# **Understanding Self-Supervised** Learning Dynamics without **Contrastive Pairs<sup>1</sup>** Mehmet F. Demirel

[1] Yuandong Tian, Xinlei Chen, Surya Ganguli Proceedings of the 38th International Conference on Machine Learning, PMLR 139:10268-10278, 2021.

July 29, 2021

## **Contrastive vs Non-contrastive Learning**

- Contrastive: Minimize the difference between positive pairs, and constrast negative pairs.
  - Former encourages modeling invariances while the latter prevents collapse.
- Non-contrastive: No negative pair contrasting, e.g., SimSiam and BYOL
  - Pair of Siamese networks: Want the views from the online + predictor network and the target network to match.

# **Non-contrastive Learning**

- Non-contrastive: No negative pair contrasting, e.g., SimSiam and BYOL.
  - Pair of Siamese networks: Want the views from the online + predictor network and the target network to match.
  - Target network is not trained with GD, i.e., a direct copy or a momentum encoder.
  - Don't require large batch size, memory queue, or negative pairs.



## Why do these models not collapse?

- Predictor and stop-gradient is essential in non-contrastive SSL. BYOL and SimSiam collapse without either of these.
- EMA or momentum encoder is not necessary in BYOL and SimSiam.
- Both BYOL and SimSiam say that the predictor must be optimal in achieving minimal error when predicting the target network's outputs from the online network's.
- BYOL suggests that no weight decay leads to unstable results.



online

## A simple model

Simple, bias-free, linear BYOL model.

Minimize 
$$J(W, W_p) = \frac{1}{2} \mathbb{E}_{x_1, x_2} \left[ \left\| W_p f_1 - \text{StopGrad} \right\| \right]$$

W: Linear online network

 $W_a$ : Linear target network

 $W_p$ : Linear predictor network





## **Training Dynamics in Closed Form**

Lemma 1. BYOL learning dynamics following Eqn. 1:  $\dot{W}_p = \alpha_p \left(-W_p W(X + X') + W_a X\right) W^{\intercal} - \eta W_p (2)$   $\dot{W} = W_p^{\intercal} \left(-W_p W(X + X') + W_a X\right) - \eta W \quad (3)$   $\dot{W}_a = \beta (-W_a + W) \quad (4)$ 

Here,  $X := \mathbb{E}[\bar{x}\bar{x}^{\mathsf{T}}]$  where  $\bar{x}(x) := \mathbb{E}_{x' \sim p_{aug}}(\cdot|x) [x']$  is the average augmented view of a data point x and X' := $\mathbb{E}_{x} [\mathbb{V}_{x'|x}[x']]$  is the covariance matrix  $\mathbb{V}_{x'|x}[x']$  of augmented views x' conditioned on x, subsequently averaged over the data x. Note that  $\alpha_{p}$  and  $\beta$  reflect *multiplicative learning rate ratios* between the predictor and target networks relative to the online network. Finally, the terms involving  $\eta$  reflect weight decay.



## **Exploration 1** Balancing of W and $W_p$ that comes from weight decay

- In BYOL and SimSiam, the match between the representations produced by online and target networks cannot be explained solely by the predictor weights.
- A non-zero weight decay, parametrized by  $\eta$ , will remove the second term on the RHS  $\implies$  more balance between online and predictor networks.

**Theorem 1** (Weight decay promotes balancing of the predictor and online networks.). Completely independent of the particular dynamics of  $W_a$  in Eqn. 4, the update rules (Eqn. 2 and Eqn. 3) possess the invariance  $W(t)W^{\mathsf{T}}(t) = \alpha_p^{-1}W_p^{\mathsf{T}}(t)W_p(t) + e^{-2\eta t}C,$ (5) initialization of W and  $W_p$ .

where C is a symmetric matrix that depends only on the



### **Exploration 2** No predictor or no stop-gradient = collapse

- Shown to be true in empirical studies various times, but there is no theoretical explanation.
- When there is no EMA, meaning that  $W_a = W$ ,

$$\frac{d}{dt}vec(W) = -H(t) \cdot vec(W)$$

where H(t) is a PSD matrix. This implies that if the minimal eigenvalue of H(t) is bounded below,  $W(t) \rightarrow 0$ , i.e., collapse.

Similar case for no predictior, i.e.  $W_p = I$ .

Theorem 2 (The stop-gradient signal is essential for success.). With  $W_{\rm a} = W$  (SimSiam case), removing the stop-gradient signal yields a gradient update for W given by positive semi-definite (PSD) matrix  $H(t) := X' \otimes$  $(W_{p}^{\intercal}W_{p} + I_{n_{2}}) + X \otimes \tilde{W}_{p}^{\intercal}\tilde{W}_{p} + \eta I_{n_{1}n_{2}}$  (here  $\tilde{W}_{p} :=$  $W_p - I_{n_2}$  and  $\otimes$  is the Kronecker product):

$$\frac{\mathrm{d}}{\mathrm{d}t}\mathrm{vec}(W) = -H(t)\mathrm{vec}(W). \tag{6}$$

If the minimal eigenvalue  $\lambda_{\min}(H(t))$  over time is bounded below,  $\inf_{t>0} \lambda_{\min}(H(t)) \geq \lambda_0 > 0$ , then  $W(t) \to 0$ .







## Assumptions

- covariance.
- of W(t), i.e.  $W_a$  points to the same direction as the online weight.

