$. Consider a corpus of *r* documents, consisting of terms from a vocabulary of size *m*. Let **D** be an $\mathbb{R}^{m \times r}$ matrix representing the underlying corpus. Each column of **D** represents a document, and each row represents a term in the vocabulary. Each element $[\mathbf{D}]_{i,j}$ represents the frequency of term *i* in document *j*. Topic modeling [8, 31, 32, 4] considers each document to be composed as a mixture of *k* topics. Each topic prescribes a frequency with which the terms in the vocabulary occur given that topic. Further, the proportion of each topic varies across documents with the frequency of terms given topic remaining invariant.

We can state the topic modelling problem as: $\mathbf{D} = \mathbf{CW}$, where \mathbf{C} is an $\mathbb{R}^{m \times k}$ matrix, $[\mathbf{C}]_{i,j}$ represents the frequency of term *i* given topic *j*, and \mathbf{W} is an $\mathbb{R}^{k \times r}$ matrix, where $[\mathbf{W}]_{i,j}$ represents the proportion of topic *i* in document *j*. Note that all three matrices are column normalized. The isomorphism is then between document and domain, topic and label, term and input sample, i.e., $\mathbf{D} = \mathbf{CW} \equiv \mathbf{Q}_{X|D} = \mathbf{Q}_{X|Y}\mathbf{Q}_{Y|D}$. In both the cases, we are interested in decomposing a known matrix into two unknown matrices. This formulation is examined as a non-negative matrix factorization problem with an added simplicial constraint on the columns (columns sum to 1) [3, 27].

Identifiability of the topic modeling problem is well-established [20, 4, 27, 32, 12]. We leverage
 the isomorphism to topic modeling to extend this identifiability condition to our LLS setting. We
 formalize the adaption here:

Theorem 1. (adapted from Proposition 1 in Huang et al. [32]) Assume A.1, A.2 and A.4 hold (A.4 in the discrete setting is referred to as the anchor word condition). Then the solution to $\mathbf{Q}_{X|D} = \mathbf{Q}_{X|Y}\mathbf{Q}_{Y|D}$ is uniquely identified.

We refer readers to Huang et al. [32] for a proof of this theorem. Intuitively, Theorem 1 states that if each label y has at least one token in the input space that has support only in y, and A.1, A.2 hold, then the solution to $\mathbf{Q}_{X|Y}$, $\mathbf{Q}_{Y|D}$ is unique. Furthermore, under this condition, there exist algorithms that can recover $\mathbf{Q}_{X|Y}$, $\mathbf{Q}_{Y|D}$ within some permutation [32, 27, 3, 4].

Extensions to the continuous case We will prove identifiability in the continuous setting, when 200 $\mathcal{X} = \mathbb{R}^p$ for some $p \ge 1$. In addition to A.1–A.4, we make an additional assumption that we have 201 oracle access to q(d|x), i.e., the true domain discriminator for mixture distribution Q. This is implied 202 by assuming access to the marginal q(x, d) from which we observe our samples. Formally, we define 203 a push forward function f such that $[f(x)]_d = q(d|x)$, then push the data forward through f to obtain 204 outputs in Δ^{r-1} . In the proof of Theorem 2, we will show that these outputs can be discretized in a 205 fashion that maps anchor subdomains to anchor words in a tabular, discrete setting. We separately 206 remark that the anchor word outputs are in fact extreme corners of the convex polytope in Δ^{r-1} which 207 encloses all f(x) mass; we discuss this geometry further in App. F. After constructing the anchor 208 word discretization, we appeal to Theorem 1 to recover $\mathbf{Q}_{Y|D}$. Given $\mathbf{Q}_{Y|D}$, we show that we can 209 use Bayes' rule and the LLS condition (Definition 1) to identify the distribution $q(y|x, d) = p_d(y|x)$ 210 over latent variable y. We formalize this in the following theorem: 211

Theorem 2. Let the distribution Q over random variables X, Y, D satisfy Assumptions A.1–A.4. Assuming access to the joint distribution q(x, d), we show that the following quantities are identifiable: (i) $\mathbf{Q}_{Y|D}$, (ii) q(y|X = x), for all $x \in \mathcal{X}$ that lies in the support (i.e. q(x) > 0); and (iii) q(y|X = x, D = d), for all $x \in \mathcal{X}$ and $d \in \mathcal{R}$ such that q(x, d) > 0.

Before presenting a proof sketch for Theorem 2, we first present key lemmas (we include their proofs in App. B).

Lemma 1. Under the same assumptions as Theorem 2, the matrix $\mathbf{Q}_{Y|D}$ and f(x) = q(d|x) uniquely determine q(y|x) for all $y \in \mathcal{Y}$ and $x \in \mathcal{X}$ such that q(x) > 0.

Lemma 1 states that given matrix $\mathbf{Q}_{Y|D}$ and oracle domain discriminator, we can uniquely identify q(y|x). In particular, we show that for any $x \in \mathcal{X}$, q(d|x) can be expressed as a convex combination of the k columns of $\mathbf{Q}_{D|Y}$ (which is computed from $\mathbf{Q}_{Y|D}$ and is column rank k) and the coefficients of the combination are q(y|x). Combining this with the linear independence of the columns of $\mathbf{Q}_{D|Y}$, we show that these coefficients are unique. In the following lemma, we show how the identified q(y|x) can then be used to identify q(y|x, d):

Lemma 2. Under the same assumptions as Theorem 2, for all $y \in \mathcal{Y}$, and $x \in \mathcal{X}$ such that q(x, d) > 0. the matrix $\mathbf{Q}_{Y|D}$ and q(y|x) uniquely determine q(y|x, d).

To prove Lemma 2, we show that we can combine the conditional distribution over the labels given a sample $x \in \mathcal{X}$ with the prior distribution of the labels in each domain to determine the posterior distribution over labels given the sample x and the domain of interest. Next, we introduce a key property of the domain discriminator classifier f:

Lemma 3. Under the same assumptions as Theorem 2, for all x, x' in anchor sub-domain, i.e., $x, x' \in A_y$ for a given label $y \in \mathcal{Y}$, we have f(x) = f(x'). Further, for any $y \in \mathcal{Y}$, if $x \in A_y, x' \notin A_y$, then $f(x) \neq f(x')$.

Lemma 3 implies that the oracle domain discriminator f maps all points in an anchor subdomain, and only those points in that anchor subdomain to the same point in f(x) = q(d|x) space. We can now present a proof sketch for Theorem 2 (full proof in App. B):

Proof sketch of Theorem 2. The key idea of the proof lies in proposing a discretization such that some 238 subset of anchor subdomains for each label y in the continuous space map to distinct anchor words in 239 discrete space. In particular, if there exists a discretization of the continuous space \mathcal{X} that for any 240 $y \in \mathcal{Y}$, maps all $x \in A_y$ to the same point in the discrete space, but no $x \notin A_y$ maps to this point, then 241 this point serves as an anchor word. From Lemma 3, we know that all the $x \in A_y$ and only the $x \in A_y$ 242 get mapped to specific points in the f(x) space. Pushing all the $x \in \mathcal{X}$ through f, we know from A.4 243 that there exists k point masses of size ϵ , one for each $f(A_y)$ in the f(x) = q(d|x) space. We can now 244 inspect this space for point masses of size at least ϵ to find at most $\mathcal{O}(1/\epsilon)$ such point masses among 245 which are contained the k point masses corresponding to the anchor subdomains. Discretizing this 246 space by assigning each point mass to a group (and non-point masses to a single additional group), 247 we have k groups that have support only in one y each. Thus, we have achieved a discretization 248 with anchor words. Further, since the discrete space arises from a pushforward of the continuous 249 space through f, the discrete space also satisfies the latent label shift assumption A.1. We now use 250 Theorem 1 to claim identifiability of $\mathbf{Q}_{Y|D}$. We then use Lemmas 1 and 2 to prove parts (ii) and (iii). 251

252 5 DDFA Framework

Motivated by our identifiability analysis, in this section, we present an algorithm to estimate 253 $\mathbf{Q}_{Y|D}$, q(y|x), and q(y|x, d) when X is continuous by exploiting domain structure and approximat-254 ing the true domain discriminator f. Intuitively, q(y|x, d) is the domain specific classifier $p_d(y|x)$ 255 and q(y|x) is the classifier for data from aggregated domains. $\mathbf{Q}_{Y|D}$ captures label marginal for indi-256 vidual domains. A naive approach would be to aggregate data from different domains and exploit 257 recent advancements in unsupervised learning [55, 45, 9, 10]. However, aggregating data from multi-258 ple domains loses the domain structure that we hope to leverage. We highlight this failure mode of 259 the unsupervised clustering method in Sec. 6. 260

Discriminate We begin Algorithm 1 by creating a split of the unlabeled samples into the training 261 and validation sets. Using the unlabeled data samples and the domain that each sample originated 262 from, we first train a domain discriminative classifier f. The domain discriminative classifier outputs 263 a distribution over domains for a given input. This classifier is trained with cross-entropy loss to 264 predict the domain label of each sample on the training set. With unlimited data, the minimizer of 265 this loss is the true f, as we prove in App. C. To avoid overfitting, we stop training f when the cross-266 267 entropy loss on the validation set stops decreasing. Note that here the validation set also only contains 268 domain labels (and no information about true labels).

Discretize We now push forward all the samples from the training and validation sets through the 269 domain discriminator to get vector $f(x_i)$ for each sample x_i . In the proof of Theorem 2, we argue 270 that when working with true f, and the entire marginal q(x), we can choose a discretization satisfying 271 272 the anchor word assumption by identifying point masses in the distribution of f(x) and filtering to 273 include those of at least ϵ size. In the practical setting, because we have only a finite set of data points and a noisy f, we use clustering to approximately find point masses. We choose $m \ge k$ and recover 274 m clusters with any standard clustering procedure (e.g. K-means). If the noise in f is sufficiently 275 small and the clustering sufficiently granular, we intuit that our m discovered clusters include k pure 276 clusters, each of which only contains data points from a different anchor subdomain which are tightly 277 arranged around the true $f(A_y)$ for the corresponding label y. This clustering is superior to a naive 278 clustering on the input space because close proximity in this space indicates similarity in q(d|x). 279

Let us denote the learned clustering function as c, where c(x) is the cluster assigned to a datapoint x. We now leverage the cluster id $c(x_i)$ of each sample x_i to discretize sample into a finite discrete

