

A. Deferred proofs

Proof of Lemma 1. For the sake of the proof, assume that we have two labeled set of samples of size m and m_L from $p_{\mathcal{X}\mathcal{Y}}$, call them respectively S and S_L . The set S represents our unlabeled sample, and the set S_L represents the labeled sample. For any $\delta \in (0, 1)$, we would like to find a $\gamma > 0$ such that with probability $1 - \delta$, for all $i \in 1, \dots, n$,

$$\left| \frac{1}{m} \sum_{(x, \mathbf{y}) \in S} \ell(\phi_i(x), \mathbf{y}) - \frac{1}{m_L} \sum_{(x, \mathbf{y}) \in S_L} \ell(\phi_i(x), \mathbf{y}) \right| \leq \gamma. \quad (6)$$

The sample S represents the unlabeled data x_1, \dots, x_m we have access to. In fact, $\frac{1}{m} \sum_{(x, \mathbf{y}) \in S} \ell(\phi_i(x), \mathbf{y}) = \hat{R}(\phi_i; X, \mathbf{Y}^*)$. The inequality (6) implies that for the true labeling of the unlabeled data x_1, \dots, x_m , for any $i \in 1, \dots, n$, it holds that:

$$\hat{R}(\phi_i; X, \mathbf{Y}^*) \in [\hat{\mu}_i - \gamma, \hat{\mu}_i + \gamma]$$

where $\hat{\mu}_i = \frac{1}{m_L} \sum_{(x, \mathbf{y}) \in S_L} \ell(\phi_i(x), \mathbf{y})$ is the empirical mean computed from the labeled sample S_L .

By using Hoeffding's inequality, we have that for a fixed i , it holds that

$$\begin{aligned} & \mathbb{P}_{S, S_L} \left(\left| \frac{1}{m} \sum_{(x, \mathbf{y}) \in S} \ell(\phi_i(x), \mathbf{y}) - \frac{1}{m_L} \sum_{(x, \mathbf{y}) \in S_L} \ell(\phi_i(x), \mathbf{y}) \right| > \gamma \right) \\ & \leq 2 \exp \left(\frac{-2\gamma^2}{\sum_{j=1}^m \left(\frac{B}{m}\right)^2 + \sum_{j=1}^{m_L} \left(\frac{B}{m_L}\right)^2} \right) \\ & = 2 \exp \left(\frac{-2m_L m \gamma^2}{B^2(m + m_L)} \right) = \frac{\delta}{n}. \end{aligned} \quad (7)$$

By taking a union bound and solving (7) with respect to γ , the statement follows.

Proof of Lemma 4. By invoking Lemma 3, it is easy to see that the function $\hat{R}(\mathbf{h}_\theta; X, \mathbf{Y}')$ is a convex combination of convex functions with respect to θ , hence it is also convex in θ . Let $\mathbf{v} \in \partial \hat{R}(\mathbf{h}_\theta; X, \mathbf{Y}')$. If a function is convex, then there exists at least one subgradient for each point of its domain, so \mathbf{v} is well defined. Then, we have that for any $\theta'' \in \Theta$, it holds that

$$\hat{R}(\mathbf{h}_{\theta''}, \mathbf{Y}') - \hat{R}(\mathbf{h}_{\theta'}, \mathbf{Y}') \geq \mathbf{v}^T (\theta'' - \theta').$$

As $f(\theta'') \geq \hat{R}(\mathbf{h}_{\theta''}; X, \mathbf{Y}')$, we have that

$$f(\theta'') - f(\theta') \geq \mathbf{v}^T (\theta'' - \theta'),$$

which implies that \mathbf{v} is a subgradient of f at θ' .

