

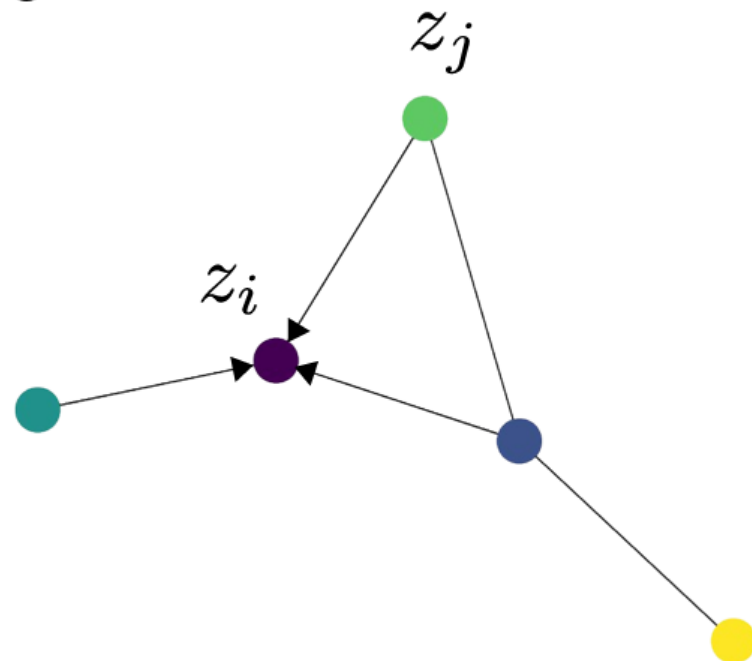
# **Not too little, not too much: a theoretical analysis of graph (over)smoothing**

Khushee Kapoor

# Message Passing Neural Networks

Graph Neural Networks (GNNs) work mostly by **Message-Passing**:

$$z_i^{(k)} = \text{AGG}_{\theta_k}(z_i^{(k-1)}, \{z_j^{(k-1)}\}_{j \in \mathcal{N}_i})$$

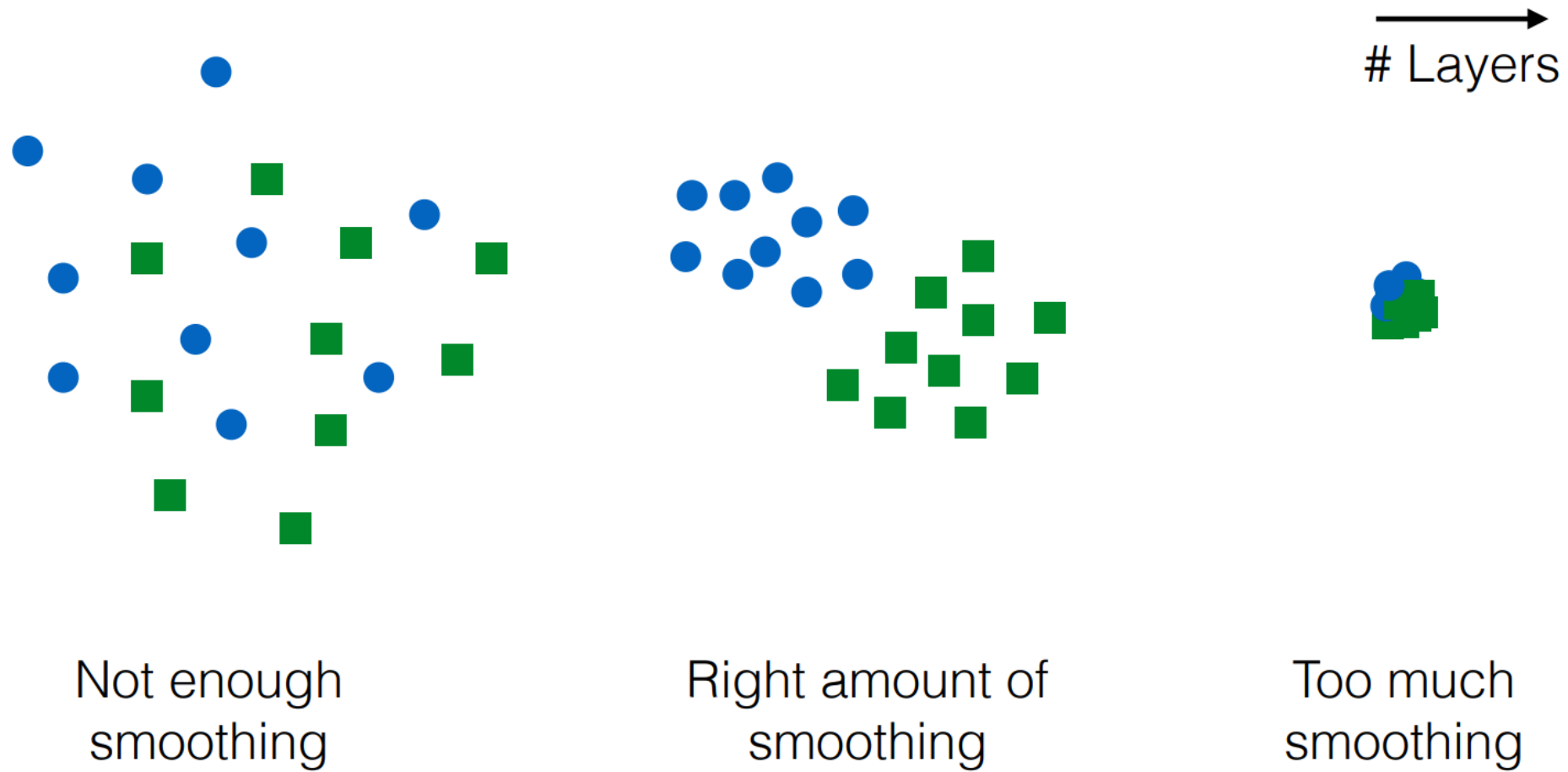


Here we use classic **mean aggregation**:

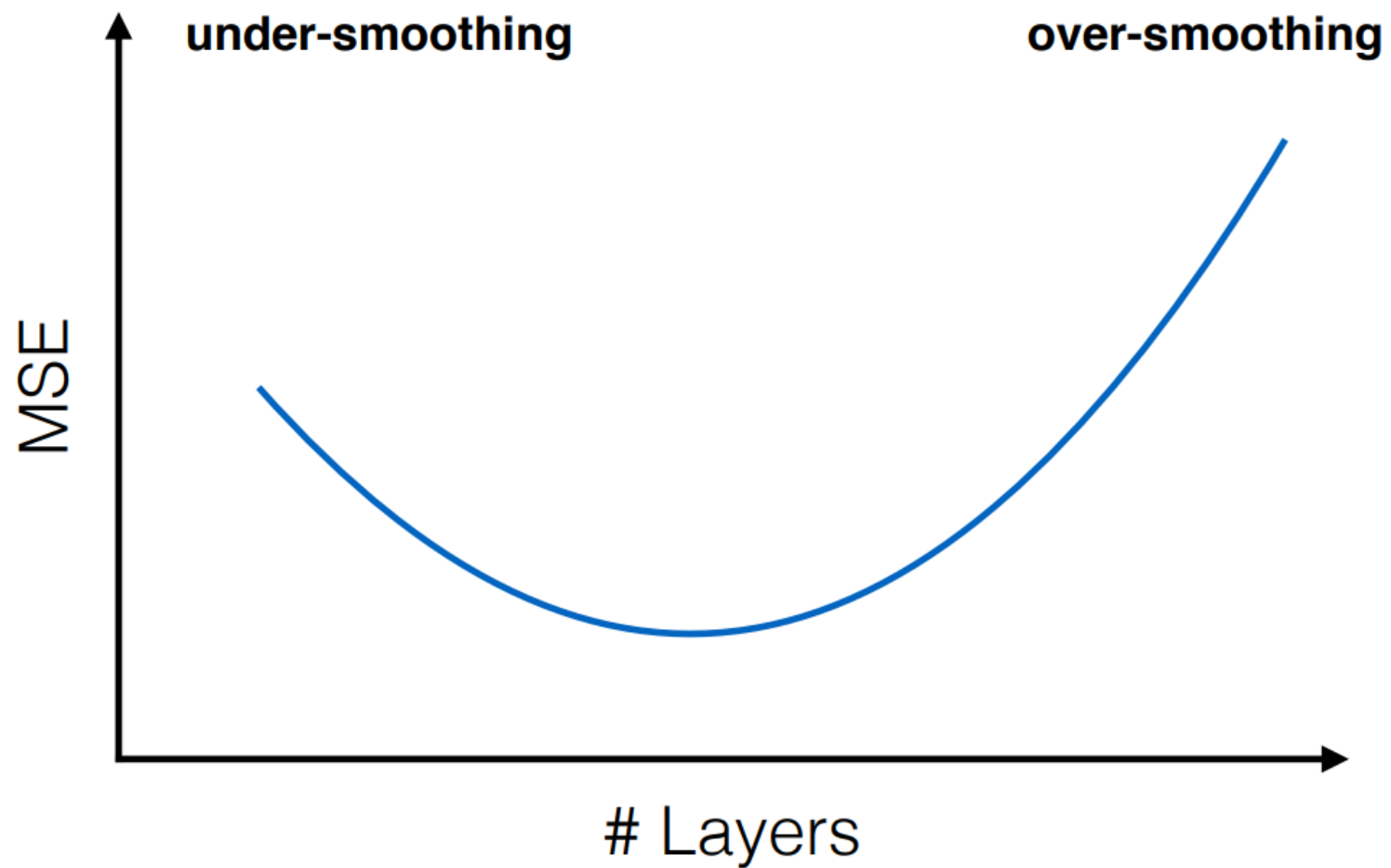
$$z_i^{(k)} = \frac{1}{\sum_j a_{ij}} \sum_j a_{ij} \Psi_{\theta_k}(z_j^{(k-1)})$$

Note that this is just  $Z^{(k)} = L \Psi_{\theta_k}(Z^{(k-1)})$  with  $L = D^{-1}A$

# Smoothing Effect of Graph Convolution

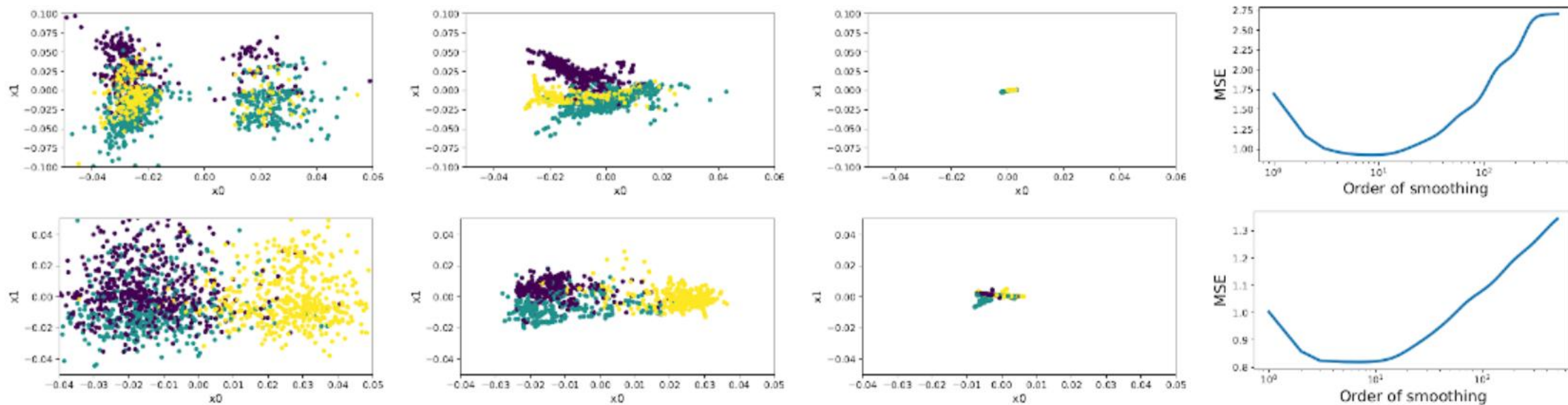


# Smoothing Effect of Graph Convolution



# Oversmoothing vs Sufficient Depth

- **Oversmoothing** restricts GNNs depth: node representations often **provably converge to a constant**... But **some** smoothing helps in practice!



