

Label Supervised Contrastive Learning for Imbalanced Text Classification in Euclidean and Hyperbolic Embedding Spaces

Anonymous ACL submission

Abstract

Text classification is an important problem with a wide range of applications in NLP. However, naturally occurring data is imbalanced which can induce biases when training classification models. In this work, we introduce a novel contrastive learning (CL) approach to help with imbalanced text classification task. CL has an inherent structure which pushes similar data closer in embedding space and vice versa using data samples anchors. However, in traditional CL methods text embeddings are used as anchors, which are scattered over the embedding space. We propose a CL approach which learns key anchors in the form of label embeddings and uses them as anchors. This allows our approach to bring the embeddings closer to their labels in the embedding space and divide the embedding space between labels in a fairer manner. We also introduce a novel method to improve the interpretability of our approach in a multi-class classification scenario. This approach learns the inter-class relationships during training which provide insight into the model decisions. Since our approach is focused on dividing the embedding space between different labels we also experiment with hyperbolic embeddings since they have been proven successful in embedding hierarchical information. Our proposed method outperforms several state-of-the-art baselines by an average 11% F1. Our interpretable approach highlights key data relationships and our experiments with hyperbolic embeddings give us important insights for future investigations. We will release the implementation of our approach with the publication.

1 Introduction

A common way of approaching the text classification problem is training a model using pre-trained text embeddings as language features (Mikolov et al., 2013; Pennington et al., 2014; Devlin et al., 2018). These embeddings can be fine-tuned using the signals from an objective function to improve

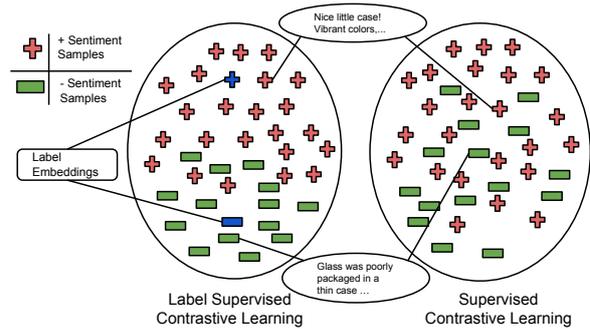


Figure 1: SCL can cause the embeddings for positive and negative sentiment text samples to be dispersed together in the embedding space (right illustration). Our approach in contrast utilizes the embedding space more effectively (left illustration). This is also shown in the form of Euclidean distance between embeddings of text samples of opposite sentiment. Our approach embeds these samples farther away from each other than SCL in terms of Euclidean distance: 13.2 vs. 3.2.

their efficacy for the classification task at hand. However, a common impediment to training a robust classifier is the fact that naturally occurring data is imbalanced. Since classifier predictions reflect the distribution of the training data, they can induce bias. There are many approaches proposed to address this issue, such as oversampling, under-sampling, using weighted objective functions or using situation/domain specific methods to improve the robustness of classification models (Chawla et al., 2002a; Tahir et al., 2012). Our work focuses on introducing a novel algorithm to deal with the challenges of imbalanced data.

Recent research shows an increasing use of contrastive learning (CL) to solve different problems in areas of computer vision and NLP (Gao et al., 2021a; Hénaff et al., 2019; Jaiswal et al., 2021). In this work, we explore CL to address the problem of imbalanced text classification. In general, CL uses anchors to embed similar samples closer in the embedding space while pushing dissimilar examples away. Unsupervised CL (Tian et al., 2019) tries to

066 contrast a data sample, called anchor, with every
067 sample in the batch while supervised CL (SCL)
068 (Khosla et al., 2020a) tries to utilize label informa-
069 tion and embed samples from the same class as the
070 anchor closer to each other. However, these CL
071 approaches rely on utilizing data embeddings as
072 anchors which are scattered over the embedding
073 space. We hypothesize that label embeddings, rep-
074 resenting a label category in the embedding space,
075 can be utilized as key anchors in CL. This allows a
076 model to embed data samples closer to their cate-
077 gory representations and results in a model learning
078 better embedding representations for the data. We
079 present an illustration in the figure 1 where we com-
080 pare how SCL divides embedding space in compar-
081 ison to our approach utilizing label embeddings
082 in a binary classification task. This shows that our
083 approach is able to achieve better class separa-
084 tion between data belong to different labels. This
085 is highlighted by the fact that distance between a
086 positive and negative text embedding pair is larger
087 when our approach is utilized in comparison to the
088 SCL.

