O 1 Unsupervised Visual Representation Learning by Synchronous Momentum Grouping Bo Pang¹, Yifan Zhang¹, Yaoyi Li², Jia Cai², and Cewu Lu^{1*}

¹ Shanghai Jiao Tong University {pangbo, zhangyf_sjtu, lucewu}@sjtu.edu.cn ² HuaWei Technologies Co., Ltd. {liyaoyi, caijiai1}@huawei.com

Reporter: Hongguang Zhu

1.supervised learning:

- simple and efficient
- expensive and time-consuming annotating costs
- non-ideal generalization (annotation bias)

2. unsupervised learning:

- Instance-level framework
- Clustering-based scheme

Instance-level framework

(1) contrastive learning methods:

- reduce the distance between two augmented views
- push apart view from negative samples





SimCLR, MOCOs...

(2) asymmetric network methods

- get rid of negative sample
- introduce predictor network
- stop-gradient operation

idea: only adopt positive samples \leftrightarrow mode collapse

- encode pervious sample information in weight of predictor
- serve as negative gradient during back-propagation







DINO: centering replaces the projection and predictor

Understanding Self-Supervised Learning Dynamics without Contrastive Pairs

(3) feature decorrelation methods

reduce the redundancy between differnet feature dimensions





Algorithm 1 PyTorch-style pseudocode for Barlow Twins. f: encoder network lambda: weight on the off-diagonal terms N: batch size D: dimensionality of the embeddings mm: matrix-matrix multiplication off diagonal: off-diagonal elements of a matrix eye: identity matrix for x in loader: # load a batch with N samples # two randomly augmented versions of x $y_a, y_b = augment(x)$ # compute embeddings z = f(y = a) # NxDz b = f(y b) # NxD# normalize repr. along the batch dimension z_a_norm = (z_a - z_a.mean(0)) / z_a.std(0) # NxD z b norm = (z b - z_b.mean(0)) / z_b.std(0) # NxD # cross-correlation matrix c = mm(z a norm.T, z b norm) / N # DxD # loss c diff = (c - eve(D)) . pow(2) # DxD# multiply off-diagonal elems of c diff by lambda off_diagonal(c_diff).mul_(lambda) loss = c diff.sum() # optimization step loss.backward()

optimizer.step()

Clustering-based scheme

- clustering: K-means / optimal transport
- but some restriction:
 - asynchronous two-satge training (DeepClustering)
 - no gradient propagate through group features
 - local equipartition constrainty to aviod degeneracy (SwAV)
 - not show an advantage in performance

instance-level discriminarion:

- supervisory signal:every sample as a category
- False negatives / Sampling Bias
- limited resource (large batch/ memory bank)

SwAV: Swapping Assignments between multiple Views



Group-level discrimination:

- supervisory signal hysteresis
- local equipartition decrease the validity of grouping
- not show an advantage in performance

This paper: SMoG

Novel:

- integrate the instance contrastive and clustering based methods
- contrastive group : allow gradient propagating through groups
- synchronous: instance grouping and representation learning

Achieve:

- avoid the false negatives and wipes off limiting factors of clustering methods
- surpasses the vanilla supervised-level performance for the first time
- multiple downstream tasks , CNN and Transformer backbones



Fig. 1: Group contrastive vs. instance contrastive learning. Compared with instance contrasting (red parts), group contrasting (blue part) learns representations through higher-level semantics, which can reduce the chance that the already learned similar instances are still treated as negative pairs (false negatives). The significant reduction of contrast elements after the grouping process also makes global contrasting easier to calculate. The colorized instances or groups denote the positive pairs and the grey ones are negatives.

gap : 3% -> surpass surpervise

instance contrastive learning :

(1) define group assigning function:

• update the groups {g}

• generate instance group c_i

$$L_i = -\log \frac{\exp(\sin(f_{\theta}(\mathbf{x}_i^a), \hat{f}_{\eta}(\mathbf{x}_i^b)/\tau)}{\sum_{\mathbf{x}_j \neq \mathbf{x}_i^a, \mathbf{x}_j \in \mathbf{Y} \subseteq \mathbf{X}} \exp(\sin(f_{\theta}(\mathbf{x}_i^a), \hat{f}_{\eta}(\mathbf{x}_j))/\tau)},$$

$$\mathbf{c}_i^a, \{\mathbf{g}\} \leftarrow \phi(f_\theta(\mathbf{x}_i^a)|\{\mathbf{g}\})$$

$$L_i = -\log \frac{\exp(\sin(\mathbf{c}_i^a, \mathbf{c}_i^b)/\tau)}{\sum_{\mathbf{g}_j \in \mathbf{G}} \exp(\sin(\mathbf{c}_i^a, \mathbf{g}_j)/\tau)},$$

How to design this function?

(1) c_i can replace the lastest $f_{\theta}(x_i)$ to participate in contrasting.