Proof of Theorem 5. We need to show that $f(\theta)$ is convex and L -Lipschitz continuous with respect to θ to apply the standard convergence result for constant step size subgradient optimization (Bertsekas, 2015), which yields

$$f(\tilde{\theta}) - f(\hat{\theta}) \leq \frac{\text{diameter}(\Theta)^2 + L^2 h^2 T}{2hT}. \quad (8)$$

To show that $f(\theta)$ is convex it is straightforward to see that $\hat{R}(\mathbf{h}_\theta; X, \mathbf{Y})$ is convex in θ as it is the convex combination of convex functions in θ . For any $\lambda \in [0, 1]$, we have that

$$\begin{aligned} f(\lambda \theta_1 + (1 - \lambda) \theta_2) &= \max_{\mathbf{Y} \in \mathbb{Y}^\circ} \hat{R}(h_{\lambda \theta_1 + (1 - \lambda) \theta_2}; X, \mathbf{Y}) \\ &\leq \max_{\mathbf{Y} \in \mathbb{Y}^\circ} \left[\lambda \hat{R}(\mathbf{h}_{\theta_1}; X, \mathbf{Y}) + (1 - \lambda) \hat{R}(\mathbf{h}_{\theta_2}; X, \mathbf{Y}) \right] \\ &\leq \lambda \max_{\mathbf{Y} \in \mathbb{Y}^\circ} \hat{R}(\mathbf{h}_{\theta_1}; X, \mathbf{Y}) + (1 - \lambda) \max_{\mathbf{Y} \in \mathbb{Y}^\circ} \hat{R}(\mathbf{h}_{\theta_2}; X, \mathbf{Y}) \\ &= \lambda f(\theta_1) + (1 - \lambda) f(\theta_2). \end{aligned}$$

Also, $f(\theta)$ is L -Lipschitz continuous with respect to θ . In fact, it is straightforward to see that $\hat{R}(\mathbf{h}_\theta; X, \mathbf{Y})$ is also L -Lipschitz continuous with respect to θ . For any $\theta_1, \theta_2 \in \Theta$, we have that

$$\begin{aligned} |f(\theta_1) - f(\theta_2)| &\leq \max_{\mathbf{Y} \in \mathbb{Y}^\circ} |\hat{R}(\mathbf{h}_{\theta_1}; X, \mathbf{Y}) - \hat{R}(\mathbf{h}_{\theta_2}; X, \mathbf{Y})| \\ &\leq L \|\theta_1 - \theta_2\|_2. \end{aligned}$$

The subgradient of $f(\theta)$ in θ is computed by using Lemma 4. The last part of the Theorem immediately follows by substituting h and T in (8) as in the Theorem statement.

Proof of Lemma 6. For any $i \in 1, \dots, n$, we have that

$$\frac{\partial}{\partial \theta_i} \ell(\mathbf{h}_\theta(x), \mathbf{e}) = 2 \left(\phi_i(x)^T \cdot \mathbf{h}_\theta(x) - \phi_i(x)^T \cdot \mathbf{e} \right).$$

Therefore, we can bound the norm of the gradient of ℓ as

$$\begin{aligned} \|\nabla \ell(\mathbf{h}_\theta(x), \mathbf{e})\|_2 &= 2 \sqrt{\sum_{i=1}^n (\phi_i(x)^T \cdot (\mathbf{h}_\theta(x) - \mathbf{e}))^2} \\ &\leq 2 \sqrt{\sum_{i=1}^n (1)^2} \\ &\leq 2\sqrt{n}. \end{aligned}$$

The first inequality is an application of Hölder's Inequality, as $\|\phi_i(x)^T\|_1 = 1$ and $\|\mathbf{h}_\theta(x) - \mathbf{e}\|_\infty \leq 1$. This implies that the function $\ell(\mathbf{h}_\theta(x), \mathbf{e})$ is $2\sqrt{n}$ -Lipschitz continuous with respect to θ .

