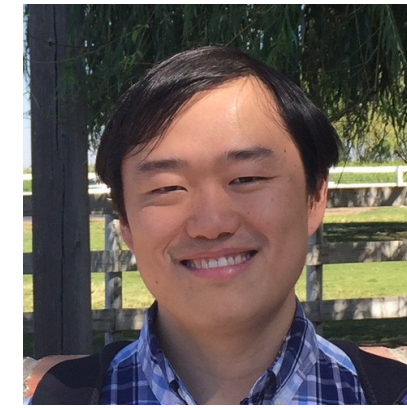
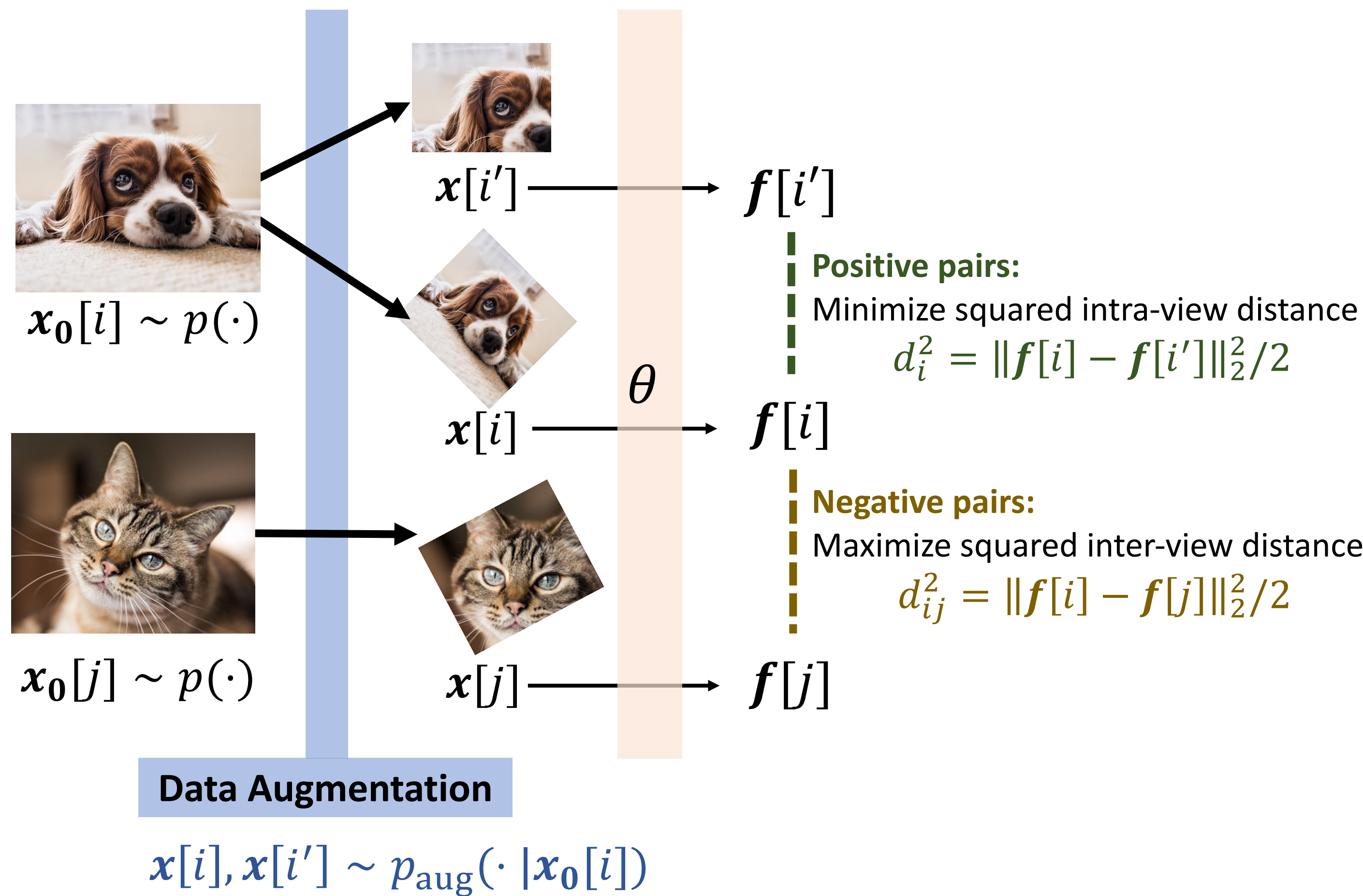


