# Unsupervised Label Refinement Improves Dataless Text Classification

**Zewei Chu,**[1]   **Karl Stratos,**[2]  **Kevin Gimpel**[3]

[1]The University of Chicago, 5730 S Ellis Ave, Chicago, IL 60637, USA
[2]Rutgers University, 110 Frelinghuysen Road Piscataway, NJ 08854, USA
[3]Toyota Technological Institute at Chicago, 6045 S Kenwood Ave, Chicago, IL 60637, USA
zeweichu@gmail.com, karlstratos@gmail.com, kgimpel@ttic.edu

## Abstract

Dataless text classification is capable of classifying documents into previously unseen labels by assigning a score to any document paired with a label description. While promising, it crucially relies on accurate descriptions of the label set for each downstream task. This reliance causes dataless classifiers to be highly sensitive to the choice of label descriptions and hinders the broader application of dataless classification in practice. In this paper, we ask the following question: how can we improve dataless text classification using the inputs of the downstream task dataset? Our primary solution is a clustering based approach. Given a dataless classifier, our approach refines its set of predictions using $k$-means clustering. We demonstrate the broad applicability of our approach by improving the performance of two widely used classifier architectures, one that encodes text-category pairs with two independent encoders and one with a single joint encoder. Experiments show that our approach consistently improves dataless classification across different datasets and makes the classifier more robust to the choice of label descriptions. [1]

## 1   Introduction

Dataless text classification aims at classifying text into categories without using any annotated training data from the task of interest. Prior work (Chang et al. 2008; Song and Roth 2014) has shown that with effective ways to represent texts and labels, dataless classifiers can perform text classification on unbounded label sets if suitable descriptions of the labels are provided.

There have been many previous efforts in dataless or zeroshot text classification (Dauphin et al. 2013; Nam, Mencía, and Fürnkranz 2016; Li et al. 2016; Ma, Cambria, and Gao 2016; Shu, Xu, and Liu 2017; Fei and Liu 2016; Zhang, Lertvittayakumjorn, and Guo 2019; Yogatama et al. 2017; Mullenbach et al. 2018; Rios and Kavuluru 2018; Meng et al. 2019). Many different settings have been considered across this prior work, and some have used slightly different definitions of dataless classifiers. In this paper, we use the term "dataless text classification" to refer to methods that: (1) can assign scores to any document-category pair, and (2) do not require any annotated training data from downstream tasks. A dataless classifier can therefore be immediately adapted to a particular label set in a downstream task

dataset by scoring each possible label for a document and returning the label with the highest score. Dataless classifiers are typically built from large-scale freely available text resources such as Wikipedia (Chang et al. 2008; Yin, Hay, and Roth 2019).

A well known problem with dataless classifiers is that the choice of label names has a significant impact on performance (Chang et al. 2008). As dataless classifiers rely purely on the label descriptions in a downstream task, there is typically no tailoring or fine-tuning of the classifier for a given dataset. A poor choice of label descriptions could jeopardize the performance of dataless classifiers on a particular text classification task, so prior work has addressed this with modifications to label descriptions. Chang et al. (2008) manually expand the label names of the 20 newsgroups dataset. Yin, Hay, and Roth (2019) expand labels by their WordNet definitions.

To illustrate the problem, Table 1 shows various choices of label names when applying a dataless classifier to the 4-class AG News dataset. When we change the descriptions of the four labels, performance of our dataless text classifier[3] changes drastically. The broader application of dataless text classifiers is hindered by their fragility caused by the choice of label descriptions. It is unclear how practitioners should choose label descriptions for practical use.

In this paper, we ask the following question: how can we improve dataless text classification provided the *unlabeled* set of input texts for the downstream task in addition to its label descriptions? Our approach, which we refer to as unsupervised label refinement (ULR), is based on $k$-means clustering. We develop variations of our approach so that it can be applied to different styles of dataless text classifiers to improve their performance. Table 1 shows results when applying ULR to our dataless text classifier. In all cases, accuracies improve after applying ULR, with larger gains when using weaker label descriptions.

Our contributions in this paper are:

- We propose unsupervised label refinement (ULR), a $k$-means clustering based approach to improve dataless classifiers.
- We develop variations of ULR that can be applied to different model architectures of dataless classifiers. Exper-

---

[1]Code and data available at https://github.com/ZeweiChu/ULR.

[3]ROBERTA dual encoder architecture as described in Section 4

| choice of label names | **world** **sports** **business** **science technology** | international health finance technology | world politics sports business and finance science and technology | world news health business science and technology | world health and sports commerce science technology | world news sports finance science |
|---|---|---|---|---|---|---|
| dataless | 69.9 | 55.9 | 71.6 | 49.9 | 55.0 | 68.8 |
| dataless + ULR | 70.7 | 78.3 | 78.4 | 70.2 | 70.9 | 76.1 |

Table 1: Accuracy (%) of a ROBERTA dual encoder dataless classifier[2] on AG News with different choices of label names. The original label set (boldfaced) is "world", "sports", "business", and "sci/tech". The "dataless + ULR" row shows accuracies after applying unsupervised label refinement (details in Section 3).

iments on dual encoder and single encoder architectures show that ULR almost always improves performance.

- Experiments show that ULR improves robustness of dataless classification against choices of label names, making dataless classifiers more practically useful.