Algorithm 1 DDFA Training

- **input** $k \ge 1, r \ge k, \{(x_i, d_i)\}_{i \in [n]} \sim q(x, d), \text{ A class of functions } \mathcal{F} \text{ from } \mathbb{R}^p \to \mathbb{R}^r$ 1: Split into train set T and validation set V
- 2: Train $\hat{f} \in \mathcal{F}$ to minimize cross entropy loss for predicting d|x on T with early stopping on V
- 3: Push all $\{x_i\}_{i \in [n]}$ through \widehat{f}
- 4: Train clustering algorithm on the n points $\{\hat{f}(x_i)\}_{i \in [n]}$, obtain m clusters.
- 5: $c(x_i) \leftarrow \text{Cluster id of } \widehat{f}(x_i)$ 6: $\widehat{q}(c(X) = a | D = b) \leftarrow \frac{\sum_{i \in [n]} \mathbb{I}[c(x_i) = a, d_i = b]}{\sum_{j \in [n]} \mathbb{I}[d_j = b]}$ 7: Populate $\widehat{\mathbf{Q}}_{c(X)|D}$ as $[\widehat{\mathbf{Q}}_{c(X)|D}]_{a,b} \leftarrow \widehat{q}(c(X) = a|D = b)$ 8: $\hat{\mathbf{Q}}_{c(X)|D}, \hat{\mathbf{Q}}_{Y|D} \leftarrow \text{NMF}(\hat{\mathbf{Q}}_{c(X)|D})$ output $\widehat{\mathbf{Q}}_{Y|D}, \widehat{f}$

Algorithm 2 DDFA Prediction

input $\widehat{\mathbf{Q}}_{Y|D}, \widehat{f}, (x', d') \sim q(x, d)$ 1: Populate $\hat{\mathbf{Q}}_{D|Y}$ as $[\hat{\mathbf{Q}}_{D|Y}]_{d,y} \leftarrow \frac{[\hat{\mathbf{Q}}_{Y|D}]_{y,d}}{\sum_{d''=1}^{d''=r} [\hat{\mathbf{Q}}_{Y|D}]_{y,d''}}$ 2: Assign $\hat{q}(y|X = x') \leftarrow \left[\left(\hat{\mathbf{Q}}_{D|Y} \right)^{\dagger} \hat{f}(x') \right]_{z}$ 3: Assign $\hat{q}(y|X = x', D = d') \leftarrow \frac{[\hat{\mathbf{Q}}_{D|Y}]_{d',y}\hat{q}(y|X = x')}{\sum\limits_{y'' \in [k]} [\hat{\mathbf{Q}}_{D|Y}]_{d',y''}\hat{q}(y''|X = x')}$ 4: $y_{\text{pred}} \leftarrow \arg \max_{y \in [k]} \widehat{q}(y|X = x', D)$ **output** : $\hat{q}(y|X = x', D = d') = \hat{p}_{d'}(y|x'), \hat{q}(y|X = x'), y_{\text{pred}}$

space [m]. Combining cluster id with the domain source d_i for each sample, we estimate $\mathbf{Q}_{c(X)|D}$ 282 by simply computing, for each domain, the fraction of its samples assigned to each cluster. 283

Factorize We apply an NMF algorithm to $\hat{\mathbf{Q}}_{c(X)|D}$ to obtain our estimates of $\hat{\mathbf{Q}}_{c(X)|Y}$ and $\hat{\mathbf{Q}}_{Y|D}$. 284

Adjust We begin Algorithm 2 by considering a test point (x', d'). To make a prediction, if we had 285 access to oracle f and true $\mathbf{Q}_{Y|D}$, we could precisely compute q(y|x') (Lemma 1). However, in place of these true quantities, we plug in the estimates \hat{f} and $\hat{\mathbf{Q}}_{Y|D}$. Since our estimates contain noise, the 286 287 estimate $\hat{q}(y|x')$ is found by left-multiplying $\hat{f}(x')$ with the pseudo-inverse of $\hat{\mathbf{Q}}_{D|Y}$, as opposed 288 to solving a sufficient system of equations. As our estimates \widehat{f} and $\widehat{\mathbf{Q}}_{D|Y}$ approach the true values, 289 the projection of $\hat{f}(x')$ into the column space of $\hat{Q}_{D|Y}$ tends to $\hat{f}(x')$ itself, so the pseudo-inverse 290 approaches the true solution. Now we can use the constructive procedure introduced in the proof of 291 Lemma 2 to compute the plug-in estimate $\hat{q}(y|x', d') = \hat{p}_{d'}(y|x')$. 292

Experiments 6 293

Baselines We select the unsupervised classification method SCAN as a state-of-the-art baseline [55]. 294 SCAN pretrains a ResNet [29] backbone using SimCLR [11] and MoCo [30] setups (pretext tasks). 295 SCAN then trains a clustering head to minimize the SCAN loss (refer [55] for more details)¹. We 296 make sure to evaluate SCAN on the same potentially class-imbalanced test subset we create for each 297 experiment. Since SCAN is fit on a superset of the data DDFA sees, we believe this gives a slight 298 data advantage to the SCAN baseline (although we acknowledge that the class balance for SCAN 299 training is also potentially different from its evaluation class balance). To evaluate SCAN, we use 300 the public pretrained weights available for CIFAR-10, CIFAR-20, and ImageNet-50. We also train 301 SCAN ourselves on the train and validation portions of the FieldGuide2 and FieldGuide28 datasets 302 with a ResNet18 backbone and SimCLR pretext task. We replicate the hyperparameters used for 303 CIFAR training. 304

¹SCAN code: https://github.com/wvangansbeke/Unsupervised-Classification

Datasets First we examine standard multiclass image datasets CIFAR-10, CIFAR-20 [36], and ImageNet-50 [17] containing images from 10, 20, and 50 classes respectively. Images in these datasets typically focus on a single large object which dominates the center of the frame, so unsupervised classification methods which respond strongly to similarity in visual space are well-suited to recover true classes up to permutation. These datasets are often believed to be separable (i.e., single true label applies to each image), so every example falls in an anchor subdomain (satisfying A.4).

Motivated by the application of LLS problem, we consider the FieldGuide dataset 2 , which contains 311 images of moths and butterflies. The true classes in this dataset are species, but each class contains 312 images taken in immature (caterpillar) and adult stages of life. Based on the intuition that butterflies 313 from a given species look more like butterflies from other species than caterpillars from their own 314 species, we hypothesize that unsupervised classification will learn incorrect class boundaries which 315 distinguish caterpillars from butterflies, as opposed to recovering the true class boundaries. Due 316 to high visual similarity between members of different classes, this dataset may indeed have slight 317 overlap between classes. However, we hypothesize that anchor subdomain still holds, i.e., there 318 exist some images from each class that could only come from that class. Additionally, if we have 319 access to data from multiple domains, it is natural to assume that within each domain the relative 320 distribution of caterpillar to adult stages of each species stay relatively constant as compared to 321 prevalence of different species. We create two subsets of this dataset: FieldGuide2, with two species, 322 and FieldGuide28, with 28 species. 323

LLS Setup The full sampling procedure for semisynthetic experiments is described in App. D. Roughly, we sample $p_d(y)$ from a symmetric Dirichlet distribution with concentration α/k , and enforce maximum condition number κ on $\mathbf{Q}_{\mathbf{Y}|\mathbf{D}}$. Small α and small κ encourages sparsity in $\mathbf{Q}_{\mathbf{Y}|\mathbf{D}}$, so each label tends to only appear in a few domains. Larger parameters encourages $p_d(y)$ to tend toward uniform. We draw from test, train, and valid datasets without replacement to match these distributions, but discard some examples due to class imbalance.

Training and Evaluation The algorithm uses train and validation data consisting of pairs of images 330 and domain indices. We train ResNet50 [29] (with added dropout) on images x_i with domain indices 331 d_i as the label, choose best iteration by valid loss, pass all training and validation data through \hat{f} , and 332 cluster pushforward predictions $\hat{f}(x_i)$ into $m \ge k$ clusters with Faiss K-Means [35]. We compute the 333 $\hat{\mathbf{Q}}_{c(X)|D}$ matrix and run NMF to obtain $\hat{\mathbf{Q}}_{c(X)|Y}, \hat{\mathbf{Q}}_{Y|D}$. To make columns sum to 1, we normalize 334 columns of $\widehat{\mathbf{Q}}_{c(X)|Y}$, multiply each column's normalization coefficient over the corresponding row of 335 $\widehat{\mathbf{Q}}_{Y|D}$ (to preserve correctness of the decomposition), and then normalize columns of $\widehat{\mathbf{Q}}_{Y|D}$. Some 336 NMF algorithms only output solutions satisfying the anchor word property [3, 37, 27]. We found the 337 strict requirement of an exact anchor word solution to lead to low noise tolerance. We therefore use 338 the Sklearn implementation of standard NMF [13, 53, 46]. 339

We instantiate the domain discriminator as ResNet18, and preseed its backbone with SCAN [55] pre-trained weights or [55] contrastive pre-text weights. We denote these models DDFA (SI) and DDFA (SPI) respectively. We predict class labels with Algorithm 2. With the Hungarian algorithm, implemented in [14, 56], we compute the highest true accuracy among any permutation of these labels (denoted "Test acc"). With the same permutation, we reorder rows of $\hat{P}_{Y|D}$, then compute the average absolute difference between corresponding entries of $\hat{Q}_{Y|D}$ and $Q_{Y|D}$ (denoted " $Q_{Y|D}$ err").

Results On CIFAR-10, we observe that DDFA alone is incapable of matching highly competitive 346 state-of-the-art baseline SCAN performance-however, in suitably sparse problem settings (small 347 α), it comes substantially close, indicating good recovery of true classes. Due to space constraints, 348 we include CIFAR-10 results in App. E. DDFA (SI) combines SCAN's strong pretrain with domain 349 350 discrimination fine-tuning to outperform SCAN in the easiest, sparsest setting and *certain* denser settings. On CIFAR-20, baseline SCAN is much less competitive, so our DDFA(SI) dominates 351 baseline SCAN in all settings except the densest (Table 1). These results demonstrate how adding 352 domain information can dramatically boost unsupervised baselines. 353

On FieldGuide-2, DDFA (SPI) outperforms SCAN baselines across all problem settings and domain counts (Table 2); in sparser settings, the accuracy gap is 20-30%. In this dataset, SCAN performs only slightly above chance, reflecting perhaps a total misalignment of learned class distinctions with true species boundaries. We do not believe that SCAN is too weak to effectively detect image groupings

²FieldGuide: https://sites.google.com/view/fgvc6/competitions/butterflies-moths-2019

Table 1: *Results on CIFAR-20.* With DDFA (RI) we refer to DDFA with randomly initialized backbone. With DDFA (SI) we refer to DDFA's backbone initialized with SCAN. Note that in DDFA (SI), we do not leverage SCAN for clustering. α is the Dirichlet parameter used for generating label marginals in each domain, κ is the maximum allowed condition number of the generated $\mathbf{Q}_{Y|D}$ matrix.

r	Approaches	$lpha: 0.5, \; \kappa: 8$		$\alpha:3, \kappa:12$		$\alpha:10,\;\kappa:20$	
-		Test acc	$\mathbf{Q}_{Y D}$ err	Test acc	$\mathbf{Q}_{Y D}$ err	Test acc	$\mathbf{Q}_{Y D}$ err
20	SCAN DDFA (RI) DDFA (SI)	0.4241 0.4872 0.7507	0.0452 0.0248	0.4313 0.3047 0.503	0.043 0.0284	0.4362 0.1036 0.3252	0.0727 0.0348
25	SCAN DDFA (RI) DDFA (SI)	0.437 0.4678 0.8364	0.0505 0.0201	0.4631 0.3157 0.7403	0.0483 0.0195	0.4434 0.0788 0.4819	0.0802 0.0322
30	SCAN DDFA (RI) DDFA (SI)	0.4387 0.5406 0.8158	0.0453 0.0219	0.4514 0.287 0.7247	0.05 0.0222	0.4214 0.091 0.5473	0.0789 0.0258

Table 2: *Results on FieldGuide*. With DDFA (SPI) we refer to DDFA's backbone initialized with pretext training done adopted by SCAN. Note that only the backbone of DDFA is initialized with SCAN-pretext weights and not the final layers. α is the Dirichlet parameter used for generating label marginals in each domain, κ is the maximum allowed condition number of the generated $\mathbf{Q}_{Y|D}$ matrix.

