# Understanding Self supervised learning without negative

Two layer model example:

setup
$$\begin{cases} \text{online Net weight } w \in \mathbb{R}^{n_2 \times n_1} \\[6pt] \text{predictor} : w_p \in \mathbb{R}^{n_2 \times n_2} \\[6pt] \text{target Net} : \bar{w}_a \in \mathbb{R}^{n_2 \times n_1} \\[6pt] x \in \mathbb{R}^{n_1} \quad // \text{ input data} \end{cases}$$

Two augmentation: $z: x_1, x_2 \sim P_{aug}(\cdot | x)$

$$f_1 = w x_1 \in \mathbb{R}^{n_2} \quad // \text{ online rep.}$$

$$f_{2a} = w_a x \in \mathbb{R}^{n_2} \quad // \text{ target rep.}$$

① BYOL objective

$$J(w, w_p) := \frac{1}{2} \mathbb{E}_{x_1, x_2} \left[ \| w_p f_1 - \text{stopgrad}(f_{2a}) \|_2^2 \right]$$

⑪ whereas $w_a \simeq \text{Exponential MA}(w)$

## BYOL learning Dynamics:

gradient {

$$\dot{W_P} = \frac{\partial J}{\partial w_p} = \alpha_p \left( -w_p w(x + x') + w_a x \right) \bar{w}^T - \eta \, w_p$$

learning rate ratio ↗

weight decay ↗

$$\dot{W} = w_p^T \left( -w_p w(x + x') + w_a x \right) - \eta \, w$$

$$\dot{w_a} = \hat{\beta} \left( -w_a + w \right)$$

}

[Every thing is in here]

→ Expectation of outere product matrix.

$$\boxed{x := E\left[ \bar{x}_o \bar{x}^T \right]} \quad \bar{x}(x) := E_{\nu \sim pang(\eta)}\left[ \hat{x} \right]$$

mean

Average augmented view of data point

Expected cov mat. [

$$x' = E\left[ \nabla_{\hat{x} | x}\left[ \hat{x}' \right] \right] \quad E(cov)$$

↳ covariance matrix of aug view $\hat{x}'$

]

Requires Simplified Assumption: fore analysis

→ ⓥ Proportial EMA

$$w_a(t) = c(t) \, w(t)$$

→ ⓝ Isotopic Data augmentation:

Avg data covariance : $X = I$

data Avg avg cov : $X' = \sigma^2 I$

$\rightarrow$ $W_p$ is symmetric

$\rightarrow$ Eigen decomposition.

$F$ : correlation matrix of output of $W$

$$F = W X W^T$$

# Findings ①

Eigen space of $W_p$ aligns to $F$

[we can approximate $W_p$ from $F$]

### Theorm 3:

under some condition $F W_p - W_p F \rightarrow 0$

## Direct Pred Method:

Estimate $\hat{F} \rightarrow$ Get $W_p$ by the following

$$\hat{F} = \hat{U} \wedge_F \hat{U}^T$$

$$\wedge_F = \text{diag} [s_1, s_2 \cdots s_d]$$

then, define

$$P_j = \sqrt{s_j} + \epsilon \max_j s_j$$

$$\Rightarrow \quad w_P = \hat{U} \, \text{diag} [P_j] \, \hat{U}^T$$

To estimate   correlation   matrix

$$\hat{F} = \rho \hat{F} + (1-\rho) \, \mathbb{E}_B [f f^T]$$

$$\Downarrow$$

Expectation over batch.