Proof of Lemma 7. First, we will prove that $\ell(\mathbf{h}_\theta(x), \mathbf{e})$ is bounded. Without loss of generality, suppose that $e_i = 1$. We have that

$$\ell(\mathbf{h}_\theta(x), \mathbf{e}) = -\ln \left(\frac{\exp(\mathbf{w}_i^T \cdot \mathbf{x})}{\sum_{c=1}^k \exp(\mathbf{w}_c^T \cdot \mathbf{x})} \right).$$

It is easy to see that $\ell(\mathbf{h}_\theta(\mathbf{x}), \mathbf{e}) \geq 0$. By using the Cauchy-Schwarz inequality, we have that

$$\begin{aligned} \ell(\mathbf{h}_\theta(\mathbf{x}), \mathbf{e}) &= -\ln \left(\frac{\exp(\mathbf{w}_i^T \cdot \mathbf{x})}{\sum_{c=1}^k \exp(\mathbf{w}_c^T \cdot \mathbf{x})} \right) \\ &\leq -\ln \left(\frac{\exp(-B_w B_x)}{k \exp(B_w B_x)} \right) \\ &\leq 2B_w B_x + \ln k . \end{aligned}$$

Now, we prove that $\ell(\mathbf{h}_\theta(\mathbf{x}), \mathbf{e})$ is Lipschitz continuous with respect to θ . For a fixed $(\mathbf{x}, \mathbf{e}) \in \mathcal{X} \times \mathcal{Y}$, consider the function $\omega(\mathbf{p}) : \mathbb{R}^k \rightarrow \mathcal{Y}^\circ$, defined as

$$\omega(\mathbf{p}) \doteq - \sum_{c=1}^k e_c \cdot \ln(p_c) ,$$

and let

$$\mathbf{h}(\theta) \doteq \left(\frac{\exp(\mathbf{w}_1^T \cdot \mathbf{x})}{\sum_{c=1}^k \exp(\mathbf{w}_c^T \cdot \mathbf{x})}, \dots, \frac{\exp(\mathbf{w}_k^T \cdot \mathbf{x})}{\sum_{c=1}^k \exp(\mathbf{w}_c^T \cdot \mathbf{x})} \right)^T ,$$

where $\theta = (\mathbf{w}_1 \dots \mathbf{w}_k)^T$, and observe that $\ell(\mathbf{h}_\theta(\mathbf{x}), \mathbf{e}) = \omega \circ \mathbf{h}(\theta)$.

It is well known that ℓ is $L_\omega L_h$ -Lipschitz continuous with respect to θ , where L_ω and L_h are the Lipschitz constants respectively of ω and \mathbf{h} . It is also a known result that $L_\omega \leq 1$ (see for example Proposition 4 of (Gao & Pavel, 2018)).

We now want to compute L_h . We will use the fact that $\max_{\theta \in \Theta} \|J_h(\theta)\|_F \leq L_h$, where J_h denotes the Jacobian matrix of \mathbf{h} and $\|\cdot\|_F$ denotes the Frobenius norm.

For ease of notation, let $\mathbf{h}(\theta) = \mathbf{p} = (p_1, \dots, p_k)^T$. We have that for any $i \in 1, \dots, k$, it holds that

$$\begin{aligned} \frac{\partial[\mathbf{h}(\theta)]_i}{\partial \mathbf{w}_j} &= p_i p_j \mathbf{x} \quad \text{for } j \neq i , \text{ and} \\ \frac{\partial[\mathbf{h}(\theta)]_i}{\partial \mathbf{w}_i} &= (p_i - p_i^2) \mathbf{x} . \end{aligned}$$

Therefore, we can bound the square of the Frobenius norm of the Jacobian matrix of \mathbf{h} with

$$\begin{aligned} \|J_h(\theta)\|_F^2 &= \sum_{i,j} \left\| \frac{\partial[\mathbf{h}(\theta)]_i}{\partial \mathbf{w}_j} \right\|_2^2 \\ &\leq \|\mathbf{x}\|_2^2 \left(\sum_i [p_i(1-p_i)]^2 + \sum_{i \neq j} [p_i p_j]^2 \right) \\ &\leq \|\mathbf{x}\|_2^2 (k + k^2/2) \leq \|k\mathbf{x}\|_2^2 . \end{aligned}$$

We can conclude that \mathbf{h} is kB_x -Lipschitz continuous, and the statement follows.