FieldGuide-2			FieldGuide-28					
Approaches	r	$\alpha:3,\ \kappa:5$		r	$\alpha: 0.5, \ \kappa: 12$		$\alpha:3,\ \kappa:20$	
pp:outlies	-	Test acc	$\mathbf{Q}_{Y D}$ err	-	Test acc	$\mathbf{Q}_{Y D}$ err	Test acc	$\mathbf{Q}_{Y D}$ err
SCAN DDFA (SPI)	2	0.6000 0.8576	- 0.1211	28	0.3634 0.6002	0.0309	0.3042 0.3062	0.0268
SCAN DDFA (SPI)	3	0.6123 0.8478	- 0.0351	37	0.3253 0.7337	0.0263	0.2709 0.4987	0.0327
SCAN DDFA (SPI)	5	0.5892 0.6827	- 0.3395	42	0.358 0.6068	0.0377	0.2956 0.4407	0.0341
SCAN DDFA (SPI)	10	0.5865 0.8094	- 0.2799	47	0.3897 0.5624	0.0414	0.3187 0.4573	0.0341

on this data; instead we acknowledge that the domain information available to DDFA (SPI) (and not to SCAN) is informative for ensuring recovery of the true class distinction between species as opposed to the visually striking distinction between adult and immature life stages. Results from more domains are available in App. E. We also show the results on FieldGuide-28 over a range of domains in (Table 2). Our method outperforms SCAN on all settings with the highest observed accuracy difference ranging up to 30%.

364 7 Conclusion

Our theoretical results demonstrate that under LLS, we can leverage shifts among previously seen 365 domains to recover labels in a purely unsupervised manner, and our practical instantiation of the 366 DDFA framework demonstrates both (i) the practical efficacy of our approach; and (ii) that generic 367 unsupervised methods can play a key role both in clustering discriminator outputs, and providing 368 weights for initializing the discriminator. We believe that this work is just the first step in a new 369 370 direction for leveraging structural assumptions together with distribution shift to perform unsupervised learning. Within the LLS setup, several components of the DDFA framework warrant further 371 investigation: (i) the deep domain discriminator can be enhanced in myriad ways; (ii) for clustering 372 discriminator outputs, we might develop methods specially tailored to our setting; (iii) clustering 373 might be replaced altogether with geometrically informed methods that directly identify the corners 374 of the polytope; (iv) the theory of LLS can be extended beyond identification to provide statistical 375 results that might hold when q(d|y) can only be noisily estimated, and when only finite samples are 376 available for the induced topic modeling problem. 377

378 **References**

- [1] Amr Alexandari, Anshul Kundaje, and Avanti Shrikumar. Adapting to label shift with biascorrected calibration. In *International Conference on Machine Learning (ICML)*, 2019.
- [2] Martin Arjovsky, Léon Bottou, Ishaan Gulrajani, and David Lopez-Paz. Invariant risk minimization. *arXiv preprint arXiv:1907.02893*, 2019.
- [3] Sanjeev Arora, Rong Ge, Ravindran Kannan, and Ankur Moitra. Computing a nonnegative matrix factorization – provably. In *Symposium on Theory of Computing (STOC)*, 2012. doi: 10.1145/2213977.2213994. URL https://doi.org/10.1145/2213977.2213994.
- [4] Sanjeev Arora, Rong Ge, and Ankur Moitra. Learning topic models–going beyond svd. In
 Foundations of Computer Science (FOCS). IEEE, 2012.
- [5] Kamyar Azizzadenesheli, Anqi Liu, Fanny Yang, and Animashree Anandkumar. Regularized
 learning for domain adaptation under label shifts. In *International Conference on Learning Representations (ICLR)*, 2019.
- [6] Jessa Bekker and Jesse Davis. Estimating the class prior in positive and unlabeled data through
 decision tree induction. In Association for the Advancement of Artificial Intelligence (AAAI), vol ume 32, 2018. URL https://ojs.aaai.org/index.php/AAAI/article/view/11715.
- [7] Jessa Bekker and Jesse Davis. Learning from positive and unlabeled data: a survey. *Machine Learning*, 109(4):719–760, apr 2020. doi: 10.1007/s10994-020-05877-5. URL https://doi.org/10.1007%2Fs10994-020-05877-5.
- [8] David M Blei, Andrew Y Ng, and Michael I Jordan. Latent dirichlet allocation. *Journal of Machine Learning research*, 3(Jan):993–1022, 2003.
- [9] Mathilde Caron, Piotr Bojanowski, Armand Joulin, and Matthijs Douze. Deep clustering
 for unsupervised learning of visual features. In *Proceedings of the European conference on computer vision (ECCV)*, pages 132–149, 2018.
- [10] Mathilde Caron, Piotr Bojanowski, Julien Mairal, and Armand Joulin. Unsupervised pretraining of image features on non-curated data, 2019. URL https://www.cv-foundation.
 org/openaccess/content_iccv_2015/papers/Doersch_Unsupervised_Visual_
 Representation_ICCV_2015_paper.pdf.
- [11] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework
 for contrastive learning of visual representations. In *International conference on machine learning (ICML)*, pages 1597–1607. PMLR, 2020.
- [12] Yinyin Chen, Shishuang He, Yun Yang, and Feng Liang. Learning topic models: Identifiability
 and finite-sample analysis. *arXiv preprint arXiv:2110.04232*, 2021.
- [13] Andrzej Cichocki and Anh-Huy Phan. Fast local algorithms for large scale nonnegative matrix
 and tensor factorizations. *IEICE Transactions*, 92-A:708–721, 03 2009. doi: 10.1587/transfun.
 E92.A.708.
- [14] David F Crouse. On implementing 2d rectangular assignment algorithms. *IEEE Transactions* on Aerospace and Electronic Systems, 52(4):1679–1696, 2016.
- [15] Luke N. Darlow, Elliot J. Crowley, Antreas Antoniou, and Amos J. Storkey. Cinic-10 is not
 imagenet or cifar-10, 2018. URL https://arxiv.org/abs/1810.03505.
- [16] Thiago de Paulo Faleiros and Alneu de Andrade Lopes. On the equivalence between algorithms
 for non-negative matrix factorization and latent dirichlet allocation. In *European Symposium on Artificial Neural Networks (ESANN)*, 2016.
- [17] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale
 hierarchical image database. In *IEEE Conference on Computer Vision and Pattern Recognition* (CVPR). Ieee, 2009.

- [18] Chris Ding, Tao Li, and Wei Peng. On the equivalence between non-negative matrix factorization
 and probabilistic latent semantic indexing. *Computational Statistics & Data Analysis*, 52(8):
 3913–3927, 2008.
- [19] Carl Doersch, Abhinav Gupta, and Alexei A. Efros. Unsupervised visual representation learning by context prediction. In *International Conference on Computer Vision (ICCV)*, 2015. URL https://www.cv-foundation.org/openaccess/content_iccv_2015/
 papers/Doersch_Unsupervised_Visual_Representation_ICCV_2015_paper.pdf.
- [20] David Donoho and Victoria Stodden. When does non-negative matrix factorization give a correct
 decomposition into parts? *Advances in Neural Information Processing Systems (NeurIPS)*, 16,
 2003.
- [21] Marthinus C. du Plessis, Gang Niu, and Masashi Sugiyama. Class-prior estimation for learning
 from positive and unlabeled data. *Machine Learning*, 106(4):463–492, nov 2016. doi: 10.1007/
 s10994-016-5604-6. URL https://doi.org/10.1007%2Fs10994-016-5604-6.
- [22] Charles Elkan and Keith Noto. Learning classifiers from only positive and unlabeled data. In
 SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD), pages 213–220, 2008.
- [23] Saurabh Garg, Yifan Wu, Sivaraman Balakrishnan, and Zachary Lipton. A unified view of label
 shift estimation. In Advances in Neural Information Processing Systems (NeurIPS), 2020. URL
 https://arxiv.org/abs/2003.07554.
- [24] Saurabh Garg, Yifan Wu, Alex Smola, Sivaraman Balakrishnan, and Zachary Lipton. Mixture
 proportion estimation and PU learning: A modern approach. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2021. URL https://arxiv.org/abs/2111.00980.
- [25] Eric Gaussier and Cyril Goutte. Relation between plsa and nmf and implications. In *Conference on Research and Development in Information Retrieval (SIGIR)*, 2005.
- [26] Spyros Gidaris, Praveer Singh, and Nikos Komodakis. Unsupervised representation learning by
 predicting image rotations, 2018. URL https://arxiv.org/abs/1803.07728.
- [27] Nicolas Gillis and Stephen A. Vavasis. Fast and robust recursive algorithms for separable non negative matrix factorization. *IEEE Transactions on Pattern Analysis and Machine Intelligence*,
 36(4):698–714, apr 2014. doi: 10.1109/tpami.2013.226. URL https://doi.org/10.1137%
 2F130946782.
- [28] Mark Girolami and Ata Kabán. On an equivalence between plsi and lda. In *Conference on Research and Development in Information Retrieval (SIGIR)*, 2003.
- [29] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
 recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recogni- tion (CVPR)*, 2016.
- [30] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for
 unsupervised visual representation learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9729–9738, 2020.
- 462 [31] Thomas Hofmann. Probabilistic latent semantic indexing. In *Conference on Research and* 463 *Development in Information Retrieval (SIGIR)*, 1999.
- Kejun Huang, Xiao Fu, and Nikolaos D Sidiropoulos. Anchor-free correlated topic modeling:
 Identifiability and algorithm. In *Advances in Neural Information Processing Systems (NeurIPS)*,
 2016.
- [33] Aapo Hyvärinen and Erkki Oja. Independent component analysis: algorithms and applications.
 Neural networks, 13(4-5):411–430, 2000.
- [34] Dmitry Ivanov. DEDPUL: Difference-of-estimated-densities-based positive-unlabeled learning.
 arXiv preprint arXiv:1902.06965, 2019.