(2) group feature c_i should update synchronously with instance feature $f_{\theta}(x_i)$

 c_i^a and c_i^b tend to be same group feature -> fail to gather similar instances and get them closer



group assigning function:

$$\mathbf{c}_i^a, \{\mathbf{g}\} \leftarrow \phi(f_\theta(\mathbf{x}_i^a)|\{\mathbf{g}\})$$

global clustering each iteration: computing cost -> momentum grouping scheme

momentum grouping scheme:

$$\mathbf{c}_{i} = \operatorname{argmin}_{\mathbf{g}_{k}}(\operatorname{sim}(f_{\theta}(\mathbf{x}_{i}), \mathbf{g}_{k}))$$
$$\mathbf{g}_{k} \leftarrow \beta * \mathbf{g}_{k} + (1 - \beta) * \operatorname{mean}_{\mathbf{c}_{t} = \mathbf{g}_{k}} f_{\theta}(\mathbf{x}_{t}),$$

collapse:

- the scale of the groups are unbalanced
- the instance collapse into few groups

-> periodically grouping on the cached large feature set to relocate the groups





 (x_P^a) , new group features) $\leftrightarrow (x_E^b)$, old group feature)

(4) replace group features:

$$c_i^b = \operatorname{argmin}(\operatorname{sim}(x_E^b, g_k))$$

(2) χ_E^a update the groups \rightarrow new group features: $\mathbf{c}_i = \operatorname{argmin}_{\mathbf{g}_k}(\operatorname{sim}(f_{\theta}(\mathbf{x}_i), \mathbf{g}_k))$ $\mathbf{g}_k \leftarrow \beta * \mathbf{g}_k + (1 - \beta) * \operatorname{mean}_{\mathbf{c}_t = \mathbf{g}_k} f_{\theta}(\mathbf{x}_t),$

Table 1: Linear protocol results on ImageNet. ResNet50 is adopted. † denotes the model adopts the multi-crop training strategy. "acc" means accuracy.

model	epoch batchsize top1 acc top5 acc			
supervised	100ep	256	76.1	92.7
instance contrastin	ve metho	od		
ReSSL [74]	200	256	69.6	
MoCoV2 23	800	256	71.1	90.1
SimSiam 9	800	256	71.3	-
InfoMin Aug. 60	800	256	73.0	91.1
MoCoV3 10	400	4096	73.1	-
MoCoV3	800	4096	73.8	13 <u>1</u> 1
BYOL 21	400	4096	73.2	2 <u>-</u> 2
BYOL	800	4096	74.3	91.6
Barlow Twins 70	1000	2048	73.2	91.0
RELIC 47	1000	4096	74.8	92.2
SSL-HSIC 37	1000	4096	74.8	92.2
group-based metho	d			
DeepCluster	400	256	52.2	-
ODC [71]	400	256	57.6	
PCL [36]	200	256	67.6	0-
DeepClusterV2	400	4096	70.2	-
SwAV	400	4096	70.1	<u></u>
CoKe 54	800	1024	72.2	<u></u>
SMoG	400	2048	73.6	91.3
SMoG	800	4096	74.5	91.9
with multi-crop				
SwAV†	400	4096	74.6	8779
SwAV [†]	800	4096	75.3	0.70
DC-v2†	800	4096	75.2	-
DINO † 7	800	4096	75.3	-
UniGrad † 59	800	4096	75.5	-
NNCLR † 15	1000	4096	75.6	92.4
SMoGt	400	4096	76.4	93.1

Linear Evaluation:

nodel	backbone	throughput	param	top1 acc
upervised	SwinT	808	28	77.8 / 81.3
Supervised	DeiT-S/16	1007	21	77.5 / 79.8
MoBY	DeiT-S/16	1007	21	72.8
MoBY	SwinT	808	28	75.0
MoCoV3	DeiT-S/16	1007	21	72.5
MoCoV3	ViT-B/16	312	85	76.7
DINO	DeiT-S/16	1007	21	72.5
DINO †	DeiT-S/16	1007	21	77.0
EsViT 35	SwinT	808	28	70.5
EsViT †	SwinT	808	28	77.0
SMoG	SwinT	808	28	74.5
SMoG †	SwinT	808	28	77.7



larger Resnet: not an increasing superiority over the supervised one

model	backbone	param	top1 acc	top5 acc
Currentiand	Res50(x2)	188	77.8	93.8
Supervised	Res50(x4)	375	78.9	94.5
DVOI	Res50(x2)	188	77.4	93.6
BIOL]	Res50(x4)	375	78.6	94.2
C AV	Res50(x2)	188	77.3	-
SWAV F	Res50(x4)	375	77.9	-
CM-C	Res50(x2)	188	78.0	93.9
SMOG	Res50 (x4)	375	79.0	94.4

Table 4: Semi-supervised results on ImageNet.