089 Our proposed approach uses two embedding
090 modules to embed text and labels individually
091 which are fine-tuned using label-supervised CL
092 (LSCL). The text embedding module utilizes a self-
093 attention mechanism to make use of token-level
094 relations while the label embedding module is an
095 embedding layer. Both of these are fine-tuned using
096 our proposed CL objective. In addition, we also
097 experiment with hyperbolic embeddings, where
098 pre-trained model (e.g. BERT), provides repre-
099 sentations with hyperbolic structure (Chen et al.,
100 2021). We show that our approach outperforms sev-
101 eral SOTA and CL baselines in both Euclidean and
102 hyperbolic spaces. Finally, we also try to improve
103 the interpretability of our model by proposing a
104 modification to our approach which allows it to
105 represent inter-class relationships in an intuitive
106 manner for a multi-class classification task. We
107 structure our contributions in the following man-
108 ner: section 2 focuses on related work, section 3
109 highlights the background formulations on CL and
110 section 4 details our approach. Section 5 presents
111 generalization of our model to hyperbolic spaces,
112 section 6 describes our experiment setup followed
113 by the performance analysis in section 7. Finally,
114 we conclude by presenting the limitations and con-
115 clusion of our work.

2 Related Work 116

117 Data imbalance is a common problem and clas-
118 sification literature has adopted a variety of ap-
119 proaches to deal with the biases it might introduce.
120 One of these ways is oversampling of less frequent
121 data. SMOTE is the first minority oversampling
122 method (Chawla et al., 2002b). Iglesias et al. (2013)
123 presents a hidden markov model which generates
124 data from minority distribution. Other works focus
125 on the use of oversampling on the basis of sample
126 difficulty (Tian et al., 2021). Song et al. (2016)
127 combines the SMOTE technique with a K-Means
128 based undersampling algorithm to try and improve
129 classifier performance on an imbalanced dataset.
130 Some methods undersample the majority class sam-
131 ples to create a balanced data distribution for the
132 training process. Smith et al. (2013); Anand et al.
133 (2010) both present methods which use a notion of
134 sample difficulty to undersample the majority class
135 samples.

136 Some works rely on weighing the objective func-
137 tion to deal with data imbalance. The idea is to in-
138 crease the loss contribution for the minority classes
139 during the training. Cao et al. (2019); Chen et al.
140 (2016); Park et al. (2021) each presents a different
141 way of weighing the label-specific loss.

142 There is a third class of works which tries to
143 introduce novel algorithms focused on the data
144 imbalance problem. These methods avoid induc-
145 ing biases that might arise because of distribution
146 changes in data. An example is (Gao et al., 2021c)
147 which introduces a convolution based algorithm to
148 handle the class imbalance problem in data. Our
149 work fits in this category as we explore the use of
150 label-supervised CL to address this problem. An-
151 other example is Díaz-Vico et al. (2018), which
152 uses cost-sensitive learning to regularize the poste-
153 rior distributions for a given sample. This relies on
154 domain specific information which can be hard to
155 obtain in realistic scenarios (Krawczyk, 2016).

156 Lately, contrastive learning is being used in a
157 variety of tasks due to its effective utilization of
158 embedding space. Kang et al. (2021) present KCL
159 which is a variation of SCL algorithm (Khosla et al.,
160 2020b) and explores the use of contrastive learning
161 for learning balanced embedding spaces in the area
162 of computer vision. Lopez-Martin et al. (2022);
163 Zhang et al. (2022) present label-centered varia-
164 tions of CL methods but do not explore the data-
165 imbalance effects or the effect of computational
166 spaces on the model performance.

Hyperbolic spaces are becoming well-known for their superiority in embedding hierarchical information like WordNet graphs (Nickel and Kiela, 2017, 2018). This is because of their natural hierarchical structure. We view the classification task as a sub-class of hierarchical problem where a label embedding represents a category and each data sample is near to its label embedding. This is why we try to assess the performance of our model in the hyperbolic space as well. Another motivation for our work comes from Chen et al. (2021), which show that BERT embeddings contain hyperbolic structure between tokens by probing BERT embedding in hyperbolic spaces.

3 Contrastive Learning Overview

Contrastive learning tries to embed similar samples closer in the embedding space by trying to make the samples closer to their anchors. Formally, CL can be expressed as (Tian et al., 2019; Khosla et al., 2020b):

3.1 Contrastive Learning

We can define $\{(t_1, y_1), (t_2, y_2), \dots, (t_N, y_N)\} = D$ as a dataset consisting of a set of text $t_i = \{w_{i1}, w_{i2}, \dots, w_{is_n}\}$ and label pairs y_i , where s_n is the length of the text sample t_i and w_{ij} is the token representation corresponding to the j^{th} token in the text sample t_i . Given an embedding representation x_i for the text sample t_i , we can define contrastive learning objective L for mini-batches $B_k \subset D$ of size b_n as:

$$\frac{-1}{b_n} \sum_{x_i \in X_k} \log \frac{\exp(\text{sim}(x_i, x_i^+))}{\sum_{x_j \in \{x_i^+ \} \cup A(i)} \exp(\text{sim}(x_i, x_j))} \quad (1)$$

where sim is a similarity function (usually the dot product), $A(i) = \{x_j | x_j \neq x_i, x_j \in X_k\}$, X_k is set of text representations in the mini-batch B_k and x_i^+ is an augmented representation of the text sample t_i . This objective causes a model to learn embedding for x_i which are closer to its augmentation and pushes it away from other examples in the mini-batch.

