

Graph neural networks as dynamical systems

Francesco Di Giovanni

Twitter

First Italian School in GDL: July 25–28, Pescara

Presentation outline

- | Graph preliminaries
- | Spectral analysis and Dirichlet energy on graphs
- | Dynamical systems on graphs
- | MPNNs as multi-particle systems and the gradient flow framework (GRAFF)
- | Presentation of *Graph Neural Networks as Gradient Flows*

Introduction

Preliminaries on graph operators

- | $G = (V, E)$ is an *undirected* graph with $|V| = n$ and $i \sim j$ if $(i, j) \in E$
- | \mathbf{A}, \mathbf{D} are $n \times n$ adjacency and (diagonal) degree matrices
- | The *normalized* adjacency is $\bar{\mathbf{A}} = \mathbf{D}^{-\frac{1}{2}} \mathbf{A} \mathbf{D}^{-\frac{1}{2}}$
- | The **Laplacian** $\mathbf{L} = \mathbf{D} - \bar{\mathbf{A}}$ is an operator acting on signals $\mathbf{f} : V \rightarrow \mathbb{R}$ as

$$(\mathbf{L} \mathbf{f})_i = f_i - \sum_{j \sim i} \frac{f_j}{d_i d_j}$$

Preliminaries on graph operators

- | $G = (V, E)$ is an *undirected* graph with $|V| = n$ and $i \sim j$ if $(i, j) \in E$
- | \mathbf{A}, \mathbf{D} are $n \times n$ adjacency and (diagonal) degree matrices
- | The *normalized* adjacency is $\bar{\mathbf{A}} = \mathbf{D}^{-\frac{1}{2}} \mathbf{A} \mathbf{D}^{-\frac{1}{2}}$
- | The **Laplacian** $\mathbf{L} = \mathbf{D} - \bar{\mathbf{A}}$ is an operator acting on signals $\mathbf{f} : V \rightarrow \mathbb{R}$ as

$$(\mathbf{L}\mathbf{f})_i = f_i - \sum_{j \sim i} \frac{f_j}{d_i d_j}$$

The Laplacian \mathbf{L} eigenvalues satisfy $0 = \lambda_0 < \lambda_1 < \dots < \lambda_{n-2} < \lambda_{n-1} = 2n - 2$, with $\lambda_0 = 0$ and $\lambda_{n-1} = 2n - 2$, and are called (graph) *frequencies*, eigenvectors are denoted by $\{\mathbf{v}_k\}_{k=0}^{n-1}$

Signal on graphs: Dirichlet energy and smoothness

Consider a signal (feature) $f : V \rightarrow \mathbb{R}$ e.g. temperature of each node

We write $\mathbf{f} = (f_1, \dots, f_n)$ $\mathbf{f} = c$

can be used to measure smoothness of f : the **Dirichlet energy**^[1] E^{Dir} is defined by

$$E^{\text{Dir}}(\mathbf{f}) := \frac{1}{4} \sum_{i,j} \left\| \frac{f_i}{d_i} - \frac{f_j}{d_j} \right\|^2 = \frac{1}{2} \mathbf{f}^T \mathbf{L} \mathbf{f} = \frac{1}{2} \sum_{i,j} c_{ij} (f_i - f_j)^2$$

[1] Zhou and Schölkopf (2005)

Signal on graphs: Dirichlet energy and smoothness

Consider a signal (feature) $f : V \rightarrow \mathbb{R}$ e.g. temperature of each node

We write $\mathbf{f} = (f_1, \dots, f_n)$ $\mathbf{f} = c$

can be used to measure smoothness of f : the **Dirichlet energy**^[1] E^{Dir} is defined by

$$E^{\text{Dir}}(\mathbf{f}) := \frac{1}{4} \sum_{i,j} \left\| \frac{f_i}{d_i} - \frac{f_j}{d_j} \right\|^2 = \frac{1}{2} \mathbf{f}^T \mathbf{L} \mathbf{f} = \frac{1}{2} \sum_{i,j} \mathbf{L}_{ij} f_i f_j$$

the frequency components of \mathbf{f} determine the variations of the signal along edges

The quantity $f_i/d_i - f_j/d_j := \mathbf{g}(i,j)$ is the **gradient** of f along (i,j) E

^[1] Zhou and Schölkopf (2005)

A rough picture: low-pass vs high-pass filtering

Consider a dynamical process $t \mapsto \mathbf{f}(t) \in \mathbb{R}^n$ starting at \mathbf{f}_0 . $\dot{\mathbf{f}}(t) = \mathbf{c}(\mathbf{f}(t))$

A rough picture: low-pass vs high-pass filtering

Consider a dynamical process $t \mapsto \mathbf{f}(t) \in \mathbb{R}^n$ starting at \mathbf{f}_0 $\mathbf{f}(t) = \mathbf{c}(t)$

If the high-frequency components $|c(t)|$, with $\gamma \gg 0$, decrease with time, then the process acts as '**low-pass filtering**' — smooths the signal out

A rough picture: low-pass vs high-pass filtering

Consider a dynamical process $t \mapsto \mathbf{f}(t) \in \mathbb{R}^n$ starting at \mathbf{f}_0 $\mathbf{f}(t) = e^{-\lambda t} \mathbf{c}(t)$

If the high-frequency components $|c_i(t)|$, with $\lambda_i \gg 0$, decrease with time, then the process acts as '**low-pass** filtering' — smooths the signal out

If the low-frequency components $|c_i(t)|$, with $\lambda_i \approx 0$, decrease with time, then the process acts as '**high-pass** filtering' — sharpens the signal

Figure 1: First four Laplacian eigenvectors of Minnesota Road graph. Figure taken from [Bronstein et al. \(2017\)](#)

A prototypical low-pass filtering: the graph heat equation

Consider an input signal $f_0 : V \rightarrow \mathbb{R}$ and recall that $\mathbf{f} \in E^{\text{Dir}}(\mathbf{f}) = \frac{1}{2} \mathbf{f}, \quad \mathbf{f}$

If we want to *minimize* E^{Dir} take infinitesimal steps in the direction of steepest descent

$$\text{Heat equation : } \dot{\mathbf{f}}(t) = - \mathbf{f} E^{\text{Dir}}(\mathbf{f}(t)) = - \mathbf{f}(t), \quad \mathbf{f}(0) = \mathbf{f}_0.$$

This is a **gradient flow**: $E^{\text{Dir}}(\mathbf{f}(t)) = 0$ and $\mathbf{f}(t) = \mathbf{f}$ s.t. $\mathbf{f} = \mathbf{0}$ i.e.
 $\mathbf{f} \in \text{span}(\bar{d}_1, \dots, \bar{d}_n)$

Low-pass dynamics ‘features become indistinguishable’ when $t \gg 1$

Multiple channels

Consider $F : V \rightarrow \mathbb{R}^d$ with matrix representation $F \in \mathbb{R}^{n \times d}$. E^{Dir} can be extended as

$$E^{\text{Dir}}(F) = \frac{1}{4} \sum_{(i,j) \in E} \left\| \frac{\mathbf{f}_i}{d_i} - \frac{\mathbf{f}_j}{d_j} \right\|^2 = \frac{1}{2} \text{trace}(F^T F)$$

The gradient flow of E^{Dir} yields heat equation in each feature channel^[2]:

$$\dot{\mathbf{f}}^r(t) = -\mathbf{f}^r(t), \quad 1 \leq r \leq d$$

^[2] 'Channels' = 'feature components' = 'feature coordinates'

The formalism

We can vectorize a matrix $F \in \mathbb{R}^{n \times d}$! $\text{vec}(F) \in \mathbb{R}^{nd}$

We use the Kronecker product $I_d \otimes I_n \in \mathbb{R}^{nd \times nd}$ to rewrite E^{Dir} as

$$E^{\text{Dir}}(F) = \frac{1}{2} \text{vec}(F); (I_d \otimes I_n) \text{vec}(F)$$

The heat equation can also be rewritten by 'stacking the columns as'

$$\text{vec}(E(t)) = (I_d \otimes I_n) \text{vec}(F(t))$$

Upshot: this formalism reduces matrix ODE to a vector ODE! vectorized ODEs are much easier to deal with

A motivating example

How to determine if a dynamical process on a graph is dominated by the low or high frequencies?