- [35] Jeff Johnson, Matthijs Douze, and Hervé Jégou. Billion-scale similarity search with GPUs.
 IEEE Transactions on Big Data, 7(3):535–547, 2019.
- [36] Alex Krizhevsky and Geoffrey Hinton. Learning Multiple Layers of Features from Tiny Images.
 Technical report, Citeseer, 2009.
- [37] Abhishek Kumar, Vikas Sindhwani, and Prabhanjan Kambadur. Fast conical hull algorithms
 for near-separable non-negative matrix factorization. In *International Conference on Machine Learning*, pages 231–239. PMLR, 2013.
- [38] Zachary Lipton, Yu-Xiang Wang, and Alexander Smola. Detecting and correcting for label shift
 with black box predictors. In *International Conference on Machine Learning (ICML)*. PMLR,
 2018.
- [39] Tongliang Liu and Dacheng Tao. Classification with noisy labels by importance reweighting.
 IEEE Transactions on Pattern Analysis and Machine Intelligence (PAMI), 38(3):447–461, mar
 2016. doi: 10.1109/tpami.2015.2456899. URL https://doi.org/10.1109%2Ftpami.2015.
 2456899.
- [40] James MacQueen et al. Some methods for classification and analysis of multivariate observations.
 In *Proceedings of the fifth Berkeley symposium on mathematical statistics and probability*, volume 1, pages 281–297. Oakland, CA, USA, 1967.
- [41] Krikamol Muandet, David Balduzzi, and Bernhard Schölkopf. Domain generalization via invariant feature representation. In *International Conference on Machine Learning (ICML)*.
 PMLR, 2013.
- [42] Mehdi Noroozi and Paolo Favaro. Unsupervised learning of visual representations by solving
 jigsaw puzzles. In *European conference on computer vision*, pages 69–84. Springer, 2016.
- [43] Pentti Paatero and Unto Tapper. Positive matrix factorization: A non-negative factor model with
 optimal utilization of error estimates of data values. *Environmetrics*, 5(2):111–126, 1994.
- [44] Christos H Papadimitriou, Prabhakar Raghavan, Hisao Tamaki, and Santosh Vempala. Latent
 semantic indexing: A probabilistic analysis. *Journal of Computer and System Sciences*, 61(2):
 217–235, 2000.
- [45] Sungwon Park, Sungwon Han, Sundong Kim, Danu Kim, Sungkyu Park, Seunghoon Hong,
 and Meeyoung Cha. Improving unsupervised image clustering with robust learning. *CoRR*,
 abs/2012.11150, 2020. URL https://arxiv.org/abs/2012.11150.
- [46] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel,
 P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher,
 M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research (JMLR)*, 12, 2011.
- [47] Kedar S Prabhudesai, Boyla O Mainsah, Leslie M Collins, and Chandra S Throckmorton.
 Augmented latent dirichlet allocation (lda) topic model with gaussian mixture topics. In 2018
 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pages 2451–2455. IEEE, 2018.
- [48] Douglas A Reynolds. Gaussian mixture models. *Encyclopedia of biometrics*, 741(659-663),
 2009.
- [49] Marco Saerens, Patrice Latinne, and Christine Decaestecker. Adjusting the Outputs of a Classifier to New a Priori Probabilities: A Simple Procedure. *Neural Computation*, 2002.
- [50] Clayton Scott. A Rate of Convergence for Mixture Proportion Estimation, with Application to
 Learning from Noisy Labels. In *Artificial Intelligence and Statistics (AISTATS)*, 2015. URL
 https://proceedings.mlr.press/v38/scott15.html.
- [51] D Seung and L Lee. Algorithms for non-negative matrix factorization. In *Advances in neural information processing systems (NeurIPS)*, 2001.

- [52] Amos Storkey. When training and test sets are different: characterizing learning transfer.
 Dataset shift in machine learning, 30:3–28, 2009.
- [53] Vincent YF Tan and Cédric Févotte. Automatic relevance determination in nonnegative matrix
 factorization with the/spl beta/-divergence. *IEEE Transactions on Pattern Analysis and Machine Intelligence (PAMI)*, 35(7), 2012.
- [54] Dongping Tian. Research on plsa model based semantic image analysis: A systematic review. J.
 Inf. Hiding Multim. Signal Process., 9(5):1099–1113, 2018.
- [55] Wouter Van Gansbeke, Simon Vandenhende, Stamatios Georgoulis, Marc Proesmans, and Luc
 Van Gool. Scan: Learning to classify images without labels, 2020. URL https://arxiv.
 org/abs/2005.12320.
- [56] Pauli Virtanen, Ralf Gommers, Travis E. Oliphant, Matt Haberland, Tyler Reddy, David 528 Cournapeau, Evgeni Burovski, Pearu Peterson, Warren Weckesser, Jonathan Bright, Stéfan J. 529 van der Walt, Matthew Brett, Joshua Wilson, K. Jarrod Millman, Nikolay Mayorov, Andrew 530 R. J. Nelson, Eric Jones, Robert Kern, Eric Larson, C J Carey, İlhan Polat, Yu Feng, Eric W. 531 Moore, Jake VanderPlas, Denis Laxalde, Josef Perktold, Robert Cimrman, Ian Henriksen, E. A. 532 Quintero, Charles R. Harris, Anne M. Archibald, Antônio H. Ribeiro, Fabian Pedregosa, Paul 533 van Mulbregt, and SciPy 1.0 Contributors. SciPy 1.0: Fundamental Algorithms for Scientific 534 Computing in Python. Nature Methods, 17:261–272, 2020. doi: 10.1038/s41592-019-0686-2. 535
- [57] Kun Zhang, Bernhard Schölkopf, Krikamol Muandet, and Zhikun Wang. Domain adaptation
 under target and conditional shift. In *International Conference on Machine Learning (ICML)*.
 PMLR, 2013.

539	1.	For a	all authors
540 541		(a)	Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope? [Yes]
542		(b)	Did you describe the limitations of your work? [Yes]
543 544 545 546 547		(c)	Did you discuss any potential negative societal impacts of your work? [N/A] We believe that this work, which proposes a novel unsupervised learning problem does not present a significant societal concern. While this could potentially guide practitioners to improve classification, we do not believe that it will fundamentally impact how machine learning is used in a way that could conceivably be socially salient.
548 549		(d)	Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]
550	2.	If yo	bu are including theoretical results
551 552 553		(a) (b)	Did you state the full set of assumptions of all theoretical results? [Yes] See Sec. 3 Did you include complete proofs of all theoretical results? [Yes] See Sec. 4 and appendices
554	3.	If yo	ou ran experiments
555 556 557 558		(a)	Did you include the code, data, and instructions needed to reproduce the main experi- mental results (either in the supplemental material or as a URL)? [Yes] We include all the necessary details to replicate our experiments in appendices. We plan to release code with an updated version of the draft.
559 560 561		(b)	Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes] Yes, we describe crucial details in Sec. 6 and defer precise details to appendices.
562 563		(c)	Did you report error bars (e.g., with respect to the random seed after running experi- ments multiple times)? [Yes] We include results with multiple seeds in appendices
564 565		(d)	GPUs, internal cluster, or cloud provider)? [Yes] Refer to experimental setup in App. D.
566	4.	If yo	bu are using existing assets (e.g., code, data, models) or curating/releasing new assets
567 568		(a)	If your work uses existing assets, did you cite the creators? [Yes] Refer to experimental setup in App. D.
569		(b)	Did you mention the license of the assets? [Yes] Refer to experimental setup in App. D.
570 571		(c)	Did you include any new assets either in the supplemental material or as a URL? [Yes] Refer to experimental setup in App. D.
572 573		(d)	Did you discuss whether and how consent was obtained from people whose data you're using/curating? [N/A]
574 575		(e)	Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [N/A]
576	5.	If yo	ou used crowdsourcing or conducted research with human subjects
577 578		(a)	Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]
579 580		(b)	Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
581 582		(c)	Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? $[N/A]$

583 A Proofs of Lemmas

In this section, we present several new lemmas which are required to prove Theorem 2, and provide proofs. We also provide proofs for Lemmas 1, 2, and 3, all presented in Section 4.

Fix $y \in \mathcal{Y}$. Whenever assumption A.4 is satisfied, we define a "maximal" A_y as the union of every A'_y which is an ϵ -anchor sub-domain for label y. This maximal A_y , an ϵ -anchor sub-domain for label y in itself, must contain every point $x \in \mathcal{X}$ such that q(x) > 0, q(x|y) > 0, and for all $y'' \in \mathcal{Y} \setminus \{y\}$, q(x|y'') = 0. We assume throughout this section, without restatement, that A_y refers to this maximal A_y .

Lemma 4. Let distribution Q over random variables X, Y, D satisfy A.1–A.4. Then for all $y \in \mathcal{Y}$, q(y) > 0. That is, all labels have nonzero probability under Q.

593 Proof of Lemma 4. Proof by contradiction. Let $y \in \mathcal{Y}$. Assume q(y) = 0.

594 We can write
$$q(y) = \sum_{d \in \mathcal{R}} q(d)q(y|D = d)$$
.

595 Then
$$q(y) = \sum_{d \in \mathcal{R}} \gamma_y q(y|D = d) = \sum_{d \in \mathcal{R}} \frac{1}{r} q(y|D = d)$$

596 Then
$$q(y) = \frac{1}{r} \sum_{d \in \mathcal{R}} q(y|D = d)$$

- Since $q(y|D = d) \ge 0$ for all $d \in \mathcal{R}$, we see that $q(y) = 0 \implies q(y|D = d) = 0$ for all $d \in \mathcal{R}$.
- Then $[\mathbf{Q}_{Y|D}]_{y,d} = 0$ for all $d \in \mathcal{R}$. Then there is a row (row d) in the matrix $\mathbf{Q}_{Y|D}$ in which every entry is 0, so $\mathbf{Q}_{Y|D}$ cannot be rank k. This violates assumption A.2.

- 600 Then by contradiction we have shown q(y) > 0.
- **Lemma 5.** Let distribution Q over random variables X, Y, D satisfy Assumptions A.1–A.4.

602 Let
$$x \in \mathcal{X}$$
 such that $q(x) > 0$.

Then if $x \in A_y$ for some $y \in \mathcal{Y}$, we have that q(y|X = x) = 1, and for all $y' \in \mathcal{Y} \setminus \{y\}$, q(y'|X = x) = 0.

The converse is also true: if q(y|X = x) = 1 for some $y \in \mathcal{Y}$ and $q(y'|X = x) = 0 \quad \forall y' \in \mathcal{Y} \setminus \{y\}$, then we know that $x \in A_y$.

607 Proof of Lemma 5. We prove directions one at a time.

Forward direction.

609 We seek to show that:

608

610
$$x \in A_y \implies [q(y|X=x)=1, \forall y' \in \mathcal{Y} \setminus \{y\}, q(y'|X=x)=0].$$

Assume
$$x \in A_y$$
. We recall that we earlier assumed $q(x) > 0$.

612
$$q(x) = \sum_{y'' \in \mathcal{Y}} q(y'')q(x|Y = y'') = q(y)q(x|Y = y) + \sum_{y' \in \mathcal{Y} \setminus \{y\}} q(y')q(x|Y = y')$$

613 Now
$$q(x) = q(y)q(x|Y = y) + \sum_{y' \in \mathcal{Y} \setminus \{y\}} q(y')(0) = q(y)q(x|Y = y)$$

Because
$$q(x|y) > 0$$
 (by A.4) and $q(y) > 0$ (by Lemma 4), we know that $q(x) = q(y)q(x|Y = y) > 0$.