*Learning on increasingly smoothed features on **Cora** and **Pubmed**, illustrated by principal components.*

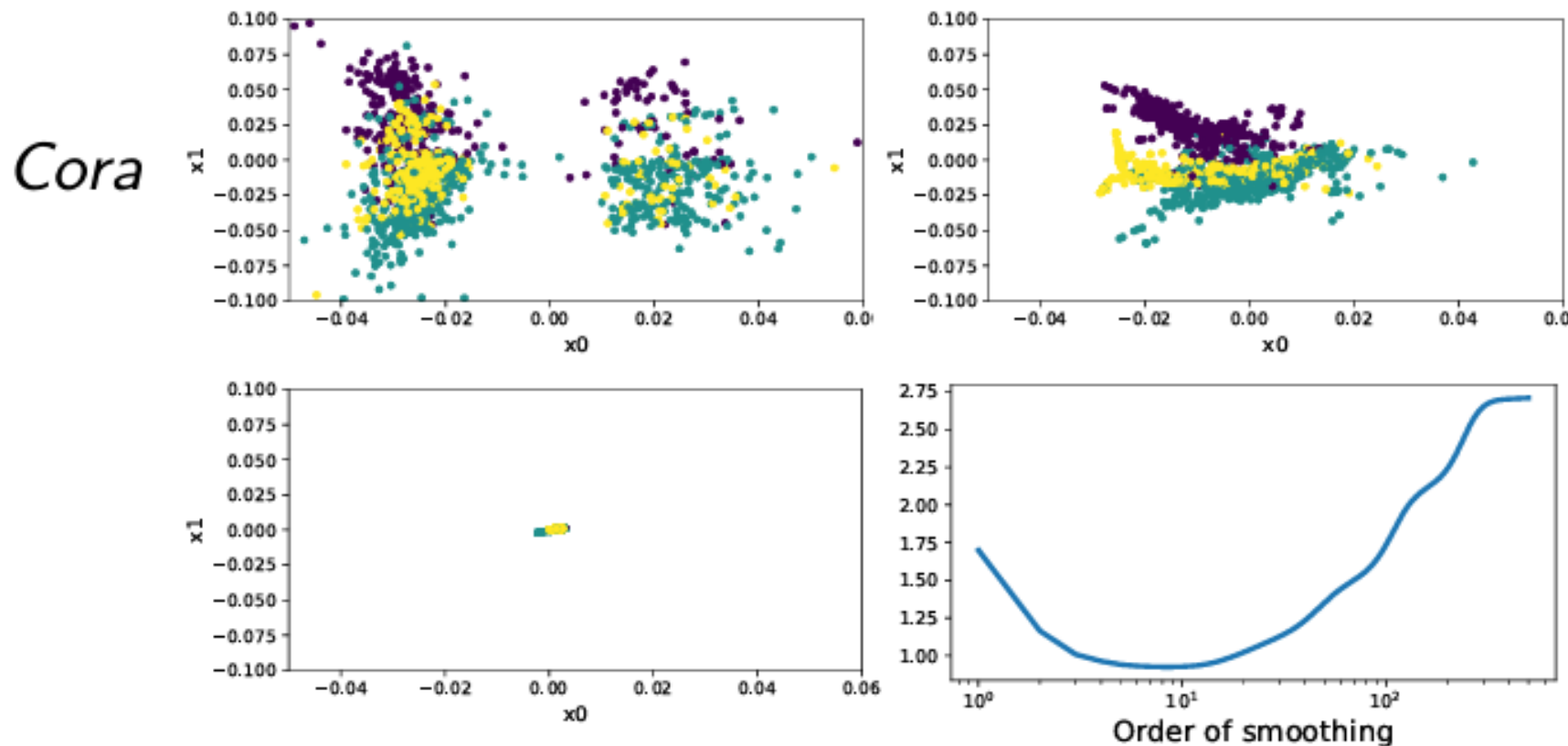
# Oversmoothing vs Sufficient Depth

- ▶ **Theoretical** analyses of GNNs are generally **"in the deep limit"**
- ▶ **Sufficiently deep** GNNs can be as powerful as 1-WL [1]
- ▶ "Well-engineered" GNNs can model useful **diffusion processes** [2]
- ▶ Of course, **only valid when oversmoothing does not occur...**
  - ▶ i.e., not with usual aggregation functions like average over neighbors



# Oversmoothing vs Sufficient Depth

**Oversmoothing** is a well-studied phenomenon “preventing” GNNs from being “too deep” in practice. E.g., for mean aggregation: 
$$L^k Z \xrightarrow[k \rightarrow \infty]{} c1_n$$





# Oversmoothing vs Sufficient Depth

But... most analyses showing the power of GNNs **take the limit  $k \rightarrow \infty$  !**

(*not* for mean aggregation, obviously)

- sufficiently deep GNNs are “**Weisfeiler-Lehman**” powerful [Xu et al. 2019]

- some GNNs model a **diffusion process** that separates well data, etc

[Bodnar et al. 2022]