A motivating example

How to determine if a dynamical process on a graph is dominated by the low or high frequencies? Use E^{Dir} after normalization

A motivating example

How to determine if a dynamical process on a graph is dominated by the low or high frequencies? Use E^{Dir} after normalization

Consider $\dot{F}(t) = AF(t)$ $\vec{f}(F(t)) = (I_d \quad A)\vec{f}(F(t))$, with $F(0) = F_0$

Recall that $A = I$ so we can solve as

$$f^r(t) = e^{At} f^r(0) = e^{(I - A)t} f^r(0); \quad 1 \leq r \leq d$$

A motivating example

How to determine if a dynamical process on a graph is dominated by the low or high frequencies? Use E^{Dir} after normalization

Consider $\dot{F}(t) = AF(t)$ $\vec{f}^r(t) = (I_d - A)\vec{f}^r(t)$, with $F(0) = F_0$

Recall that $A = I - \dots$ so we can solve as

$$\vec{f}^r(t) = e^{A t} \vec{f}^r(0) = e^{-(I - A)t} \vec{f}^r(0); \quad 1 \leq r \leq d$$

Expand each channel in the basis $\{g_i\}$ satisfying $A g_i = (1 - \lambda_i) g_i$:

$$\vec{f}^r(t) = \sum_i e^{-(1 - \lambda_i)t} \langle g_i, \vec{f}^r(0) \rangle g_i$$

A motivating example

Recall that \mathbf{v}_0 is the smoothest eigenvector i.e. $\mathbf{v}_0 = \mathbf{0}$

The projection along \mathbf{v}_0 is the one growing the fastest^[3] since

$$\mathbf{h}^r(t); \mathbf{v}_0 = e^{(1-\lambda)t} \mathbf{h}^r(0); \mathbf{v}_0$$

The dynamics are 'dominated' by the low-frequencies: $\mathbf{F}^D(\mathbf{F}(t)) \rightarrow \mathbf{0}$?

[3] Unless $\mathbf{h}^r(0); \mathbf{v}_0 = 0$ which is only true in a smaller subspace of \mathbb{R}^n

[4] Unless $\mathbf{h}^r(0); \mathbf{v}_i = 0$ for all $\lambda_i > 0$

A motivating example

Recall that \mathbf{v}_0 is the smoothest eigenvector i.e. $\mathbf{v}_0 = \mathbf{0}$

The projection along \mathbf{v}_0 is the one growing the fastest^[3] since

$$\langle \mathbf{f}^r(t); \mathbf{v}_0 \rangle = e^{(1-\lambda_0)t} \langle \mathbf{f}^r(0); \mathbf{v}_0 \rangle$$

The dynamics are 'dominated' by the low-frequencies: $E^{\text{Dir}}(\mathbf{f}^r(t)) \rightarrow \mathbf{0}$? No!^[4]

$$E^{\text{Dir}}(\mathbf{f}^r(t)) = \frac{1}{2} \langle \mathbf{f}^r(t); \mathbf{f}^r(t) \rangle = \sum_{\lambda > 0} e^{(1-\lambda)t} \langle \mathbf{f}^r(0); \mathbf{v}_\lambda \rangle^2 \neq 1$$

[3] Unless $\langle \mathbf{f}^r(0); \mathbf{v}_0 \rangle = 0$ which is only true in a smaller subspace \mathbb{R}^d

[4] Unless $\langle \mathbf{f}^r(0); \mathbf{v}_\lambda \rangle = 0$ for all $\lambda > 0$

A motivating example

Looking at E^{Dir} is not enough! we should normalize first: in fact we have

$$E^{\text{Dir}}(\|F(t)\|) \neq \|E^{\text{Dir}}(F(t))\|; \quad t \neq 1$$

A motivating example

Looking at E^{Dir} is not enough! we should normalize \mathbf{r} : in fact we have

$$E^{\text{Dir}}(\mathbf{F}(t) = \|\mathbf{F}(t)\| \mathbf{r}) \neq 0; \quad t \neq 1$$

and for each channel r $\exists f_1^r$ s.t.

$$f^r(t) = \|\mathbf{f}^r(t)\| \mathbf{r} \neq f_1^r; \quad \mathbf{f}_1^r = 0$$

A motivating example

Looking at E^{Dir} is not enough! we should normalize first: in fact we have

$$E^{\text{Dir}}(F(t)) = \int |F(t)|^2 dt \neq 0; \quad \int |F(t)|^2 dt = 1$$

and for each channel r $d \geq 1$ f_1^r s.t.

$$f_1^r(t) = \int |f_1^r(t)|^2 dt = f_1^r; \quad \int f_1^r = 0$$

Upshot: Analyse $F(t)$ via $E^{\text{Dir}}(F(t)) = \int |F(t)|^2 dt$! **Rayleigh quotient** of d

Definition

A dynamical system $F(t)$ initialized at $F(0)$ is Low-Frequency-Dominant (LFD) if $E^{\text{Dir}}(F(t) = \|F(t)\|) \rightarrow 0$ for $t \rightarrow \infty$.

Definition

A dynamical system $F(t)$ initialized at $F(0)$ is Low-Frequency-Dominant (LFD) if $E^{\text{Dir}}(F(t) = \text{jj} F(t) \text{jj}) \rightarrow 0$ for $t \rightarrow \infty$.

Does it make sense?

Lemma

A dynamical system is LFD iff for each sequence $t_k \rightarrow \infty$ there exist a subsequence $t_{j_k} \rightarrow \infty$ and F_1 s.t. $F(t_{j_k}) \rightarrow \text{jj} F(t_{j_k}) \text{jj} \rightarrow F_1$ and $f_{j_k} \rightarrow 0$.

LFD dynamics: numerical example

A numerical example of LFD dynamics $\bar{x} = 4:0$, $\bar{y} = 0:1$

$$E(t) = AF(t) ; \quad = \begin{matrix} & " & \# \\ & 1 & 0 \\ & 0 & 0 \end{matrix}$$

LFD dynamics: numerical example

A numerical example of LFD dynamics: $\lambda_1 = 4:0$, $\lambda_2 = 0:1$

$$\dot{F}(t) = AF(t) ; \quad A = \begin{matrix} & \text{"} & \text{\#} \\ \begin{matrix} 1 & 0 \\ 0 & 0 \end{matrix} \end{matrix}$$

In both cases the eigenvector e_1 dominates the dynamics

- | Top: solution becomes unbounded
- | Bottom: evolution of $\|F(t)\| = \|F(0)\| e^{-\lambda_2 t}$
 - ! convergence $t_{ker}(\epsilon)$ where
 - we only distinguish nodes based on their degree

High-frequency-dominant: HFD

Note that $E^{\text{Dir}}(F) = \frac{1}{2} \sum_{j,j} F_{jj}^2 \leq E^{\text{Dir}}(F) = \frac{1}{2} \sum_{j,j} F_{jj}^2$

High-frequency-dominant: HFD

Note that $E^{\text{Dir}}(F) = \frac{1}{2} \|F\|^2 = E^{\text{Dir}}(F = \|F\|)$

Definition

A dynamical system $F(t)$ initialized at $F(0)$ is High-Frequency-Dominant (HFD) if $E^{\text{Dir}}(F(t) = \|F(t)\|) \approx 2$ for $t \gg 1$.

High-frequency-dominant: HFD

Note that $E^{\text{Dir}}(F) = \frac{1}{2} \sum_j |F_j|^2 \leq E^{\text{Dir}}(F = \sum_j F_j) = \frac{1}{2}$

Definition

A dynamical system $F(t)$ initialized at $F(0)$ is High-Frequency-Dominant (HFD) if $E^{\text{Dir}}(F(t) = \sum_j F_j(t)) \leq \frac{1}{2}$ for $t \geq 1$.

Does it make sense?

Lemma

A dynamical system is HFD iff for each sequence $t_k \geq 1$ there exist a subsequence $t_{j_k} \geq 1$ and F_1 s.t. $F(t_{j_k}) = \sum_j F_j(t_{j_k}) \leq F_1$ and $f_1^r = f_1^r$.

Why do we need HFD?

Consider $F(t) = A F(t)$! eigenvector dominates the dynamics

- | Evolution of $F(t) = \|F(t)\|$! convergence $t \rightarrow \infty$ where we distinguish nodes based on the largest frequency eigenvector (right figure)

Homophily vs heterophily aka short vs long range interactions

Semi-supervised setting V_{tr} V labelled! predict labels on V_{test}

Homophily: Neighbours often share labels labels are smooth i.e. low-pass is 'good'

Heterophily: 1 homophily! labels are not smooth i.e. low-pass is 'bad'

Homophily vs heterophily aka short vs long range interactions

Semi-supervised setting V_{tr} V labelled! predict labels on V_{test}

Homophily: Neighbours often share labels labels are smooth i.e. low-pass is 'good'

Heterophily: 1 homophily! labels are not smooth i.e. low-pass is 'bad'

Dual perspective: short-range relations vs long-range relations relevant for graph classification and regression tasks on molecules

A layer of Graph Convolutional Network (GCN)^[5] is defined by:

$$F(t+1) = \text{ReLU}(AF(t)W(t))$$

A is the message-passing matrix and W(t) is the 'channel-mixing'

[5] Kipf and Welling (2017)

[6] Nt and Maehara (2019); Oono and Suzuki (2020); Cai and Wang (2020)

A layer of Graph Convolutional Network (GCN)^[5] is defined by:

$$F(t+1) = \text{ReLU}(AF(t)W(t))$$

A is the message-passing matrix and W(t) is the 'channel-mixing'

- | Poor performance on heterophilic graphs
- | Degradation when increasing depth (over-smoothing)^[6]

[5] Kipf and Welling (2017)