## 2 Background: Dataless Text Classification

Dataless text classification (Chang et al. 2008; Chen et al. 2015; Song and Roth 2014; Yin, Hay, and Roth 2019) aims at building a single, universal text classifier that can be applied to any text classification task with a given set of label descriptions. Dataless classifiers can be used on an unbounded set of categories. There is typically no tailoring or fine-tuning of the classifier for a dataset other than through specifying the label descriptions.

Since annotated data in the target task is not available for training, the choice of label descriptions plays a critical role in the performance of dataless classifiers (Chang et al. 2008; Yin, Hay, and Roth 2019). With a dataless classifier, a score is produced for each text-category pair, indicating their semantic relatedness. Text classification then becomes a ranking problem, i.e., picking the category that has the highest semantic relatedness with the text.

Several have used EXPLICIT SEMANTIC ANALYSIS (ESA) (Gabrilovich and Markovitch 2007) as text representations in dataless text classification (Chang et al. 2008; Song and Roth 2014; Wang, Domeniconi, and Hu 2009). Both label descriptions and text are encoded into ESA vectors. Cosine similarity is used to compute scores between text and categories. Yin, Hay, and Roth (2019) directly compute text-category relatedness with a single BERT (Devlin et al. 2019) model. Chang et al. (2008) and Yin, Hay, and Roth (2019) exemplify two typical modeling choices for dataless classifiers, namely dual encoder and single encoder architectures, respectively. We will introduce them briefly and consider both types in our experiments.

**Dual encoder model.** With the dual encoder model, the category and text are fed into the encoder separately, each producing a vector representation. The text and category encoders could have shared or independent parameters. In our experiments, we always share parameters, i.e., we use the same encoder for both the categories and texts. A distance function takes both the category and text vectors and produces a scalar value. In our experiments, this scoring function can be either cosine distance or Euclidean (L2) distance.
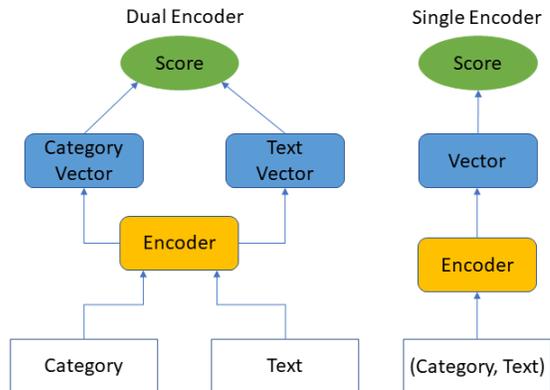


Figure 1: Dual encoder and single encoder architectures

**Single encoder model.** With a single encoder model, the category is combined with the text as a single sequence and fed into an encoder. The output of the encoder is a single vector that contains the information from both the category and the text. This vector can pass through a linear layer and produce a score for this particular document-category pair.

Figure 1 demonstrates the architectures of typical dual and single encoder models for text classification.

## 3 Unsupervised Label Refinement (ULR)

In this section, we introduce unsupervised label refinement (ULR). ULR uses the components of a dataless classifier and refines representations of labels with a modified $k$-means clustering algorithm. While dataless text classifiers are designed to handle an unbounded set of categories, they are used and evaluated on a particular set of documents with a set of labels. The idea of our approach is to leverage the assumption that the documents in a text classification dataset are separable according to the accompanying set of labels. That is, given a strong document encoder, the documents should be separable by label in the encoded space. This assumption is similarly made when performing clustering for unsupervised document classification (Liang and Klein 2009).

We use the set of unlabeled input texts to refine the predictions of our dataless classifiers via clustering. To better inform the algorithm, we initialize the clusters by using our dataless classifiers run on the provided label set for each

task. The algorithm takes on different forms for the dual and single encoder models. Details are provided below.

## ULR for Dual Encoder Architectures

In the setting of a dual encoder model, we propose to perform $k$-means clustering among text representations, i.e., of vectors produced by the text encoder. The assumption is that texts falling under the same category will be close in the semantic space. We want to adapt our dataless classifier's predictions based on the natural clustering structure in the encodings of the texts.

We show the $k$-means algorithm for dual encoder architectures in Algorithm 1. We use one encoder enc to encode

---

**Algorithm 1:** Unsupervised label refinement for dual encoder architectures

**Data:** documents $\mathcal{T}$, categories $\mathcal{C}$, encoder enc, scoring function score
initialize the centroids as $r_c = \text{enc}(c) \, \forall c \in \mathcal{C}$;
**while** *not converged* **do**
    **for** $t \in \mathcal{T}$ **do**
        **for** $c \in \mathcal{C}$ **do**
            $s_{tc} = \text{score}(\text{enc}(t), r_c)$;
        **end**
        $\text{pred}_t = \text{argmin}_c s_{tc}$;
    **end**
    **for** $c \in \mathcal{C}$ **do**
        $r_c = \left( \frac{\sum_{t:\text{pred}_t=c} \text{enc}(t)}{\text{count}(\{t:\text{pred}_t=c\})} + \text{enc}(c) \right)/2$;
    **end**
    $\text{objective} = \sum_t s_{t\text{pred}_t}$;
**end**
**Result:** $s_{tc}$, objective, $\text{pred}_t$, and $r_c$ of all iterations

---

texts and categories. We link centroids to categories and initialize the centroids using the encodings of the corresponding categories. The algorithm converges when no data point (text representation) updates its cluster assignment, i.e., the centroids stop updating. In our experiments, we run a maximum of 100 iterations. We perform model selection ("early stopping") based on the minimum value of objective among iterations.

