Exploring the gap between collapsed and whitened features in self-supervised learning

Bobby He and Mete Ozay ICML 2022



Self-supervised learning (SSL)

- Obtaining labelled data can be expensive.
- Self-supervised learning (SSL) learns from unlabelled data $X = \{x_n\}_{n=1}^N$ producing useful encoder features, $h_{\theta}(X) \in \mathbb{R}^{N \times d}$ for downstream tasks.
- **Popular SSL strategy** (SimCLR, BYOL, SwAV, Barlow Twins + more): encourage feature invariance to some input transformations, T_1, T_2 , i.e. $h_{\theta}(T_1(x)) \approx h_{\theta}(T_2(x))$



Feature collapse vs. whitening in SSL

- Common goal in SSL: avoiding *collapsed* NN features $h_{m{ heta}}(m{x}') = m{c}, \; orall m{x}'$.
- Some approaches: contrastive learning (SimCLR) or clustering (SwAV).
- *Whitening* or *decorrelating* features has been proposed as a sufficient condition to avoid collapse.^{1,2,3,4}
- **Barlow Twins**² adds a regulariser to the SSL loss function to encourage decorrelated features, whilst still achieving invariance to transformations:

Barlow Twins objective = Transformation Invariance loss + Feature Decorrelation reg.

[1] Ermolov et al, Whitening for Self-Supervised Representation Learning, 2020
[2] Zbontar et al, Barlow Twins: Self-Supervised Learning via Redundancy Reduction, 2021
[3] Hua et al, On Feature Decorrelation in Self-Supervised Learning, 2021
[4] Bardes et al, VICReg: Variance-Invariance-Covariance Regularization for Self-Supervised Learning, 2021

The spectrum of feature eigenspectra

- Whitened & Collapsed features represent two extremes of a gap:
 - Exactly whitened features have identity covariance: $h_{m{ heta}}(m{X})^T h_{m{ heta}}(m{X}) = I_d$
 - OTOH, collapsed features have (at most) a rank-one covariance.
- This gap can be characterised by rate of decay of feature eigenvalues $\{\lambda_i\}_{i=1}^d$:
 - If we have **power-law** decay, $\lambda_i = \frac{1}{i^{\beta}}$, $\forall 1 \le i \le d$ with exponent $\beta \ge 0$ then:

eta=0 denotes whitened features and $\ eta=\infty$ denotes collapsed features.



Factors affecting degree of collapse/whitening

- Using our insights from power-law decay, we prove (for linear NNs) that:
 - (Commonly found) projection layers enable encoders to observe decaying eigenspectra.
 - Deeper projection layers lead to more collapsed encoder features.
 - Stronger Barlow Twins regularisation encourages more whitened features.



The effect of whitening on generalisation

- By varying projection depth and regularisation strength, we empirically show that more whitened features do not necessarily generalise better (above right).
- We prove that more collapsed features may generalise better in regimes of **low-labelled data**; a common setting in SSL.



Method: Post-hoc Manipulation of the Principle Axes and Trace

- Propose **PostMan-Pat** (PMP) to control rate of feature eigenvalue decay, and deliver improved label efficiency in SSL.
 - Given a pretrained encoder $h_{\theta}(x) \in \mathbb{R}^d$ and power-law exponent $\beta \ge 0$ PMP computes an eigenspectrum-dependent rescaling matrix $W_{\text{PMP}} \in \mathbb{R}^{d \times d}$
 - PMP defines a new encoder: $h_{\rm PMP}({m x}) \leftarrow h_{m heta}({m x}) W_{\rm PMP}$ for linear evaluation.
 - $h_{\rm PMP}(\boldsymbol{x})$ is constructed to have power-law behaviour in eigenvalues, with exponent β acting as a hyperparameter that controls feature eigenvalue decay.