[6] Nt and Maehara (2019); Oono and Suzuki (2020); Cai and Wang (2020)

Theorem (Cai and Wang)

Let $(1 - \sigma_T)^2 := \max_t (1 - \sigma_t)^2$ and $s_T = \max_t \sigma_t \sin(W(t))$. Then the solution $F(T)$ of GCN satisfies

$$E^{\text{Dir}}(F(T)) = (s_T (1 - \sigma_T))^{2T} E^{\text{Dir}}(F(0)):$$

Theorem (Cai and Wang)

Let $(1 - s_T)^2 := \max_t (1 - s_T(t))^2$ and $s_T = \max_t s_T(t) \sin(W(t))$. Then the solution $F(T)$ of GCN satisfies

$$E^{\text{Dir}}(F(T)) = (s_T(1 - s_T))^2 E^{\text{Dir}}(F(0)):$$

- | If singular values of $W(t)$ are controlled in terms of the spectrum of L , solution of GCN becomes increasingly smoother

Theorem (Cai and Wang)

Let $(1 - \sigma_T)^2 := \max_t (1 - \sigma_t)^2$ and $s_T = \max_t \sigma_t \sin(W(t))$. Then the solution $F(T)$ of GCN satisfies

$$E^{\text{Dir}}(F(T)) = (s_T (1 - \sigma_T))^{2T} E^{\text{Dir}}(F(0)):$$

- | If singular values of $W(t)$ are controlled in terms of the spectrum of L , solution of GCN becomes increasingly smoother
- | GCN should succeed with homophily but fail with heterophily

Theorem (Cai and Wang)

Let $(1 - \sigma)^2 := \max_{\lambda} (1 - \lambda)^2$ and $s_T = \max_t \tau \text{sing}(W(t))$. Then the solution $F(T)$ of GCN satisfies

$$E^{\text{Dir}}(F(T)) = (s_T (1 - \sigma))^{2T} E^{\text{Dir}}(F(0)):$$

- | If singular values of $W(t)$ are controlled in terms of the spectrum of L , solution of GCN becomes increasingly smoother
- | GCN should succeed with homophily but fail with heterophily
- | If $T \gg 1$, we converge to $\text{ker}(L)$ i.e. only information to separate nodes by degree

Are graph convolutional models doomed?

- I What if the singular values of W are not bounded by $(1 - \epsilon)^2$?

Are graph convolutional models doomed?

- | What if the singular values of W are not bounded by $(1 - \epsilon)^2$?
- | Can we require more structure on W ?

Are graph convolutional models doomed?

- | What if the singular values of W are not bounded by $(1 - \epsilon)^2$?
- | Can we require more structure on W ?
- | What is the interpretation of W ?

Are graph convolutional models doomed?

- | What if the singular values of W are not bounded by $(1 - \epsilon)^2$?
- | Can we require more structure on W ?
- | What is the interpretation of W ?
- | What is the 'minimal requirement' for a graph convolutional framework to be useful?

Figure 2: ActuaGRAFF dynamics: attractive and repulsive forces lead to a non-smoothing process able to separate labels

Joint w/ J. Rowbottom, B. Chamberlain, T. Markovich, M. Bronstein (2022)

- I We propose a **gradient flow framework (GRAFF)** for MPNNs where the equations follow the direction of steepest descent of learnable energy

Outline of the contributions

- I We propose a **gradient flow framework (GRAFF)** for MPNNs where the equations follow the direction of steepest descent of learnable energy
- I We show how the channel-mixing **W** can learn to induce either **FD** or **HFD** dynamics via its spectrum

Outline of the contributions

- I We propose a **gradient flow framework (GRAFF)** for MPNNs where the equations follow the direction of steepest descent of learnable energy
- I We show how the channel-mixing can learn to induce either FD or HFD dynamics via its spectrum
- I This allows us to interpret MPNNs as multi-particle dynamics with attractive and repulsive forces generated by positive and negative eigenvalues of

Outline of the contributions

- | We propose a **gradient flow framework (GRAFF)** for MPNNs where the equations follow the direction of steepest descent of learnable energy
- | We show how the channel-mixing can learn to induce either LFD or HFD dynamics via its spectrum
- | This allows us to interpret MPNNs as multi-particle dynamics with attractive and repulsive forces generated by positive and negative eigenvalues of \mathbf{W}
- | Show that LFD \Rightarrow HFD dynamics induced by this framework adapt to the underlying homophily/heterophily

Residual networks as discrete ODEs

A ResNet $F(t + \Delta t) = F(t) + \Delta t \text{ResNet}(F(t))$ is the Euler discretization of an ODE (as the step-size $\Delta t \rightarrow 0$)

$$F'(t) = \text{ResNet}(F(t))$$

ODE theory! analysing and improving ResNets

Figure 3: Dynamics of ResNet vs ODE. Figure taken from [Chen et al. \(2018\)](#)

[7] [Haber and Ruthotto \(2018\)](#); [Chen et al. \(2018\)](#)

Residual networks as discrete ODEs

A ResNet $F(t + \Delta t) = F(t) + \Delta t \text{ResNet}(F(t))$ is the Euler discretization of an ODE (as the step-size $\Delta t \rightarrow 0$)

$$F'(t) = \text{ResNet}(F(t))$$

ODE theory! analysing and improving ResNets

What about residual MPNNs?

$$F(t + \Delta t) = F(t) + \Delta t \text{MPNN}(G; F(t)) \quad F'(t) = \text{MPNN}(G; F(t))$$

Figure 3: Dynamics of ResNet vs ODE. Figure taken from [Chen et al. \(2018\)](#)

[7] [Haber and Ruthotto \(2018\)](#); [Chen et al. \(2018\)](#)

The linear GCN^[8] system

$$F(t+1) = AF(t)W(t) \quad E(t) = AF(t)W(t) \quad F(t)$$

[8] Wu et al. (2019)

The linear GCN^[8] system

$$F(t+1) = AF(t)W(t) \quad E(t) = AF(t)W(t) - F(t)$$

If we use the τ -formalism: GCN is the unit step-size discretization of

$$\text{vec}(E(t)) = (W(t)^T A - I)\text{vec}(F(t))$$

! we'll see that the dampening term I is responsible for FD dynamics

[8] Wu et al. (2019)

Continuous Graph Neural Network (CGNN) $\text{set } W = W^> !$

$$E(t) = F(t) + F(t)W + F(0)$$

^[9] Xhonneux et al. (2020)

Continuous Graph Neural Network (CGNN) $\text{set } W = W^> !$

$$E(t) = F(t) + F(t)W + F(0)$$

- | CGNN is a gradient flow
- | We'll prove that this is never HFD
- | Source term $F(0)$ increases expressive power

^[9] Xhonneux et al. (2020)

Graph Neural Diffusion (GRAND)^[10] is the 'continuous' version of GAT^[11]

$$F(t) = (I - A(F(t))) F(t)$$

- | $A(F(t))$ is an attention matrix over the edge set
- | (Linear) GRAND is a diffusion process with maximum principle low-pass filter and over-smoothing

^[10] Chamberlain et al. (2021)

^[11] Velicković et al. (2018)

PDE-GCN_D^[12] is a diffusion process given by

$$F(t) = F(0)W(t)^{\top}W(t)$$

! We'll prove that this is a smoothing process and hence suitable for heterophilic graphs

^[12] Eliasof et al. (2021)

Second-order variants^[13] ! by design they prevent over-smoothing

$$F(t) = \text{MPNN}(G; F(t)) \quad F(t) \quad F(t)$$

However, why oscillatory behaviour? Do we need them?

^[13] Eliasof et al. (2021); Rusch et al. (2022)

Second-order variants^[13] ! by design they prevent over-smoothing

$$F(t) = \text{MPNN}(G; F(t)) \quad F(t) \quad F(t)$$

However, why oscillatory behaviour? Do we need them?

The actual equations are parametric how to choose them?

^[13] Eliasof et al. (2021); Rusch et al. (2022)

Second-order variants^[13] ! by design they prevent over-smoothing

$$F(t) = \text{MPNN}(G; F(t)) \quad F(t) \quad F(t)$$

However, why oscillatory behaviour? Do we need them?

The actual equations are parametric how to choose them?

Upshot: Learn an energy rather than the equations!

^[13] Eliasof et al. (2021); Rusch et al. (2022)

Dynamical systems as gradient flows

Dynamical systems are **gradient flows** when $E : \mathbb{R}^N \rightarrow \mathbb{R}$:

$$\dot{F}(t) = -\text{ODE}(\dot{F}(t)) = -\text{grad}_F E(F(t)) \quad E(F(t)) \rightarrow 0:$$

Gradient flows are easier to analyze and interpret since the solution $F(t)$ is minimizing E .

What if we parametrize an energy rather than the MPNN equations?

Dynamical systems as gradient flows

Dynamical systems are gradient flows when $E : \mathbb{R}^N \rightarrow \mathbb{R}$:

$$\dot{F}(t) = -\text{ODE}(\dot{F}(t)) = -\text{grad}_F E(F(t))$$

Gradient flows are easier to analyze and interpret since the solution $F(t)$ is minimizing E

What if we parametrize an energy rather than the MPNN equations?