Our $k$-means algorithm differs from standard $k$-means as our updated centroids are interpolated with the initial category embeddings. In standard $k$-means, the centroids are typically initialized randomly. In our case, since we link centroids to categories and use our encoder to provide initial centroids, we want to leverage the information in our category embeddings across clustering iterations. Therefore, we average the "new centroid" with the original category vector, which serves as a kind of regularization. In preliminary experiments we found this modification to stabilize performance so we use it in our experiments reported below.

## ULR for Single Encoder Architectures

In our single encoder architecture, a score is produced for each document-category pair indicating its relatedness. For each document, after exponentiating and normalizing the

score over all categories, we obtain a distribution indicating the probability of the document belonging to each category. For text $t$ and category $c$, we have a probability $p_{tc}$, where $\sum_c p_{tc} = 1$.

A straightforward way of classifying each document is to pick the category of which it has the highest probability score, i.e., $\text{argmax}_c p_{tc}$. Another way of interpreting this classification rule is to define $|\mathcal{C}|$ centroid vectors, each being an identity distribution one–hot$(c)$,[4] and pick the category having the minimum Jensen-Shannon Divergence (Lin 1991) with the document probability vector, i.e., $\text{argmin}_c \text{JS}(p_t, \text{one–hot}(c))$. Here JS is the function to compute the Jensen-Shannon Divergence between two distributions. Some prior work has explored using probability distribution based distance metric for clustering (Banerjee et al. 2005).

---

**Algorithm 2:** $k$-means algorithm for single encoder architectures

**Data:** documents $\mathcal{T}$, categories $\mathcal{C}$, scoring function score, function to compute JS divergence JS
initialize the centroids as $r_c = \text{one–hot}(c) \, \forall c \in C$;
compute each document's probability distribution over categories as
$[p_t]_c \propto \exp(\text{score}(t, c)) \, \forall t \in \mathcal{T}, c \in \mathcal{C}$;
**while** *not converged* **do**
    **for** $t \in \mathcal{T}$ **do**
        **for** $c \in \mathcal{C}$ **do**
            $s_{tc} = \text{JS}(p_t, r_c)$;
        **end**
        $\text{pred}_t = \text{argmin}_c s_{tc}$;
    **end**
    **for** $c \in \mathcal{C}$ **do**
        $r_c = \frac{\sum_{t:\text{pred}_t=c} p_t}{\text{count}(\{t:\text{pred}_t=c\})}$;
    **end**
**end**
**Result:** $\text{pred}_t$ of the last iteration

---

It is natural to represent each category as a distribution by a one-hot vector. However, in a real text classification problem, the semantics of a category is affected by how the annotators view it. For instance, a news document could relate to both "business" and "science & technology", though it will only have a single annotated category in the downstream task dataset. With this intuition, we propose to represent each category by a soft distribution over all the categories, but not necessarily a one-hot vector.

Algorithm 2 describes our $k$-means clustering approach applied on the single encoder model. In this algorithm our predicted categories will be the clustering assignments of the last iteration. Unlike with the dual encoder model, we do not do early stopping. We also do not use interpolated centroids as one-hot vectors may not necessarily be good category embeddings in this setting.

---

[4] Here we abuse the notation of one–hot$(c)$ to represent a one-hot vector that has a "1" at the index of category $c \in \mathcal{C}$.

# 4 Experimental Setup

In this section, we introduce the datasets we used for evaluation and the dual encoder and single encoder dataless classifiers we used to run ULR experiments.

## Evaluation

We use four text classification datasets spanning different domains for evaluation. They are: AG News[5] (AG), which uses 4 classes and covers the newswire domain; DBpedia (DBP; Lehmann et al. 2015), which has 14 classes and is from the domain of encyclopedias; Yahoo (Zhang, Zhao, and Lecun 2015), which has 10 classes and addresses categorizing questions in online question fora; and 20 newsgroups (20NG; Lang 1995), with 20 classes which are types of newsgroups.

We do not use any data or labels from the training set in our main experiments, but only rely on label descriptions. We use the official label names from these datasets, and only expand them if the original label name is provided as abbreviations such as "sci_tech". The exact label names we used are in the appendix.

## Dataless Classifiers

We experiment with multiple dataless text classifiers that vary in terms of their complexity. Our simplest classifier uses an encoder that averages pretrained GloVe (Pennington, Socher, and Manning 2014) word embeddings. We also fine-tune a ROBERTA model (Liu et al. 2019) in both single and dual encoder settings, using ROBERTA-base (110M parameters). We choose ROBERTA instead of BERT as ROBERTA outperforms BERT in a variety of text classification tasks (Liu et al. 2019). We do not run experiments with traditional dataless classifiers such as ESA (Chang et al. 2008; Gabrilovich and Markovitch 2007) as ESA vectors are of extremely high dimension, making it computationally difficult to apply ULR. GloVe and ROBERTA produce lower dimensional vectors that are more computationally amenable to refinement. Also, we experimented with ESA and found that it does not perform as well as our ROBERTA based models.[6] Next we describe the details of the three dataless classifiers that we use in our experiments.

**GloVe Dual Encoder.** We use GloVe (Pennington, Socher, and Manning 2014) in the dual encoder setting, simply averaging word vectors to represent both the categories and the documents. We use the 300 dimensional GloVe vectors[7] trained on Common Crawl.[8] We experiment with two distance functions when using GloVe: cosine and L2.