616 Then
$$q(y|X = x) = \frac{q(y)q(x|Y = y)}{q(x)} = \frac{q(x)}{q(x)} = 1$$
 (Bayes' rule).

By the properties of a probability distribution, $1 = q(y|X = x) + \sum_{y' \in \mathcal{Y} \setminus \{y\}} q(y'|X = x)$.

Then because q(y|X = x) = 1, we know that:

619
$$1 = 1 + \sum_{y' \in \mathcal{Y} \setminus \{y\}} q(y'|X = x)$$

Then for all $y' \in \mathcal{Y} \setminus \{y\}$, it must be that q(y'|X = x) = 0620 Then we have shown q(y|X = x) = 1, and for all $y' \in \mathcal{Y} \setminus \{y\}$, q(y'|X = x) = 0. 621 • Converse. 622 We seek to show that: 623 $[q(y|X = x) = 1, \forall y' \in \mathcal{Y} \setminus \{y\}, q(y'|X = x) = 0] \implies x \in A_y.$ 624 Assume q(y|X = x) = 1 and for all $y' \in \mathcal{Y} \setminus \{y\}, q(y'|X = x) = 0$. We recall that we 625 earlier assumed q(x) > 0. 626 q(y) > 0 by Lemma 4. 627 Then $q(x|Y = y) = \frac{q(y|X = x)q(x)}{q(y)} = \frac{(1)q(x)}{q(y)} > 0.$ 628 Let $y' \in \mathcal{Y} \setminus \{y\}$. Then $q(x|Y = y') = \frac{q(y'|X = x)q(x)}{q(y')} = \frac{(0)q(x)}{q(y')} = 0$. 629 Then because q(x|Y = y) > 0 and $\forall y' \in \mathcal{Y} \setminus \{y\}, q(x|Y = y) = 0$, we see that $x \in A_y$. 630 631

Lemma 6. Let random variables X, Y, D and distribution Q satisfy Assumptions A.1–A.4. Then, the matrix $\mathbf{Q}_{D|Y}$ defined as an $r \times k$ matrix whose elements $[\mathbf{Q}_{D|Y}]_{i,j} = Q(D = i|Y = i)$

j), and each column is a conditional distribution over the domains given a label , has linearly independent columns.

- 636 Furthemore, $\mathbf{Q}_{D|Y}$ can be computed directly from only $\mathbf{Q}_{Y|D}$.
- Proof of Lemma 6. Let random variables X, Y, D and distribution Q satisfy Assumptions A.1–A.4.

Each
$$[\mathbf{Q}_{D|Y}]_{d,y} = q(d|Y=y) = \frac{q(y|D=d)q(d)}{q(y)} = \frac{q(y|D=d)\gamma_d}{q(y)} = \frac{q(y|D=d)}{rq(y)}$$

Since each column of $\mathbf{Q}_{D|Y}$ is a probability distribution that sums to 1, and rq(y) is constant down each column y, we can obtain $\mathbf{Q}_{D|Y}$ by simply taking $\mathbf{Q}_{Y|D}^{\top}$, in which each $[\mathbf{Q}_{Y|D}^{\top}]_{d,y} = [\mathbf{Q}_{Y|D}]_{y,d} = q(y|D = d)$, and normalizing the columns so they sum to 1.

The matrix $\mathbf{Q}_{Y|D}$ has linearly independent rows by Assumption A.2. Then $\mathbf{Q}_{Y|D}^{\top}$ has linearly independent columns. Scaling these columns by a nonzero value does not change their linear independence, so the columns of $\mathbf{Q}_{D|Y}$ are also linearly independent.

Then matrix $\mathbf{Q}_{D|Y}$ has linearly independent columns, and can be computed by taking $\mathbf{Q}_{Y|D}^{\top}$ and normalizing its columns.

647

Lemma 7. Let random variables
$$X, Y, D$$
 and distribution Q satisfy Assumptions A.1–A.4

- 649 Let $d \in \mathcal{R}, x \in \mathcal{X}, y \in \mathcal{Y}$.
- 650 Then q(d|X = x, Y = y) = q(d|Y = y).
- 651 Proof of Lemma 7. $q(d|X = x, Y = y) = \frac{q(x|D = d, Y = y)q(D = d|Y = y)}{q(x|Y = y)}$

652
$$= \frac{p_d(x|Y=y)q(d|Y=y)}{q(x|Y=y)}$$

653
$$= \frac{p(x|Y=y)q(d|Y=y)}{q(x|Y=y)}$$

654
$$= \frac{q(x|Y=y)q(d|Y=y)}{q(x|Y=y)}$$

655
$$= q(d|Y=y)$$

656

657 Proof of Lemma 1. Let distribution Q over random variables X, Y, D satisfy Assumptions A.1-A.4.

- Let $x \in \mathcal{X}$ with q(x) > 0, and $y \in \mathcal{Y}$.
- Assume we know $\mathbf{Q}_{Y|D}$ and $[f(x)]_d = q(d|X = x)$.
- 660 With $\mathbf{Q}_{Y|D}$, we know $q(y|D = d) \forall y, d$.
- Also, with the oracle f, we are able to obtain $q(d|X = x) \forall x, d$.
- For all $d \in \mathcal{R}$, we can write that $q(d|X = x) = \sum_{y' \in \mathcal{Y}} q(d|X = x, Y = y')q(y'|X = x) = \sum_{y' \in \mathcal{Y}} q(d|Y = y')q(y'|X = x)$, using Lemma 7.
- ⁶⁶⁴ Define the vector-valued function $g: \mathcal{X} \to \mathbb{R}^k$ such that $[g(x)]_y = q(y|X = x)$ for all $x \in \text{supp}(X)$.
- $\mathbf{Q}_{D|Y}$ is a matrix of shape $r \times k$, with $[\mathbf{Q}_{D|Y}]_{i,j} = Q(D = i|Y = j)$. It can be computed from $\mathbf{Q}_{Y|D}$ and has linearly independent columns—both facts shown in Lemma 6.
- Then $[f(x)]_d = q(d|X = x) = \mathbf{Q}_{D|Y}[d, :]g(x)$, a product between the *d*th row vector of $\mathbf{Q}_{D|Y}$ and the column vector g(x).

669 Then
$$f(x) = \mathbf{Q}_{D|Y}g(x)$$
.

- This system is a linear system with $r \ge k$ equations. Recalling that $\mathbf{Q}_{D|Y}$ has k linearly independent
- columns, we can select any k linearly independent rows of $\mathbf{Q}_{D|Y}$ to solve the equation for the true, unique solution for the unknown vector g(x).
- Another way to describe this is with the pseudo-inverse: $g(x) = (\mathbf{Q}_{D|Y})^{\dagger} f(x)$.
- 674 Then we have $[g(x)]_y = q(y|X = x)$ for all $y \in \mathcal{Y}$.

675

- 676 Proof of Lemma 2. Let distribution Q over random variables X, Y, D satisfy Assumptions A.1-A.4.
- 677 Let $x \in \mathcal{X}, d \in \mathcal{R}$ with q(x, d) > 0, and $y \in \mathcal{Y}$.
- Assume we know matrix $\mathbf{Q}_{Y|D}$ and $q(y'|X = x), \ \forall y' \in \mathcal{Y}$.
- 679 We can compute $\mathbf{Q}_{D|Y}$ from $\mathbf{Q}_{Y|D}$ via Lemma 6.

680
$$q(y|X = x, D = d) = \frac{q(y, x, d)}{q(x, d)} = \frac{q(d|X = x, Y = y)q(y|X = x)q(x)}{q(d|X = x)q(x)}$$

681 Using Lemma 7,
$$q(d|X = x, Y = y) = q(d|Y = y)$$
.

682 Then
$$q(y|X = x, D = d) = \frac{q(d|Y = y)q(y|X = x)q(x)}{q(d|X = x)q(x)} = \frac{q(d|Y = y)q(y|X = x)}{q(d|X = x)}$$

The denominator q(d|X = x) is constant across all values of y, so we can write that $q(y|X = x, D = d) \propto q(d|Y = y)q(y|X = x)$ and normalize to find the probability:

685
$$q(y|X = x, D = d) = \frac{q(d|Y = y)q(y|X = x)}{\sum_{y' \in \mathcal{Y}} q(d|Y = y')q(y'|X = x)}$$

- 686 We know q(d|Y = y) as $[\mathbf{Q}_{D|Y}]_{d,y}$, and every q(d|Y = y'), where $y' \in \mathcal{Y}$, as $[\mathbf{Q}_{D|Y}]_{d,y'}$.
- We also know q(y|X = x) and every q(y'|X = x) where $y' \in \mathcal{Y}$, by the lemma assumptions.
- 688 Then we can compute q(y|X = x, D = d).

Proof of Lemma 3. Let distribution Q over random variables X, Y, D satisfy Assumptions A.1-A.4. Recall $f : \mathbb{R}^p \to \mathbb{R}^r$ is a vector-valued oracle function such that $[f(x)]_d = q(d|X = x)$ for all $x \in \operatorname{supp}_Q(X)$.

Also let us recall that $\mathbf{Q}_{D|Y}$ is defined as an $r \times k$ matrix whose elements $[\mathbf{Q}_{D|Y}]_{i,j} = Q(D = i|Y = j)$, and each column is a conditional distribution over the domains given a label. It has linearly independent columns by Lemma 6.

First recognize that $\forall d \in \mathcal{R}, x \in \mathcal{X} : q(x) > 0, \ [f(x)]_d = q(d|X = x) = \sum_{y'' \in \mathcal{Y}} q(d, y''|X = x)$

697 $= \sum_{y'' \in \mathcal{Y}} q(d|Y = y'', X = x)q(y''|X = x) = \sum_{y'' \in \mathcal{Y}} q(d|Y = y'')q(y''|X = x), \text{ by Lemma 7.}$

Then recognize that we can write $f(x) = \sum_{y'' \in \mathcal{Y}} q(y''|X = x) \mathbf{Q}_{D|Y}[:, y'']$, where $\mathbf{Q}_{D|Y}[:, y'']$ is the

699 y''th column of $\mathbf{Q}_{D|Y}$.

Now we could also rewrite
$$f(x) = \mathbf{Q}_{D|Y} \left[Q(Y=1|X=x) \dots Q(Y=k|X=x) \right]^{\top}$$

- 701 We now begin the bulk of the proof.
- 702 Let $y \in \mathcal{Y}$.
- 703 Let $x \in A_y : q(x) > 0$.