4 Proposed Approach

We propose a supervised CL approach which uses label embeddings as anchors and causes the model to learn representations which are closer to their respective label representations or key anchors in the embedding space. An architecture diagram for our approach, Label Supervised Contrastive

Learning (LSCL), is presented in the figure 2 and its formulation L_{LSCL} is given as follows:

$$L_{LSCL} = \frac{1}{b_n} \sum_{x_i \in X_k} -\log \frac{\exp(\text{sim}(x_i, l_i))}{\sum_{l_j \in L} \exp(\text{sim}(x_i, l_j))} \quad (2)$$

where L is the set of all label representations. This approach embeds the text samples closer to their label embeddings in the embedding space. Labels for each text embedding can be predicted by choosing the label whose embedding is closest.

4.0.1 Increasing Interpretability Through Learning Inter-Class Relationships

In a multi-class classification scenario, sometimes label categories are related to each other, e.g. emotions *love* and *joy* are likely to be expressed in similar ways in many cases. In such cases it is hard to interpret how model embedded certain text samples in certain parts of the embedding space. Considering this we modify our approach to learn interpretable inter-class relationships, in form of a weight matrix, so these could be used to highlight the reasoning behind model decisions. This variation L_{LSCL-W} can be formulated as follows:

$$\frac{-1}{b_n} \sum_{x_i \in X_k} \log \frac{\exp(\text{sim}(x_i, l_i))}{\sum_{l_j \in L-l_i} w_{ij} \exp(\text{sim}(x_i, l_j))} \quad (3)$$

where $w_{ij} \in W^{|L| \times |L|}$ is a weight matrix we learn during the training process and $w_{ij} = 1$ when $i = j$. A problem here is that a learning method would just take the weight matrix W to zero. To prevent that, we add a Shannon Entropy (Shannon, 1948) regularization term to the objective which ensure that there is a relative difference in the magnitude of weights so the new objective L'_{LSCL-W} becomes:

$$L'_{LSCL-W} = L_{LSCL-W} + \lambda H(W_i) \\ H(W_i) = - \sum_{w_{ij} \in W_i} w_{ij} \log(w_{ij}) \quad (4)$$

where w_{ij} is the relation between labels l_i and l_j . The greater the weight the more difficult to separate data belonging to these two labels which is why the model assigns a higher weight to the contrastive weight of these labels. The λ is a term between 0 and 1 to control the contribution of entropy objective. W is not symmetric because of the data imbalance.

4.0.2 How Our Approach Helps with Data Imbalance

Our approach tries to bring the data samples in the closer to their respective labels and push the other

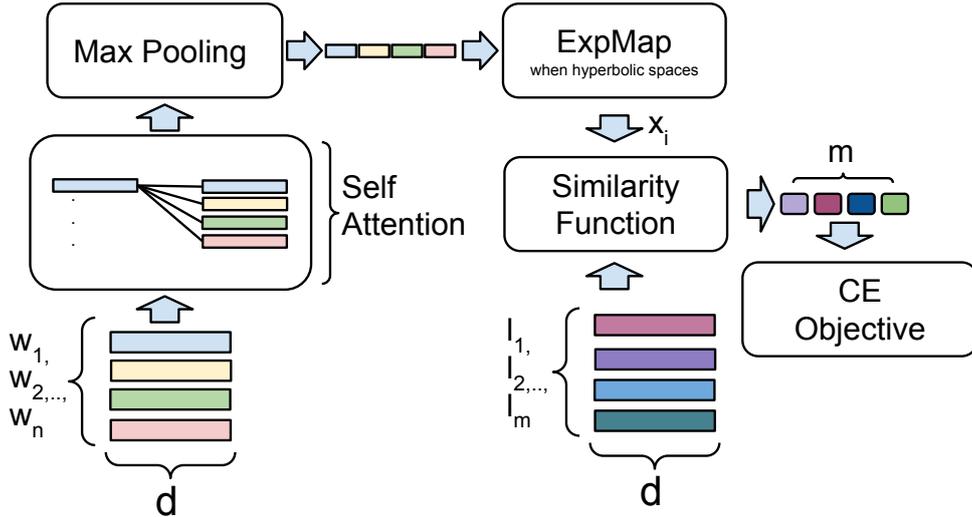


Figure 2: A batch of utterances is passed through a self-attention encoder to obtain text embeddings. These embeddings may be passed through an exponential map function to obtain embeddings in hyperbolic plane. Label embeddings are obtained by passing the input labels through a label embedding layer. These label embeddings are used as anchors in the CL objective which outputs loss signals for fine-tuning both the text encoder and label embedding layer together.

label embeddings away. This creates a push-pull effect for data samples w.r.t to the label embeddings. Both of these effects help improve the model performance. The data samples belonging to the majority class help improve the performance for the minority classes in this way as these samples push the minority label embeddings away as well while trying to get close to their respective label embeddings.