---

[5]https://www.di.unipi.it/~gulli/AG_corpus_of_news_articles.html

[6]As a comparison to the results of the ROBERTA dual encoder in Table 3, ESA accuracies (%) are 71.2 for AG, 62.5 for DBP, 29.7 for Yahoo, and 25.1 for 20NG.

[7]http://nlp.stanford.edu/data/glove.840B.300d.zip

[8]https://commoncrawl.org/

**ROBERTA Dual Encoder.** The category $c$ and text $t$ are fed separately to ROBERTA using the formatting "[CLS] $c$ [SEP]" and "[CLS] $t$ [SEP]". We use the average of the final-layer hidden states produced by ROBERTA as category and text vectors. A scoring function takes both the category and text vectors to produce a scalar value. We experiment with dot product and L2 distance as scoring functions.

**ROBERTA Single Encoder.** The category $c$ and text $t$ are combined in the form "[CLS] $c$ [SEP] $t$ [SEP]", tokenized, and encoded using ROBERTA. We truncate $t$ to ensure the category-document pair is within 128 tokens. The vector representation of the "[CLS]" token (after a linear transformation and non-linear activation) is then passed to a linear layer to produce a score.

## Fine-tuning ROBERTA models

We use the NATCAT dataset (Anonymous, 2020) to fine-tune the ROBERTA models. NATCAT comprises document-category pairs from three resources: Wikipedia, Stack Exchange, and Reddit. Wikipedia documents are paired with their annotated categories and ancestor categories. Stack Exchange question and question descriptions are paired with their corresponding question domain. Reddit post titles are paired with their subreddit names. Additional details and sample instances from NATCAT are provided in the supplementary material. Each NATCAT document comes with positive and negative categories. A positive category describes the document, and a negative category is randomly sampled and is irrelevant to the document. The ROBERTA models are fine-tuned as binary classifiers to indicate whether a category is positive for a document.

We use the Huggingface framework (Wolf et al. 2019) to fine-tune all ROBERTA models, using 300k instances from NATCAT. The peak learning rate is set to be 0.00002, and we perform learning rate warmup for 10% of the training steps and then linearly decay the learning rate. We set the random seed to be 1 in all experiments.

**Fine-tuning ROBERTA dual encoder.** While many other combinations of score and loss function could be considered, we report results with two particular combinations: dot product paired with binary cross entropy and L2 distance paired with a contrastive hinge loss. When dot product is used, ROBERTA is fine-tuned to minimize binary cross entropy between $\text{score}(\text{enc}(c), \text{enc}(t))$ and a binary label $y \in \{0, 1\}$.

When using L2 distance, we fine-tune ROBERTA to minimize a contrastive hinge loss:

$$x_p = \text{score}(\text{enc}(c_p(t)), \text{enc}(t)) \tag{1}$$

$$x_n = \text{score}(\text{enc}(c_n(t)), \text{enc}(t)) \tag{2}$$

$$\text{loss} = \max(x_p + \gamma - x_n, 0) \tag{3}$$

where $\text{score}$ is the squared L2 distance, $c_p(t)$ is a positive category for text $t$, $c_n(t)$ is a negative category, and $\gamma$ is a parameter indicating the margin. This loss aims to make negative category-text pairs have higher squared L2 distance than negative pairs by the margin.

|        |          | AG   | DBP  | Yahoo | 20NG | avg  |
|--------|----------|------|------|-------|------|------|
| cosine | baseline | 66.1 | 43.5 | **29.4** | 29.8 | 42.2 |
|        | + ULR    | **78.3** | **46.4** | 27.9 | **30.8** | **45.9** |
| L2     | baseline | 40.5 | 15.7 | 14.5 | 19.5 | 22.6 |
|        | + ULR    | **58.7** | **35.7** | **23.3** | **25.9** | **35.9** |

Table 2: Accuracies (%) when applying ULR to the GloVe dual encoder architecture ("baseline"), for two different choices of distance function (cosine and L2).

|        |          | AG   | DBP  | Yahoo | 20NG | avg  |
|--------|----------|------|------|-------|------|------|
| cosine | baseline | **74.0** | 84.6 | 52.3 | 36.4 | 61.8 |
|        | + ULR    | 73.7 | **93.3** | **54.3** | **40.3** | **65.4** |
| L2     | baseline | 69.9 | 78.8 | 55.7 | **37.8** | 60.6 |
|        | + ULR    | **70.7** | **90.5** | **61.3** | 37.4 | **65.0** |

Table 3: Accuracies (%) when applying ULR to the ROBERTA dual encoder architecture ("baseline").

**Fine-tuning ROBERTA single encoder.** When used as a single encoder, ROBERTA is fine-tuned to minimize binary cross entropy between the predicted score $\text{score}(\text{enc}(c, t))$ and a binary label $y$, where $\text{score}$ is a linear function that transforms the vector into a scalar score.

## 5 Experimental Results of ULR

**Dual Encoder Models.** With the dual encoder models, we used two sets of distance measures. The first distance is the cosine distance of two vectors. In this case, we always normalize the vector representations before applying ULR. The second distance measure is the L2 distance.

Table 2 shows results for the GloVe dual encoder model with two distance functions. Except for the single case of cosine distance with Yahoo, all accuracies improve, with some improving by large amounts (up to 20% absolute).