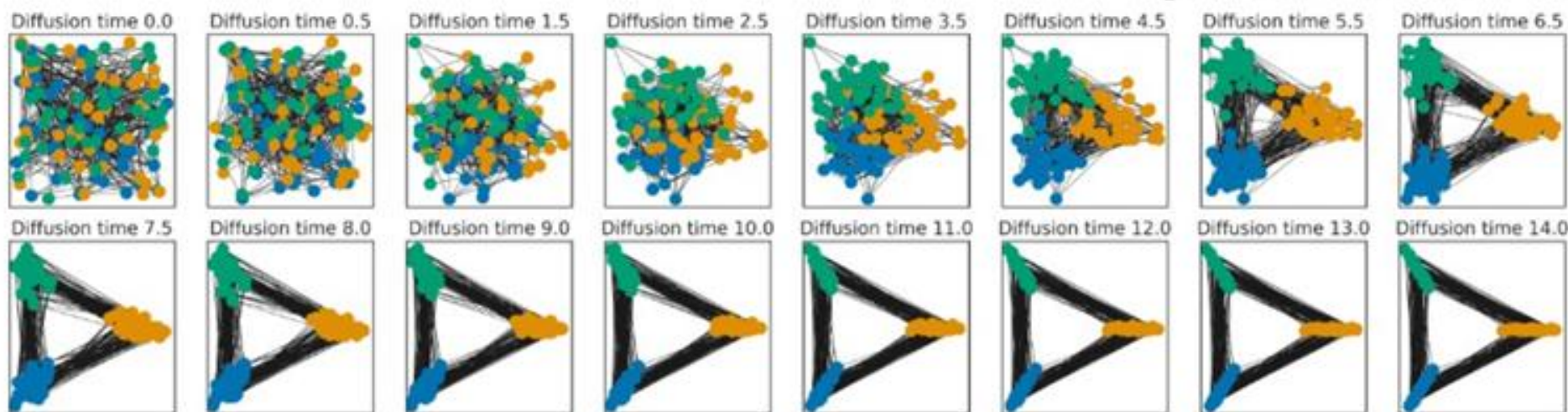
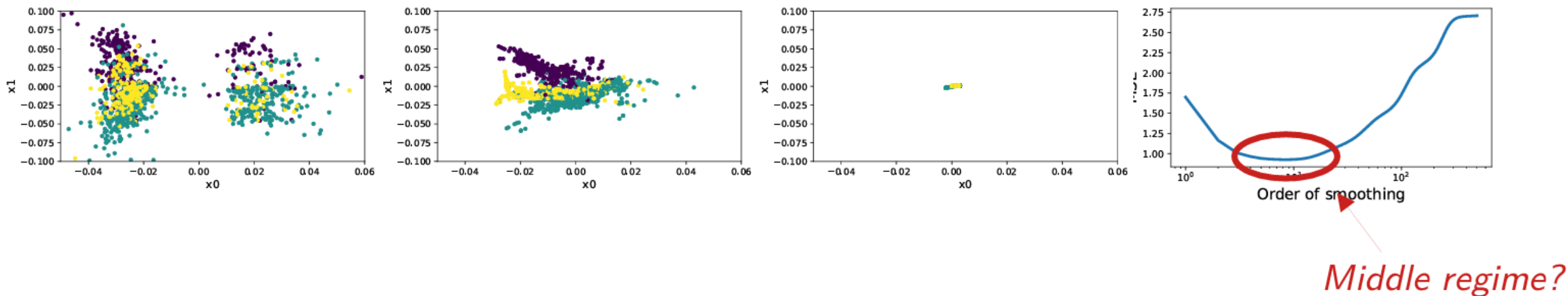


Figure 7. Sheaf diffusion process disentangling the  $C = 3$  classes over time. The nodes are coloured by their class.



# Oversmoothing vs Sufficient Depth

Can “good smoothing” and oversmoothing co-exist? *Why?*



*Take-home message: smoothing collapses node features, but **not everything collapses at the same speed***

## Latent Space Random Graph Model

Each node  $i$  is associated with a **latent variable**  $x_i \in \mathbb{R}^d$  with  $d \gg p$ .

A complete graph with **edge weights**

$$a_{ij} = W(x_i, x_j) = \epsilon + e^{-\frac{1}{2}\|x_i - x_j\|_2^2}$$

Latent variables and node labels  $(x_i, y_i)$  are drawn iid from some distribution.

**Node features** are a linear projection of the latent variables to a lower dimension:

$$z_i = M^T x_i$$

for some unknown  $M \in \mathbb{R}^{d \times p}$  and  $M^T M = I$ .

## Over-smoothing as $k \rightarrow \infty$

$L = D^{-1}A$  is a stochastic matrix,

$$L^k \rightarrow 1_n \bar{d}^T, \text{ where } \bar{d}_i = \frac{\text{degree of node } i}{\text{sum of all degrees}}.$$

Therefore  $Z^{(k)} = L^k Z \rightarrow 1_n \bar{d}^T Z$ , i.e., each node has identical representation.

Since the test nodes have identical representations, the predictions will be the same:  $\hat{y}_{te}^{(k)} = Z^{(k)} \hat{\beta}^{(k)} \rightarrow c 1_{n_{te}}$  for some constant  $c$ .

Using the closed-form solution for the ridge regression can get

$$c = \frac{1}{n_{tr}} \left( \frac{\|Z^T \bar{d}\|_2^2}{\|Z^T \bar{d}\|_2^2 + \lambda} \right) \sum_{i=1}^{n_{tr}} y_i \quad \text{average training labels}$$

# Model of Random Graph

We model **both phenomena** at once:

**We give simple examples in which finite smoothing provably helps learning, before oversmoothing kicks in.**

- ▶ Smoothing collapses features, but **not everything collapse at the same speed**
- ▶ **Some subspaces may collapse faster**, which helps regression
- ▶ **Communities collapse onto themselves**, which helps classification

- ▶ Unknown **latent variables**  $x_i \in \mathbb{R}^d$ , labels  $y_i \in \mathbb{R}$ :

$$(x_i, y_i) \stackrel{iid}{\sim} P, \quad 1 \leq i \leq n$$

- ▶ **Graph structure**: Gaussian kernel

$$a_{ij} = W(x_i, x_j), \quad W(x, x') = e^{-\|x-x'\|^2} + \epsilon$$

- ▶ For simplicity: no “Bernoulli edges” & small  $\epsilon > 0$
- ▶ **Node features**: **partial observation**