# On the Role of Nonlinearity in Training Dynamics of Contrastive Learning on 1-layer Network

Yuandong Tian  
yuandong@meta.com

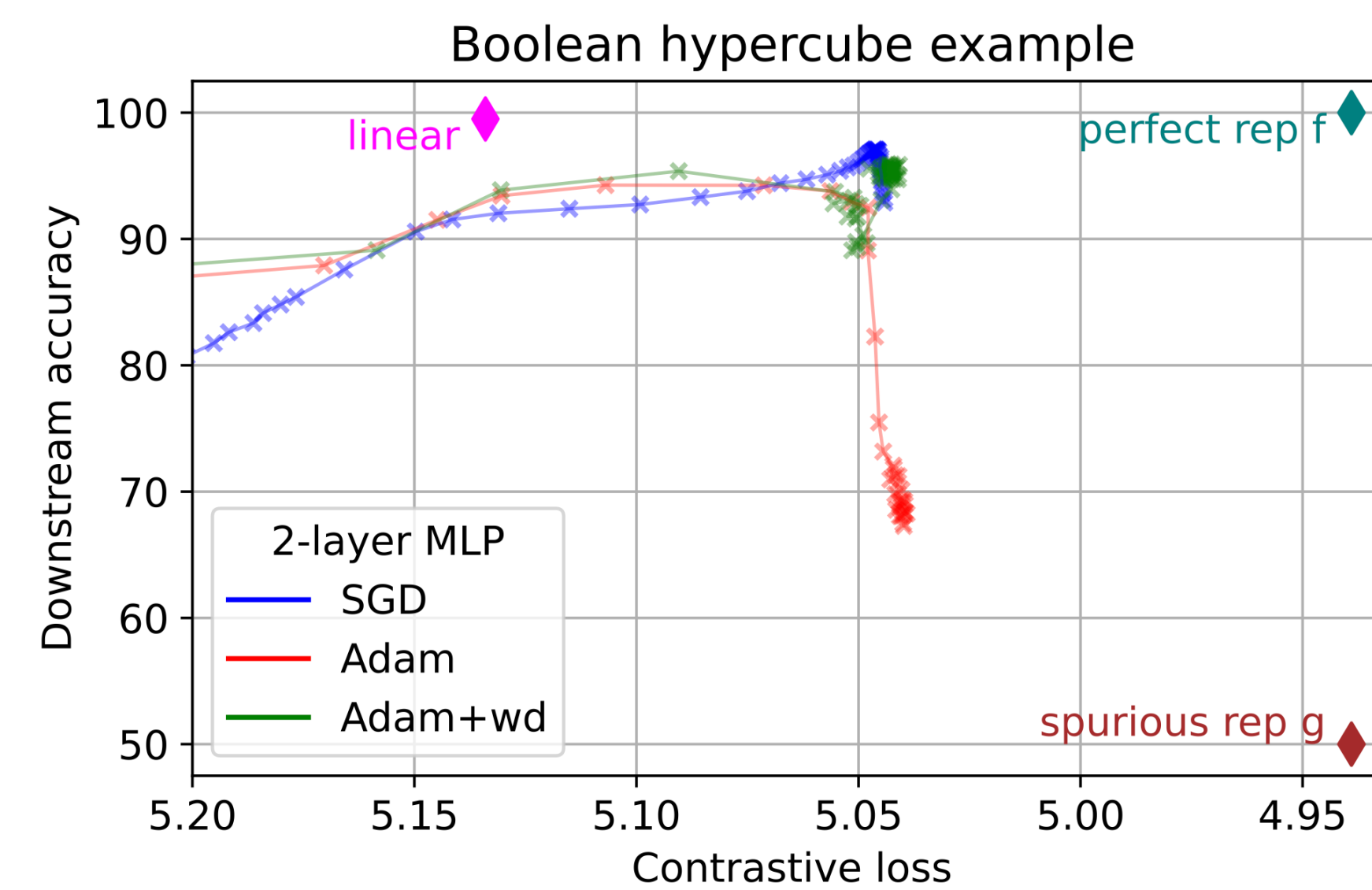


## Background



E.g., InfoNCE loss:  $\mathcal{L}_{nce} := -\tau \sum_{i=1}^N \log \frac{e^{-d_i^2/\tau}}{\epsilon e^{-d_i^2/\tau} + \sum_{j \neq i} e^{-d_{ij}^2/\tau}}$

## Is CL just loss + blackbox function family?

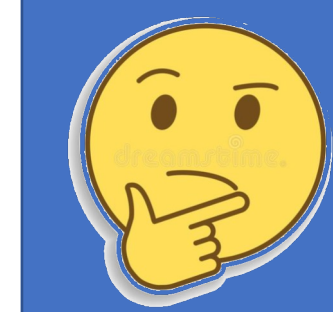


[Tables from N. Saunshi et al, Understanding Contrastive Learning Requires Incorporating Inductive Biases, ICML'22]

### Inductive bias matters! (e.g., architecture, optimizer)



When the network is linear, CL can be shown to perform like PCA ( $\alpha$ -CL)



How to understand the properties of CL with nonlinear network?

## $\alpha$ -CL: a unified framework [Tian, NeurIPS'22]

$$\theta_{t+1} = \theta_t + \eta \nabla_{\theta} \mathcal{E}_{\text{sg}(\alpha_t)}(\theta_t)$$

Pairwise importance  $\alpha_t = \alpha(\theta_t)$

The pairwise importance  $\alpha$  can be

- determined by  $\alpha(\theta) = \arg \min_{\alpha \in \mathcal{A}} \mathcal{E}_{\alpha}(\theta) - \mathcal{R}(\alpha)$ , with different regularization  $\mathcal{R}(\alpha)$ .
