

Amortized Eigendecomposition for Neural Networks

Tianbo Li, Zekun Shi, Jiayi Zhao, Min Lin



SEA AI Lab
National University of Singapore

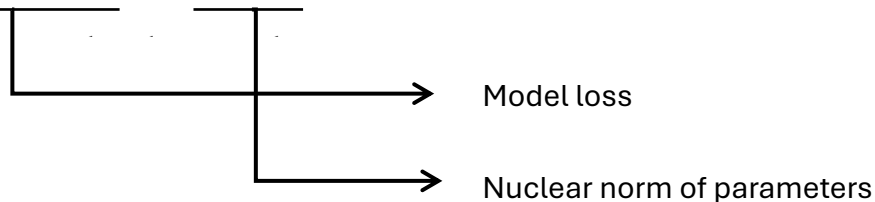
An Optimization Problem Involving Eigendecomposition

$$\min_{\boldsymbol{\theta}} f(h_{\boldsymbol{\theta}}(\mathbf{X}), \mathbf{V}, \boldsymbol{\Lambda})$$

$$\text{s.t. } \mathbf{V}^{\top} \boldsymbol{\Lambda} \mathbf{V} = \mathbf{A} \quad \text{where } \mathbf{A} \text{ is constructed from } h(x)$$

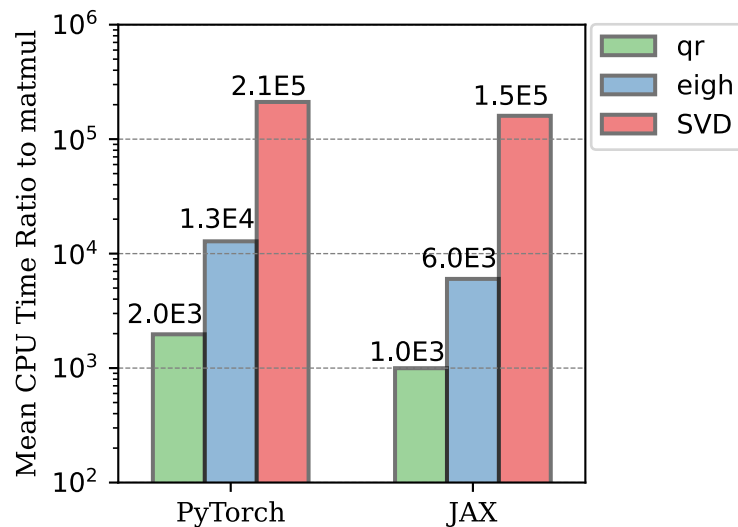
Nuclear Norm Regularization

$$\min_{\boldsymbol{\theta}} f(h_{\boldsymbol{\theta}}(\mathbf{X})) + \eta \|\boldsymbol{\theta}\|_*$$



Eigendecomposition is Slow

In both PyTorch and JAX, performing **Singular Value Decomposition** and **eigendecomposition** can be 10x to 100x slower than **matrix multiplication** for a 10000 x 10000 matrix.



The idea

Eigendecomposition itself is an optimization process;
however, it **does not need to fully converge** at each step;
it only needs to **reach the desired outcome** in the end.

Amortized Eigendecomposition

Algorithm 1 The conventional eigendecomposition in a neural network outlined in Eq. (2)

Input: Dataset \mathbf{X} , encoder h_θ , task f_ω ;

- 1: Initialize model parameter θ and ω ;
- 2: **while** not converged **do**
- 3: compute \mathbf{A} from $h_\theta(\mathbf{X})$;
- 4: $\mathbf{V}, \mathbf{\Lambda} = \text{eigh}(\mathbf{A})$;
- 5: compute $f_\omega(h_\theta(\mathbf{X}), \mathbf{V}, \mathbf{\Lambda})$;
- 6: update θ, ω by gradient descent;
- 7: **end while**

Algorithm 2 The amortized eigendecomposition technique outlined in Eq. (12)

Input: Dataset \mathbf{X} , encoder h_θ , task f_ω , and η ;

- 1: Initialize model parameter θ, ω and \mathbf{W} ;
- 2: **while** not converged **do**
- 3: compute \mathbf{A} from $h_\theta(\mathbf{X})$;
- 4: $\mathbf{U} = \text{QR}(\mathbf{W}), \mathbf{\Sigma} = \text{diag}(\mathbf{U}^\top \mathbf{A} \mathbf{U})$;
- 5: compute $f_\omega(h_\theta(\mathbf{X}), \mathbf{U}, \mathbf{\Sigma}) + \eta g(\mathbf{A}, \mathbf{U})$;
- 6: update $\theta, \omega, \mathbf{W}$ by gradient descent;
- 7: **end while**

Table 1: Evaluation of execution times per iteration on three tasks.

Task	Dimension	Backbone	Backbone+	Backbone+	Speed-up
		time (s/iter)	eigh/svd	our method	
		t_0	t_1	t_2	$\frac{t_1-t_0}{t_2-t_0}$
Nuclear norm regularization	128×128	5.275E-2	8.323E-2	6.025E-2	4.06×
	256×256	5.600E-2	1.209E-1	6.080E-2	13.5×
	512×512	7.186E-2	2.616E-1	7.366E-2	105.4×
Latent-space PCA	256×2	4.178E-3	1.446E-2	1.117E-2	1.47×
	512×2	6.792E-3	2.918E-2	2.224E-2	1.45×
	1028×2	1.434E-2	7.018E-2	5.467E-2	1.39×
Low-rank GCN	2708×16	1.021E-3	1.769E-2	1.732E-3	23.4×
	3312×16	1.367E-3	2.825E-2	2.498E-3	23.7×
	19717×16	1.931E-2	4.941E+0	2.731E-2	615.2×