5 Generalization to Multiple Computational Spaces

Hyperbolic models show great promise for embedding hierarchical or graph structures (Nickel and Kiela, 2017, 2018). Our CL approach treats the classification problem as a hierarchical task by trying to learn the embedding regions for their respective labels. In addition, Chen et al. (2021) shows that pre-trained text embeddings contain hyperbolic structure. Due to these reasons we explore the effect of hyperbolic embeddings on our approach and show that these models perform competitively to their Euclidean counterparts and outperform all the baselines.

5.1 Manifold Centric Label Embeddings

We wanted to make use of the information encoded in the pre-trained textual representations and they are usually trained in Euclidean space. Due to this reason, we make use of hyperbolic exponential map

to obtain hyperbolic textual embeddings. However, label embeddings need not have any such restriction so we embed the labels in a manifold specific representation space. This entails that hyperbolic versions of our approach embed labels directly in the hyperbolic space so there is no need to use exponential map to obtain label embeddings.

5.2 Notion of Similarity

Contrastive learning uses a measure of similarity to embed similar examples closer to each other in a higher dimensional space. We generalize the notion of similarity between two vectors h and h' across Euclidean and hyperbolic manifolds, in an intuitive manner, as follows:

$$sim_{manifold}(h, h') = -d_{manifold}(h, h') \quad (5)$$

where $d_{manifold}$ represents the manifold specific distance function.

5.2.1 Vector Similarity in Euclidean Space

Following the formulation specified above the similarity function can be defined as:

$$sim_{eucl}(h, h') = - \sum_{i <= d} \sqrt{(h_i - h'_i)^2} \quad (6)$$

5.2.2 Vector Similarity in Hyperbolic Space

For our hyperbolic model variation, we use Lorentz formulation of Riemannian Manifolds because (Nickel and Kiela, 2018) suggests that Lorentz formulation of hyperbolic space is numerically stable

310 compared to the *Poincare*' formulation. The simi- 355
311 larity function for the hyperbolic variation is thus 356
312 given as: 357

$$313 \begin{aligned} \text{sim}_{\text{lorentz}}(h, h') &= -\text{arcosh}(-\langle h, h' \rangle_L) \\ \langle h, h' \rangle_L &= -h_0 h'_0 + \sum_{1 \leq i \leq d} h_i h'_i \end{aligned} \quad (7) \quad 358$$

314 6 Experiment Setup 359

315 We conduct experiments to compare the perfor- 358
316 mance of our approach on two classification tasks 359
317 with several baselines. In addition, we conduct 360
318 experiments to compare the performance of our 361
319 approach between hyperbolic and Euclidean em- 362
320 beddings. We rely on 256 dimensional variation of 363
321 BERT (Turc et al., 2019) to obtain the seed embed- 364
322 dings for our text encoder. 365

323 6.1 Datasets 366

324 We rely on two datasets for the purpose of our eval- 365
325 uation: 1) Amazon Reviews Sentiment Classifica- 366
326 tion (Keung et al., 2020) 2) Twitter Emotion Clas- 367
327 sification dataset¹. We create a binary sentiment 368
328 classification task from the former by splitting the 369
329 the review ratings into positive and negative classes. 370
330 Reviews with rating greater ≥ 4 are categorized as 371
331 positive and reviews with rating ≤ 2 are considered 372
332 negative. We induce a data imbalance of 9:1 for 373
333 positive and negative classes respectively to obtain 374
334 an imbalanced dataset containing a total of 15000 375
335 reviews.

336 Twitter emotion dataset is a multi-class data with 376
337 six emotions: sadness, joy, love, anger, fear, sur- 377
338 prise, contains a total of 20000 tweets and is natu- 378
339 rally imbalanced. Class ratios for both datasets are 379
340 given in the tables 1 and 2. 380

341 6.2 Model Parameters 381

342 Figure 2 shows the architecture of the text encoder 382
343 we use for CL. We utilize a self-attention layer 383
344 to embed the text embeddings. When we need to 384
345 obtain the hyperbolic embeddings we utilize the 385
346 exponential map operation to project the euclidean 386
347 embeddings into the hyperbolic space. We seed 387
348 our text embedding layer with BERT embeddings 388
349 which improves the training time of the model dur- 389
350 ing fine-tuning with CL. The right side of the archi- 390
351 tecture diagram shows the label embeddings which 391
352 are used to computer similarity with the text embe- 392
353 ddings. These embeddings are fine-tuned using the 393
354 LSCL training objective shown in the section 4. 394

¹<https://huggingface.co/datasets/emotion>

355 We use a prefix of E or H to indicate whether the 356
357 model utilizes euclidean embeddings or hyperbolic 358
359 ones respectively. 360

361 When using euclidean embeddings we fine-tune 362
363 our model using the Adam optimizer (Kingma and 364
365 Ba, 2014) while we use Reimannian SGD² to opti- 366
367 mize the hyperbolic weights as it relies on the ex- 368
369 ponential map to update the weights using Reiman- 370
371 nian gradients. Inspired from (Gao et al., 2021b), 372
373 we use a dropout layer (*rate: 0.1*) to obtain the 374
375 augmented representations when needed. We use 376
377 a learning rate of 10^{-3} for Adam and a learning 378
379 rate of 10^{-1} for Reimannian optimizer with a batch 379
380 size of 64. 380