Proof of Theorem 8. From Chapter 14 of Mitzenmacher & Upfal (2017), we know that

$$R(\mathbf{h}_{\hat{\theta}}) \leq \hat{R}(\mathbf{h}_{\hat{\theta}}; X, \mathbf{Y}^*) + 2\hat{\mathfrak{R}}_m(\mathcal{L}; X, \mathbf{Y}^*) + O \left(B \sqrt{\frac{\ln \frac{1}{\delta}}{m}} \right) .$$

By definition of $f(\cdot)$, it holds that $\hat{R}(\mathbf{h}_{\hat{\theta}}; X, \mathbf{Y}^*) \leq f(\hat{\theta})$. As $\hat{\theta}$ is the optimal solution of (2), we have that $f(\hat{\theta}) \leq f(\theta^*)$. Let $\mathbf{Y}' \doteq \arg \max_{\mathbf{Y} \in \mathbb{Y}^\circ} \hat{R}(\mathbf{h}_{\theta^*}; X, \mathbf{Y})$. It holds that

$$\begin{aligned} f(\theta^*) &= \hat{R}(\mathbf{h}_{\theta^*}; X, \mathbf{Y}') \\ &= \hat{R}(\mathbf{h}_{\theta^*}; X, \mathbf{Y}') + \hat{R}(\mathbf{h}_{\theta^*}; X, \mathbf{Y}^*) - \hat{R}(\mathbf{h}_{\theta^*}; X, \mathbf{Y}^*) \\ &= \hat{R}(\mathbf{h}_{\theta^*}; X, \mathbf{Y}^*) + |\hat{R}(\mathbf{h}_{\theta^*}; X, \mathbf{Y}') - \hat{R}(\mathbf{h}_{\theta^*}; X, \mathbf{Y}^*)| . \end{aligned}$$

By using the fact that ℓ is bounded, and the definition of diameter $D_{\mathbb{Y}^\circ}$, we have that

$$\begin{aligned} &|\hat{R}(\mathbf{h}_{\theta^*}; X, \mathbf{Y}') - \hat{R}(\mathbf{h}_{\theta^*}; X, \mathbf{Y}^*)| \\ &= \left| \frac{1}{m} \sum_{j=1}^m \sum_{c=1}^k \ell(\mathbf{h}_{\theta^*}(x_j), \mathbf{e}_c) (y'_{jc} - y_{jc}^*) \right| \\ &\leq B \frac{1}{m} \sum_{j=1}^m \sum_{c=1}^k |y'_{jc} - y_{jc}^*| \leq B D_{\mathbb{Y}^\circ} . \end{aligned}$$

To wrap it up, it results that

$$\begin{aligned} R(\mathbf{h}_{\hat{\theta}}) &\leq \hat{R}(\mathbf{h}_{\theta^*}; X, \mathbf{Y}^*) + B D_{\mathbb{Y}^\circ} + 2\hat{\mathfrak{R}}_m(\mathcal{L}; X, \mathbf{Y}^*) \\ &\quad + O \left(B \sqrt{\frac{\ln \frac{1}{\delta}}{m}} \right) \\ &\leq R(\mathbf{h}_{\theta^*}) + B D_{\mathbb{Y}^\circ} + 4\hat{\mathfrak{R}}_m(\mathcal{L}; X, \mathbf{Y}^*) \\ &\quad + O \left(B \sqrt{\frac{\ln \frac{1}{\delta}}{m}} \right) \\ &\leq R(\mathbf{h}_{\theta^*}) + B D_{\mathbb{Y}^\circ} \\ &\quad + \sup_{\mathbf{Y} \in \mathbb{Y}^\circ} 4\hat{\mathfrak{R}}_m(\mathcal{L}; X, \mathbf{Y}) + O \left(B \sqrt{\frac{\ln \frac{1}{\delta}}{m}} \right) . \end{aligned}$$

Proof of Lemma 9 The proof is along the same lines of the proof of Lemma 1, but we take a union bound with respect to all the nK intervals $\Delta_{i,c,\hat{c}}$ for $i \in 1, \dots, n$, $c \in 1, \dots, k$, and $\hat{c} = 1, \dots, k_i$. Moreover, as for any $j \in 1, \dots, m$, we have that $y_{j,c}[\phi_i(x_j)]_{\hat{c}} \leq 1$, we take $B = 1$ during the proof (as defined in Lemma 1).