Goal: Learn E generalizing E^{Dir} and right notion of smoothness for the problem

$$E(F) = \text{MPNN}(G; F) = \text{grad}_F E(G; F)$$

GNNs as Gradient Flows part 1: taking inspiration from harmonic maps

Harmonic map flow in continuous space

$f : \mathbb{R}^n \rightarrow \mathbb{R}^d$ smooth with h a constant metric. The Dirichlet energy of f is

$$E(f; h) = \frac{1}{2} \int_{\mathbb{R}^n} \| \nabla f \|^2_h dx = \frac{1}{2} \int_{\mathbb{R}^n} \sum_{q,r=1}^d \sum_{j=1}^n h_{qr} \partial_j f^q \partial_j f^r(x) dx$$

! measures the smoothness of f wrt h

[14] Kimmel et al. (1997); Perona and Malik (1990)

Harmonic map flow in continuous space

$f : \mathbb{R}^n \rightarrow (\mathbb{R}^d; \mathbf{h})$ smooth with \mathbf{h} a constant metric. The Dirichlet energy of f is

$$E(f; \mathbf{h}) = \frac{1}{2} \int_{\mathbb{R}^n} \| \nabla f \|^2_{\mathbf{h}} dx = \frac{1}{2} \int_{\mathbb{R}^n} \sum_{q,r=1}^d \sum_{j=1}^n h_{qr} \partial_j f^q \partial_j f^r(x) dx$$

! measures the smoothness of f wrt \mathbf{h}

Eells and Sampson (1964) studied gradient flow of E given by $f_{-t} = \text{grad}_{-t} E(f(t))$ to find minimizers of E ! extended to manifolds **harmonic map flow**

For PDE-based image processing, gradient flows of E recover the Perona-Malik diffusion^[14]

^[14] Kimmel et al. (1997); Perona and Malik (1990)

Extending the formalism to graphs

$f : \mathbb{R}^n \rightarrow \mathbb{R}$ smooth with h a constant metric The Dirichlet energy of f is

$$E(f; h) = \frac{1}{2} \int_{\mathbb{R}^n} \| \nabla f \|_h^2 dx = \frac{1}{2} \int_{\mathbb{R}^n} \sum_{q,r=1}^d h_{qr} \frac{\partial f}{\partial x^q} \frac{\partial f}{\partial x^r} (x) dx$$

Extending the formalism to graphs

$f : \mathbb{R}^n \rightarrow \mathbb{R}^d$; \mathbf{h} smooth with \mathbf{h} a constant metric. The Dirichlet energy of f is

$$E(f; \mathbf{h}) = \frac{1}{2} \int_{\mathbb{R}^n} \|\nabla f\|_{\mathbf{h}}^2 dx = \frac{1}{2} \sum_{q,r=1}^d \sum_{j=1}^n \int_{\mathbb{R}^n} h_{qr} \partial_j f^q \partial_j f^r(x) dx$$

! Replace $\int_{\mathbb{R}^n}$ with $\sum_{i,j \in V}$ and ∂_j with $\mathbf{r}_{(i,j)}$:

$$E_W^{\text{Dir}}(F) := \frac{1}{4} \sum_{q,r=1}^d \sum_{j \in E} h_{qr} (\mathbf{r}_{(i,j)} F^q)_{ij} (\mathbf{r}_{(i,j)} F^r)_{ij} = \frac{1}{4} \sum_{(i,j) \in E} \mathbf{W}_{(i,j)} (\mathbf{r}_{(i,j)} F)_{ij}^2$$

with $\mathbf{H} = \mathbf{W}^T \mathbf{W}$ with $\mathbf{W} \in \mathbb{R}^{d \times d}$

Extending the formalism to graphs

$f : \mathbb{R}^n \rightarrow \mathbb{R}^d$; \mathbf{h} smooth with \mathbf{h} a constant metric. The Dirichlet energy of f is

$$E(f; \mathbf{h}) = \frac{1}{2} \int_{\mathbb{R}^n} \|\nabla f\|_{\mathbf{h}}^2 dx = \frac{1}{2} \sum_{q,r=1}^d \sum_{j=1}^n \int_{\mathbb{R}^n} h_{qr} \frac{\partial f^q}{\partial x_j} \frac{\partial f^r}{\partial x_j}(x) dx$$

! Replace $\int_{\mathbb{R}^n}$ with $\sum_{i,j \in V}$ and $\frac{\partial}{\partial x_j}$ with $\mathbf{r}_{(i,j)}$:

$$E_W^{\text{Dir}}(F) := \frac{1}{4} \sum_{q,r=1}^d \sum_{i,j \in V} h_{qr} (\mathbf{r}_{(i,j)} F^q)_{ij} (\mathbf{r}_{(i,j)} F^r)_{ij} = \frac{1}{4} \sum_{(i,j) \in E} \mathbf{J}^T W(\mathbf{r}_{(i,j)} F)_{ij} \mathbf{J}^2:$$

with $\mathbf{H} = \mathbf{W}^T \mathbf{W}$ with $\mathbf{W} \in \mathbb{R}^{d \times d}$

If we minimize E_W^{Dir} we expect $\mathbf{J}(\mathbf{r}_{(i,j)} F)_{ij} \mathbf{J}$ to shrink 'except' when inside $\ker(\mathbf{H})$

Generalized harmonic flow on graphs is smoothing

We treat W as learnable weights and study the gradient flow of E_W^{Dir} :

$$\dot{F}(t) = -r_F E_W^{\text{Dir}}(F(t)) = -F(t)W^>W$$

Generalized harmonic flow on graphs is smoothing

We treat W as learnable weights and study the gradient flow of E_W^{Dir} :

$$\dot{F}(t) = -\text{r}_F E_W^{\text{Dir}}(F(t)) = -F(t)W^>W$$

Proposition (Di G., Rowbottom, et al.)

The dynamics is smoothing. Let P_W^{ker} be the projection onto $\ker(W^>W)$, then

$$\|E^{\text{Dir}}(F(t))\| \leq e^{-2t \text{gap}(W^>W) \text{gap}(V)} \|F(0)\|^2 + E^{\text{Dir}}(\|(P_W^{\text{ker}} - I_n) \text{vec}(F(0))\|); \quad t \geq 0:$$

$\forall \epsilon > 0 \exists \delta > 0$: for each $i \in V$ we have $f_i(t) \leq \delta + P_W^{\text{ker}} f_i(0)$.

Generalized harmonic flow on graphs is smoothing

We treat W as learnable weights and study the gradient flow of E_W^{Dir} :

$$\dot{F}(t) = -\text{r}_F E_W^{\text{Dir}}(F(t)) = -F(t)W^>W$$

Proposition (Di G., Rowbottom, et al.)

The dynamics is smoothing. Let P_W^{ker} be the projection onto $\ker(W^>W)$, then

$$\|E^{\text{Dir}}(F(t))\| \leq e^{-2t \text{gap}(W^>W) \text{gap}(F)} \|F(0)\|^2 + E^{\text{Dir}}((P_W^{\text{ker}} - I_n) \text{vec}(F(0))); \quad t \geq 0:$$

$\forall i \in \{1, 2, \dots, d\}$: for each $i \in V$ we have $f_i(t) \leq \frac{1}{d_i} + P_W^{\text{ker}} f_i(0)$.

Generalized harmonic flow on graphs is smoothing

We treat W as learnable weights and study the gradient flow of E_W^{Dir} :

$$\dot{F}(t) = -\text{r}_F E_W^{\text{Dir}}(F(t)) = -F(t)W^{\top}W$$

Proposition (Di G., Rowbottom, et al.)

The dynamics is smoothing. Let P_W^{ker} be the projection onto $\ker(W^{\top}W)$, then

$$\|E^{\text{Dir}}(F(t))\| \leq e^{-2t \text{gap}(W^{\top}W) \text{gap}(V)} \|F(0)\|^2 + E^{\text{Dir}}(\|(P_W^{\text{ker}} - I_n) \text{vec}(F(0))\|); \quad t \geq 0:$$

$\forall i \in [1, 2] \subset \mathbb{R}^d$: for each $i \in V$ we have $f_i(t) \leq \frac{1}{d_i} + P_W^{\text{ker}} f_i(0)$.

Generalized harmonic flow on graphs is smoothing

We treat W as learnable weights and study the gradient flow of E_W^{Dir} :

$$\dot{F}(t) = -\text{r}_F E_W^{\text{Dir}}(F(t)) = -F(t)W^>W$$

Proposition (Di G., Rowbottom, et al.)