369 6.3 Baselines 369

370 We compare our proposed approach with several 370
371 baselines. We divide the baselines in two groups: 371
372 1) SOTA baselines – baselines designed to help 372
373 with data imbalance in classification task; and 2) 373
374 CL baselines – baselines utilizing other versions of 374
375 contrastive. 375

376 6.3.1 Baselines for Imbalanced Classification 376

377 We use the following baselines to indicate the ad- 377
378 vantages of using a label-supervised CL approach 378
379 to deal with the problem of class imbalance in a 379
380 classification task. 380

381 **SetConv:** Gao et al. (2021c) presents a convolu- 381
382 tion based method to learn better representations 382
383 for the minority class samples. It utilizes a minor- 383
384 ity class representative as anchor to learn kernel 384
385 weights during the training process. 385

386 **GILE:** Pappas and Henderson (2019) uses joint 386
387 embeddings obtained using a dimension-wise prod- 387
388 uct of text and label embeddings. Their approach 388
389 uses a fully-connected layer to score these joint 389
390 embeddings and makes use of binary cross-entropy 390
391 objective to train the model. 391

392 **BertGCN:** Lin et al. (2021) treats the textual 392
393 data as a graph of token and document representa- 393
394 tions. The graph encodes token-level information 394
395 using measures like tf-idf and documents using 395
396 BERT representations. The approach utilizes a 396
397 graph convolution operation to obtain a vector rep- 397
398 resentation for a given text document. 398

399 6.3.2 Contrastive Learning Baselines 399

400 We utilize the following CL approaches to highlight 400
401 the advantages of utilizing our CL approach in a 401

²<https://github.com/facebookresearch/poincare-embeddings>

classification task.

K-Contrastive Learning: Kang et al. (2021) presents KCL, a variation of supervised contrastive learning in the domain of computer vision which learns balanced features spaces. Instead of using batch data samples as positive and negative anchors their approach samples k samples for each class from training data.

Supervised Contrastive Learning: SCL (Khosla et al., 2020a) is a CL approach which tries to contrast data samples from one class with data samples belonging to other classes while trying to bring the data samples from same classes closer to each other. As highlighted by the results presented below, this is a poor choice for imbalanced classification as skew in data distribution will create a bias in favor of majority class data when the model tries to bring samples from same class together.

7 Performance Analysis

We evaluate the performance of our approach on two tasks: binary sentiment classification and multi-class emotion classification. Both tasks highlight different aspects of our approach as a binary classification task with sufficient disparity in labels might be easier than a multi-class classification task which requires a model to learn inter-class relationships. For all our experiments, we measure the overall performance of a model using macro F1 score average because it equally weighs the model performance of the minority classes; hence reflects effect of data imbalance. Our key insights are:

- Our proposed CL approach is able to outperform the baselines in both computational spaces as shown in the tables 1 and 2).
- Euclidean version of our approach achieves the best overall performance as shown in the tables 1 and 2.
- We can improve model-decision interpretability by learning inter-class relationship weights. This is highlighted in the figure 4.
- Visualizing our approach in a 2-dimensional setting shows that hyperbolic version of our approach divides the embedding space fairly in the binary setting. This is highlighted in the figure 3.

7.1 Baseline Performance Comparison

We compare the performance of our approach with several contrastive learning and SOTA baselines

Model	Macro F1	Positive Class F1	Negative Class F1
Class Ratios		0.9	0.1
SOTA Baselines			
SetConv	0.682	0.888	0.476
GILE	0.706	0.951	0.462
BertGCN	0.702	0.948	0.455
Contrastive Learning Baselines			
SCL	0.594	0.95	0.237
KCL(k=5)	0.646	0.944	0.346
Our Approach			
HLSCL	0.72	0.930	0.511
ELSCL	0.779	0.959	0.6

Table 1: This table shows the per class F1 scores achieved by our model and their corresponding macro averages on Amazon Reviews Sentiment classification task. We show the results of both hyperbolic and euclidean models. The bold numbers represent the best performing model.

as stated in the section 6.3. In short, our approach outperforms the best SOTA baseline by a margin of 7% and 14% in the tasks of binary sentiment classification and multi-class emotion classification, respectively. These results are shown in the tables 1 and 2 respectively. In addition our approach does not sacrifice the majority class performance for a gain in minority class performance. This can be observed in both the binary and multi-class classification settings as our model consistently outperforms all the baselines in both overall and per-class performance, as highlighted in the table 1.

In the multi-class classification setting, the best performing baseline for the minority emotion *surprise* is BertGCN with a macro F1 of 38%. Our approach utilizing hyperbolic embeddings outperforms BertGCN by 7% in the minor class while achieving better performance in the majority classes – sadness and joy, as shown in the table 2.