704	Points in same anchor sub-domain map together.
705	Let $x' \in A_y$ such that $q(x') > 0$. We now seek to show that $f(x) = f(x')$.
706 707	Recall that $x, x' \in A_y$. By Lemma 5, $q(y X = x) = q(y X = x') = 1$. Also, $\forall y'' \in \mathcal{Y} \setminus \{y\}, q(y'' X = x) = q(y'' X = x') = 0$.
708	Then for all $y'' \in \mathcal{Y}$, $q(y'' X = x) = q(y'' X = x')$.
709	Therefore, $\forall d \in \mathcal{R}, [f(x)]_d = q(d X = x) = \sum_{y'=0}^{\infty} q(d Y = y'')q(y'' X = x) =$
710	$\sum_{y'' \in \mathcal{Y}} q(d Y = y'')q(y'' X = x') = q(d X = x') = [f(x')]_d.$
711	Then $f(x) = f(x')$.
712 713	• Point outside of the anchor sub-domain does not map with points in the anchor sub- domain.
714	Let $x_0 \notin A_y$ such that $q(x_0) > 0$. We now seek to show that $f(x) \neq f(x_0)$.
715 716	Because $x_0 \notin A_y$ with $q(x_0) > 0$, and because A_y is maximal, then by Lemma 5, it must be that one of the following cases is true:
717 718	- Case 1: $q(y X = x_0) \neq 1$ - Case 2: $q(y' X = x_0) > 0$ for some $y' \in \mathcal{Y} \setminus \{y\}$.
719	In all circumstances, $\exists y'' \in \mathcal{Y} : q(y'' x_0) \neq Q(Y = y'' x).$
720	$[Q(Y = 1 X = x)Q(Y = k X = x)]^{\top} \neq [Q(Y = 1 X = x_0)Q(Y = k X = x_0)]^{\top}.$
721	Because $\mathbf{Q}_{D Y}$ has linearly independent columns (shown in Lemma 6), we now know that
722	$f(x) = \mathbf{Q}_{D Y} [Q(Y = 1 X = x) \dots Q(Y = k X = x)]^{\top}$
723	$\neq \mathbf{Q}_{D Y} \left[Q(Y=1 X=x_0) \dots Q(Y=k X=x_0) \right]^{\top} = f(x_0).$
724	So $f(x) \neq f(x_0)$.
725	

726 **B Proof of Theorem 2**

Fix $y \in \mathcal{Y}$. Whenever assumption A.4 is satisfied, we define a "maximal" A_y as the union of every A'_y which is an ϵ -anchor sub-domain for label y. This maximal A_y , an ϵ -anchor sub-domain for label y in itself, must contain every point $x \in \mathcal{X}$ such that q(x) > 0, q(x|y) > 0, and for all $y'' \in \mathcal{Y} \setminus \{y\}$, q(x|y'') = 0. We assume throughout this section that A_y refers to this maximal A_y .

Proof of Theorem 2. Let distribution Q over random variables X, Y, D satisfy Assumptions A.1-A.4.

Recall $f : \mathcal{X} \to \mathbb{R}^r$ is a vector-valued oracle function such that $[f(x)]_d = q(d|X = x) \forall x \in$ supp_Q(X). It is known because we know the marginal q(x, d).

734 Let $y \in \mathcal{Y}$.

Then by Lemma 3, f sends every $x \in A_y$ (and no other $x \notin A_y$) to the same value. We overload notation to denote this as $f(A_y)$.

737 Then $Q(f(X) = f(A_y)) = Q(X \in A_y) \ge \epsilon$.

Then in the marginal distribution of f(X) with respect to distribution Q, there is a distinct point mass on each $f(A_y)$, with mass at least ϵ .

Because we know the marginal q(x, d), we know the marginal q(x), so we can obtain the distribution of f(X) with respect to distribution Q.

If we analyze the marginal distribution of f(X) with respect to distribution Q, and recover all point masses with mass at least ϵ , we can recover no more than $\mathcal{O}(1/\epsilon)$ such points. We set $m \in \mathbb{Z}^+$ so that the number of points we recovered is m - 1.

We denote a mapping $\psi : \mathbb{R}^r \to [m]$. This mapping sends each value of f(x) corresponding to a point mass with mass at least ϵ to a unique index in $\{1, ..., m-1\}$. It sends any other value in \mathbb{R}^p to m. We note that the ordering of the point masses might have (m-1)! permutations.

- Notice that the point mass on each $f(A_y)$ must be recovered among these m-1 masses.
- 749 Recall that $\forall y \in \mathcal{Y}, [f(x) = f(A_y) \iff x \in A_y].$

Then for each $y \in \mathcal{Y}$, $[\psi(f(X)) = \psi(f(A_y)) \iff X \in A_y]$, because ψ does not send any other value in \mathbb{R}^r besides $f(A_y)$ to $\psi(f(A_y))$.

- For convenience, we now define a mapping $c: \mathcal{X} \to [m]$ such that $c = \psi \circ f$.
- We will also abuse notation here to denote $c(A_y) = \psi(f(A_y)) = \psi(f(x)), \forall x \in A_y$

Then c(X) is a discrete, finite random variable that takes values in [m]. As c is a pushforward function on X, c(X) satisfies the label shift assumption because X does.

- We might now define a matrix $\mathbf{Q}_{c(X)|D}$ in which each entry $[\mathbf{Q}_{c(X)|D}]_{i,d} = Q(c(X) = i|D = d)$.
- Because c(X) satisfies label shift, we know that we can decompose $\mathbf{Q}_{c(X)|D} = \mathbf{Q}_{c(X)|Y}\mathbf{Q}_{Y|D}$.

758
$$Q(c(x) = c(A_y)|Y = y) = Q(X \in A_y|Y = y) > 0.$$

- 759 $Q(c(x) = c(A_y)|Y \neq y) = Q(X = A_y|Y \neq y) = 0.$
- Then for each $y \in \mathcal{Y}$, the row with row index $c(A_y)$ is positive in the *y*th column, and zero everywhere else.
- Restated, for each $y \in \mathcal{Y}$, there is some row with positive entry exactly in *y*th column. This is precisely the anchor word assumption for a discrete, finite random variable.
- We already know that $\mathbf{Q}_{Y|D}$ is full row-rank, so because $\mathbf{Q}_{c(X)|Y}$ satisfies the anchor word assumption, we can identify $\mathbf{Q}_{Y|D}$ up to permutation of rows by Theorem 1.

766

767 C Minimizing Cross-Entropy Loss yields Domain Discriminator

- Let distribution Q over random variables X, Y, D satisfy Assumptions A.1-A.4.
- We here examine the behavior of the cross-entropy loss, in the infinite data case (when we can work
- with expectations over the entire distribution instead of empirical expectations over a finite set ofdatapoints).
- Define the vector-valued function $z : \mathcal{R} \to \mathbb{R}^r$ such that z(d) is a one-hot vector of length r, such that $[z(d)]_i = 1$, iff d = i.
- Then we write the cross-entropy loss with targets as true domains as

775
$$\mathcal{L}_{CE} = \mathbb{E}_{(X,D)\sim Q} \left[-\sum_{i=1}^{i=r} [z(D)]_i \log([f(X)]_i) \right]$$

776
$$\mathcal{L}_{CE} = \mathbb{E}_X \mathbb{E}_{D|X} [-\sum_{i=1}^{i=r} [z(D)]_i \log([f(X)]_i)]$$

- 777 $\mathcal{L}_{CE} = \mathbb{E}_X \left[-\sum_{i=1}^{i=r} \mathbb{E}_{D|X} \left[[z(D)]_i \log([f(X)]_i) \right] \right]$
- 778 $\mathcal{L}_{CE} = \mathbb{E}_X \left[-\sum_{i=1}^{i=r} \log([f(X)]_i) \mathbb{E}_{D|X}[[z(D)]_i] \right]$

779
$$\mathcal{L}_{CE} = \mathbb{E}_X \left[-\sum_{i=1}^{i=r} \log([f(X)]_i) (1 \times Q(D=i|X) + 0 \times (1 - Q(D=i|X))) \right]$$

780
$$\mathcal{L}_{CE} = \mathbb{E}_X \left[-\sum_{i=1}^{i=r} \log([f(X)]_i) Q(D=i|X) \right]$$

In order to find the minimizer of the cross entropy loss over the class of all functions from $R^p \rightarrow [0, 1]^r$, we formulate the following objective with the Lagrange constraint:

783
$$J = \min_{[f(X)]_1 \dots [f(X)]_r} \mathbb{E}_X \left[-\sum_{i=1}^{i=r} \log([f(X)]_i) Q(D=i|X) \right] + \lambda \left(\sum_{i=1}^{i=r} [f(X)]_i - 1 \right)$$

Setting partial derivative with respect to $[f(X)]_r$ to 0

785
$$-\frac{Q(D=i|X)}{[f^{\star}(X)]_i} + \lambda = 0$$

786
$$[f^{\star}(X)]_i = \frac{1}{\lambda}Q(D=i|X)$$

787 From KKT condition, the optimal solution lies on constraint surface, giving:

788
$$\sum_{i=1}^{i=r} [f^{\star}(X)]_i = 1$$

- 789 $\sum_{i=1}^{i=r} \frac{1}{\lambda} Q(D=i|X) = 1$
- 790 $\frac{1}{\lambda} \sum_{i=1}^{i=r} Q(D=i|X) = 1$
- 791 $\frac{1}{\lambda} = 1$
- 792 $\lambda = 1$
- Finally, we get $[f^*(X)]_i = Q(D = i|X)$, so the optimal f^* by the cross entropy loss as defined will in fact recover the oracle domain discriminator.

795 **D** Additional Experimental Details

- 796 Our code is available at
- 797 https://github.com/latentlabelshift-anonymous/latentlabelshift
- ⁷⁹⁸ Here we present the full generation procedure for semisynthetic example problems, and discuss the ⁷⁹⁹ parameters.
- 1. Choose a Dirichlet concentration parameter $\alpha > 0$, maximum condition number $\kappa \ge 1$ (with respect to 2-norm), and domain count $r \ge k$.
- 802 2. For each $y \in [k]$, sample $p_d(y) \sim \text{Dir}(\frac{\alpha}{k} \mathbf{1}_k)$.
- 803 3. Populate the matrix $\mathbf{Q}_{Y|D}$ with the computed $p_d(y)$ s. If $\operatorname{cond}(\mathbf{Q}_{D|Y}) \ge k$, return to step 2 and re-sample.
- 4. Distribute examples across domains according to $\mathbf{Q}_{Y|D}$, for each of train, test, and valid sets. This procedure entails creating a quota number of examples for each (class, domain) pair, and drawing datapoints without replacement to fill each quota. We must discard excess examples from some classes in the dataset due to class imbalance in the $\mathbf{Q}_{D|Y}$ matrix. Due to integral rounding, domains may be *slightly* imbalanced.
- 5. Conceal true class information and return (x_i, d_i) pairs.

It is important to note the role of κ and α in the above formulation. Although they are unknown 811 parameters to the classification algorithm, they affect the sparsity of the $\mathbf{Q}_{Y|D}$ and difficulty of the 812 problem. Small α encourages high sparsity in $p_d(y)$, and large α causes $p_d(y)$ to tend towards a 813 uniform distribution. κ has a strong effect on the difficulty of the problem. Consider the case when 814 k = 2. When $\kappa = 1$, the only potential $P_{Y|D}$ matrices are \mathbf{I}_2 up to row permutation (which means 815 that domains and classes are exactly correlated, so the domain indicates the class and the problem 816 is supervised). In the other limit, if $\kappa \to +\infty$, we may generate $P_{Y|D}$ matrices that are singular, 817 breaking needed assumptions for domain discriminator output to uniquely identify true class of 818 anchor subdomains. κ also helps control the class imbalance (if a row of $\mathbf{Q}_{Y|D}$ is small, indicating 819 that the class is heavily under-represented across all domains, the condition number will increase). 820

821 D.1 FieldGuide-2 and FieldGuide-28 Datasets

The dataset and description is available at https://sites.google.com/view/fgvc6/ competitions/butterflies-moths-2019. For the purpose of our experiments. From this data we create two datasets FieldGuide-2 and FieldGuide-28. For FieldGuide-28 we select the 28 classes which have 1000 datapoints in the training file. Since the test set provided in the website does not have annotations, we manually create a test set by sampling 200 datapoints from each of the 28 classes. The FieldGuide-2 dataset is created by considering two classes from the created FieldGuide-28 dataset.