$$z_i = Mx_i \in \mathbb{R}^p, \quad p < d$$

- ▶ Does **not** satisfy JL lemma: **loss of information**

**Can smoothing (and only smoothing, here) help recover lost information before oversmoothing dominates?**

# Model of Random Graph

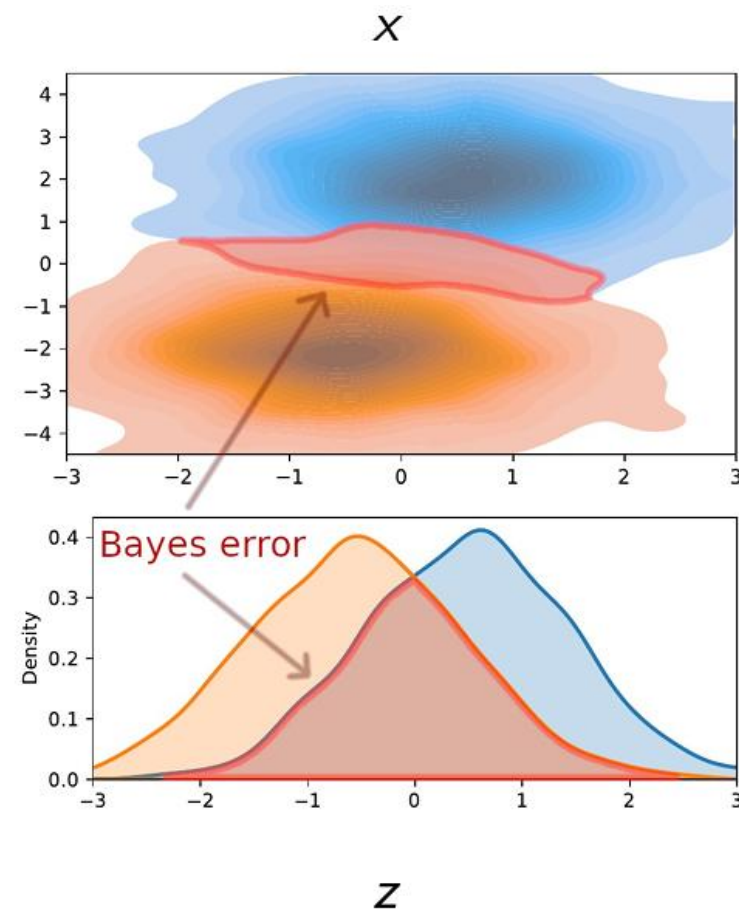
## Random graph model:

$$(x_i, y_i) \sim P, \quad a_{ij} = W(x_i, x_j), \quad z_i = Mx_i$$

With  $M \in \mathbb{R}^{p \times d}$ ,  $p < d$   $W(x, x') = e^{-\|x - x'\|^2} + \epsilon$

No Johnson-Lindenstrauss here. There **is loss of information** in the node features.

Can **mean aggregation** recover some of the information **before oversmoothing occurs** ?



We will focus on **linear GCN** with a Mean Square Error (MSE) loss.

Given an input matrix  $Z \in \mathbb{R}^{n \times p}$ , the output after  $k$  rounds of mean aggregation is

$$Z^{(k)} = L^k Z$$

where  $L = D^{-1}A$  is the normalized adjacency matrix.

Consider learning with MSE loss and ridge regularization

$$\min_{\beta} \frac{1}{2n_{tr}} \left\| y_{tr} - Z_{tr}^{(k)} \beta \right\|_2^2 + \lambda \|\beta\|_2^2$$

where subscript *tr* means training.

Denote

$$\hat{\beta}^{(k)} = \operatorname{argmin}_{\beta} \frac{1}{2n_{tr}} \left\| y_{tr} - Z_{tr}^{(k)} \beta \right\|_2^2 + \lambda \|\beta\|_2^2$$

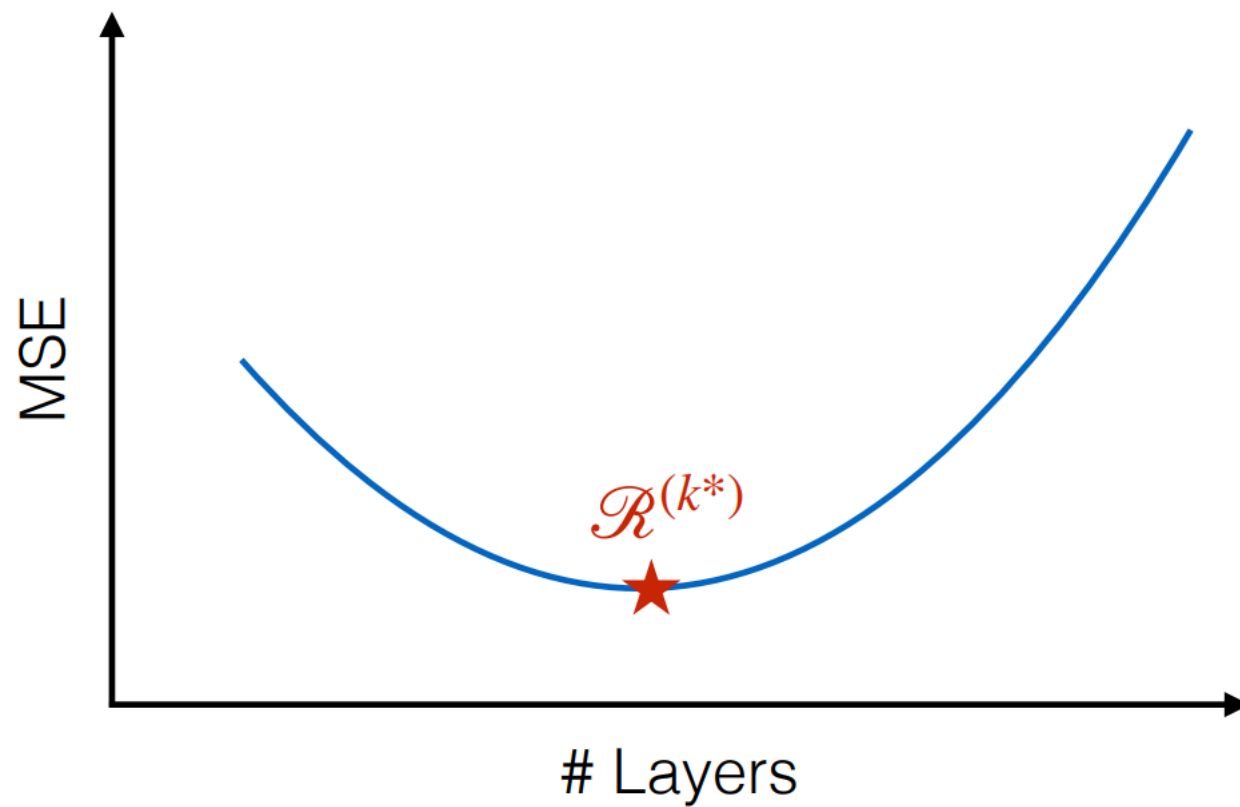
The test risk is

$$\mathcal{R}^{(k)} = \frac{1}{n_{te}} \left\| y_{te} - Z_{te}^{(k)} \hat{\beta}^{(k)} \right\|_2^2$$