Comparing the performance of our approach with CL baselines in the tables 1 and 2, specially SCL, shows the our approach to CL outperforms the other approaches in the task of imbalanced text classification.

7.2 Performance Comparison Among Computational Spaces

As described earlier, our formulation of the classification problem inspires us to test the performance of hyperbolic space embeddings in the tasks of binary and multi-class text classification tasks. In both cases, euclidean embeddings are better at embedding the text samples in the hidden space. However, hyperbolic variant of our approach still

Model	Macro F1	Sadness	Joy	Love	Anger	Fear	Surprise
Class Ratios		0.292	0.335	0.0815	0.135	0.121	0.0357
SOTA Baselines							
SetConv	0.361	0.425	0.469	0.297	0.314	0.378	0.283
GILE	0.401	0.607	0.675	0.242	0.42	0.325	0.138
BertGCN	0.554	0.712	0.778	0.330	0.571	0.55	0.383
Contrastive Learning Baselines							
SCL	0.285	0.555	0.646	0.0523	0.213	0.243	0.0
KCL(k=5)	0.299	0.508	0.63	0.0971	0.219	0.295	0.047
Our Approach							
HLSCl	0.621	0.757	0.774	0.553	0.597	0.595	0.451
ELSCl	0.695	0.793	0.836	0.611	0.704	0.637	0.591

Table 2: This table shows the per class and macro F1 scores achieved by our model on the task of emotion classification. We present both the hyperbolic and euclidean versions of our approach. The best performance numbers have been made bold.

outperforms all the baselines. This is evident from the results in the tables 1 and 2. In the case, of binary classification task, the highest performance difference between models in both spaces is minor, approximately 2% macro F1 score, but this difference increases in the case of multi-class sentiment classification task to approximately 8% macro F1. This shows that Euclidean models are better at the task of imbalanced classification even though hyperbolic models are effective classifiers.

7.3 Analyzing Embedding Space

We train our approach in both euclidean and hyperbolic spaces with 2-dimensional embeddings to visualize how our approach divides the embedding space. We find that hyperbolic variation of our approach divides the space more fairly between the minority and majority class in the binary classification case. This is interesting and may require further investigation in future work, as we fail to observe such a result when it comes to the multi-classification task. This could be because of data characteristics or may point to an innate trait of hyperbolic embeddings.

7.4 Interpreting Model Decisions Using Inter-Class Relationships

As described in the section 4, we proposed an approach to make model decisions interpretable by learning the inter-class relationships in the form of weights between 0 and 1. We train a model with the weighted variation of our approach and results highlight that model tries to distance embeddings which belong to similar emotions more than those belonging to different ones. This is apparent by looking at the weights in the figure 4 which shows that relationship weight between the positive labels *love* and *joy* (0.540) is higher in con-

trast to the weight between opposite ones *joy* and *sadness* (0.186). Similarly, weight between correlated emotions like *anger* and *surprise* (0.447) is higher than between emotions which are not correlated like *anger* and *love* (0.0558). This is interesting as this shows that model is capturing the fact that some emotions even though not similar are correlated. Another interesting insight is that the relationship between non-opposite categories like *anger* and *surprise* or *surprise* and *joy* are comparatively higher. This may point to an interesting characteristic of the data and alludes the fact that text expressing surprise can both be positive or negative. These results highlight that, along with improving interpretability, our approach can be utilized to highlight data specific characteristics and relationships. These may be used in data modeling or adopting data specific approaches for implementing practical solutions.

8 Limitations and Future Work

Our current approach is limited by the architecture of the label embedding layer. In our current implementation the label embeddings are obtained using a simple embedding which is fine-tuned during training along with text embedding module. In future works, we should experiment with more sophisticated ways to obtain label embeddings to check if we can improve our approach further.

Our approach, specially with hyperbolic embeddings, may have applications in hierarchical classification tasks where classes have a hierarchy and relationships between data samples and their classes are more complex. Such a task may be able to better utilize the natural structure of hyperbolic plane more effectively. In addition, our hyperbolic models fall behind in performance to their Euclidean



Figure 3: Space division by the 2-dimensional variation of our approach with negative text-samples. The figure shows how our approach divides the embedding space when trained with hyperbolic vs. euclidean embeddings. Rectangular space shows the normalized euclidean space while the circular shows a hyperbolic disk of Poincaré radius=1.

	Sadness	Joy	Love	Anger	Fear	Surprise
Sadness	--	0.0285	0.1939	0.436	0.1273	0.2144
Joy	0.1855	--	0.1708	0.399	0.0375	0.2072
Love	0.1398	0.5399	--	0.1284	0.1027	0.0893
Anger	0.1142	0.2254	0.0558	--	0.157	0.4475
Fear	0.3165	0.2517	0.1144	0.1889	--	0.1284
Surprise	0.1326	0.5381	0.0569	0.0683	0.204	--

Figure 4: Cells with darker red colors represent that model learns to separate these pairs more.

counterparts so more investigation is needed into how can hyperbolic spaces be used to learn effective classifiers.