B. Additional Experimental Details

We provide further information specifying the experimental setup used to generate our figures.

B.1. Weak Supervision Sources

We first build the weak supervision sources on our two datasets as follows.

Animals with Attributes. Each class is annotated with a binary vector of attributes. For each attribute, we train a binary classifier by finetuning a ResNet-18 using labeled data from the seen classes. When we consider a classification task between two unseen classes, we use as weak supervision sources the classifiers for the attributes which are different between the animals of these two unseen classes. We report the results of the 4 binary classification tasks which have the lowest majority vote accuracy. We chose these particular tasks to demonstrate the abilities of our methods on the tasks that have the least accurate weak supervision sources.

DomainNet. We sample 5 of the 25 classes of DomainNet with the largest number of datapoints. For each domain, we use 60% of the available data for those classes to fine tune a pretrained ResNet-18 network. We perform this procedure on two disjoint samples of test classes to illustrate our results on two distinct multiclass classification tasks.

In our experiments, we use the pretrained ResNet-18 from PyTorch. We finetune this ResNet-18 network following the approach described in (He et al., 2016), using cross-entropy loss.

B.2. Algorithm Hyperparameters

The subgradient method (Algorithm 1) used to train AMCL-CC and AMCL-LR uses the following hyperparameters:

AMCL-CC: We set $\delta = 0.1$, and build the constraints as in Lemma 1. We use $\varepsilon = 0.1$, and define the step size h and the number of iterations T as in Theorem 5, using $L = 2\sqrt{n}$ and diameter of Θ equal to $\sqrt{2}$.

AMCL-LR: In this case, the loss function is bounded as in Lemma 7. Since this value could be potentially very large, which in turn it would result in large intervals and number of iterations, we use the value $B = 0.1$ in the experiments. We set δ to 0.1 and build the constraints as in Lemma 1. We do not bound the set of weights Θ : in the experiments, the norm of the weights of the multinomial logistic regression model has never diverged. We run the subgradient algorithm for $T = 1000$ iterations with step size $h = 0.02$.

C. Additional Figures

C.1. Animals with Attributes

We provide the remaining figures for our experiments on the Animals with Attributes dataset. The last two binary classification tasks are bat v. rat and horse v. giraffe.

From Figure 3, we note that our methods show similar results as the figures displayed in the main body of the paper. AMCL-LR matches or outperforms all other methods on both tasks, over all ranges labeled data. AMCL-CC is within a few accuracy points of the other baselines and AMCL-LR on these tasks.

C.2. DomainNet

We provide the remaining figures for our experiments on the DomainNet dataset. We provide histograms when using the other 4 domains as the target task and also provide histograms for results on another of the samples of 5 classes. The first sample of classes as mentioned in the main body of the paper is {sea turtle, vase, whale, bird, violin}. The second sample is {tornado, trombone, submarine, feather, zebra}.

From Figures 4–8, we note that in most domains our methods perform better than or match all other approaches, namely in both samples of Clipart, Quickdraw, Painting, and the second sample of Sketch. Our methods achieve slightly lower accuracy than the best performing baseline on the Real domain and on the second sample of the Infograph domain, although they are not beaten by a single baseline in all of these tasks. We believe that the combination of our theoretical guarantees and that our methods achieve similar or sometimes better empirical performance captures the benefits of AMCL.