PMP Experiments: CIFAR-10



PMP Experiments: ImageNet-1K

1) Low-labelled data (w. ResNet-50)

Method		Top-1			Top-5		
PRETRAIN	EVAL	0.3%	1%	10%	0.3%	1%	10%
SIMCLR (69.3)	LP MLP PMP	$\begin{array}{c} 34.2_{\pm.2} \\ 31.8_{\pm.4} \\ \underline{35.9}_{\pm.2} \end{array}$	48.1 45.2 <u>50.9</u>	61.0 61.5 <u>62.5</u>	$57.2_{\pm.3} \\ 54.1_{\pm.3} \\ \underline{57.9}_{\pm.2}$	73.8 71.1 <u>76.6</u>	84.3 85.0 <u>85.2</u>
	NFT	$\textbf{39.8}_{\pm.2}$	52.5	67.5	$\textbf{65.5}_{\pm.2}$	78.9	88.7
SwAV (74.7)	LP MLP PMP	$\frac{36.5}{34.6_{\pm.4}}$ $39.3_{\pm.3}$	<u>53.8</u> 52.0 55.9	68.2 67.5 <u>68.5</u>	$\begin{array}{c} \underline{61.8}_{\pm.1} \\ 59.3_{\pm.2} \\ \textbf{64.1}_{\pm.2} \end{array}$	78.8 77.2 79.7	88.8 88.6 <u>89.1</u>
	NFT	$32.4_{\pm.3}$	53.6	70.8	$57.8_{\pm .2}$	<u>79.1</u>	90.5
BARLOW (73.5)	LP MLP PMP	$\begin{array}{c} 39.9_{\pm.1} \\ 37.6_{\pm.1} \\ \textbf{42.3}_{\pm.1} \end{array}$	55.0 53.0 56.2	63.2 66.3 <u>67.3</u>	$\begin{array}{c} 63.8_{\pm.2} \\ 61.5_{\pm.1} \\ \textbf{65.8}_{\pm.1} \end{array}$	79.0 76.9 79.7	83.4 87.2 <u>88.6</u>
	NFT	$\underline{40.7}_{\pm.1}$	<u>55.3</u>	70.0	$\textbf{65.8}_{\pm.1}$	<u>79.6</u>	89.9
SUPERVISED		—	25.4	56.4	-	48.4	80.4

2) Transfer learning

Метно	DD	TRANSFER DATASET			
PRETRAIN	EVAL	C-100	CARS	FLOWERS	
SIMCLR	LP	65.26	46.88	84.65	
	PMP	66.13	<u>47.83</u>	<u>85.88</u>	
BARLOW	LP	74.19	69.36	92.29	
	PMP	75.10	<u>69.67</u>	92.54	
SwAV	LP	75.24	63.39	90.47	
	PMP	<u>76.10</u>	<u>64.46</u>	<u>92.00</u>	

3)

Different encoder architectures (ViT)

	Top-1			Top-5		
EVAL	0.3%	1%	10%	0.3%	1%	10%
LP	$54.0_{\pm.3}$	64.5	72.3	$78.0_{\pm.3}$	86.6	91.2
PMP	$\underline{55.2}_{\pm.2}$	<u>65.1</u>	72.4	$\underline{78.7}_{\pm.2}$	86.8	91.2

Conclusion

- > We highlight the spectrum that exists between collapsed and whitened features in SSL:
 - Parameterise by power-law behaviour
 - Influential hyperparameters: projection depth and regularisation strength.
 - Implications for generalisation, particularly in low-labelled data regimes
- Propose Post-hoc Manipulation of the Principle Axes and Trace (PostMan-Pat or PMP):
 - Controls rate of decay of feature eigenvalues.
 - Consistently outperforms linear probing, and often outperforms fine-tuning on low-labelled data.

Limitations/future work:

- ★ Devising (self-)supervised learning methods that consider feature eigenspectra directly.
- \star Extending PMP to fine-tuning settings.
- ★ PMP doesn't consider the actual features themselves, just the importance of their relative weighting.
- \star Performance depends on power-law β hyperparameter, can we automatically tune this?