Another significant limitation of our approach, lie in the problem formulation. One powerful aspect of CL approaches is that they do not need label information. However, we rely on the presence of label information in the corpus to learn label embeddings. This may not always be possible. In the future, we may be able to combine our approach with traditional CL approaches. This will involve dividing the embedding space during pre-training in the first phase. Using the results from this pre-training, we may be able to obtain key anchors by averaging out the embeddings in a region. These key anchors may then be used in an approach similar to ours to reduce noise in the CL training and better split the embedding space between different distributions in the data.

Finally, weighted variation of our CL objective is successful in quantifying relationship between class pairs. This provides additional insight into how our model is making decisions and improves interpretability. It even helps decipher information which is not obvious without a detailed look at data, like relationship between correlated emotions.

However, it does not help in improving the performance. Investigation into how this information can be used to learn better classifiers is another possible venue for future work. Similarly, using this information to design data specific solutions for deployment may offer another avenue for future research.

9 Conclusion

We present a novel CL approach which uses label embeddings as anchors for the task of imbalanced text classification in both the binary and multi-class classification settings. Our approach outperforms several baselines by a margin of 7% in the binary classification task and a margin of 15% in the multi-class classification task. In addition, we extend our approach to hyperbolic spaces, show its effectiveness in the task of imbalanced data classification. We also conduct a study of how our approach utilizes embedding space and show that it may be worth for future investigation that hyperbolic models divide the embedding space in a fairer manner than euclidean counterparts. Finally, we present a interpretable variation of our approach for multi-class classification which helps us draw important conclusions about data relationships.

References

Ashish Anand, Ganesan Pugalenthi, Gary Fogel, and Ponnuthurai Suganthan. 2010. [An approach for classification of highly imbalanced data using weighting and undersampling](#). *Amino acids*, 39:1385–91.

Kaidi Cao, Colin Wei, Adrien Gaidon, Nikos Arachiga,

612	and Tengyu Ma. 2019. Learning imbalanced datasets with label-distribution-aware margin loss . In <i>Advances in Neural Information Processing Systems</i> , volume 32. Curran Associates, Inc.	665
613		666
614		667
615		668
616	N. V. Chawla, K. W. Bowyer, L. O. Hall, and W. P. Kegelmeyer. 2002a. SMOTE: Synthetic minority over-sampling technique . <i>Journal of Artificial Intelligence Research</i> , 16:321–357.	669
617		670
618		671
619		672
620	Nitesh V. Chawla, Kevin W. Bowyer, Lawrence O. Hall, and W. Philip Kegelmeyer. 2002b. Smote: Synthetic minority over-sampling technique. <i>J. Artif. Int. Res.</i> , 16(1):321–357.	673
621		674
622		675
623		676
624	Boli Chen, Yao Fu, Guangwei Xu, Pengjun Xie, Chuanqi Tan, Mosha Chen, and Liping Jing. 2021. Probing BERT in hyperbolic spaces . <i>CoRR</i> , abs/2104.03869.	677
625		678
626		679
627		680
628	Kewen Chen, Zuping Zhang, Jun Long, and Hao Zhang. 2016. Turning from tf-idf to tf-igm for term weighting in text classification . <i>Expert Systems with Applications</i> , 66:245–260.	682
629		683
630		
631		
632	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding .	684
633		685
634		686
635		
636	David Díaz-Vico, Aníbal R. Figueiras-Vidal, and José R. Dorronsoro. 2018. Deep mlps for imbalanced classification . In <i>2018 International Joint Conference on Neural Networks (IJCNN)</i> , pages 1–7.	687
637		688
638		689
639		690
640	Tianyu Gao, Xingcheng Yao, and Danqi Chen. 2021a. Simcse: Simple contrastive learning of sentence embeddings .	691
641		692
642		693
643	Tianyu Gao, Xingcheng Yao, and Danqi Chen. 2021b. SimCSE: Simple contrastive learning of sentence embeddings . In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 6894–6910, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.	694
644		695
645		696
646		697
647		698
648		
649		
650	Yang Gao, Yi-Fan Li, Yu Lin, Charu Aggarwal, and Latifur Khan. 2021c. Setconv: A new approach for learning from imbalanced data .	699
651		700
652		701
653	Olivier J. Hénaff, Aravind Srinivas, Jeffrey De Fauw, Ali Razavi, Carl Doersch, S. M. Ali Eslami, and Aaron van den Oord. 2019. Data-efficient image recognition with contrastive predictive coding .	702
654		703
655		704
656		
657	E.L. Iglesias, A. Seara Vieira, and L. Borrajo. 2013. An hmm-based over-sampling technique to improve text classification . <i>Expert Systems with Applications</i> , 40(18):7184–7192.	705
658		706
659		707
660		
661	Ashish Jaiswal, Ashwin Ramesh Babu, Mohammad Zaki Zadeh, Debapriya Banerjee, and Fillia Makedon. 2021. A survey on contrastive self-supervised learning . <i>Technologies</i> , 9(1).	708
662		709
663		710
664		711
	Bingyi Kang, Yu Li, Sa Xie, Zehuan Yuan, and Jiashi Feng. 2021. Exploring balanced feature spaces for representation learning . In <i>International Conference on Learning Representations</i> .	712
		713
		714
		715
		716
	Phillip Keung, Yichao Lu, György Szarvas, and Noah A. Smith. 2020. The multilingual amazon reviews corpus. In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing</i> .	
	Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron Maschinot, Ce Liu, and Dilip Krishnan. 2020a. Supervised contrastive learning .	
	Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron Maschinot, Ce Liu, and Dilip Krishnan. 2020b. Supervised contrastive learning . <i>CoRR</i> , abs/2004.11362.	
	Diederik P. Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization .	
	Bartosz Krawczyk. 2016. Learning from imbalanced data: Open challenges and future directions . <i>Progress in Artificial Intelligence</i> , 5.	
	Yuxiao Lin, Yuxian Meng, Xiaofei Sun, Qinghong Han, Kun Kuang, Jiwei Li, and Fei Wu. 2021. BertGCN: Transductive text classification by combining GNN and BERT . In <i>Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021</i> , pages 1456–1462, Online. Association for Computational Linguistics.	
	Manuel Lopez-Martin, Antonio Sanchez-Esguevillas, Juan Ignacio Arribas, and Belen Carro. 2022. Supervised contrastive learning over prototype-label embeddings for network intrusion detection . <i>Information Fusion</i> , 79:200–228.	
	Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. Efficient estimation of word representations in vector space .	
	Maximilian Nickel and Douwe Kiela. 2017. Poincaré embeddings for learning hierarchical representations . <i>CoRR</i> , abs/1705.08039.	
	Maximilian Nickel and Douwe Kiela. 2018. Learning continuous hierarchies in the lorentz model of hyperbolic geometry . <i>CoRR</i> , abs/1806.03417.	
	Nikolaos Pappas and James Henderson. 2019. GILE: A generalized input-label embedding for text classification . <i>Transactions of the Association for Computational Linguistics</i> , 7:139–155.	
	Seulki Park, Jongin Lim, Younghun Jeon, and Jin Young Choi. 2021. Influence-balanced loss for imbalanced visual classification . In <i>Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)</i> , pages 735–744.	