- directly specified ( $\alpha$ -CL-direct)

Define  $\mathcal{E}_{\alpha}$  as the trace of **contrastive covariance**  $\mathbb{C}_{\alpha}[\cdot]$ :

$$\mathcal{E}_{\alpha}(\theta) := \frac{1}{2} \text{tr} \mathbb{C}_{\alpha}[f_{\theta}(x)]$$

where the *contrastive covariance*

$$\mathbb{C}_{\alpha}[x] := \frac{1}{2N^2} \sum_{i,j} \alpha_{ij} [(x[i] - x[j])(x[i] - x[j])^T - (x[i] - x[i'])(x[i] - x[i'])^T]$$



**Goal: Analyze the local maxima of the energy function  $\mathcal{E}_{\alpha}$**

## Setup

### The Assumptions of Homogenous Activations

We assume the activation satisfies  $h(x) = h'(x)x$   
This includes Linear, ReLU and monomial activations (with additional constant)

### Connect $\mathbb{C}_{\alpha}[\cdot]$ with regular variance $\mathbb{V}[\cdot]$

If  $\alpha$  satisfies  $\alpha_{ij} = \mathcal{K}(x_0[i], x_0[j])$ , where  $\mathcal{K}(x, y) = \sum_{l=0}^{+\infty} \phi_l(x)\phi_l(y)$

is a kernel, then for any function  $g(\cdot)$ :

$$\mathbb{C}_{\alpha}[g(x)] \rightarrow \sum_{l=0}^{+\infty} z_l^2 \mathbb{V}_{x_0 \sim \tilde{p}_l(\cdot; \alpha)} [\mathbb{E}_{x \sim p_{\text{aug}}(\cdot | x_0)} [g(x)]]$$