- Assumption (Symmetric Predictor):  $W_p(t) = W_p^{\top}(t)$ 

• Assumption (Isotropic Data and Augmentation): X = I and  $X' = \sigma^2 I$ , which comes from the assumptions that that data distribution p(x) has zero mean and identity covariance whereas the augmentation distribution  $p_{aug}(\cdot | x)$  has x mean and  $\sigma^2 I$ 

• Assumption (Proportional EMA): EMA weight  $W_a(t) = \tau(t)W(t)$  is a linear function

## **New Dynamics**

Under these assumptions, we have the following dynamics

$$\begin{split} \dot{W}_p &= -\frac{\alpha_p}{2}(1+\sigma^2)\{Wp,F\} + \alpha_p\tau F - \eta W_p \\ \dot{F} &= -(1+\sigma^2)\{W_p^2,F\} + \tau\{W_p,F\} - 2\eta F \\ &= AB + BA \text{ and } F = \mathbb{E}[f_1f_1^{\mathsf{T}}] = WXW^{\mathsf{T}} \text{ is the correlation} \\ \text{put of the predictor } W \end{split}$$

where  $\{A, B\}$ : matrix of the input of the predictor  $W_p$ .

## **Exploration 3** Eigenspace of $W_p$ aligns with F

Theorem 3: Under certain conditions,

$$FW_p - W_p F \to 0$$

and the eigenspace of  $W_p$  and F gradually **aligns**.





## **Decoupled Dynamics**

 $\Lambda_{W_p} = diag[p_1, ..., p_d]$ , and  $F = U\Lambda_F U^T$  where  $\Lambda_F = diag[s_1, ..., s_d]$ .

$$\dot{p}_j = lpha_p s_j [ au - \dot{s}_j] = 2p_j s_j [ au - s_j \dot{ au}]$$

Invariance holds:  $s_j(t) = \alpha_p^{-1} p_j^2(t) + e^{-2\eta t} c_j$ 

Let columns of U be the common eigenvectors of  $W_p$  and F so that  $W_p = U \Lambda_{W_p} U^{+}$  where



## No collapse

1D dynamics of the eigenvalue  $p_j$  of  $W_p$ :





## Effect of Weight Decay



**BUT**, if weight decay is large, then the eigenspace alignment condition is more likely to satisfy!

## **Effect of Other Hyperparameters**



$$s_j(t) = \alpha_p^{-1} p_j^2(t) + e^{-2\eta t} c_j$$



• • • Size of triv

Condition of alignn

Training

|                    | $\alpha_p \uparrow$                                    | $\beta\downarrow$ |
|--------------------|--|-------------------|
| vial basin         |  |                   |
| eigenspace<br>nent |  |                   |
| speed              |  |                   |
|                    | Eigenvalue of F won't<br>grow (no feature<br>learning) |                   |

## DirectPred

- Hmm, it looks like it is very important that eigenspaces of  $W_p$  and F align well. Why not do this directly without relying on gradient descent?
- Directly set linear  $W_p$ , so no optimization.

1. Estimate 
$$\hat{F} = \rho \hat{F} + (1 - \rho) \mathbb{E}[f$$

- 2. Eigen-decompose  $\hat{F} = \hat{U}\Lambda_F \hat{U}^T, \Lambda_F = diag[s_1, \dots, s_d]$



3. Set  $W_p = \hat{U} \operatorname{diag}[p_j] \hat{U}^{\mathsf{T}}$  following the invariance  $p_j = \sqrt{s_j} + \epsilon \max_i s_j$ 

### Eigenspaces are always and automatically aligned!





| Number of enoche |                                    |                  |  |  |  |  |  |
|------------------|------------------------------------|------------------|--|--|--|--|--|
| Number of epochs |                                    |                  |  |  |  |  |  |
| )                | 300                                | 500              |  |  |  |  |  |
| STL-10           |                                    |                  |  |  |  |  |  |
| 0.16             | $78.77\pm0.97$                     | $78.86 \pm 1.15$ |  |  |  |  |  |
| 0.11             | $79.90 \pm 0.66$                   | $80.28 \pm 0.62$ |  |  |  |  |  |
| 0.52             | $75.25\pm0.74$                     | $75.25\pm0.74$   |  |  |  |  |  |
| IFAR-10          |                                    |                  |  |  |  |  |  |
| 0.23             | $\textbf{88.88} \pm \textbf{0.15}$ | $89.52 \pm 0.04$ |  |  |  |  |  |
| 0.29             | $88.83 \pm 0.10$                   | $89.56 \pm 0.13$ |  |  |  |  |  |
| 0.20             | $88.57\pm0.15$                     | $89.33 \pm 0.27$ |  |  |  |  |  |

## DirectPred

### Downstream classification (ImageNet):

| BYOL variants                                 | Accuracy (60 ep) |       | Accuracy (300 ep) |             |  |
|---|------------------|-------|-------------------|-------------|--|
|   | Top-1            | Top-5 | Top-1             | Top-5       |  |
| 2-layer predictor*                            | <b>64.7</b>      | 85.8  | 72.5              | 90.8        |  |
| linear predictor                              | 59.4             | 82.3  | 69.9              | 89.6        |  |
| DirectPred                                    | 64.4             | 85.8  | 72.4              | <b>91.0</b> |  |
| * 2 lorron mus distants DVOI deferste setting |                  |       |                   |             |  |

2-layer predictor is BYOL default setting.