- 717 Jeffrey Pennington, Richard Socher, and Christopher
718 Manning. 2014. [GloVe: Global vectors for word
719 representation](#). In *Proceedings of the 2014 Confer-
720 ence on Empirical Methods in Natural Language Pro-
721 cessing (EMNLP)*, pages 1532–1543, Doha, Qatar.
722 Association for Computational Linguistics.
- 723 C. E. Shannon. 1948. [A mathematical theory of commu-
724 nication](#). *Bell System Technical Journal*, 27(3):379–
725 423.
- 726 Michael R. Smith, Tony R. Martinez, and Christophe G.
727 Giraud-Carrier. 2013. An instance level analysis of
728 data complexity. *Machine Learning*, 95:225–256.
- 729 Jia Song, Xianglin Huang, Sijun Qin, and Qing Song.
730 2016. [A bi-directional sampling based on k-means
731 method for imbalance text classification](#). In *2016
732 IEEE/ACIS 15th International Conference on Com-
733 puter and Information Science (ICIS)*, pages 1–5.
- 734 Muhammad Atif Tahir, Josef Kittler, and Fei Yan. 2012.
735 [Inverse random under sampling for class imbalance
736 problem and its application to multi-label classifica-
737 tion](#). *Pattern Recognition*, 45(10):3738–3750.
- 738 Jiachen Tian, Shizhan Chen, Xiaowang Zhang, Zhiyong
739 Feng, Deyi Xiong, Shaojuan Wu, and Chunliu Dou.
740 2021. [Re-embedding difficult samples via mutual
741 information constrained semantically oversampling
742 for imbalanced text classification](#). In *Proceedings of
743 the 2021 Conference on Empirical Methods in Natu-
744 ral Language Processing*, pages 3148–3161, Online
745 and Punta Cana, Dominican Republic. Association
746 for Computational Linguistics.
- 747 Yonglong Tian, Dilip Krishnan, and Phillip Isola. 2019.
748 [Contrastive multiview coding](#).
- 749 Iulia Turc, Ming-Wei Chang, Kenton Lee, and Kristina
750 Toutanova. 2019. Well-read students learn better:
751 On the importance of pre-training compact models.
752 *arXiv preprint arXiv:1908.08962v2*.
- 753 Zhenyu Zhang, Yuming Zhao, Meng Chen, and Xi-
754 aodong He. 2022. [Label anchored contrastive learn-
755 ing for language understanding](#). In *Proceedings of
756 the 2022 Conference of the North American Chap-
757 ter of the Association for Computational Linguistics:
758 Human Language Technologies*, pages 1437–1449,
759 Seattle, United States. Association for Computational
760 Linguistics.