The dynamics is smoothing. Let P_W^{ker} be the projection onto $\ker(W^>W)$, then

$$\|E^{\text{Dir}}(F(t))\| \leq e^{-2t \text{gap}(W^>W) \text{gap}(V)} \|F(0)\|^2 + E^{\text{Dir}}(\|(P_W^{\text{ker}} - I_n)\text{vec}(F(0))\|); \quad t \geq 0:$$

$\forall i \in \{1, 2\} \subset \mathbb{R}^d$: for each $i \in V$ we have $f_i(t) \leq P_{d_i} + P_W^{\text{ker}} f_i(0)$.

Generalized harmonic flow on graphs is smoothing

We treat W as learnable weights and study the gradient flow of E_W^{Dir} :

$$\dot{F}(t) = -\text{r}_F E_W^{\text{Dir}}(F(t)) = -F(t)W^>W$$

Proposition (Di G. , Rowbottom , et al.)

The dynamics is smoothing. Let P_W^{ker} be the projection onto $\ker(W^>W)$, then

$$\|E^{\text{Dir}}(F(t))\| \leq e^{-2t \text{gap}(W^>W)} \|F(0)\|^2 + E^{\text{Dir}}((P_W^{\text{ker}} - I_n)\text{vec}(F(0))); \quad t \geq 0:$$

$\forall i \in \{1, 2, \dots, d\}$: for each $i \in V$ we have $f_i(t) \leq \frac{1}{d_i} + P_W^{\text{ker}} f_i(0)$.

- I No W separates the limit embeddings of nodes with same degree and input features

[15] Similar to [Nt and Maehara \(2019\)](#); [Oono and Suzuki \(2020\)](#)

[16] This is different from [Nt and Maehara \(2019\)](#); [Oono and Suzuki \(2020\)](#); [Cai and Wang \(2020\)](#)

A few comments on the graph harmonic flow

- | No W separates the limit embeddings of nodes with same degree and input features
- | If W has zero kernel, nodes with same degrees converge to the same representation and over-smoothing occurs^[15]

[15] Similar to [Nt and Maehara \(2019\)](#); [Oono and Suzuki \(2020\)](#)

[16] This is different from [Nt and Maehara \(2019\)](#); [Oono and Suzuki \(2020\)](#); [Cai and Wang \(2020\)](#)

A few comments on the graph harmonic flow

- | No W separates the limit embeddings of nodes with same degree and input features
- | If W has zero kernel, nodes with same degrees converge to the same representation and over-smoothing occurs^[15]
- | Over-smoothing occurs independently of the spectral radius if its eigenvalues are positive— even for equations which lead to residual MPNNs when discretized^[16]

[15] Similar to [Nt and Maehara \(2019\)](#); [Oono and Suzuki \(2020\)](#)

[16] This is different from [Nt and Maehara \(2019\)](#); [Oono and Suzuki \(2020\)](#); [Cai and Wang \(2020\)](#)

GNNs as Gradient Flows part 2: multi-particle energy approach

A more general energy

We can rewrite $E_W^{\text{Dir}}(F) = \frac{1}{2} \sum_i \langle \mathbf{f}_i; W \rangle \langle W \mathbf{f}_i \rangle + \frac{1}{2} \sum_{i,j} a_{ij} \langle \mathbf{f}_i; W \rangle \langle W \mathbf{f}_j \rangle$

Replace $W \succ W$ with **symmetric** matrices $W \in \mathbb{R}^{d \times d}$!

$$E^{\text{tot}}(F) := \frac{1}{2} \sum_i \langle \mathbf{f}_i; \mathbf{f}_i \rangle + \frac{1}{2} \sum_{i,j} a_{ij} \langle \mathbf{f}_i; W \mathbf{f}_j \rangle = E^{\text{ext}}(F) + E_W^{\text{pair}}(F)$$

A more general energy

We can rewrite $E_W^{\text{Dir}}(F) = \frac{1}{2} \sum_i \langle f_i; W \rangle \langle f_i; W \rangle + \frac{1}{2} \sum_{i,j} a_{ij} \langle f_i; W \rangle \langle f_j; W \rangle$

Replace $W \succ W$ with **symmetric** matrices $W \in \mathbb{R}^{d \times d}$!

$$E^{\text{tot}}(F) := \frac{1}{2} \sum_i \langle f_i; f_i \rangle + \frac{1}{2} \sum_{i,j} a_{ij} \langle f_i; W f_j \rangle = E^{\text{ext}}(F) + E_W^{\text{pair}}(F)$$

The gradient flow of E^{tot} is

$$\dot{F}(t) = -\nabla_F E^{\text{tot}}(F(t)) = -F(t) + AF(t)W:$$

Node-features $\$$ particles in \mathbb{R}^d with energy E^{tot}

- | E^{ext} is independent of the graph topology external field
- | E_W^{pair} potential energy, with W defining pairwise interactions of adjacent nodes

Node-features \mathbf{x} particles in \mathbb{R}^d with energy E^{tot}

| E^{ext} is independent of the graph topology external field

| E_W^{pair} potential energy, with W defining pairwise interactions of adjacent nodes

Decompose $W = W^+ + W^-$ into positive and negative eigenvalues

Attraction vs repulsion

$$W = \sum_i \dots + \sum_j \dots$$

$$E^{\text{tot}}(F) = \frac{1}{2} \sum_i f_i^2 + \frac{1}{4} \sum_{ij} (r_{ij} - F)_{ij}^2 + \frac{1}{4} \sum_{ij} (r_{ij} + F)_{ij}^2$$

Attraction vs repulsion

$$W = \sum_{i,j} w_{ij} > 0$$

$$E^{\text{tot}}(F) = \frac{1}{2} \sum_i |f_i|^2 + \sum_{i,j} w_{ij} (f_i - f_j)^2 = \frac{1}{2} \sum_i |f_i|^2 + \sum_{i,j} w_{ij} (f_i^2 - 2f_i f_j + f_j^2)$$

The gradient flow minimizes E^{tot} ! W encodes..

- | **attraction** via its **positive eigenvalues** since $\sum_{i,j} w_{ij} (f_i - f_j)^2$ decreases edge-wise
- | **repulsion** via its **negative eigenvalues** since $\sum_{i,j} w_{ij} (f_i - f_j)^2$ increases edge-wise

Consider $F(t) = AF(t)W$ \Rightarrow $\text{vec}(F(t)) = (W \quad A)\text{vec}(F(t))$

Spectrum of W induces LFD or HFD

Consider $F(t) = AF(t)W$ $\vec{\text{vec}}(F(t)) = (W \quad A)\vec{\text{vec}}(F(t))$

Write the spectrum of W as λ_r^W with $\lambda_+^W = (\max \lambda_r^W)_+$ and $\lambda_-^W = (\min \lambda_r^W)$

Spectrum of W induces LFD or HFD

Consider $F(t) = AF(t)W$ $\vec{\text{vec}}(F(t)) = (W \quad A)\vec{\text{vec}}(F(t))$

Write the spectrum of W as λ_i^W with $\lambda_+^W = (\max \lambda_i^W)_+$ and $\lambda_-^W = (\min \lambda_i^W)_-$

Any eigenvalue of $W \quad A$ can be written as $\lambda_i^W \quad \lambda_j^A = \lambda_i^W (1 \quad \lambda_j^A)$

Spectrum of W induces LFD or HFD

Consider $F(t) = AF(t)W$ $\vec{\text{vec}}(F(t)) = (W \ A)\vec{\text{vec}}(F(t))$

Write the spectrum of W as λ_i^W with $\lambda_+^W = (\max \lambda_i^W)_+$ and $\lambda_-^W = (\min \lambda_i^W)_-$

Any eigenvalue of $W \ A$ can be written as $\lambda_i^W \lambda_j^A = \lambda_i^W (1 - \lambda_j^W)$

Let P_W be the projection onto the eigenspace of $W \ A$ associated with $\lambda_j^W := \lambda_j^W (1 - \lambda_j^W)$. Recall that λ_+^W is the largest eigenvalue of $W = I - A$

Proposition (Di G. , Rowbottom , et al.)

If $\alpha > \frac{W}{+}$, then $F(t) = AF(t)W$ is HFD for a.e. $F(0)$: there exists ϵ_{HFD} such that ^[17]

$$E^{\text{Dir}}(F(t)) = e^{2t} \left(\frac{1}{2} \|P_W F(0)\|^2 + O(e^{-2t \epsilon_{\text{HFD}}}) \right); \quad t \geq 0;$$

and $\|F(t)\|$ converges to $\|F_1\| \in \mathbb{R}^n$ such that $f_1^r = f_1^r$, for $1 \leq r \leq d$.

^[17] We have an explicit formula depending on 'spectral gaps' of W

Proposition (Di G. , Rowbottom , et al.)

If $\lambda > \frac{W}{+}$, then $F(t) = A F(t) W$ is HFD for a.e. $F(0)$: there exists ϵ_{HFD} such that ^[18]

$$E^{\text{Dir}}(F(t)) = e^{2t} \left(\frac{1}{2} \|P_W F(0)\|^2 + O(e^{-2t \epsilon_{\text{HFD}}}) \right); \quad t \geq 0;$$

and $\|F(t)\|$ converges to $\|F_1\| \in \mathbb{R}^n$ such that $f_1^r = f_1^r$, for $1 \leq r \leq d$.