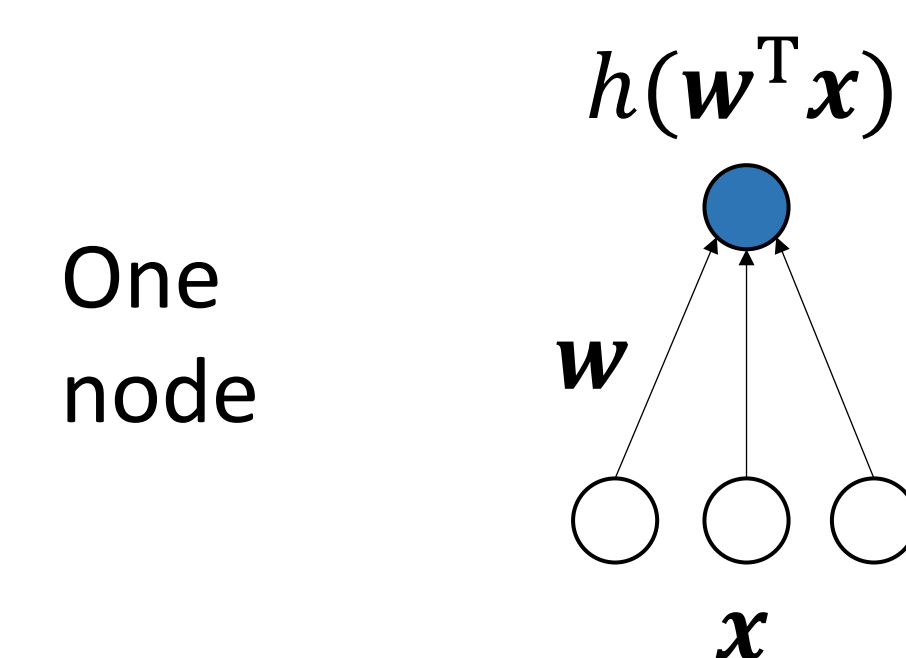
Where  $\tilde{p}_l(x; \alpha) := \frac{1}{z_l(\alpha)} p_D(x)\phi_l(x; \alpha)$  is adjusted density of the data, and  $z_l(\alpha)$  is the normalization constant.

### Example of Kernel-like $\alpha$

Uniform  $\alpha_u := 1$

Gaussian  $\alpha_g := \exp\left(-\frac{\|h(w^T x_0[i]) - h(w^T x_0[j])\|_2^2}{2\tau}\right)$