828 D.2 Hyperparameters and Implementation Details: SCAN baseline

In all cases, we initialize the SCAN [55] network with the clustering head attached, sample data according to the $\mathbf{Q}_{D|Y}$ matrix, and predict classes.

With the Hungarian algorithm, implemented in [14, 56], we compute the highest true accuracy among any permutation of these labels (denoted "Test acc").

• CIFAR-10 and CIFAR-20 Datasets [36]

- We use ResNet-18 [29] backbone with weights trained by SCAN-loss and obtained from the SCAN repo https://github.com/wvangansbeke/Unsupervised-Classification.
- We use the same transforms present in the repo for test data.

• ImageNet-50 Dataset [17]

- We use ResNet-50 backbone with weights trained by SCAN-loss and obtained from the SCAN repo.
- 840 We use the same transforms present in the repo for test data.

FieldGuide-2 and FieldGuide-28 Datasets

For each of the two datasets, we pretrain a different SCAN baseline network (including pretext and SCAN-loss steps) on all available data from the dataset. The backbone for each is ResNet-18.

For training the pretext task, we use the same transform strategy used in the repo for CIFARfor training the pretext task, we use the same transform strategy used in the repo for CIFARdata (with mean and std values as computed on the Fieldguide-28 dataset, and crop size 224). For training SCAN, we resize the smallest image dimension to 256, perform a random horizontal flip and random crop to size 224. We also normalize. For validation we resize smallest image dimension to 256, center crop to 224, and normalize.

850 D.3 Hyperparameters and Implementation Details: DDFA (RI)

- ⁸⁵¹ This is the DDFA procedure with random initialization.
- The bulk of this procedure is described in Section 6, but for completeness we reiterate here.
- We train ResNet50 [29] (with random initialization and added dropout) based on the implementation
- from https://github.com/kuangliu/pytorch-cifar on images x_i with domain indices d_i as
- the label, choose best iteration by valid loss, pass all training and validation data through \hat{f} , and
- cluster pushforward predictions $\hat{f}(x_i)$ into $m \ge k$ clusters with Faiss K-Means [35]. We compute the
- ⁸⁵⁷ $\widehat{\mathbf{Q}}_{c(X)|D}$ matrix and run NMF to obtain $\widehat{\mathbf{Q}}_{c(X)|Y}, \widehat{\mathbf{Q}}_{Y|D}$. To make columns sum to 1, we normalize
- columns of $\widehat{\mathbf{Q}}_{c(X)|Y}$, multiply each column's normalization coefficient over the corresponding row
- of $\widehat{\mathbf{Q}}_{Y|D}$ (to preserve correctness of the decomposition), and then normalize columns of $\widehat{\mathbf{Q}}_{Y|D}$.

Some NMF algorithms only output solutions satisfying the anchor word property [3, 37, 27]. We

- found the strict requirement of an exact anchor word solution to lead to low noise tolerance. We
- therefore use the Sklearn implementation of standard NMF [13, 53, 46].
- We predict class labels with Algorithm 2. With the Hungarian algorithm, implemented in [14, 56],
- we compute the highest true accuracy among any permutation of these labels (denoted "Test acc").
- With the same permutation, we reorder rows of $\hat{P}_{Y|D}$, then compute the average absolute difference
- between corresponding entries of $\widehat{\mathbf{Q}}_{Y|D}$ and $\mathbf{Q}_{Y|D}$ (denoted " $\mathbf{Q}_{Y|D}$ err").
- ⁸⁶⁷ Hyperparameters were tuned by repeatedly consulting validation domain discrimination loss and final
- classification task accuracy (valid and test) on CIFAR-10 and CINIC-10 (similar to an extension of
- CIFAR-10) [15]. For this reason, we acknowledge that the hyperparameters may be overfit to CIFAR-
- 10 in particular. Final evaluation runs used the following fixed hyperparameters:

871 Common Hyperparameters

- 872 Architecture: ResNet-50 with added dropout
- Faiss KMeans number of iterations (niter): 100
- 874 Faiss Kmeans number of clustering redos (nredo): 5
- 875 Learning Rate: 0.001
- ⁸⁷⁶ Learning Rate Decay: Exponential, parameter 0.97
- 877 SKlearn NMF initialization: random

878 Dataset-Specific Hyperparameters

879 •	CIFAR-10 Dataset
880	Training Epochs: 100
881	Number of Clusters (m): 30
882 •	CIFAR-20 Dataset
883	Training Epochs: 100
884	Number of Clusters (m) : 60
885 •	ImageNet-50 Dataset
886	Not evaluated.

• FieldGuide-2 and FieldGuide-28 Datasets

888 Not evaluated.

889 D.4 Hyperparameters and Implementation Details: DDFA (SI) and DDFA (SPI)

890 This is the DDFA procedure with SCAN initialization.

The procedure is identical to the standard DDFA procedure, except that SCAN [55] pre-trained weights or SCAN [55] contrastive pre-text weights are used to initialize the domain discriminator

⁸⁹³ before it is fine-tuned on the domain discrimination task. Hyperparameters used also differ.

When SCAN pretrained weights are available, we use those. When they are not, we train SCAN ourselves.

Hyperparameters were tuned by repeatedly consulting validation domain discrimination loss and final classification task accuracy (valid and test) on CIFAR-10—as well as validation domain discrimination loss alone on CIFAR-20. For this reason, we acknowledge that the hyperparameters may be overfit to CIFAR-10 in particular (and to a lesser extent, to CIFAR-20). Final evaluation runs used the following fixed hyperparameters:

901 **Common Hyperparameters**

- Faiss KMeans number of iterations (niter): 100
- ⁹⁰³ Faiss Kmeans number of clustering redos (nredo): 5
- 904 Learning Rate: 0.00001
- ⁹⁰⁵ Learning Rate Decay: Exponential, parameter 0.97
- 906 SKlearn NMF initialization: random

907 Dataset-Specific Hyperparameters

908 • CIFAR-10 Dataset

- 909 Architecture: ResNet-18
- Pre-seed: Weights trained with SCAN pretext and SCAN-loss on entirety of CIFAR-10
- 911 (from SCAN repo).
- 912 Training Epochs: 25
- 913 Number of Clusters (m): 10
- 914 Transforms used: Same as SCAN repo.

915 • CIFAR-20 Dataset

- 916 Architecture: ResNet-18
- Pre-seed: Weights trained with SCAN pretext and SCAN-loss on entirety of CIFAR-20
- 918 (from SCAN repo).
- 919 Training Epochs: 25
- 920 Number of Clusters (*m*): 20
- ⁹²¹ Transforms used: Same as SCAN repo.
- ImageNet-50 Dataset
- 923 Architecture: ResNet-50
- Pre-seed: Weights trained with SCAN pretext and SCAN-loss on entirety of ImageNet-50
- 925 (from SCAN repo).
- 926 Training Epochs: 25
- 927 Number of Clusters (m): 50
- ⁹²⁸ Transforms used: Same as SCAN repo.
- 929 FieldGuide-2 Dataset
- 930 Architecture: ResNet-18
- Pre-seed: Weights trained with SCAN pretext on entirety of FieldGuide-2 (trained by us).
- 932 Training Epochs: 30
- Number of Clusters (*m*): 2

934	Transforms used for pretext: Same strategy as CIFAR-10 in SCAN repo with appropriate
935	mean, std, and crop size 224.
936	Transform used for SCAN: Resize to 256, Random horizontal flip, Random crop to 224,
937	normalize
938	Learning rate used for SCAN: 0.001 (other hyperparameters were same as in SCAN repo
939	for CIFAR-10)
940	Note: During one of the random seeds of training, test data transforms were mismatched with
941	train transforms (specifically, missing the Resize(256) transform on test only). We consider
942	this to disadvantage our approach for that random seed as compared to the SCAN baseline,
943	which uses the proper transforms in all seeds. We report these results regardless due to the
944	fact that our approach still competes effectively even despite the transform disadvantage.
945	This random seed is the one displayed in Section 6 of the main paper. The results shown
946	in App. E are different results, with the Resize(256) included (and are therefore the best,
947	fair-footing evaluation comparison between our method and baseline).
948	 FieldGuide-28 Dataset
949	Architecture: ResNet-18
950	Pre-seed: Weights trained with SCAN pretext on entirety of FieldGuide-28 (trained by us).
951	Training Epochs: 60
952	Number of Clusters (m): 28
953	Transforms used for pretext: Same strategy as CIFAR-10 in SCAN repo with appropriate
954	mean, std, and crop size 224.
955	Transform used for SCAN: Resize to 256, Random horizontal flip, Random crop to 224,
956	normalize
957	Learning rate used for SCAN: 0.01 (other hyperparameters were same as in SCAN repo for
958	CIFAR-10)
959	Note: During one of the random seeds of training, test data transforms were mismatched with
960	train transforms (specifically, missing the Resize(256) transform on test only). We consider
961	this to disadvantage our approach for that random seed as compared to the SCAN baseline,
962	which uses the proper transforms in all seeds. We report these results regardless due to the
963	fact that our approach still competes effectively even despite the transform disadvantage.
964	This random seed is the one displayed in Section 6 of the main paper. The results shown
965	in App. E are different results, with the Resize(256) included (and are therefore the best,
966	tair-tooting evaluation comparison between our method and baseline).