With ROBERTA dual encoder model, the choice of distance function matches the scoring function we used when fine-tuning ROBERTA, i.e., cosine distance is used with dot product scoring and L2 distance is used with the contrastive hinge loss. Table 3 shows results of ULR in the dual encoder setting with ROBERTA encoders. ULR improves accuracies by 3.6% to 4.4% on average, and the improvements are consistent across distance functions and datasets, except for the cases of 20 NEWSGROUPS with L2 distance and AG with cosine distance, which show slight degradations.

**Single Encoder Model.** Table 4 summarizes the results of applying ULR to the ROBERTA single encoder architecture. ULR improves performance across all four datasets, ranging from 0.5% for 20NG up to 6.8% on DBP.

**Label Ensembles.** Finding the best choice of label names for dataless classifiers can be difficult without labeled data.

|          | AG   | DBP  | Yahoo | 20NG | avg  |
|----------|------|------|-------|------|------|
| baseline | 72.6 | 81.8 | 59.3 | 36.0 | 62.4 |
| + ULR    | **75.1** | **88.6** | **60.0** | **36.5** | **65.1** |

Table 4: Accuracies (%) when applying ULR to the ROBERTA single encoder architecture ("baseline").

|        |          | AG   | DBP  | Yahoo |
|--------|----------|------|------|-------|
| cosine | ensemble | 79.1 | 84.6 | **54.5** |
|        | + ULR    | **81.1** | **93.3** | 54.4 |
| L2     | ensemble | 72.4 | 79.0 | 56.5 |
|        | + ULR    | **73.4** | **90.7** | **61.7** |

Table 5: Accuracies (%) of ensemble predictions using 10 choices of label names using the ROBERTA dual encoder architecture ("ensemble"), and results when combining it with ULR.

Therefore, it is easier to supply multiple choices of label names for a given task.

We manually pick 10 different sets of label names for a particular task, generate category and text representations with our ROBERTA dual encoder model, and perform ULR. The exact choices of label names are in the appendix. The predicted scores of such 10 different settings are summed up as the final ensemble predictions. Table 5 presents the results on such ensemble predictions. Compared to Table 3, all accuracies on all tasks are improved. Even with a stronger starting point from ensembling, ULR still yields consistent improvements in accuracy, with the single case of Yahoo and cosine distance being the only one that does not improve.

## 6 Robustness to label noise and random label initialization

**Different choices of label names.** Dataless text classification tasks are known to suffer from high variance due to label descriptions. Performance can vary dramatically across different choices of category names. Table 1 shows some examples of how different choices of labels names drastically change the performance of our ROBERTA dual encoder model fine-tuned with dot product scoring function.

One advantage of ULR is that it is robust to label noise. That is, even given a poor choice of label names, ULR can help the classifier to partially recover some of the accuracy, as shown in the lower portion of Table 1.

To demonstrate that ULR for dataless text classification is robust to label noise, we now describe similar experiments on a larger scale. In particular, we test several hundred sets of label names for each of AG, DBP, and Yahoo. For each category, we randomly assign different label names to it, all having similar meaning. The exact combinations of label names are in the appendix.

We then perform dataless text classification and ULR with our ROBERTA dual encoder model, either under cosine distance or L2 distance. Table 6 shows the results of the average

|  |  | AG News | DBP | Yahoo |
|---|---|---|---|---|
| cos. | baseline | 69.6 | 79.4 | 43.8 |
|  | + ULR | **77.5** | **91.6** | **45.3** |
|  | # imp. | $\frac{213}{240} = 88.8\%$ | $\frac{3867}{3867} = 100\%$ | $\frac{732}{1015} = 72.1\%$ |
| L2 | baseline | 62.2 | 71.7 | 49.5 |
|  | + ULR | **75.2** | **85.2** | **59.3** |
|  | # imp. | $\frac{237}{240} = 98.8\%$ | $\frac{2960}{2963} = 99.9\%$ | $\frac{947}{947} = 100\%$ |

Table 6: Robustness analysis when varying choices of label names. "baseline" and "ULR" are average accuracy (%) of the ROBERTA dual encoder architecture among all category naming choices. "# imp." are the numbers and percentages of cases where the performance improves after ULR.
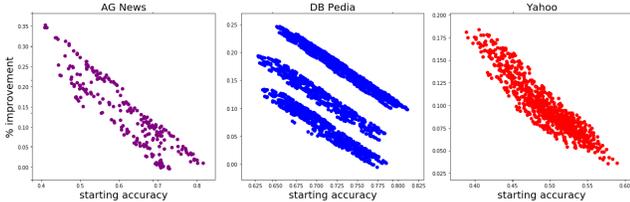


Figure 2: Accuracy (%) improvement when applying ULR compared to the ROBERTA dual encoder model (with Euclidean distance). The horizontal axis is the initial accuracy and the vertical axis is the absolute accuracy improvement after ULR.

performance gains before and after ULR. We also report the number and percentage of cases in which ULR improves accuracy. ULR improves the performance on average for all three tasks, and improves individual accuracies in the vast majority of cases. These results show that ULR is not only effective at improving the accuracy of dataless text classifiers across a wide range of label name sets, but also can help to mitigate harmful effects due to suboptimal label names.

We next study the relationship between the ROBERTA dual encoder model's initial accuracies and its accuracy gains after applying ULR. Figure 2 plots accuracy improvements vs. initial accuracies for the three datasets. We find that ULR gives larger gains when the initial accuracies are lower.