## 1-layer nonlinear network

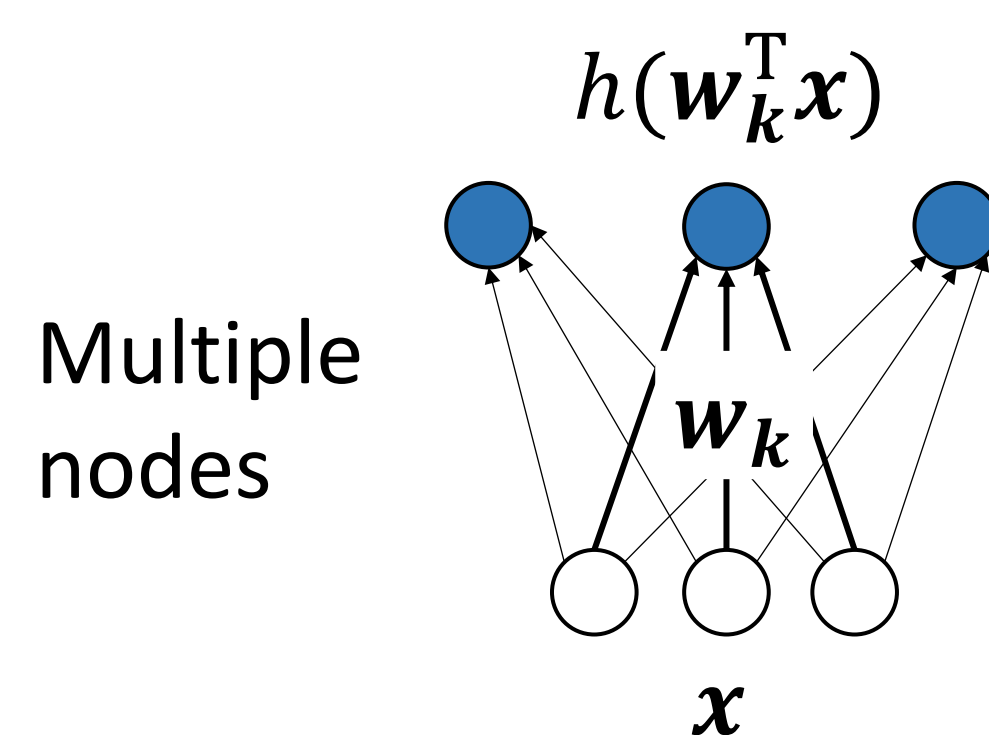


One node

$$\max_{\|w\|_2=1} \mathbb{C}_{\alpha}[f_{\theta}] = \mathbb{C}_{\alpha}[h(w^T x)] = w^T A(w) w$$

where  $A(w) := \mathbb{C}_{\alpha}[\tilde{x}^w]$

$\tilde{x}^w := x \cdot h'(w^T x)$  is the **gated** data point

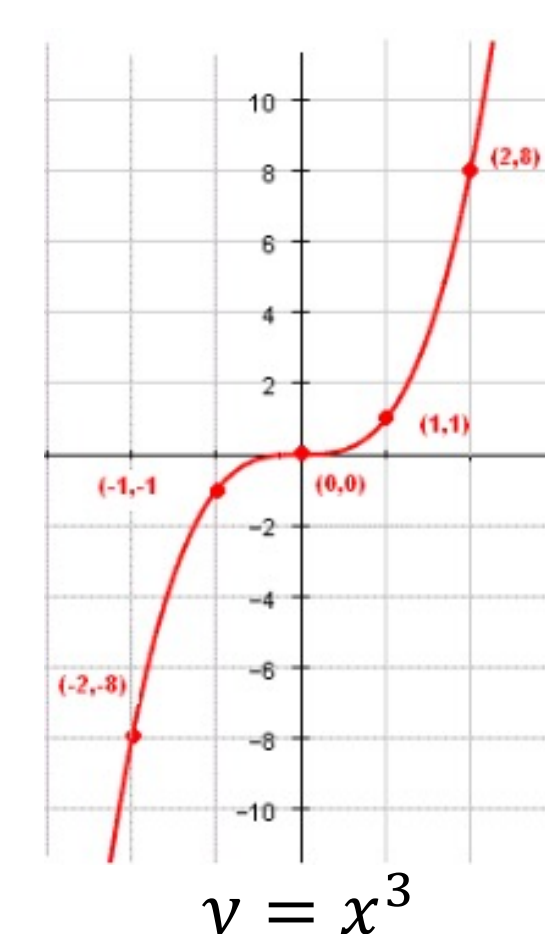


Multiple nodes

$$\max_{\|w_k\|_2=1, 1 \leq k \leq K} \text{tr} \mathbb{C}_{\alpha}[f_{\theta}] = \sum_{k=1}^K \max_{\|w_k\|_2=1} w_k^T A(w_k) w_k$$