mathad	Top-	1 Acc	c (%)	Top-	5 Ac	c (%)
method	1%	10%	100%	1%	10%	100%
$resnet-50 \times 1$						
Supervised	25.4	56.4	76.1	48.4	80.4	92.9
SimCLR	48.3	65.6	76.0	75.5	87.8	93.1
BYOL	53.2	68.8	77.7	78.4	89.0	93.9
SwAV	53.9	70.2	-	78.5	89.9	-
Barlow Twins	55.0	69.7	-	79.2	89.3	3 -
SSL-HSIC	52.1	67.9	77.2	77.7	88.6	93.6
NNCLR	56.4	69.8	-	80.7	89.3	-
SMoG	58.0	71.2	78.3	81.6	90.5	94.2
$resnet-50 \times 2$						
Supervised	223		77.8	2	127	93.8
SimCLR	58.5	71.7	-	83.0	91.2	0.0
BYOL	62.2	73.5	. 8	84.1	91.7	3
SMoG	63.6	74.4	80.2	85.6	92.4	95.2

Table 5: Transfer learning results on semantic segmentation task. We fine-tune the representations on VOC2012 and Cityscapes dataset. The segmentation model is FCN with ResNet-50.

model	Citys	capes	VOC-2012		
moder	mIoU	mAcc	mIoU	mAcc	
Supervised	73.83	82.56	73.59	83.74	
MoCoV2	74.30	83.37	70.86	80.37	
SwAV	74.80	83.01	74.97	84.27	
BYOL	74.90	83.73	74.76	84.37	
SMoG	76.03	83.97	76.22	85.01	

Table 6: Transfer learning results on object detection and instance segmentation tasks. we adopt COCO as the fine-tuning dataset. Mask RCNN with ResNet-50-FPN is the detection and segmentation model. We report the AP metrics.

	C	OCO d	et	COC) instand	ce seg.
method	AP^{bb}	AP_{50}^{bb}	AP_{75}^{bb}	AP^{mk}	AP $\frac{mk}{50}$	AP_{75}^{mk}
Rand Init	31.0	49.5	33.2	28.5	46.8	30.4
Supervised	38.9	59.6	42.7	35.4	56.5	38.1
InsDis 67	37.4	57.6	40.6	34.1	54.6	36.4
PIRL 45	38.5	57.6	41.2	34.0	54.6	36.2
MoCoV2	39.4	59.9	43.0	35.8	56.9	38.4
SwAV	38.5	60.4	41.4	35.4	57.0	37.7
DC-v2 6	38.3	60.3	41.3	35.4	56.7	38.0
BYOL	39.4	59.9	43.0	35.8	56.8	38.5
Barlow Twins	39.2	58.7	42.6	34.3	55.4	36.5
SMoG	40.1	61.6	43.7	36.9	58.7	39.3

53.7

53.4

67.2

30.2

42.1

65.8

67.2

Grouping quality:



each group, its entropy is $\sum_{i} -p_i * \log(p_i)$ where p_i is the ratio of instances that belong to class *i* (data annotation) in this group. A lower entropy means a group has more unitary semantics and is more meaningful. Thus, more groups with low

Fig. 5: Entropy of each groups. SMoG produces groups with much lower entropy which represents a better grouping quality in terms of high-level semantics.

Table 7: Ablation study on SMoG. We adopt ResNet50 as the backbone and train the unsupervised algorithm on ImageNet for 100 epochs without multicrop training strategy. We report the linear evaluation results (Top-1 accuracy).

(a) Momentum ratio β . Linearly decreasing schedule performs better.

Schedule	Top1 acc
fixed 0.99 β	65.9
$1.0 \; \beta ightarrow 0.9 \; \beta$	67.0
1.0 $\beta \rightarrow 0.99 \beta$	67.2
$1.0 \ \beta \rightarrow 0.999 \ \beta$	67.1

clustering is necessary. Tricks Top1 acc None 0.1

(b) Tricks dealing with collapse. The periodical

+ periodical clustering (pd)

+ pd & reset \hat{f}_n periodically

+ reset \hat{f}_n periodically

(c) Number of groups. 3k is	(d) Grouping algo. Momentum update
enough for SMoG.	better balances grouping and learning.

method # Groups Top1 acc Top1 acc 300 65.2 Randomly select Adopt latest 1000 67.0 67.2 Averaging update 3000 67.2 Momentum update 10000

Periodical clustering is important!

Grouping number is not sensitive Lastest feature as grouping feature is not good choice.

(RS): randomly select a group of latest instance features

(AL):
$$\mathbf{g}_k \leftarrow \operatorname{mean}_{\mathbf{c}_t = \mathbf{g}_k} f_{\theta}(\mathbf{x}_t)$$

(AU): $\mathbf{g}_k \leftarrow \mathbf{g}_k + (1/n) * (\operatorname{mean}_{\mathbf{c}_t = \mathbf{g}_k} f_{\theta}(\mathbf{x}_t) - \mathbf{g}_k)$