Among all the combinations of label names that we tried, we also report the oracle accuracies corresponding to the best combinations we found (without ULR). With cosine distance, they are 81.9% for AG, 87.5% for DBP, and 53.3% for Yahoo. With L2 distance, they are 81.5% for AG, 81.1% for DBP, and 58.9% for Yahoo. Directly applying ULR with the ROBERTA dual encoder classifier provides better performance on DBP and Yahoo, and competitive results on AG, as shown in Tables 3 and 5.

**Clustering with random initialization.** We also investigate the impact of our initialization in ULR. We consider a variation in which we randomly initialize centroids. Since we can no longer link clusters and categories, this task be-

|  |  | cosine | L2 |
|---|---|---|---|
| AG News | baseline | 37.3 | 39.9 |
|  | + ULR | **61.4** | **75.8** |
|  | # improved | $\frac{239}{240} = 99.6\%$ | $\frac{240}{240} = 100\%$ |

Table 7: Average 1-to-1 accuracy (%) of ULR with random initialization of centroids. "baseline" and "ULR" are average accuracy (%) of the ROBERTA dual encoder architecture among all random initialized centroids. "# imp." are the numbers and percentages of cases where the performance improves after ULR.

comes unsupervised text classification. A standard $k$-means algorithm is applied to update the centroids. Unlike Algorithm 1, in this experiment we do not interpolate the updated centroids with the initial centroids, as the initial centroids are randomly generated and do not have any meaning. The accuracy is calculated based on the oracle one-to-one mapping between final centroids and categories. This is often referred to as "one-to-one accuracy".

Table 7 presents results on AG with this experiment. We performed 240 trials with different random initialization of the centroids, and unsurprisingly, accuracy is improved in all cases. In this unsupervised setting, since we always choose the mapping between centroids and categories which maximizes accuracy, the initial ("baseline" in Table 7) centroid-category mapping may be different from the end ("+ ULR") mapping. With L2 distance, ULR with random initialization of centroids even outperform centroids initialized as category embeddings, which shows that unsupervised clustering is powerful in text classification with good text representations. However, we do not have the one-to-one mapping between category and final centroids in this pure unsupervised setting, so the results in Table 7 and Table 3, 5 are not directly comparable.

# 7 Unsupervised Label Refinement with Few Shot Learning

In this section, we apply ULR in the few shot learning setting. In particular, we draw 30 labeled instances for each category from the training split in the original datasets. We then further fine-tune our ROBERTA dual encoder model with the labeled instances. We adopt the hyperparameters of fine-tuning text classifiers from Wolf et al. (2019) and fine-tune for 3 epochs on these 30 labeled instances for each category. Then we apply ULR on the unlabeled test set in addition to the 30 labeled instances from each category.

Unlike Algorithm 1, these 30 labeled instances from each category are fixed to be assigned to their corresponding clusters. At every iteration, the update rule of the centroids will include the text vector representations of these 30 labeled instances.

Table 8 summarizes the results of combining 30 labeled instances for each category and ULR. Unsurprisingly, adding labeled instances improves accuracy by a large margin. Even so, ULR yields an additional improvement of

| | AG | DBP | Yahoo | 20NG | avg |
|---|---|---|---|---|---|
| baseline | 74.0 | 84.6 | 52.3 | 36.4 | 61.8 |
| + 30 labeled ins. | 81.5 | 97.1 | 66.0 | 58.1 | 75.7 |
| + ULR | **85.2** | **97.4** | **66.7** | **58.2** | **76.9** |

Table 8: Accuracies (%) when adding 30 labeled instances for each category with the ROBERTA dual encoder model ("baseline").

---

**Algorithm 3:** $k$-means algorithm for dual encoder models, with 30 labeled instances

**Data:** categories $\mathcal{C}$, unlabeled documents $\mathcal{T}$, documents $\mathcal{L}_c$ labeled with $c \; \forall c \in \mathcal{C}$, encoder enc, scoring function score

initialize the centroids as $r_c = \text{enc}(c) \; \forall c \in C$;
compute the labeled document centroids as
$r_c^l = \frac{\sum_{t \in \mathcal{L}_c} \text{enc}(t)}{|\mathcal{L}_c|} \; \forall c \in \mathcal{C}$;

**while** *not converged* **do**
  **for** $t \in \mathcal{T}$ **do**
    **for** $c \in \mathcal{C}$ **do**
      $s_{tc} = \text{score}(\text{enc}(t), \text{enc}(c))$;
    **end**
    $\text{pred}_t = \text{argmin}_c s_{tc}$;
  **end**
  **for** $c \in \mathcal{C}$ **do**
    $r_c = \left( \frac{\sum_{t:\text{pred}_t=c} \text{enc}(t)}{\text{count}(\{t:\text{pred}_t=c\})} + r_c^l + 2 \times \text{enc}(c) \right)/4$;
  **end**
  $\text{objective} = \sum_t s_{t\text{pred}_t}$;
**end**
**Result:** $s_{tc}$, objective, $\text{pred}_t$, and $r_c$ of all iterations

---

3.7% for AG, and also improves on the other datasets.

## 8  Augmenting Categories and Text

In this section, we ask the following questions: will more unlabeled text inputs from the evaluation task improve ULR? Will adding more category names benefit ULR?