Independent one node objective

### Critical Points != Local optima



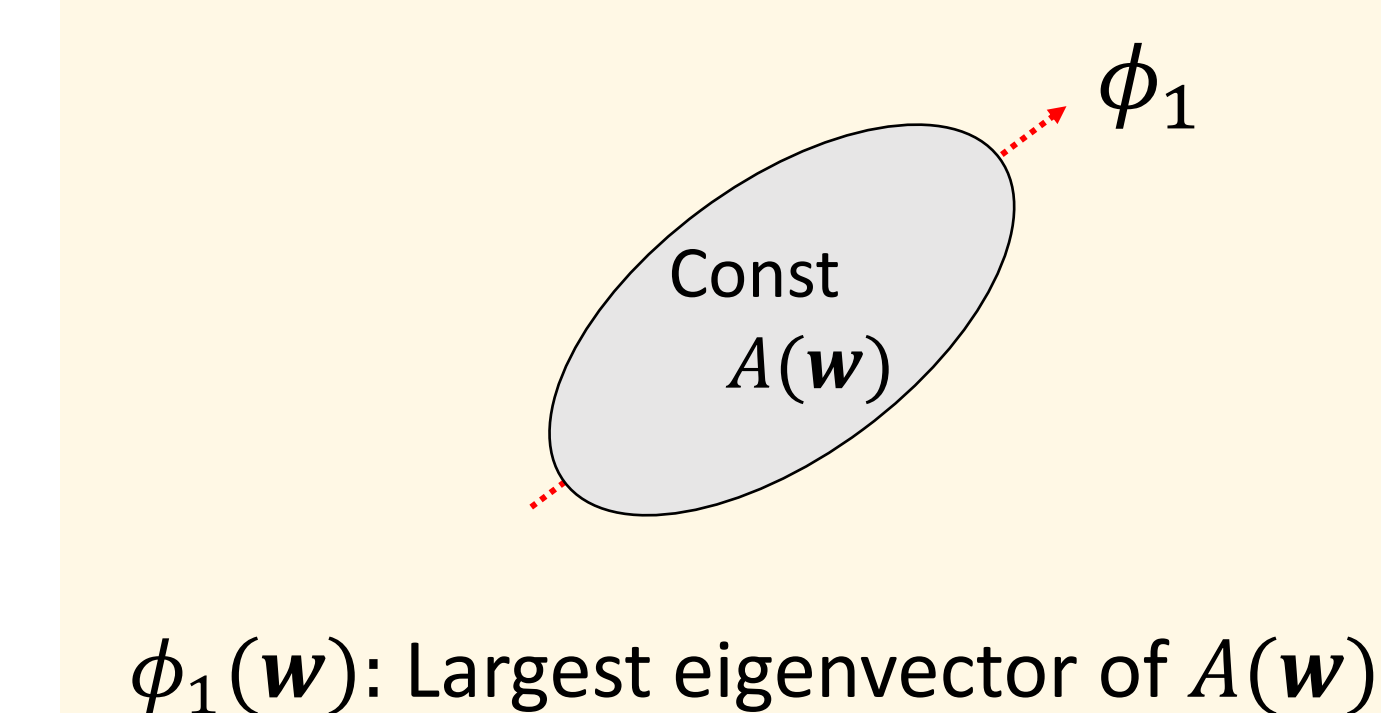
Local roughness  $\rho(w)$ :

