How much do Label Representations Matter for Image Classification?

What's in a Label?

Recent work (Chen et al, 2021) argues that the *representation of the label space* plays a critical role in a supervised learning problem in addition to the representation of the input space. They investigate this in context of image classification and endorse (1) high dimensional and (11) high entropy labels for:

1. Strong adversarial robustness of representation

sty dimension[®]

semantic embedding

2. Quicker convergence of models with fewer samples

We systematically attempt to reproduce their work and reach the conclusions:

- I. Label space manipulation seems to have *tangible* impact on training dynamics 2. Claimed dependence is inconsistent as multiple parameters are co-dependent

As such we note this area as a promising direction for further investigation, and propose a systematic study of label representations to stress-test and expand the claims of (Chen et al, 2021) on the quicker convergence with fewer samples claim. As such we will train only on 10% of the CIFAR-10 dataset throughout.



Gaussian Random Projections

Theorem: Johnson-Lindenstrauss: Given $0 < \varepsilon < 1$, a set X of *m* points in \mathbb{R}^N , and a number $n > 8 \ln(m) / \varepsilon^2$, there is a linear map $f : \mathbb{R}^N \to \mathbb{R}^n$ such that:

$$(1-\varepsilon) \|u-v\|^2 \le \|f(u)-f(v)\|^2 \le (1+\varepsilon) \|u-v\|^2$$

for all $u, v \in X$.

Dimensionality: We sample from $\sim N(0, 0.1*I)$ Gaussians of various dimensions.

Rotations: We apply gaussian projections to the standard one-hot label system to get random rotations. We use both the same dimension and the min. JL dimension.

Туре	Best Val. 600	Best Val. 100
Category (Regression)	40.64	12.52
Same Dim.	53.16	13.15
Min. JL Dim.	12.93	11.97



Label Distribution Entropies

We draw 512 dimensional labels with each element sampled independently to measure the impact of changing differential entropy. Contrary to [Chen et. al], we notice a negative trend:

Distribution	D.E. Formula
Exp.	1 - $\ln(\lambda)$
Laplace	$1 + \ln(2\sigma)$
Gaussian	$\ln(\sigma\sqrt{2\pi e})$





Language Models for Labels

If the label associated with a classification task carries a semantic meaning, then embeddings trained on this semantic meaning can be used as a potential label space. **Chen et al, 2021** experiment with BERT and Glove embeddings.

We briefly consider predicting a set of such self-supervised pre-text task labels as a vector and note that they tend to have strong performance. Please note that the BERT and CLIP embeddings are zero-centered and scaled with std. deviation.



A more general framework to preserve label-label relationships will make use of the implicit semantic distance between features rather than rely on the label to have a reasonable language realization.

Here we use contrastive learning to do "double work" to test this: in order to see if it helps for a label space to be organized as such, we: Learn a image representation where image embeddings are separated by label Form a label embedding for each label using the corresponding image embeddings and solve a supervised learning problem from the image itself

To learn the image representation, we can simply use data augmentation to get positives and treat all others as negatives, or we can use the labels to mine positives and negatives. We choose to use SimCLR and Supervised Contrastive Learning (SupCon), each from the respective family of methods.

$$l(i,j) = -\log \frac{\exp(z_i, z_j/\tau)}{\sum_{k=1, k \neq j}^{2N} \exp(z_i, z_k/\tau)} \qquad l^{out}(i) = \frac{-1}{|P(i)|} \sum_{p \in P(i)} \log \frac{\exp(z_i \cdot z_p/\tau)}{\sum_k \in A(i) \exp(z_i \cdot z_k/\tau)}$$

Then with the learnt encoder f_{i} , the corresponding label representation becomes:

We then calculate a label-label representation matrix and further use this with MDS (metric and non-metric) to learn embeddings of the labels in lower dimensions while preserving the relative distances between the labels. These sets of labels are then used for 10% data CIFAR-10 prediction as before. **Original v. Non Metric MDS**



Original v. Metric MDS



Ishaan Chandratreya, Katon Luaces

Columbia University

Protoxt datails	Dimonsions	Boot Val 600
Fretext details	Dimensions	best val. 600
Matrix Fact. on Word-Context	50	76.92
CBOW, Skip-Gram	100	76.84
Multimodal Retrieval	512	77.92
Language <i>cloze</i> , Sentence Prediction	768	78.14

Contrastive (Metric) Learning

$$l_i = \frac{1}{\sum_{j=1}^N \mathbf{1}[y_j = i]} \sum_{j=1, y_j = i}^N f(x_j)$$

Method





Method



enforced via a max-margin push-pull loss:

Method	Energy
OE (Vendrov et al.2016)	$ \max(0, v - i) $
EC (Ganea et al. 2018)	$\max(0,\epsilon(u,v) -$

for OE / EC

Geometry Probabilistic Embeddings

Metric Learning on labels





Figure from Khosla et al. (2020)



Current Deficit	Solution
High variance in low-dim. exp.	More seeds with more compute
CIFAR 10 Restriction	CIFAR100/Mini-ImageNet
Euclidean restrictive for hierarchy	Extension Survey in Paper
Space "Crowding"	Extension Survey in Paper
Require access to labels	Learn from label space