When applying ULR with dataless text classifiers, more augmented categories can be added to be the centroids of new clusters. The goal of adding such augmented categories is to enrich the semantic space of all possible categories. For example, a document belonging to "science & technology" may be further classified into "physics", "math", or other scientific subjects. Adding such augmented categories into Algorithm 1 makes the category embedding cover broader and finer semantic spaces. To accommodate augmented categories in Algorithm 1, we add more centroids which are updated across iterations. However, at prediction time, we only predict within the categories of the original label set.

We also consider adding extra unlabeled text inputs from the domain of the evaluation task. The assumption is that by adding extra text, the $k$-means clustering approach would gain more knowledge from the task domain, and the centroids would capture semantics of larger scale data.

| | | AG | DBP | Yahoo | 20NG | avg |
|---|---|---|---|---|---|---|
| cosine | baseline | 74.0 | 89.4 | 53.5 | 37.0 | 63.5 |
| | ULR | 73.7 | **93.3** | 54.3 | 40.3 | 65.4 |
| | + aug cats | **74.9** | 92.6 | **57.3** | **41.6** | **66.6** |
| | + aug text | 73.8 | 93.2 | 54.2 | 40.3 | 65.4 |
| L2 | baseline | 69.9 | 78.8 | 55.7 | 37.8 | 60.6 |
| | ULR | 70.7 | **90.5** | **61.3** | 37.4 | 65.0 |
| | + aug cats | **71.9** | 90.4 | 59.4 | **39.8** | **65.4** |
| | + aug text | 70.8 | 90.4 | **61.3** | 37.6 | 65.0 |

Table 9: Accuracies (%) when adding augmented text or categories to ULR with the ROBERTA dual encoder architecture ("baseline").

**Experiments and results.** Starting from our ROBERTA based dual encoder model as a baseline, we run ULR with augmented categories and text. Our augmented categories are the names of the 30 most popular Stack Exchange sites (Chu et al. 2020). These category names cover a broad range of topics. The exact list of augmented category names are in the appendix. As for augmented text, we pick at most 100k instances without labels from the training sets of the evaluation tasks. The experimental results are summarized in Table 9. Augmented categories give the most improvement to ULR. Adding augmented text to ULR has little impact, indicating that having more text does not always help in unsupervised label refinement.

## 9  Related Work

We discussed prior work on dataless and zero-shot text classification in Sections 1 and 2. We briefly introduce more related work in this section. Song et al. (2015) provide a text label refinement algorithm to adjust the label set with noisy and missing labels. There is also a wealth of prior work in semi-supervised text classification: using unlabeled text to improve classification performance (Nigam et al. 2000; Mukherjee and Awadallah 2020; Xie et al. 2019). These methods typically learn generally useful text representations from a large corpus of unlabeled text and use them for a specific target task with limited supervision (Howard and Ruder 2018; Devlin et al. 2019; Liu et al. 2019; Lan et al. 2020; Peters et al. 2018). Finally, in supervised classification, label descriptions are also exploited to improve text classifiers (Chai et al. 2020; Wang et al. 2019; Sun, Huang, and Qiu 2019).

## 10  Conclusion

In this paper, we have shown that our proposed $k$-means clustering based unsupervised label refinement is a simple but effective approach to improve the performance of dataless text classifiers. ULR can be applied in both single encoder or dual encoder architectures. This approach is robust against the choices of label names, making dataless text classification more useful for practitioners.

# References

Banerjee, A.; Merugu, S.; Dhillon, I. S.; and Ghosh, J. 2005. Clustering with Bregman divergences. *Journal of machine learning research* 6(Oct): 1705–1749.

Chai, D.; Wu, W.; Han, Q.; Wu, F.; and Li, J. 2020. Description based text classification with reinforcement learning. *ICML* .

Chang, M.-W.; Ratinov, L.; Roth, D.; and Srikumar, V. 2008. Importance of Semantic Representation: Dataless Classification. In *AAAI*.

Chen, X.; Xia, Y.; Jin, P.; and Carroll, J. 2015. Dataless text classification with descriptive LDA. In *AAAI*.

Chu, Z.; Chen, J.; Gimpel, K.; Faruqui, M.; Wang, M.; Chen, M.; and Si, X. 2020. How to Ask Better Questions? A Large-Scale Multi-Domain Dataset for Rewriting Ill-Formed Questions. In *AAAI*.

Dauphin, Y. N.; Tur, G.; Hakkani-Tur, D.; and Heck, L. 2013. Zero-shot learning for semantic utterance classification. *arXiv preprint arXiv:1401.0509* .

Devlin, J.; Chang, M.-W.; Lee, K.; and Toutanova, K. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In *NAACL*.

Fei, G.; and Liu, B. 2016. Breaking the closed world assumption in text classification. In *NAACL*.

Gabrilovich, E.; and Markovitch, S. 2007. Computing Semantic Relatedness using Wikipedia-based Explicit Semantic Analysis. In *IJCAI*.

Howard, J.; and Ruder, S. 2018. Universal Language Model Fine-tuning for Text Classification. In *ACL*.

Lan, Z.; Chen, M.; Goodman, S.; Gimpel, K.; Sharma, P.; and Soricut, R. 2020. ALBERT: A Lite BERT for Self-supervised Learning of Language Representations. In *International Conference on Learning Representations*. URL https://openreview.net/forum?id=H1eA7AEtvS.