$$\|(A(v) - A(w))w\|_2 \leq \rho(w)\|v - w\|_2 + \mathcal{O}(\|v - w\|_2^2)$$

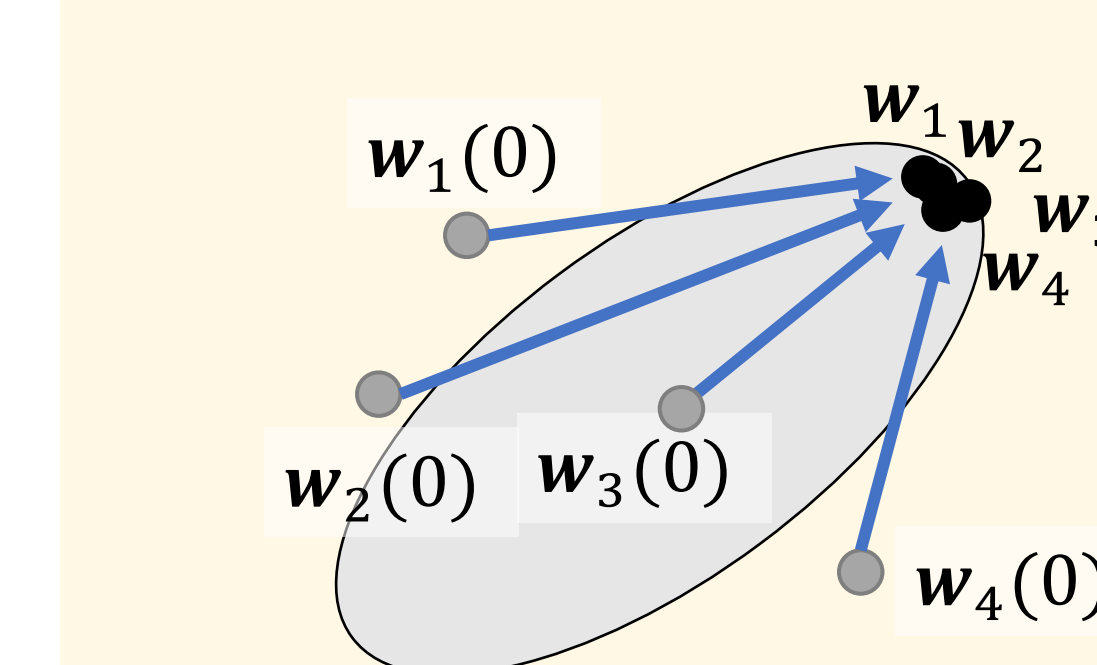
(for any  $v$  in the local neighborhood of  $w$ )

**[Theorem]** if  $A(w_*)w_* = \lambda_* w_*$ , and  $\lambda_{\text{gap}}(w_*) > \rho(w_*)$ , Then  $w_*$  is stable (i.e., local maximum)