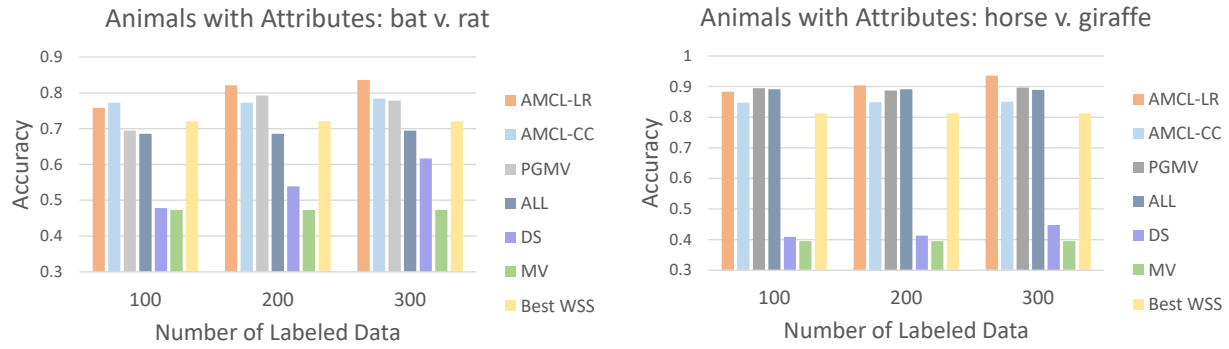


Figure 3. Experimental results on the Animals with Attributes dataset for the binary classification tasks of bat vs. rat and horse vs. giraffe as we vary the amount of labeled data. Each method uses 347 unlabeled data for bat vs. rat and 1424 unlabeled data for horse vs. giraffe.

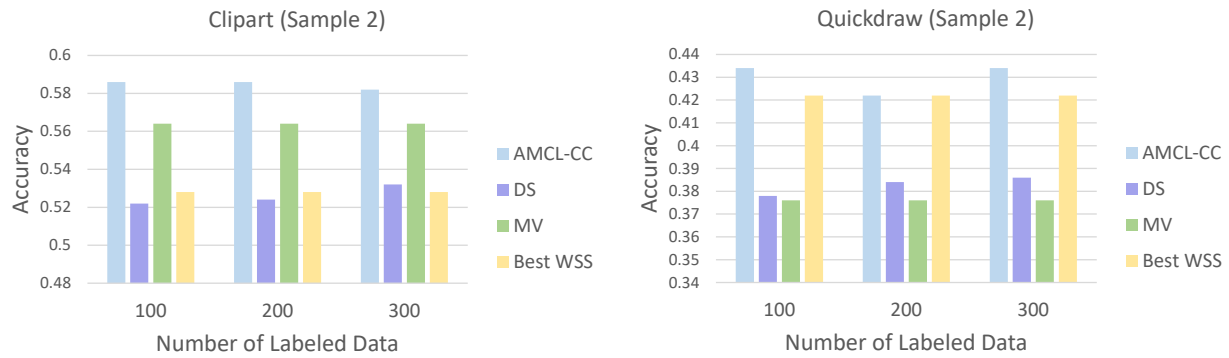


Figure 4. Experimental results on the second sample of Domain Net for the clipart and quickdraw domains as we vary the amount of labeled data. Each method uses 500 unlabeled data.