- $\mathcal{R}^{(0)}$  is the test risk without any GCN layer
- $\mathcal{R}^{(\infty)}$  denotes the asymptotic test risk as  $k \rightarrow \infty$
- **Over-smoothing:**  $\mathcal{R}^{(0)} < \mathcal{R}^{(\infty)}$
- **Key result in this paper:**  $\exists k^* \geq 1$  such that  $\mathcal{R}^{(k^*)} < \min\{\mathcal{R}^{(0)}, \mathcal{R}^{(\infty)}\}$



# Theoretical Risk



# Settings: Ridge Regression and SSL

- **Linear GNN** (also called SGC [Wu et al. 2019])

$$\hat{Y} = Z^{(k)} \beta \text{ with } Z^{(k)} = L^k Z$$

- **Semi-Supervised Learning**  $n_{tr}, n_{te} \sim n$

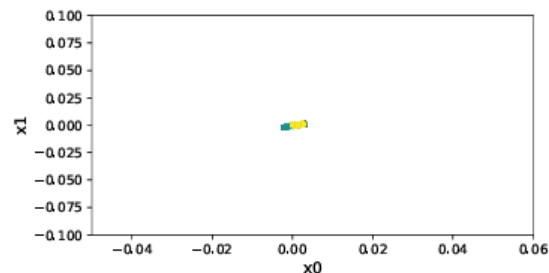
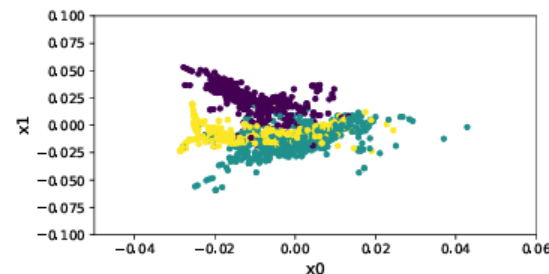
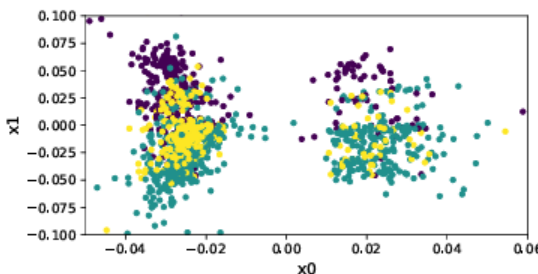
- **Ridge Regression**

$$\beta^{(k)} = \arg \min_{\beta} \frac{1}{n_{tr}} \|Z_{tr}^{(k)} \beta - Y_{tr}\|^2 + \lambda \|\beta\|^2$$

- **Test risk**

$$\mathcal{R}^{(k)} = \frac{1}{n_{te}} \|Y_{te} - Z_{te}^{(k)} \beta^{(k)}\|^2$$

**Thm: Oversmoothing**  $Z_{te}^{(k)} \beta^{(k)} \xrightarrow{k \rightarrow \infty} C 1_{n_{te}}$



**Goal:** show there is  $k^*$  s.t.

$$\mathcal{R}^{(k^*)} < \min(\mathcal{R}^{(0)}, \mathcal{R}^{(\infty)})$$

# Regression

Regression settings:  $x \sim \mathcal{N}(0, \Sigma)$ ,  $y = x^\top \beta^*$

**Thm:** if  $\Sigma, \beta^*, M$  are “well-aligned” and  $n$  is large enough,  $k^*$  exists.

**Intuition:**  $L^k X$  behaves “almost” as  
 $\mathcal{N}(0, (\text{Id} + \Sigma^{-1})^{-k} \Sigma)$

- The small eigenvalues shrink **faster** than the large ones  $\lambda_i \leftarrow \lambda_i / (1 + 1/\lambda_i)^k$
- If well-aligned (“*homophily*”), smoothing helps
- If inversely aligned (“*heterophily*”), smoothing never helps
- Proof not that simple: for  $k > 0$ , **dependent rows of  $Z$**

# Regression

Assume the latent variable  $x \sim \mathcal{N}(0, \Sigma)$  and the label  $y = x^T \beta^\star$ .  
We observe node features  $z = M^T x \in \mathbb{R}^p$ . Recall  $x \in \mathbb{R}^d$  and  $d \gg p$ .

**Example:**  $\mathcal{R}^{(\infty)}$

$y$  follows a mean 0 normal distribution.

Recall that as  $k \rightarrow \infty$ , the predictions are  $\hat{y}_{te}^{(k)} = c 1_{n_{te}}$ , where  $c = c' \sum_{i=1}^{n_{tr}} y_i$ .

In the infinite sample limit, as  $n \rightarrow \infty$ , we have  $c \rightarrow \mathbb{E}[y] = 0$ .

This means that, as  $n, k \rightarrow \infty$ , the predictions  $\hat{y}_{te}^{(k)} \rightarrow 1_{n_{te}} \mathbb{E}[y] = 0$ .

Consequently, as  $n \rightarrow \infty$ ,  $\mathcal{R}^{(\infty)} \rightarrow \mathbb{E}|y|^2 = \text{Var}(y) = \beta^{\star T} \Sigma \beta^\star$ .

What about  $\mathcal{R}^{(0)}$  and more generally  $\mathcal{R}^{(k)}$ ?