^[18] We have an explicit formula depending on 'spectral gaps' of W

Proposition (Di G. , Rowbottom , et al.)

If $\lambda > \frac{W}{+}$, then $F(t) = AF(t)W$ is HFD for a.e. $F(0)$: there exists ϵ_{HFD} such that ^[19]

$$E^{\text{Dir}}(F(t)) = e^{2t} \left(\frac{1}{2} \|P_W F(0)\|^2 + O(e^{-2t \epsilon_{\text{HFD}}}) \right); \quad t \geq 0;$$

and $\|F(t)\|$ converges to $\|F_1\| \in \mathbb{R}^n$ such that $f_1^r = f_1^r$, for $1 \leq r \leq d$.

^[19] We have an explicit formula depending on 'spectral gaps' of W

Proposition (Di G. , Rowbottom , et al.)

If $\lambda > \frac{W}{+}$, then $F(t) = AF(t)W$ is HFD for a.e. $F(0)$: there exists ϵ_{HFD} such that ^[20]

$$E^{\text{Dir}}(F(t)) = e^{2t} \left(\frac{1}{2} \|P_W F(0)\|^2 + O(e^{-2t \epsilon_{\text{HFD}}}) \right); \quad t \geq 0;$$

and $\|F(t)\|$ converges to $\|F_1\| \in \mathbb{R}^n$ such that $f_1^r = f_1^r$, for $1 \leq r \leq d$.

^[20] We have an explicit formula depending on 'spectral gaps' of W

Proposition (Di G. , Rowbottom , et al.)

If $\lambda_{\min}(W) > \frac{1}{2}$, then $F(t) = A F(0) W^t$ is HFD for a.e. $F(0)$: there exists α_{HFD} such that ^[21]

$$E^{\text{Dir}}(F(t)) = e^{2t} \left(\frac{1}{2} \|P_W F(0)\|^2 + O(e^{-2t \alpha_{\text{HFD}}}) \right); \quad t \geq 0;$$

and $\|F(t)\|$ converges to $\|F_1\|$ such that $f_1^r = f_1^r$, for $1 \leq r \leq d$.

^[21] We have an explicit formula depending on 'spectral gaps' of W

Proposition (Di G. , Rowbottom , et al.)

If $\alpha > \frac{W}{+}$, then $F(t) = AF(t)W$ is HFD for a.e. $F(0)$: there exists ϵ_{HFD} such that ^[22]

$$E^{\text{Dir}}(F(t)) = e^{2t} \left(\frac{1}{2} \|P_W F(0)\|^2 + O(e^{-2t \epsilon_{\text{HFD}}}) \right); \quad t \geq 0;$$

and $\|F(t)\|$ converges to $\|F_1\| \in \mathbb{R}^n$ such that $f_1^r = \|f_1^r\|$, for $1 \leq r \leq d$.

^[22] We have an explicit formula depending on 'spectral gaps' of W

Proposition (Di Gi., Rowbottom, et al.)

If $\mu > \frac{W}{2}$, then $F(t) = A F(0) W^t$ is HFD for a.e. $F(0)$: there exists ϵ_{HFD} such that ^[22]

$$E^{\text{Dir}}(F(t)) = e^{2t} \left(\frac{1}{2} \|P_W F(0)\|^2 + O(e^{-2t \epsilon_{\text{HFD}}}) \right); \quad t \geq 0;$$

and $\|F(t)\|$ converges to $\|F_1\|$ such that $\|f_r\| = \|f_1\|$, for $1 \leq r \leq d$.

If enough mass is distributed over the negative eigenvalues of the 'channel-mixing', graph high frequencies dominate what matters is how the spectra of A and W interact

^[22] We have an explicit formula depending on 'spectral gaps' of A and W

Source term and a more general family of energies

Equations with a source term may have better expressive power ^[23]

In our framework: add an extra energy term $E^{\text{source}}(F) := \langle F; F(0)W \rangle$!

$$E(t) = F(t) + \langle AF(t)W \rangle - \langle F(0)W \rangle$$

^[23] Xhonneux et al. (2020); Chen et al. (2020); Thorpe et al. (2021)

Source term and a more general family of energies

Equations with a source term may have better expressive power ^[23]

In our framework: add an extra energy term $E_W^{\text{source}}(F) := \langle F; F(0)W \rangle!$

$$E(t) = F(t) + A F(t)W - F(0)W :$$

We can also replace A with A satisfying $A_{ij} = 0$ if $(i; j) \not\subseteq E!$

$$E_{A;W}^{\text{pair}}(F) := \sum_{(i;j)} A_{ij} \langle f_i; W f_j \rangle :$$

^[23] Xhonneux et al. (2020); Chen et al. (2020); Thorpe et al. (2021)

Non-linear function can 'activate' the inner products in the energy:

$$E^{\text{ext}}(F) + E_W^{\text{pair}}(F) = \frac{1}{2} \sum_i^P (\mathbf{h}_i; \mathbf{f}_i) + \frac{1}{2} \sum_{i,j}^P a_{ij} (\mathbf{h}_i; \mathbf{W} \mathbf{f}_j):$$

^[24] Wu et al. (2019); Oono and Suzuki (2020); Chen et al. (2020)

Non-linear function can 'activate' the inner products in the energy:

$$E^{\text{ext}}(F) + E_W^{\text{pair}}(F) = \frac{1}{2} \sum_i (h f_i; f_i) + \frac{1}{2} \sum_{i,j} a_{ij} (h f_i; W f_j):$$

A few reasons why we keep the gradient ~~linear~~ ^{linear}

- | Perform spectral analysis in closed form ^[24]
- | We have seen no gain in performance when including non-linear activations
- | We can 'push the non-linear maps' in either the encoding block or the decoding one

[24] Wu et al. (2019); Oono and Suzuki (2020); Chen et al. (2020)

A comparison with (some) continuous GNN models

Recall the continuous models:

| LinearPDE GCN_D: $F_{\text{PDE GCN}_D}(t) = F(t)K(t)^T K(t)$

| CGNN: $F_{\text{CGNN}}(t) = F(t) + F(t) \tilde{\sim} + F(0)$ with symmetric $\tilde{\sim}$

| LinearGRAND: $F_{\text{GRAND}}(t) = \text{RW } F(t) = (I - A(F(0)))F(t)$

A comparison with (some) continuous GNN models

Recall the continuous models:

| LinearPDE GCN_D: $F_{\text{PDE GCN}_D}(t) = F(t)K(t)^{\top}K(t)$

| CGNN: $F_{\text{CGNN}}(t) = F(t) + F(t)\tilde{\sim} + F(0)$ with symmetric $\tilde{\sim}$

| LinearGRAND: $F_{\text{GRAND}}(t) = \text{RW } F(t) = (I - A(F(0)))F(t)$

Proposition (Di G. , Rowbottom , et al.)

(i) PDE GCN_D is a smoothing mode $\mathbf{E}^{\text{Dir}}(F_{\text{PDE GCN}_D}(t)) = 0$.

A comparison with (some) continuous GNN models

Recall the continuous models:

I LinearPDE GCN_D: $F_{\text{PDE GCN}_D}(t) = F(t)K(t)^T K(t)$

I CGNN: $F_{\text{CGNN}}(t) = F(t) + F(t) \tilde{\sim} + F(0)$ with symmetric $\tilde{\sim}$

I LinearGRAND: $F_{\text{GRAND}}(t) = \text{RW } F(t) = (I - A(F(0)))F(t)$

Proposition (Di G. , Rowbottom , et al.)

- (i) PDE GCN_D is a smoothing mode $E^{\text{Dir}}(F_{\text{PDE GCN}_D}(t)) = 0$.
- (ii) For a.e. $F(0)$ it holds: CGNN is never HFD and if we remove the source term, then $E^{\text{Dir}}(F_{\text{CGNN}}(t) - F_{\text{CGNN}}(t)) = e^{-\text{gap}(t)}$.

A comparison with (some) continuous GNN models

Recall the continuous models:

I LinearPDE GCN_D: $F_{\text{PDE GCN}_D}(t) = F(t)K(t)^T K(t)$

I CGNN: $F_{\text{CGNN}}(t) = F(t) + F(t) \sim + F(0)$ with symmetric \sim

I LinearGRAND: $F_{\text{GRAND}}(t) = \text{RW } F(t) = (I - A(F(0)))F(t)$

Proposition (Di G. , Rowbottom , et al.)

- (i) PDE GCN_D is a smoothing mode $E^{\text{Dir}}(F_{\text{PDE GCN}_D}(t)) = 0$.
- (ii) For a.e. $F(0)$ it holds: CGNN is never HFD and if we remove the source term, then $E^{\text{Dir}}(F_{\text{CGNN}}(t) - F_{\text{CGNN}}(t)) = e^{-\text{gap}(t)}$.
- (iii) If G is connected $F_{\text{GRAND}}(t) \rightarrow \text{ast} \rightarrow 1$, with $r = \text{mean}(f^r(0))$, $1 \leq r \leq d$.