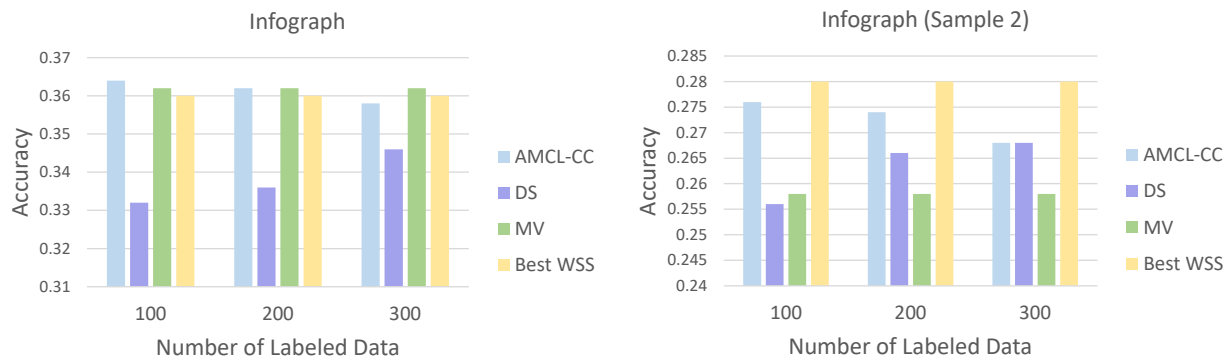


Figure 5. Experimental results on both samples of Domain Net for the Infograph domain as we vary the amount of labeled data. Each method uses 500 unlabeled data.

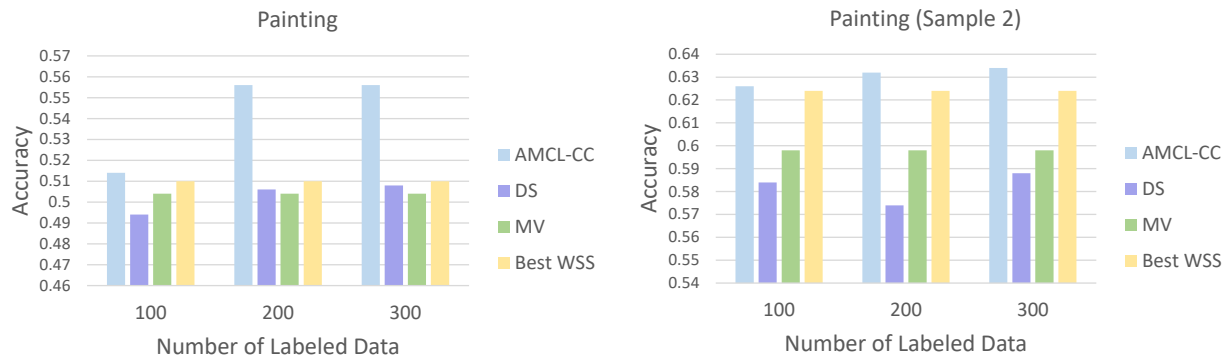


Figure 6. Experimental results on both samples of Domain Net for the Painting domain as we vary the amount of labeled data. Each method uses 500 unlabeled data.

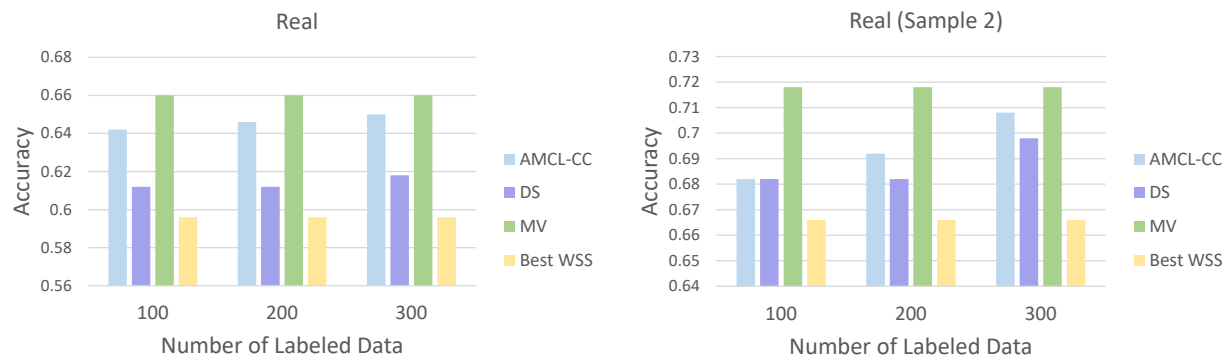


Figure 7. Experimental results on both samples of Domain Net for the Real domain as we vary the amount of labeled data. Each method uses 500 unlabeled data.

