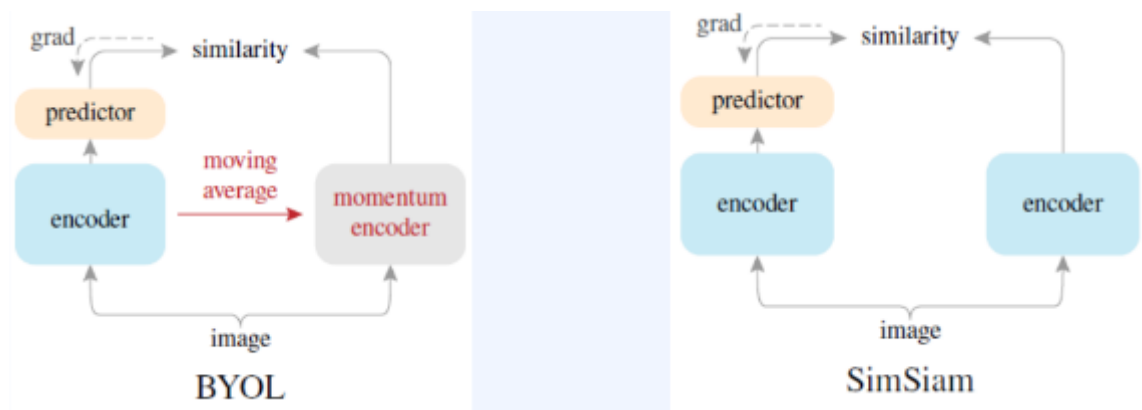


Simsiam

논문명 : Exploring Simple Siamese Representation Learning

- BYOL
 - 어떤 Transform을 적용할지 알 수 없어 평균을 적용
 - Prediction head는 모든 view의 mean값을 예측하도록 학습.
 - Projection값은 항상 다를 테니, Collapse가 일어나지 않음.
 - 사실 Momentum encoder의 역할은 없음.
 - Momentum update와 asymmetric approach를 통해서 collapse problem을 해결.
- Simsiam
 - Collapse problem을 해결하는 요소는 momentum update가 아닌, stop gradient였음을 밝힘.



BYOL이 성능이 더 좋지만 momentum encoder가 필요 없어 확장성 더 좋음

1. 학습 중단 시 encoder와 momentum encoder 둘 다 저장
2. 학습 후 중간 layer의 값 출력 시 일반 layer가 아닌 momentum encoder의 값 출력이 어려움

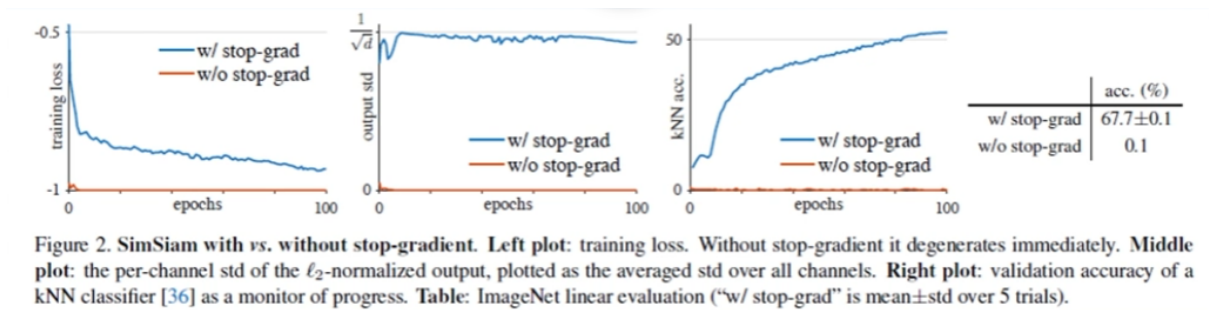
Stop-gradient

collapsing occurs when w/o stop-grad (loss = -1)

모델 출력의 표준편차가 0. (상수 출력)

KNN acc doesn't increase.