GNNs as Gradient Flows part 3: discrete setting

The requirement for symmetry

When classical MPNNs are discretized gradient flows?

The requirement for symmetry

When classical MPNNs are discretized gradient flows?

Let $A \in \mathbb{R}^{n \times n}$ be a symmetric graph vector field $(A)_{ij} = 0; (i; j) \notin E$

Consider a family of linear GNNs with shared weights of the form

$$F(t+1) = F(t) + A F(t)W + F(0)W; \quad 0 \leq t \leq T:$$

They are gradient flow of a 'multi-particle' energy iff A and W are symmetric.

Can graph convolutional models be high-frequency dominated?

Introduce step-size τ and consider gradient flow system

$$F(t + \tau) = F(t) + \tau AF(t)W; \quad W = W^T;$$

Let P_W be the projection into the eigenspace of $A = W(I - \tau A)$ associated with the eigenvalue $\lambda := j^W j(1)$ and set

$$\frac{W}{+} (1) \tau^{-1} < j^W j < 2(2) \tau^{-1} \quad (1)$$

Can graph convolutional models be high-frequency dominated?

Introduce step-size τ and consider gradient flow system

$$F(t + \tau) = F(t) + \tau AF(t)W; \quad W = W^T;$$

Let P_W be the projection into the eigenspace of $A = W(I - \tau A)$ associated with the eigenvalue $\lambda := j^W j(1 - \tau \lambda)$ and set

$$\frac{W}{+} (1) < j^W j < 2(2) < 1 \tag{2}$$

Can graph convolutional models be high-frequency dominated?

Introduce step-size τ and consider gradient flow system

$$F(t + \tau) = F(t) + \tau AF(t)W; \quad W = W^T;$$

Let P_W be the projection into the eigenspace of $A = W(I - \tau A)$ associated with the eigenvalue $\lambda_j := \frac{1}{1 + \tau \lambda_j}$ and set

$$\frac{1}{1 + \tau \lambda_j} < \lambda_j < \frac{1}{1 - \tau \lambda_j} \quad (3)$$

Can graph convolutional models be high-frequency dominated?

Theorem (Di G. , Rowbottom , et al.)

If equation 3 holds then there exists $\epsilon_{\text{HFD}} < 1$ s.t.

$$E^{\text{Dir}}(\|F(m)\|) = (1 + \epsilon_{\text{HFD}})^{2m} \frac{1}{2} \|P_W F(0)\|^2 + O\left(\frac{1 + \epsilon_{\text{HFD}}}{1 + \epsilon_{\text{HFD}}^{2m}}\right) :$$

The dynamics is HFD for a.e. $F(0)$ and $\|F(m)\| \approx \|F(0)\|$ s.t. $f_1^r = f_1^l$.

Can graph convolutional models be high-frequency dominated?

Theorem (Di G. , Rowbottom , et al.)

If equation 3 holds then there exists $\epsilon_{\text{HFD}} < \epsilon$ s.t.

$$E^{\text{Dir}}(F(m)) = (1 + \epsilon)^{2m} \frac{1}{2} \|P_W F(0)\|^2 + O\left(\frac{1 + \epsilon_{\text{HFD}}}{1 + \epsilon}\right)^{2m} :$$

The dynamics is HFD for a.e. $F(0)$ and $F(m) = \|F(m)\|_1$ s.t. $f_1^r = f_1^l$.

Conversely, if G is not bipartite, then for a.e. $F(0)$ the system $F(t+1) = AF(t)W$, with W symmetric, is LFD independent of the spectrum λ .

- ! linear discrete gradient flows can be HFD due to the negative eigenvalues
- ! Differently from previous results^[25], no bound on spectral radius coming from the graph topology as long as λ is small enough
 - ! Recall that previous over-smoothing results required to have sufficiently small singular values depending on the spectrum of
 - ! If we have symmetry and control the spectrum we can avoid over-smoothing (and in fact be HFD) in terms of positive vs negative eigenvalues

^[25] Nt and Maehara (2019); Oono and Suzuki (2020); Cai and Wang (2020)

- I Without a residual term the dynamics $\mathbf{L} \mathbf{E} \mathbf{D}$ for a.e. $F(0)$ independently of the sign and magnitude of the eigenvalues $\mathbf{W} \mathbf{f}$
 - ! provides a justification for the residual connection in terms of the spectrum $\mathbf{W} \mathbf{f}$ of
 - ! explains via induced dynamics and spectral analysis the 'expressivity' results in [Chen et al. \(2020\)](#)

Reversing time and sign of the edge weights

Let f_r^W be the spectrum of W with orthonormal eigenvectors f_r^W and $U = U^{-1}$

[26] Similar effect as in [Bo et al. \(2021\)](#); [Yan et al. \(2021\)](#)

Reversing time and sign of the edge weights

Let λ_r^W be the spectrum of A^W with orthonormal eigenvectors g_r^W and $U = [g_1^W \dots g_n^W]^T$

Introduce $z^r(t) : V \rightarrow \mathbb{R}$ defined by $z_i^r(t) = \mathbf{h}_i^r(t)$; $z^r = U^T z$, then gradient flow becomes:

$$z^r(t + \Delta t) = U (I + \Delta t \lambda_r^W (I - \text{diag}(z^r(t))) U^T z^r(t) = z^r(t) + \Delta t \lambda_r^W A z^r(t)$$

Along λ_r^W if $\lambda_r^W < 0$ then the dynamics is equivalent to flipping the sign of the edges ^[26]

[26] Similar effect as in [Bo et al. \(2021\)](#); [Yan et al. \(2021\)](#)

GNNs as Gradient Flows part 4: ablation studies and experiments

General ingredients of the framework GRAFF (Gradient Flow Framework)

- | Encoding block $E_N : \mathbb{R}^{n \times p} \rightarrow \mathbb{R}^{n \times d}$ is used to process input features $F_0 \in \mathbb{R}^{n \times p}$

General ingredients of the framework GRAFF (Gradient Flow Framework)

- | Encoding block $E_N : \mathbb{R}^{n \times p} \rightarrow \mathbb{R}^{n \times d}$ is used to process input features $F_0 \in \mathbb{R}^{n \times p}$
- | Symmetric channel-mixing matrices $W \in \mathbb{R}^{d \times d}$ that are shared across the layers

General ingredients of the framework GRAFF (Gradient Flow Framework)

- | Encoding block $E_N : \mathbb{R}^{n \times p} \rightarrow \mathbb{R}^{n \times d}$ is used to process input features $F \in \mathbb{R}^{n \times p}$
- | Symmetric channel-mixing matrices $W \in \mathbb{R}^{d \times d}$ that are shared across the layers
- | Decoding block $D_E : \mathbb{R}^{n \times d} \rightarrow \mathbb{R}^{n \times k}$, where k is the number of label classes

General ingredients of the framework GRAFF (Gradient Flow Framework)

- | Encoding block $E_N : \mathbb{R}^n \times \mathbb{P} \rightarrow \mathbb{R}^n \times \mathbb{D}$ is used to process input features $F_0 \in \mathbb{R}^n \times \mathbb{P}$
- | Symmetric channel-mixing matrices $W \in \mathbb{R}^{\mathbb{D} \times \mathbb{D}}$ that are shared across the layers
- | Decoding block $D_E : \mathbb{R}^n \times \mathbb{D} \rightarrow \mathbb{R}^n \times \mathbb{K}$, where k is the number of label classes

$$F(t + \Delta t) = F(t) + \Delta t \left(-\nabla_{F(t)} \mathcal{L}(F(t)) + A F(t) W + F(0) \right); \quad F(0) = E_N(F_0);$$

I Sumvariant: $W = W^0 + W^{\text{op}}$! 'no-control' on spectrum

[27] Provides justification to [Chen et al. \(2020\)](#)

- | Sumvariant: $W = W^0 + W^{\ominus}$! 'no-control' on spectrum
- | (Neg)-Prod $W = W^0 > W^{\ominus}$! signed eigenvalues

[27] Provides justification to [Chen et al. \(2020\)](#)

- | Sumvariant: $W = W^0 + W^0$! 'no-control' on spectrum
- | (Neg)-Prod $W = W^0 > W^0$! signed eigenvalues
- | W diagonally-dominant (DD): take W^0 symmetric with zero diagonal and $w \in \mathbb{R}^d$ defined by $w_j = \rho_j \sum_{i \neq j} |W^0_{ij}| + r_j$, and set $W = \text{diag}(w) + W^0$! by Gershgorin Theorem the model 'can' easily re-distribute mass in the spectrum [via](#) ^[27].