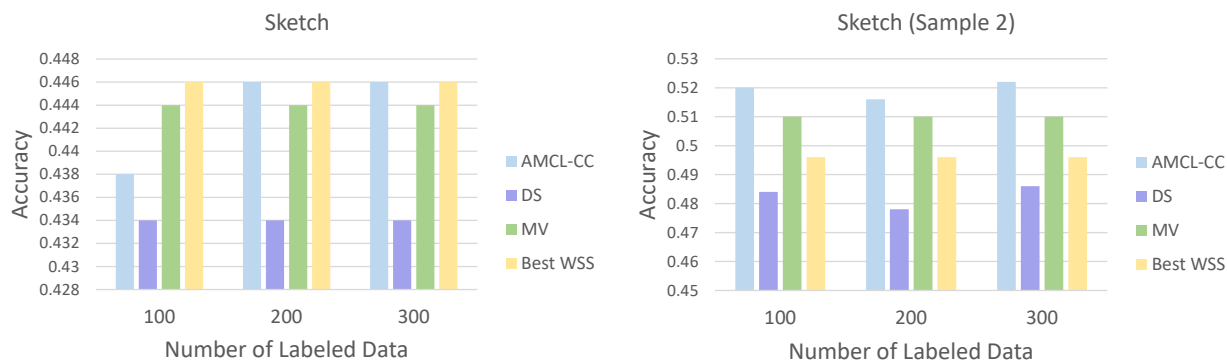


Figure 8. Experimental results on both samples of Domain Net for the Sketch domain as we vary the amount of labeled data. Each method uses 500 unlabeled data.

D. Experiments on Synthetic Data

We run synthetic experiments to show that our method is robust with respect to the addition of correlated weak supervision sources. Similar experiments have been done for ALL by [Arachie & Huang \(2019\)](#).

We consider a multiclass classification task over 5 classes, and 25 weak supervision sources ϕ_1, \dots, ϕ_{25} . In this classification task, each item of the domain \mathcal{X} has a unique true label. Given an item $x \in \mathcal{X}$, for $i \in 1, \dots, 10$, the weak supervision source ϕ_i returns the correct label with probability $1/2$, and a random label with probability $1/2$. The output of the weak supervision source ϕ_i is independent to the output of the weak supervision sources ϕ_j for $j \in \{1, \dots, 10\} \setminus \{i\}$. Therefore, the weak supervision source ϕ_i is correct with probability $\frac{1}{2}(1 + \frac{1}{k})$. For $i = 11, \dots, 25$, the weak supervision sources ϕ_i outputs the same result than the weak supervision source ϕ_1 . Note that the weak supervision sources $\phi_{11}, \dots, \phi_{25}$ do not provide any additional information with respect to the target classification task, as they add redundant constraints to the set of feasible labelings \mathbb{Y}° . The majority vote of the weak supervision sources ϕ_1, \dots, ϕ_{25} is highly affected by these dependencies, and it is very likely to provide the same answer as ϕ_1 , which is only $\frac{1}{2}(1 + \frac{1}{k})$ accurate on average. On the other hand, the majority vote of the weak supervision sources ϕ_1, \dots, ϕ_{10} would improve upon the individual accuracy of the weak supervision sources, as their output is independent.

We use 500 unlabeled examples, run experiments varying the amount of labeled data, and show that our method AMCL-CC is robust against those dependencies. For the sake of these experiments, as we want to use very small amount of labeled data, we set $\gamma = 0$ when building the constraints for \mathbb{Y}° as in Lemma 1. The experimental results are reported in Table D. The table shows that AMCL-CC is robust with respect to dependencies among weak supervision sources, whereas majority vote is greatly affected by them. In fact, in this case the majority vote does not improve upon the individual accuracy of the weak supervision sources, which is on average $\frac{1}{2}(1 + \frac{1}{k}) = \frac{3}{5}$.

Table 1. We report the experimental results on the synthetic dataset. We report the accuracy obtained by our method AMCL-CC and the majority vote, when varying the amount of labeled examples (we report the average accuracy over 3 distinct runs).

Labeled Examples	AMCL-CC	Majority Vote
100	0.902	0.595
50	0.828	0.602
25	0.819	0.598