### Linear activation

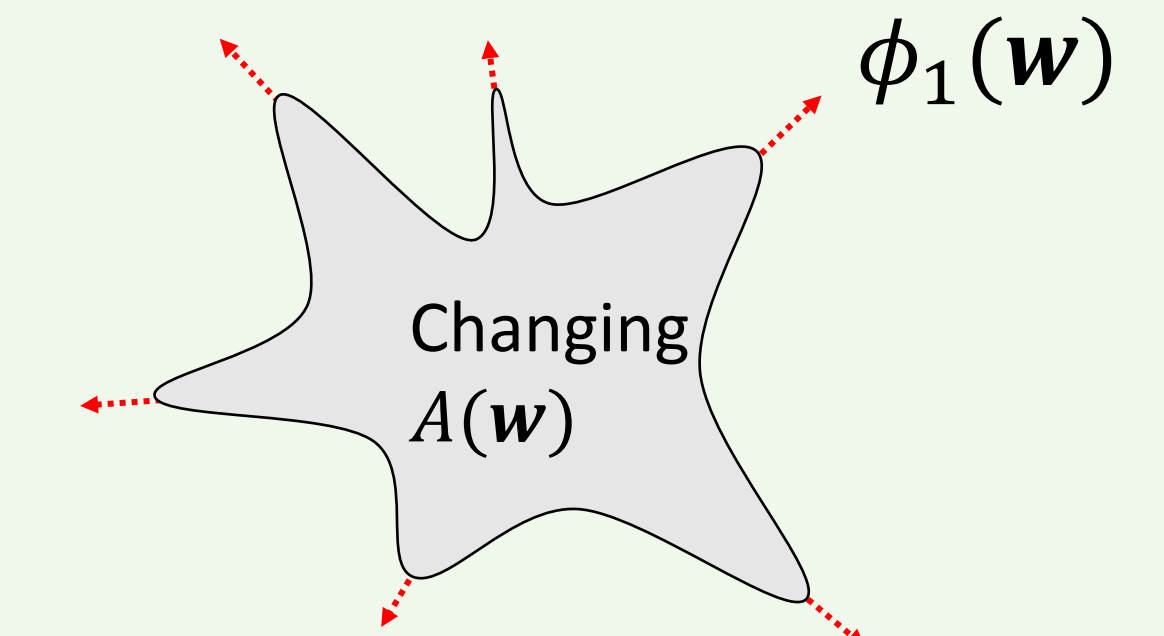


$\phi_1(w)$ : Largest eigenvector of  $A(w)$

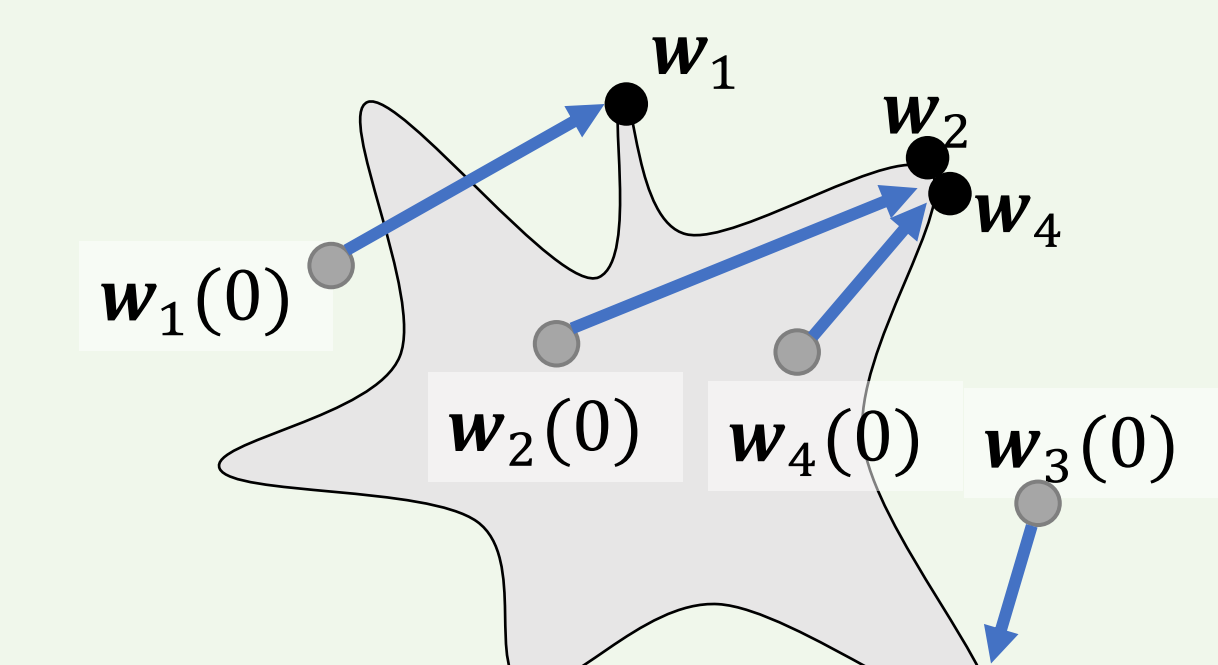


- All  $w_k \rightarrow$  global maximal eigenvector
- More nodes do NOT help.

### Homogenous nonlinear activation



**Multiple largest eigenvectors!**



- Each  $w_k$  can converge to **different** patterns
- More nodes learn **more** patterns!

**[Theorem]** Upper bound of  $\rho(w)$  in Gaussian  $\alpha_g$   
 $\ll$  Upper bound of  $\rho(w)$  in Uniform  $\alpha_u$

	CIFAR-10	STL-10
Quadratic loss (uniform $\alpha$ )	73.58 $\pm$ 0.82	67.28 $\pm$ 1.21
InfoNCE loss (normalized Gaussian $\alpha$ )	87.86 $\pm$ 0.12	83.70 $\pm$ 0.12