→ sort cut



Predictor

	pred. MLP h	acc. (%)
baseline	lr with cosine decay	67.7
(a)	no pred. MLP	0.1
(b)	fixed random init.	1.5
(c)	lr not decayed	68.1

Table 1. Effect of prediction MLP (ImageNet linear evaluation accuracy with 100-epoch pre-training). In all these variants, we use the same schedule for the encoder f (lr with cosine decay).

Batch size

batch size	64	128	256	512	1024	2048	4096
acc. (%)	66.1	67.3	68.1	68.1	68.0	67.9	64.0

Table 2. Effect of batch sizes (ImageNet linear evaluation accuracy with 100-epoch pre-training).

Similarity Function

	cosine	cross-entropy
acc. (%)	68.1	63.2

Symmetrization

	sym.	asym.	asym. 2×
acc. (%)	68.1	64.8	67.3

평가

method	batch size	negative pairs	momentum encoder	100 ep	200 ep	400 ep	800 ep
SimCLR (repro.+)	4096	✓		66.5	68.3	69.8	70.4
MoCo v2 (repro.+)	256	✓	✓	67.4	69.9	71.0	72.2
BYOL (repro.)	4096		✓	66.5	70.6	73.2	74.3
SwAV (repro.+)	4096			66.5	69.1	70.7	71.8
SimSiam	256			68.1	70.0	70.8	71.3

Table 4. **Comparisons on ImageNet linear classification.** All are based on **ResNet-50** pre-trained with **two 224×224 views**. Evaluation is on a single crop. All competitors are from our reproduction, and “+” denotes *improved* reproduction vs. original papers (see supplement).

pre-train	VOC 07 detection			VOC 07+12 detection			COCO detection			COCO instance seg.		
	AP ₅₀	AP	AP ₇₅	AP ₅₀	AP	AP ₇₅	AP ₅₀	AP	AP ₇₅	AP _{mask} ₅₀	AP _{mask}	AP _{mask} ₇₅
scratch	35.9	16.8	13.0	60.2	33.8	33.1	44.0	26.4	27.8	46.9	29.3	30.8
ImageNet supervised	74.4	42.4	42.7	81.3	53.5	58.8	58.2	38.2	41.2	54.7	33.3	35.2
SimCLR (repro.+)	75.9	46.8	50.1	81.8	55.5	61.4	57.7	37.9	40.9	54.6	33.3	35.3
MoCo v2 (repro.+)	77.1	48.5	52.5	82.3	57.0	63.3	58.8	39.2	42.5	55.5	34.3	36.6
BYOL (repro.)	77.1	47.0	49.9	81.4	55.3	61.1	57.8	37.9	40.9	54.3	33.2	35.0
SwAV (repro.+)	75.5	46.5	49.6	81.5	55.4	61.4	57.6	37.6	40.3	54.2	33.1	35.1
SimSiam, base	75.5	47.0	50.2	82.0	56.4	62.8	57.5	37.9	40.9	54.2	33.2	35.2
SimSiam, optimal	77.3	48.5	52.5	82.4	57.0	63.7	59.3	39.2	42.1	56.0	34.4	36.7

Table 5. **Transfer Learning.** All unsupervised methods are based on 200-epoch pre-training in ImageNet. *VOC 07 detection*: Faster R-CNN [30] fine-tuned in VOC 2007 trainval, evaluated in VOC 2007 test; *VOC 07+12 detection*: Faster R-CNN fine-tuned in VOC 2007 trainval + 2012 train, evaluated in VOC 2007 test; *COCO detection* and *COCO instance segmentation*: Mask R-CNN [18] (1× schedule) fine-tuned in COCO 2017 train, evaluated in COCO 2017 val. All Faster/Mask R-CNN models are with the C4-backbone [13]. All VOC results are the average over 5 trials. **Bold entries** are within 0.5 below the best.

1. Momentum encoder를 사용하지 않고, Positive only SSL 수행.
2. 어떤 요소가 Collapse problem을 해결하는 요소인지 실험적으로 증명.
3. BYOL에 비해 더 간단한 디자인으로, comparable한 성능을 도출.