# Regression

Given a psd matrix  $S \in \mathbb{R}^{d \times d}$ , define

$$R_{\text{reg}}(S) = (\Sigma^{1/2}\beta^\star)^T (I - S^{1/2}M(\lambda I + M^T S M)^{-1}M^T S^{1/2})^2 (\Sigma^{1/2}\beta^\star) \in \mathbb{R}_+$$

**[Theorem 3]**

$$\mathcal{R}^{(0)} = R_{\text{reg}}(\Sigma) + O\left(\frac{\text{poly}(\|\Sigma\|, \|\beta^\star\|, d)}{\sqrt{n}}\right)$$

$\xrightarrow{n \rightarrow \infty} 0$

➔  $\mathcal{R}^{(0)} \approx R_{\text{reg}}(\Sigma) \leq \text{Var}(y) \approx \mathcal{R}^{(\infty)}$

# Regression

**[Theorem 2]**  $\exists k^* \geq 1, \mathcal{R}^{(k^*)} < \min\{\mathcal{R}^{(0)}, \mathcal{R}^{(\infty)}\}$

Follows from Theorems 3 and 4 and **Assumption 1**:  $R_{\text{reg}}(\Sigma) > R_{\text{reg}}(\Sigma^{(1)})$

**[Theorem 3]**

$$\mathcal{R}^{(0)} = R_{\text{reg}}(\Sigma) + O\left(\frac{\text{poly}(\|\Sigma\|, \|\beta^*\|, d)}{\sqrt{n}}\right)$$

$n \rightarrow \infty \rightarrow 0$

**[Theorem 4]**

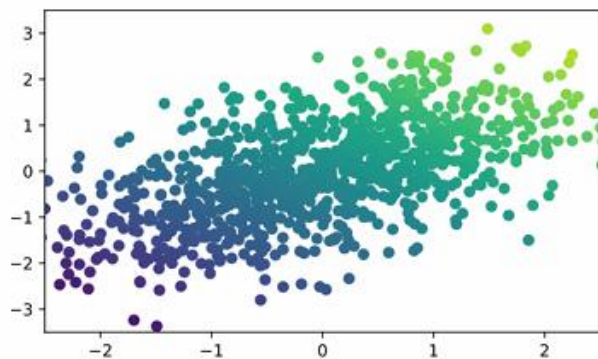
$$\mathcal{R}^{(1)} = R_{\text{reg}}(\Sigma^{(1)}) + O\left(\frac{\text{poly}(\|\Sigma\|, \|\beta^*\|, d, \epsilon^{-1})}{\sqrt{n}}\right) + O\left(\frac{C}{\epsilon^{1/5}}\right)$$

$n \rightarrow \infty \rightarrow 0$

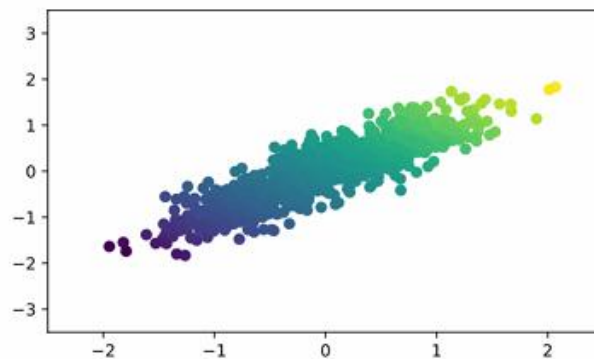
$\Sigma^{(k)} = (I + \Sigma^{-1})^{-2k} \Sigma$

# Regression

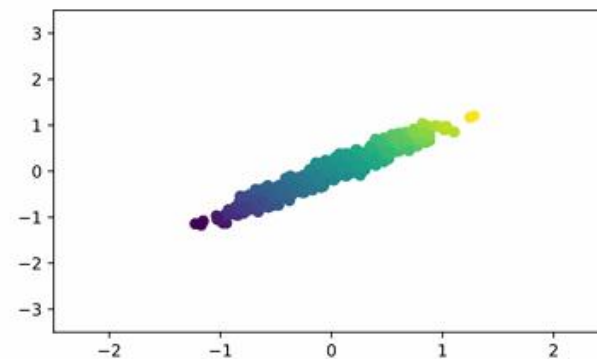
0 mean aggregation



1 mean aggregation

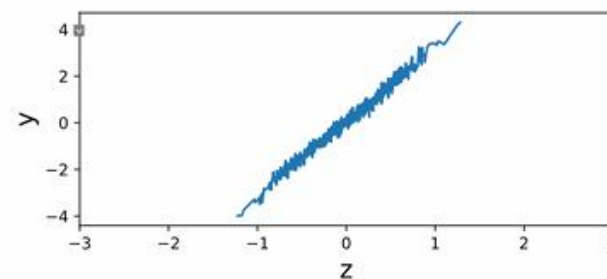
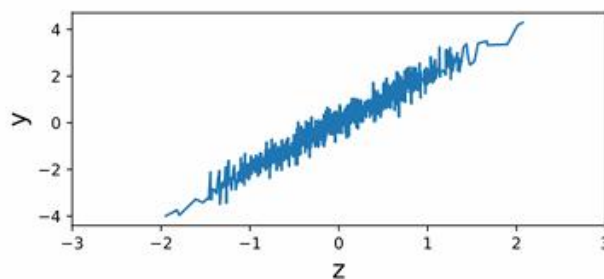
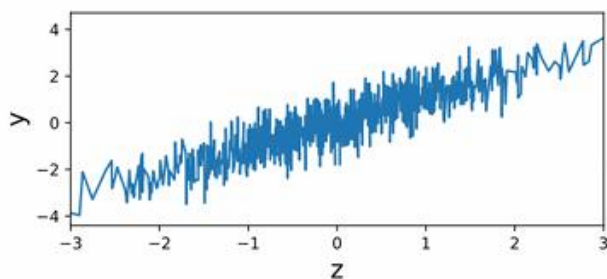