967 E Additional Experimental Results

Table 3: *Results on CIFAR-10.* Each entry is produced with the averaged result of 3 different random seeds. With DDFA (RI) we refer to DDFA with randomly initialized backbone. With DDFA (SI) we refer to DDFA's backbone initialized with SCAN. Note that in DDFA (SI), we do not leverage SCAN for clustering. α is the Dirichlet parameter used for generating label marginals in each domain, κ is the maximum allowed condition number of the generated $\mathbf{Q}_{Y|D}$ matrix, r is number of domains.

r Approache		$\alpha:0.5,\;\kappa:4$		$\alpha:3,\;\kappa:4$		$\alpha:10,\;\kappa:8$	
	- pprodentes	Test acc	$\mathbf{Q}_{Y D}$ err	Test acc	$\mathbf{Q}_{Y D}$ err	Test acc	$\mathbf{Q}_{Y D}$ err
10	SCAN DDFA (RI) DDFA (SI)	0.8205 0.7361 0.8987	0.0354 0.0233	0.8222 0.5393 0.7566	0.0481 0.0401	0.8025 0.3135 0.5359	0.0736 0.0542
15	SCAN DDFA (RI) DDFA (SI)	0.8245 0.7732 0.9605	0.0326 0.0163	0.8172 0.5319 0.8443	0.0464 0.0256	0.8107 0.275 0.7327	0.0739 0.0376
20	SCAN DDFA (RI) DDFA (SI)	0.8048 0.6883 0.9656	0.0471 0.0158	0.8109 0.5651 0.904	0.0461 0.0193	0.8164 0.2704 0.7979	0.0705 0.0298
25	SCAN DDFA (RI) DDFA (SI)	0.7974 0.7235 0.9701	0.0389 0.0132	0.8131 0.5616 0.9173	0.0444 0.0166	0.8094 0.2802 0.8204	0.0863 0.027

Table 4: *Extended Results on CIFAR-20.* Each entry is produced with the averaged result of 3 different random seeds (including the 1 seed shown in main paper results). With DDFA (RI) we refer to DDFA with randomly initialized backbone. With DDFA (SI) we refer to DDFA's backbone initialized with SCAN. Note that in DDFA (SI), we do not leverage SCAN for clustering. α is the Dirichlet parameter used for generating label marginals in each domain, κ is the maximum allowed condition number of the generated $\mathbf{Q}_{Y|D}$ matrix, r is number of domains.

r	Approaches	$\alpha: 0.5, \ \kappa: 8$		$\alpha:3, \kappa:12$		$\alpha:10, \kappa:20$	
-	- 11	Test acc	$\mathbf{Q}_{Y D}$ err	Test acc	$\mathbf{Q}_{Y D}$ err	Test acc	$\mathbf{Q}_{Y D}$ err
20	SCAN DDFA (RI) DDFA (SI)	0.4416 0.517 0.7838	0.0423 0.0225	0.4447 0.3355 0.5926	0.0451 0.0271	0.433 0.1632 0.3904	0.0568 0.0344
25	SCAN DDFA (RI) DDFA (SI)	0.4403 0.4886 0.8372	0.0487 0.0204	0.4439 0.292 0.6685	0.0485 0.0253	0.4398 0.0753 0.4869	0.0806 0.0299
30	SCAN DDFA (RI) DDFA (SI)	0.4317 0.5121 0.8197	0.0462 0.0219	0.4611 0.2992 0.7425	0.0477 0.0207	0.4327 0.0869 0.5434	0.0765 0.0283

Table 5: *Results on ImageNet-50*. Each entry is produced with the averaged result of 3 different random seeds. With DDFA (SI) we refer to DDFA's backbone initialized with SCAN. Note that in DDFA (SI), we do not leverage SCAN for clustering. α is the Dirichlet parameter used for generating label marginals in each domain, κ is the maximum allowed condition number of the generated $\mathbf{Q}_{Y|D}$ matrix, r is number of domains.

r	Approaches	$\alpha:0.5,\ \kappa:200$		$\alpha:3,\ \kappa:205$		$\alpha: 10, \ \kappa: 210$	
-	1 pprouenes	Test acc	$\mathbf{P}_{Y D}$ err	Test acc	$\mathbf{P}_{Y D}$ err	Test acc	$\mathbf{P}_{Y D}$ err
50	SCAN DDFA (SI)	0.7366 0.7204	0.0132	0.754 0.6322	- 0.0149	0.7372 0.3431	0.0217
60	SCAN DDFA (SI)	0.757 0.8179	0.0101	0.7344 0.7434	- 0.0124	0.7315 0.5784	0.0175

Table 6: *Extended Results on FieldGuide-2*. Each entry is produced with the averaged result of 3 different random seeds (this does not include the 1 seed shown in main paper results, due to different transforms—see App. D for more information). With DDFA (SI) we refer to DDFA's backbone initialized with SCAN. Note that in DDFA (SI), we do not leverage SCAN for clustering. α is the Dirichlet parameter used for generating label marginals in each domain, κ is the maximum allowed condition number of the generated $\mathbf{Q}_{Y|D}$ matrix, r is number of domains.

r	Approaches	$\alpha:0.5,\ \kappa:3$		$\alpha:3,\ \kappa:5$		$\alpha:10,\;\kappa:7$	
-		Test acc	$\mathbf{P}_{Y D}$ err	Test acc	$\mathbf{P}_{Y D}$ err	Test acc	$\mathbf{P}_{Y D}$ err
2	SCAN DDFA (SPI)	0.5784 0.7757	0.2411	0.5635 0.7729	- 0.1497	0.5731 0.6578	0.264
3	SCAN DDFA (SPI)	0.598 0.9599	0.0545	0.5889 0.8304	0.1478	0.5935 0.6933	0.2238
5	SCAN DDFA (SPI)	0.6004 0.9534	0.0934	0.5767 0.7835	0.1406	0.5792 0.6166	0.2575
7	SCAN DDFA (SPI)	0.5932 0.9037	0.1151	0.5923 0.8157	- 0.1448	0.5861 0.6608	0.1975
10	SCAN DDFA (SPI)	0.5831 0.9071	0.1547	0.5848 0.7139	0.1698	0.582 0.582	0.164

968

Table 7: *Extended Results on FieldGuide-28*. Each entry is produced with the result of 1 random seed (this does not include the 1 seed shown in main paper results, due to different transforms—see App. D for more information). With DDFA (SI) we refer to DDFA's backbone initialized with SCAN. Note that in DDFA (SI), we do not leverage SCAN for clustering. α is the Dirichlet parameter used for generating label marginals in each domain, κ is the maximum allowed condition number of the generated $\mathbf{Q}_{Y|D}$ matrix, r is number of domains.

r	Approaches	$\alpha:0.5,\ \kappa:12$		$\alpha:3, \ \kappa:20$		$\alpha: 10, \ \kappa: 28$	
	1 pprouenes	Test acc	$\mathbf{Q}_{Y D}$ err	Test acc	$\mathbf{Q}_{Y D}$ err	Test acc	$\mathbf{Q}_{Y D}$ err
28	SCAN DDFA (SPI)	0.2882 0.5472	0.0359	0.2915 0.3102	0.0338	0.314 0.3136	0.0355
37	SCAN DDFA (SPI)	0.295 0.7595	0.0276	0.3126 0.5214	0.0318	0.3145 0.3257	- 0.0409
42	SCAN DDFA (SPI)	0.2856 0.6699	0.0317	0.3177 0.4714	0.0373	0.2883 0.4082	0.0311
47	SCAN DDFA (SPI)	0.2807 0.7087	0.0352	0.3184 0.4728	0.0353	0.3044 0.2992	0.0394



Figure 3: This figure illustrates the case with 3 domains and 3 classes. The oracle domain discriminator maps points from a high-dimensional input space to a k = 3 vertex convex polytope (shaded red) embedded in Δ^{r-1} , r = 3 (shaded yellow). The anchor subdomains map to the vertices of this polytope.

970 F Discussion of Convex Polytope Geometry

The geometric properties of topic modeling for finite, discrete random variables has been explored in depth in related works (Huang et al. [32], Donoho and Stodden [20], Chen et al. [12]). The observation that columns in $\mathbf{Q}_{X|D}$ are convex combinations of columns in $\mathbf{Q}_{X|Y}$ leads to a perspective on identification of the matrix decomposition as identification of the convex polytope in \mathbb{R}^m which encloses the datapoints (the corners of which correspond to columns of $\mathbf{Q}_{X|Y}$).

Here, we briefly discuss an interesting but somewhat different application of convex polytope geometry. Instead of a convex polytope in \mathbb{R}^m with corners as columns of $\mathbf{Q}_{X|Y}$, we concern ourselves with the convex polytope in \mathbb{R}^r with corners as columns in $\mathbf{Q}_{D|Y}$, which must enclose all f(x) for $x \in \mathcal{X}, q(x) > 0$.

980 Let us assume that Assumptions A.1–A.4 are satisfied.

We recall the oracle domain discriminator f(x) = q(d|X = x). Let $x \in \mathcal{X} = \mathbb{R}^p$. Now, since the r values q(d|X = x) for $d \in \{1, 2, ..., r\}$ constitute a distribution over the random variable $d \in [r]$, each of the r values lie between 0 and 1, and also their sum adds to 1. Therefore the vector f(x) lies on the simplex Δ^{r-1} . We now express f(x) as a convex combination of the k columns of $\mathbf{Q}_{D|Y}$. We denote these column vectors $\mathbf{Q}_{D|Y}[:, y]$ for each $y \in \mathcal{Y} = [k]$. Note that each such vector also lies in the Δ^{r-1} simplex.

As an intermediate step in the proof of 3 given in App. A, we showed that each f(x) is a linear combination of these columns of $\mathbf{Q}_{D|Y}$ with coefficients q(y|X = x) for all $y \in \mathcal{Y}$.

989 That is, we can rewrite $f(x) = \mathbf{Q}_{D|Y} [Q(Y = 1|X = x) \dots Q(Y = k|X = x)]^{\top}$

Since the coefficients in the linear combination are probabilities which, taken together, form a categorical distribution, they lie between 0 and 1 and sum to 1. Thus, for all $x \in \mathcal{X}$, f(x) can be expressed as a *convex* combination of the columns of $\mathbf{Q}_{D|Y}$. Therefore, for any x, f(x) lies inside the *k*-vertex convex polytope with corners as the columns of $\mathbf{Q}_{D|Y}$ (which are linearly independent by Lemma 6). This polytope is embedded in Δ^{r-1} .

Now consider x in an anchor sub-domain, that is $x \in A_y$ for some $y \in \mathcal{Y}$. We know that if q(x) > 0, q(y|X = x) = 1, $q(y'|X = x) = 0 \forall y' \neq y$ (Lemma 5). Since the q(y|X = x) are now one-hot, we have that $f(x) = \mathbf{Q}_{D|Y}[:, y]$ for $x \in A_y$. In words, this means that f(x) is precisely the yth column of $\mathbf{Q}_{D|Y}$. It follows that the domain discriminator maps each of the k anchor sub-domains exactly to a unique vertex of the polytope.

- We could now recover the columns of $\mathbf{Q}_{D|Y}$, up to permutation, with the following procedure:
- 1001 1. Push all $x \in \mathcal{X}$ through f.
- 10022. Find the minimum volume convex polytope that contains the resulting density of points1003on the simplex. The vectors that compose the vertices of this polytope are the columns of1004 $Q_{D|Y}$, up to permutation.

Note that from Assumption A.4, we are guaranteed to have a region of the input space with at least $\epsilon > 0$ mass that gets mapped to each of the vertices when carrying out step (i). Therefore, our discovered minimum volume polytope must enclose all of these vertices. Since no mass will exist outside of the true polytope, requiring a minimum volume polytope will ensure that the recovered polytope fits the true polytope's vertices precisely (as any extraneous volume outside of the true polytope must be eliminated). Then step (ii) recovers $\mathbf{Q}_{D|Y}$, up to permutation of columns.

Having recovered $\mathbf{Q}_{D|Y}$, we can use Lemmas 1 and 2 to recover q(y|x, d).

This procedure is a geometric alternative to the clustering approach outlined in Algorithm 1. In practice, fitting a convex hull around a noisy domain discriminator may be computationally expensive, and may fail to recover the true vertices.