Lang, K. 1995. Newsweeder: Learning to filter netnews. In *Proceedings of the Twelfth International Conference on Machine Learning*, 331–339.

Lehmann, J.; Isele, R.; Jakob, M.; Jentzsch, A.; Kontokostas, D.; Mendes, P. N.; Hellmann, S.; Morsey, M.; Van Kleef, P.; Auer, S.; et al. 2015. DBpedia–a large-scale, multilingual knowledge base extracted from Wikipedia. *Semantic Web* 6(2): 167–195.

Li, Y.; Zheng, R.; Tian, T.; Hu, Z.; Iyer, R.; and Sycara, K. 2016. Joint Embedding of Hierarchical Categories and Entities for Concept Categorization and Dataless Classification. *COLING* .

Liang, P.; and Klein, D. 2009. Online EM for unsupervised models. In *Proceedings of human language technologies: The 2009 annual conference of the North American chapter of the association for computational linguistics*, 611–619.

Lin, J. 1991. Divergence measures based on the Shannon entropy. *IEEE Transactions on Information Theory* 37(1): 145–151. ISSN 1557-9654. doi:10.1109/18.61115.

Liu, Y.; Ott, M.; Goyal, N.; Du, J.; Joshi, M.; Chen, D.; Levy, O.; Lewis, M.; Zettlemoyer, L.; and Stoyanov, V. 2019. RoBERTa: A Robustly Optimized BERT Pretraining Approach. *CoRR* abs/1907.11692. URL http://arxiv.org/abs/1907.11692.

Ma, Y.; Cambria, E.; and Gao, S. 2016. Label Embedding for Zero-shot Fine-grained Named Entity Typing. *COLING* .

Meng, Y.; Shen, J.; Zhang, C.; and Han, J. 2019. Weakly-supervised hierarchical text classification. In *AAAI*.

Mukherjee, S.; and Awadallah, A. H. 2020. Uncertainty-aware Self-training for Text Classification with Few Labels. *arXiv preprint arXiv:2006.15315* .

Mullenbach, J.; Wiegreffe, S.; Duke, J.; Sun, J.; and Eisenstein, J. 2018. Explainable prediction of medical codes from clinical text. *arXiv preprint arXiv:1802.05695* .

Nam, J.; Mencía, E. L.; and Fürnkranz, J. 2016. All-in text: Learning document, label, and word representations jointly. In *Thirtieth AAAI Conference on Artificial Intelligence*.

Nigam, K.; Mccallum, A.; Thrun, S.; and Mitchell, T. 2000. Text Classification from Labeled and Unlabeled Documents using EM. In *Machine Learning*.

Pennington, J.; Socher, R.; and Manning, C. D. 2014. GloVe: Global Vectors for Word Representation. In *EMNLP*.

Peters, M.; Neumann, M.; Iyyer, M.; Gardner, M.; Clark, C.; Lee, K.; and Zettlemoyer, L. 2018. Deep Contextualized Word Representations. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, 2227–2237.

Rios, A.; and Kavuluru, R. 2018. Few-Shot and Zero-Shot Multi-Label Learning for Structured Label Spaces. In *EMNLP*.

Shu, L.; Xu, H.; and Liu, B. 2017. DOC: Deep Open Classification of Text Documents. In *EMNLP*.

Song, Y.; and Roth, D. 2014. On Dataless Hierarchical Text Classification. In *AAAI*.

Song, Y.; Wang, C.; Zhang, M.; Sun, H.; and Yang, Q. 2015. Spectral Label Refinement for Noisy and Missing Text Labels. In *AAAI*.

Sun, C.; Huang, L.; and Qiu, X. 2019. Utilizing BERT for aspect-based sentiment analysis via constructing auxiliary sentence. *arXiv preprint arXiv:1903.09588* .

Wang, J.; Sun, C.; Li, S.; Liu, X.; Si, L.; Zhang, M.; and Zhou, G. 2019. Aspect sentiment classification towards question-answering with reinforced bidirectional attention network. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, 3548–3557.

Wang, P.; Domeniconi, C.; and Hu, J. 2009. Towards a Universal Text Classifier: Transfer Learning using Encyclopedic Knowledge. In *Data Mining Workshops*.

Wolf, T.; Debut, L.; Sanh, V.; Chaumond, J.; Delangue, C.; Moi, A.; Cistac, P.; Rault, T.; Louf, R.; Funtowicz, M.; and

Brew, J. 2019. HuggingFace's Transformers: State-of-the-art Natural Language Processing. *ArXiv* abs/1910.03771.

Xie, Q.; Dai, Z.; Hovy, E.; Luong, M.-T.; and Le, Q. V. 2019. Unsupervised data augmentation for consistency training. *arXiv preprint arXiv:1904.12848* .

Yin, W.; Hay, J.; and Roth, D. 2019. Benchmarking Zero-shot Text Classification: Datasets, Evaluation and Entailment Approach. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, 3905–3914.

Yogatama, D.; Dyer, C.; Ling, C.; and Blunsom, P. 2017. Generative and Discriminative Text Classification with Recurrent Neural Networks. In *arXiv*.

Zhang, J.; Lertvittayakumjorn, P.; and Guo, Y. 2019. Integrating Semantic Knowledge to Tackle Zero-shot Text Classification. In *NAACL*.

Zhang, X.; Zhao, J.; and Lecun, Y. 2015. Character-level Convolutional Networks for Text Classification. In *NIPS*.