25th INTERNATIONAL CONFERENCE ON PATTERN RECOGNITION Milan, Italy 10 | 15 January 2021



Image Representation Learning by Transformation Regression

Xifeng Guo, **Jiyuan Liu***, Sihang Zhou, En Zhu, Shihao Dong National University of Defense Technology

Related Work



- Current self-supervised representation learning methods can be roughly separated into two categories, i.e. content based and manipulation based representation learning.
- The content based methods take advantage of data itself as the supervisory signal. They are also known as conventional unsupervised representation learning, like AE, VAE and GAN.
- The manipulation based methods define a group of manipulations (e.g., image rotation, color jitter, and inpainting) employed on images and generate corresponding labels.



Related Work • For example



Motivation



- Self-supervised learning learns representations by developing an auxiliary learning task.
- Existing methods usually model the auxiliary learning tasks as classification tasks with finite discrete labels.
- Insufficient supervisory signals restrict the learning ability.
- To make full use of the supervision from data, we design a regression model to predict the continuous parameters of a group of affine transformations.

Method



- Transform images and generate continuous labels.
- Train the neural network by the regression task.





Method





Method

$$\min_{\mathbf{W}} \frac{1}{n} \sum_{i=1}^{n} \left\| \mathcal{F} \left(\mathcal{T} \left(\mathbf{x}_{i}; y_{i} \right); \mathbf{W} \right) - y_{i} \right\|_{2}^{2}$$



Algorithm 1 Image Representation Learning Algorithm by Transformation Regression

Input: Image Dataset X; Transformation function $\mathcal{T}(\cdot, y)$ Neural Network $\mathcal{F}(\cdot, \mathbf{W})$.

Output: The parameters of the neural network **W**. Initialize the neural network parameter **W**;

for t in 1 to T do

for i in 1 to n do

Random sample a value y_i ;

Apply transformation: $\hat{\mathbf{x}}_i = \mathcal{T}(\mathbf{x}_i; y_i).$

Forward pass to get the output $\mathcal{F}(\hat{\mathbf{x}}_i; \mathbf{W})$.

end for

Compute the loss $L = \frac{1}{n} \sum_{i=1}^{n} \|\mathcal{F}(\hat{\mathbf{x}}_i; \mathbf{W}) - y_i\|_2^2$. Update the W by gradient descent.

end for

return \mathbf{W} .

• Datasets

	#examples	#training examples	#testing examples	#classes	image size
CIFAR10	60,000	50,000	10,000	10	$32 \times 32 \times 3$
CIFAR100	60,000	50,000	10,000	100	$32 \times 32 \times 3$
STL10	13,000	10,000	3,000	10	$32 \times 32 \times 3$
SVHN	99,289	$73,\!257$	26,032	10	$32 \times 32 \times 3$

• Protocol

- Train a network in a self-supervised way to learn representations.
- Extract representations from different layers.
- Use a classification model to validate the quality of representations (higher accuracy corresponds to better quality)



- Evaluation
 - Classification accuracy of a MLP on representations of different levels
- Observation
 - The method can learn good representations comparing Input with Block_i
 - The intermediate representation (Block_2) is the best.
 - More samples lead to better representation comparing STL10-10k and -100k

	Input	$Block_1$	$Block_2$	$Block_3$
CIFAR10 CIFAR100 STL 10, 10k	$66.84 \\ 39.60 \\ 54.20$	81.37 53.19	$84.88 \\ 57.06 \\ 72.67$	76.16 44.07 64.22
STL10-10k STL10-100k SVHN	$54.20 \\ 54.20 \\ 82.25$	$ \begin{array}{r} 8.53 \\ 71.00 \\ 92.29 \end{array} $	$72.07 \\77.03 \\94.28$	64.33 67.80 91.52



- Evaluation
 - Classification accuracy of a CNN on representations of different levels
- Observation
 - The method can learn good representations comparing Input with Block_i
 - The intermediate representation (Block_1) is the best.
 - More samples lead to better representation comparing STL10-10k and -100k

	Input	$Block_1$	$Block_2$	$Block_3$
CIFAR10 CIFAR100 STL10-10k STL10-100k	$93.37 \\70.84 \\77.93 \\77.93 \\06.12$	$93.65 \\71.42 \\79.93 \\81.03 \\06.54$	90.75 65.69 77.20 79.60 05.67	$81.44 \\ 49.01 \\ 65.50 \\ 67.17 \\ 02.62$

- Compare with SOTA
 - Supervised: Train the network with ground-truth labels then use a CNN to classify the learned representations of different levels.
 - Our method outperforms the other self-supervised methods and approaches the supervised method.

	CIFAR10	CIFAR100	STL10	SVHN
$DCGAN^*$ [12]	82.80	_	_	_
Split-Brain ^{\dagger} [27]	67.10	39.00	_	77.30
$Counting^{\dagger}$ [28]	50.90	18.20	—	63.40
AND^{\dagger} [21]	77.60	47.90	—	93.70
$RotNet^*$ [15]	91.16	_	_	—
TR (Ours)	93.65	71.42	79.93	96.54
Supervised	94.92	75.76	80.10	96.45



- Ablation Study
 - Use different transformations to train the network.
 - The composition of any two types of transformations leads to good performance
 - All three types of transformations achieve the best performance.

Rotation	Translation	Scaling	ACC (%)
$\begin{array}{c} 0\\ 0\\ 0\\ [-180, 180]\\ [-180, 180]\\ [-180, 180]\end{array}$	$\begin{array}{c} 0\\ [-10, 10]\\ [-10, 10]\\ 0\\ 0\\ [-10, 10]\end{array}$	$[0.5, 1.5] \\ 1.0 \\ [0.5, 1.5] \\ 1.0 \\ [0.5, 1.5] \\ 1.0 \\ 1$	90.9787.9791.4887.9992.2193.24
[-180, 180] [-180, 180]	[-10, 10] $[-10, 10]$	[0.5, 1.5]	93.65

Conclusion



- We propose a new image representation learning method by constructing a regression task whose target is to predict the continuous parameters of some transformations applied to the input image.
- Extensive experiments on various image datasets validate the effectiveness and discriminability of representation learned by our proposed transformation regression method.
- Future work
 - Exploring other types of transformations like image flipping, cropping, and color jitter.
 - Eliminating the edge effect (artifact) when applying some transformations like image rotation.

THANK YOU!