^[27] Provides justification to [Chen et al. \(2020\)](#)

Complexity and number of parameters

GRAFF scales as $\mathcal{O}(Vpd + Ed)$, where p and d are input feature and hidden dimension

! our model is faster than GCN with small number of parameters $pd + d^2 + 3d + dk$

Recall our claims about role of `channel-mixing`:

- | Positive eigenvalues of W induce attraction in a residual convolutional model

Recall our claims about role of `channel-mixing`:

- | Positive eigenvalues σ induce attraction in a residual convolutional model
- | Negative eigenvalues σ induce repulsion in a residual convolutional model

Recall our claims about role of λ channel-mixing:

- | Positive eigenvalues λ induce attraction in a residual convolutional model
- | Negative eigenvalues λ induce repulsion in a residual convolutional model
- | A non-residual convolutional model is always dominated by low-frequencies independent of the spectrum of λ

Recall our claims about role of λ channel-mixing:

- | Positive eigenvalues λ induce attraction in a residual convolutional model
- | Negative eigenvalues λ induce repulsion in a residual convolutional model
- | A non-residual convolutional model is always dominated by low-frequencies independent of the spectrum of \mathbf{W}

To investigate our claims we use the synthetic Cora dataset of [Zhu et al. \(2020\)](#)

! graphs are generated for target levels of homophily via preferential attachment: we expect LFD to be better than HFD with high homophily and vice-versa for low homophily

Goal: Explain performance wrt homophily in terms of the spectrum of M_f

- | **Neg-prod** is better than **prod** on low-homophily!
con rms HFD dynamics

Goal: Explain performance wrt homophily in terms of the spectrum of \mathbf{M}

- | **Neg-prod** is better than **prod** on low-homophily!
con rms HFD dynamics
- | **prod** (attraction-only) struggles in low-homophily
even with residual connection

Goal: Explain performance wrt homophily in terms of the spectrum λ

- | **Neg-prod** is better than **prod** on low-homophily!
con rms HFD dynamics
- | **prod** (attraction-only) struggles in low-homophily
even with residual connection
- | `neutral' variants like **sum** and **(DD)** are more
exible and outperform **GCN** con rming that
non- residual convolutional models **at** λ
irrespectively of the spectrum λ

Goal: Use homophily to assess if the evolution is smoothing compute homophily of the prediction (cross) and compare with that read from the encoding (i.e. evolution)

Goal: Use homophily to assess if the evolution is ~~smoothing~~ compute homophily of the prediction (cross) and compare with that read from the encoding (i.e. evolution)

- | **neg-prod** homophily decreases after evolution while with **prod** the prediction is smoother than the true homophily

Goal: Use homophily to assess if the evolution is ~~smoothing~~ compute homophily of the prediction (cross) and compare with that read from the encoding (i.e. evolution)

- | **neg-prod** homophily decreases after evolution while with **prod** the prediction is smoother than the true homophily
- | **(DD)** and **sum** variants adapt better to the true homophily

Goal: Use homophily to assess if the evolution is smoothing compute homophily of the prediction (cross) and compare with that read from the encoding (i.e. evolution)

- | **neg-prod** homophily decreases after evolution while with **prod** the prediction is smoother than the true homophily
- | **(DD)** and **sum** variants adapt better to the true homophily
- | The encoding compensates when the spectrum of V has a sign

Conclusions and where to next?

What was the message then?

- | Framework where the MPNNs equations minimize a multi-particle learnable energy

What was the message then?

- | Framework where the MPNNs equations minimize a multi-particle learnable energy
- | Analysis of the interaction between the spectrum of the graph and the spectrum of the Laplacian
`channel-mixing!` when and why the dynamics is low (high) frequency dominated

What was the message then?

- | Framework where the MPNNs equations minimize a multi-particle learnable energy
- | Analysis of the interaction between the spectrum of the graph and the spectrum of the channel-mixing! when and why the dynamics is low (high) frequency dominated
- | Refined existing asymptotic analysis of MPNNs to account for the role of the spectrum of the channel-mixing

What was the message then?

- | Framework where the MPNNs equations minimize a multi-particle learnable energy
- | Analysis of the interaction between the spectrum of the graph and the spectrum of the channel-mixing! when and why the dynamics is low (high) frequency dominated
- | Refined existing asymptotic analysis of MPNNs to account for the role of the spectrum of the channel-mixing
- | From a practical perspective, our framework allows for 'educated' choices resulting in a simple, more explainable convolutional model: our results refute the folklore of graph convolutional models being too 'simple' for complex benchmarks.

We restricted to a constant bilinear form W , how about non-constant alternatives $W(F; t)$ that are aware of the features? requirement for local 'heterogeneity' with efficiency

We restricted to a constant bilinear form W , how about non-constant alternatives $W(F; t)$ that are aware of the features? requirement for local 'heterogeneity' with efficiency

What can we say about dynamics that are neither LFD nor HFD?

We restricted to a constant bilinear form W , how about non-constant alternatives $W(F; t)$ that are aware of the features? requirement for local 'heterogeneity' with efficiency

What can we say about dynamics that are neither LFD nor HFD?

The energy formulation points to new models more 'physics' inspired

Thank you!

For any question/complaint/video-game recommendation do not hesitate to contact me! :

fdigiovanni (at) twitter com

[@Francesco_dgv](#)

References

- Bo, D., Wang, X., Shi, C., and Shen, H. (2021). Beyond low-frequency information in graph convolutional networks. In AAAI. AAAI Press
- Bronstein, M. M., Bruna, J., LeCun, Y., Szlam, A., and Vandergheynst, P. (2017). Geometric deep learning: going beyond euclidean data. IEEE Signal Processing Magazine, 34(4):18–42.
- Cai, C. and Wang, Y. (2020). A note on over-smoothing for graph neural networks. preprint arXiv:2006.13318
- Chamberlain, B., Rowbottom, J., Gorinova, M. I., Bronstein, M., Webb, S., and Rossi, E. (2021). Grand: Graph neural diffusion. International Conference on Machine Learning pages 1407–1418. PMLR.
- Chen, M., Wei, Z., Huang, Z., Ding, B., and Li, Y. (2020). Simple and deep graph convolutional networks. International Conference on Machine Learning, pages 1725–1735. PMLR.

- Chen, R. T., Rubanova, Y., Bettencourt, J., and Duvenaud, D. K. (2018). Neural ordinary differential equations. *Advances in neural information processing systems*, 31.
- Eells, J. and Sampson, J. H. (1964). Harmonic mappings of riemannian manifolds. *American journal of mathematics*, 86(1):109–160.
- Eliasof, M., Haber, E., and Treister, E. (2021). Pde-gcn: Novel architectures for graph neural networks motivated by partial differential equations. *Advances in Neural Information Processing Systems*, 34.
- Haber, E. and Ruthotto, L. (2018). Stable architectures for deep neural networks. *arXiv preprint arXiv:1808.08759*, 34.
- Kimmel, R., Sochen, N., and Malladi, R. (1997). From high energy physics to low level vision. In *International Conference on Scale-Space Theories in Computer Vision*, pages 236–247. Springer.
- Kipf, T. N. and Welling, M. (2017). Semi-Supervised Classification with Graph

Convolutional Networks. In *Proceedings of the 5th International Conference on Learning Representations*, ICLR '17.

Nt, H. and Maehara, T. (2019). Revisiting graph neural networks: All we have is low-pass filters. *arXiv preprint arXiv:1905.09550*.

Oono, K. and Suzuki, T. (2020). Graph neural networks exponentially lose expressive power for node classification. In *International Conference on Learning Representations*.

Perona, P. and Malik, J. (1990). Scale-space and edge detection using anisotropic diffusion. *PAMI*, 12(7):629–639.

Rusch, T. K., Chamberlain, B. P., Rowbottom, J., Mishra, S., and Bronstein, M. M. (2022). Graph-coupled oscillator networks. In *International Conference on Machine Learning*.

Thorpe, M., Nguyen, T. M., Xia, H., Strohmer, T., Bertozzi, A., Osher, S., and Wang, B. (2021). Grand++: Graph neural diffusion with a source term. In *International Conference on Learning Representations*.

- Veličković, P., Cucurull, G., Casanova, A., Romero, A., Liò, P., and Bengio, Y. (2018). Graph attention networks. In *International Conference on Learning Representations*.
- Wu, F., Souza, A., Zhang, T., Fifty, C., Yu, T., and Weinberger, K. (2019). Simplifying graph convolutional networks. In *International conference on machine learning*, pages 6861–6871. PMLR.
- Xhonneux, L.-P., Qu, M., and Tang, J. (2020). Continuous graph neural networks. In *International Conference on Machine Learning*, pages 10432–10441. PMLR.
- Yan, Y., Hashemi, M., Swersky, K., Yang, Y., and Koutra, D. (2021). Two sides of the same coin: Heterophily and oversmoothing in graph convolutional neural networks. *arXiv preprint arXiv:2102.06462*.
- Zhou, D. and Schölkopf, B. (2005). Regularization on discrete spaces. In *Joint Pattern Recognition Symposium*, pages 361–368. Springer.
- Zhu, J., Yan, Y., Zhao, L., Heimann, M., Akoglu, L., and Koutra, D. (2020). Beyond

homophily in graph neural networks: Current limitations and effective designs. *Advances in Neural Information Processing Systems*, 33:7793–7804.