2 mean aggregations



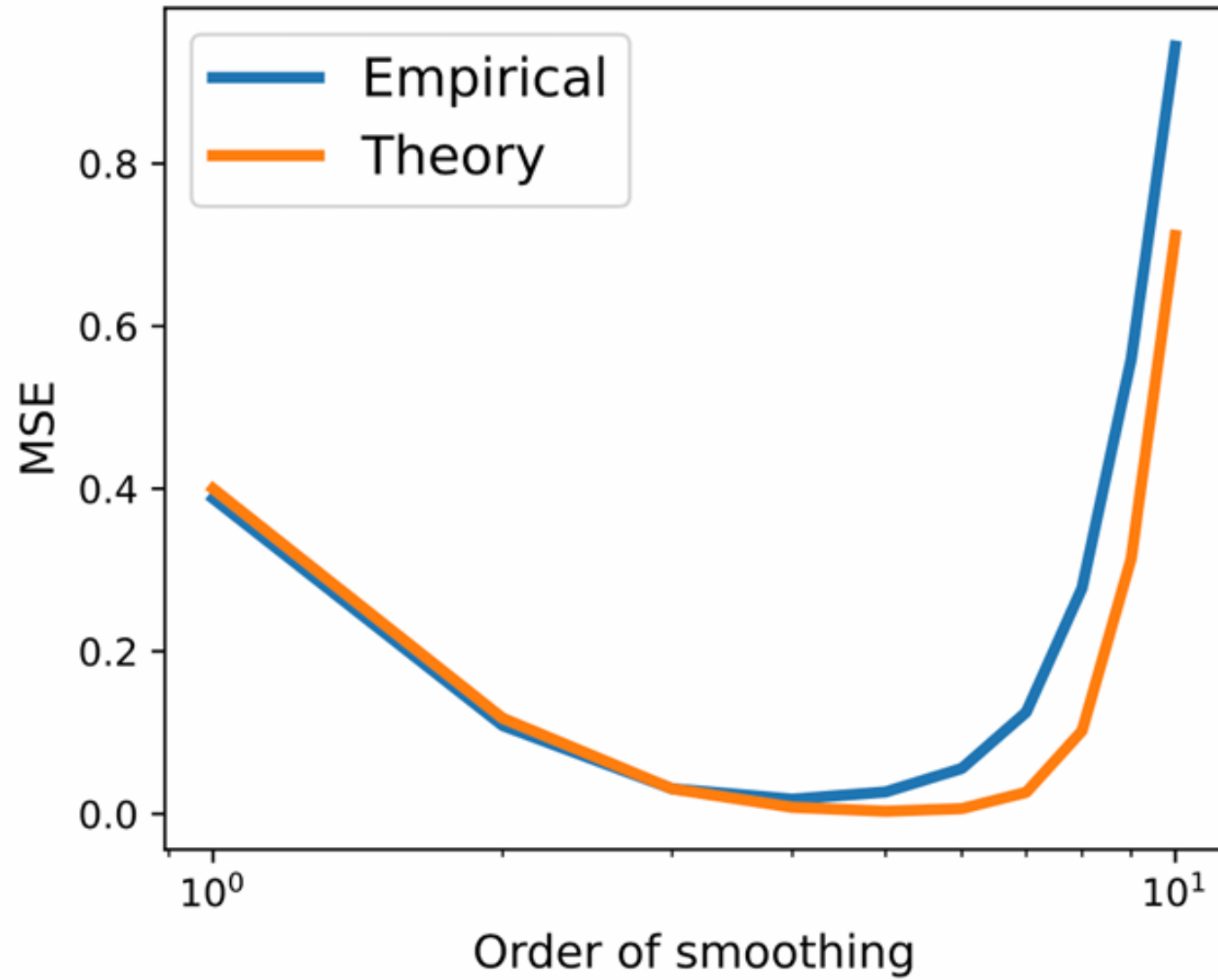
2-dim latent  
variable  $x$   
Color reflects  $y$

1-dim observed  
variable  $z$   
versus labels  $y$





# Regression



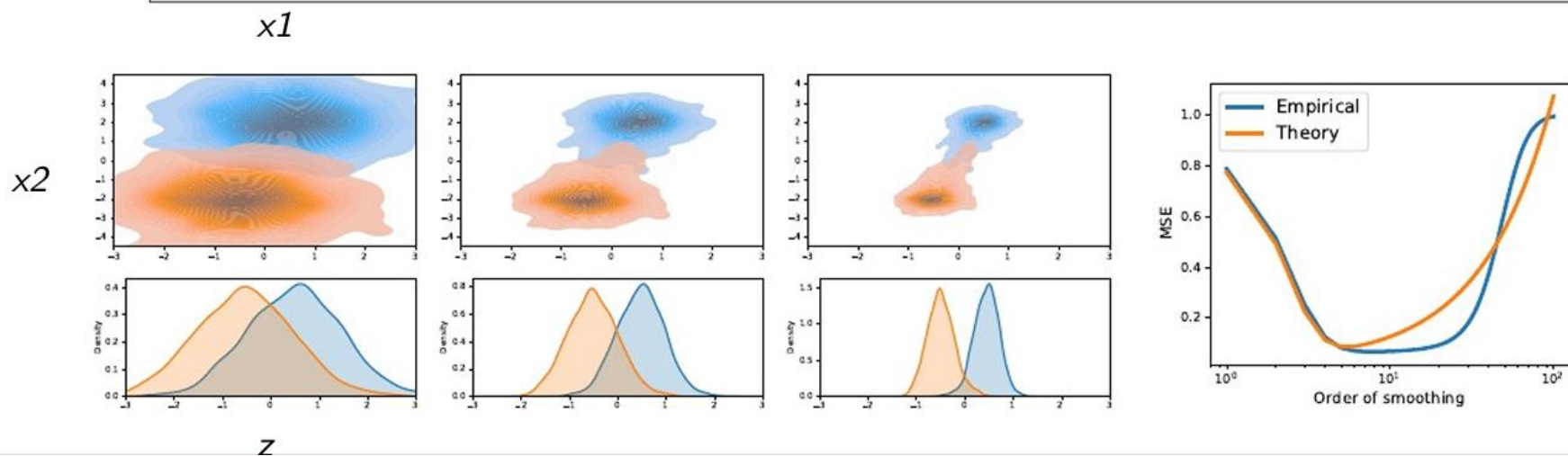
# Classification

Classif. settings:  $(x, y) \sim \frac{1}{2}\mathcal{N}(\mu, \text{Id}) \otimes \{1\} + \frac{1}{2}\mathcal{N}(-\mu, \text{Id}) \otimes \{-1\}$

**Thm:** if  $\|\mu\|, n$  are large enough and  $\|M\mu\| > 0$ ,  $k^*$  exists.

Intuition:

*The communities (initially) concentrate faster than they get close to each other.*



► This can help **separate node features  $z_i$**  before oversmoothing happens

# Conclusion

- ▶ **Beneficial smoothing and oversmoothing** generally coexist
- ▶ Theoretical analysis may give hints as to why
- ▶ There are links with the distinction **homophily/heterophily**

## Outlooks:

- ▶ Can take inspiration to **combat oversmoothing less indiscriminatively?**
  - ▶ eg, smarter normalization in GNNs to fight oversmoothing
- ▶ How to better describe and exploit the interaction between **labels, node features and graph structure?**

# Thank you!

Khushee Kapoor