# Graph Neural Ordinary Differential Equations

<u>Michael Poli</u><sup>1\*</sup>, Stefano Massaroli<sup>2\*</sup>, Junyoung Park<sup>1\*</sup> Atsushi Yamashita<sup>2</sup>, Hajime Asama<sup>2</sup>, Jinkyoo Park<sup>1</sup>

<sup>1</sup>Korea Advanced Institute of Technology (KAIST), <sup>2</sup>The University of Tokyo, \*Equal contribution authors

The First International Workshop on Deep Learning on Graphs: Methodologies and Applications (DLGMA'20) 2020-02-08





Continuous-depth Learning

Standard DL Settings  $\mathbf{h}_{s+1} = \mathbf{h}_s + \mathbf{f}(\mathbf{h}_s, \boldsymbol{\theta}_s)$ residual layer

DL on Graphs  $\mathbf{H}_{s+1} = \mathbf{H}_{s} + \sigma(\tilde{\mathbf{D}}^{-\frac{1}{2}}\tilde{\mathbf{A}}\tilde{\mathbf{D}}^{-\frac{1}{2}}\mathbf{H}_{s}\boldsymbol{\Theta}_{s})$ graph convolution layer

Objective: develop the continuous-depth paradigm for deep learning.

 $\dot{\mathbf{h}}(s) = \mathbf{f}(s, \mathbf{h}(s), \boldsymbol{\theta})$   $\dot{\mathbf{H}}(s) = \mathbf{F}_{\mathcal{G}}(s, \mathbf{H}(s), \boldsymbol{\Theta})$ 

Continuous-depth Learning

Standard DL Settings  $\mathbf{h}_{s+1} = \mathbf{h}_s + \mathbf{f}(\mathbf{h}_s, \boldsymbol{\theta}_s)$ residual layer

DL on Graphs  $\mathbf{H}_{s+1} = \mathbf{H}_{s} + \sigma(\tilde{\mathbf{D}}^{-\frac{1}{2}}\tilde{\mathbf{A}}\tilde{\mathbf{D}}^{-\frac{1}{2}}\mathbf{H}_{s}\boldsymbol{\Theta}_{s})$ graph convolution layer

Objective: develop the continuous-depth paradigm for deep learning.

 $\dot{\mathbf{h}}(s) = \mathbf{f}(s, \mathbf{h}(s), \boldsymbol{\theta})$   $\dot{\mathbf{H}}(s) = \mathbf{F}_{\mathcal{G}}(s, \mathbf{H}(s), \boldsymbol{\Theta})$ 

Continuous-depth Learning

Standard DL Settings residual layer

DL on Graphs  $\mathbf{h}_{s+1} = \mathbf{h}_s + \mathbf{f}(\mathbf{h}_s, \boldsymbol{\theta}_s)$   $\mathbf{H}_{s+1} = \mathbf{H}_s + \sigma(\tilde{\mathbf{D}}^{-\frac{1}{2}}\tilde{\mathbf{A}}\tilde{\mathbf{D}}^{-\frac{1}{2}}\mathbf{H}_s\boldsymbol{\Theta}_s)$ graph convolution layer

Objective: develop the continuous-depth paradigm for deep learning.

Neural ODEs [R. T. Chen et al., 2018]

Graph Neural ODEs  $\dot{\mathbf{h}}(s) = \mathbf{f}(s, \mathbf{h}(s), \theta)$   $\dot{\mathbf{H}}(s) = \mathbf{F}_{\mathcal{C}}(s, \mathbf{H}(s), \Theta)$ proposed approach

Graph Neural Ordinary Differential Equations (GDEs) blend discrete topological structures and differential equations.

#### GDEs blend discrete topological structures and differential equations.

#### Advantages:

- Static settings: computational advantages by incorporation of numerical methods in the forward pass.
- Dynamic settings: exploitation of the geometry of the underlying dynamics and flexibility with respect to irregular observations.

#### Notation:

set of nodes	set of edges	graph	features
$\mathcal{V} ( \mathcal{V} = n )$	$\mathcal{E} \subset \{(u,v)\}_{u,v\in\mathcal{V}}$	$\mathcal{G} := (\mathcal{V}, \mathcal{E})$	$\mathbf{X} \in \mathbb{R}^{n  imes d}$

# Graph Neural ODEs (GDEs)

Graph Neural Networks

$$\begin{cases} \mathsf{H}_{s+1} = \mathsf{H}_{s} + \mathsf{F}_{\mathcal{G}}\left(s, \mathsf{H}_{s}, \Theta_{s}\right) \\ \mathsf{H}_{0} = \mathsf{X}_{e} \end{cases}, s \in \mathbb{N}$$

where **F** is a matrix-valued nonlinear function conditioned on graph  $\mathcal{G}$  and  $\Theta_s$  is the tensor of **trainable parameters** of the *s*-th layer.

Graph Neural ODEs (GDEs) [Proposed]

$$\begin{cases} \dot{\mathbf{H}}_{s} = \mathbf{F}_{\mathcal{G}}\left(s, \mathbf{H}_{s}, \Theta\right) \\ \mathbf{H}_{0} = \mathbf{X}_{e} \end{cases}, \quad s \in \mathcal{S} \subset \mathbb{R}$$
(1)

where  $\mathbf{F}: \mathcal{S} \times \mathbb{R}^{n \times d} \times \mathbb{R}^{p} \to \mathbb{R}^{n \times d}$  is a **depth-varying vector field** defined on  $\mathcal{G}$ .

#### What do GDEs learn?

They learn a graph-conditioned vector field  $\mathbf{F}_{\mathcal{G}}$  (parametrized by a GNN) such that:

$$\mathbf{Y} := \mathbf{H}(s) = \underbrace{\mathbf{X}_{e}}_{ ext{input embedding}} + \int_{\mathcal{S}} \mathbf{F}_{\mathcal{G}}( au, \mathbf{H}( au), \mathbf{\Theta}) d au$$

with

$$\mathbf{X}_e := \mathbf{X} \mathbf{W}, \quad \mathbf{\Theta}^* = \operatorname*{arg\,min}_{\mathbf{\Theta}} \mathcal{L}$$

#### **Remark:**

The depth variable s assuming real values brings, in the limit, the map

$$\mathbf{X}\mapsto \mathbf{H}(s)=\mathbf{Y}$$

to resemble a network with infinitely dense layers  $\Rightarrow$  i.e., GDEs are the **deep limit** of GNNs.

### Static Models

Graph Convolution Neural Differential Equations (proposed) By choosing  $\mathbf{F}_{\mathcal{G}}$  as a simple GCN, we obtain the **GCDE**:

$$\frac{d}{ds}\mathbf{H}(s) = \sigma(\tilde{\mathbf{D}}^{-\frac{1}{2}}\tilde{\mathbf{A}}\tilde{\mathbf{D}}^{-\frac{1}{2}}\mathbf{H}(s)\Theta)$$

Alternative convolution filters can be employed instead.

#### Note:

- With different priors on F<sub>G</sub> we can build other types of GDEs: GAT, diffusive, message passing, etc.
- Moreover, we can include additional biases in the model structure, e.g. second-order:  $| \ddot{H}(s) = F_{\mathcal{G}}(s, H(s), \Theta)$ stochastic:  $| dH(s) = F_{\mathcal{G}}(s, H(s), \Theta_f)ds + G_{\mathcal{G}}(s, H(s), \Theta_g)dW$

### Semi-supervised Node Classification

Evaluation on Cora, Citeseer, Pubmed

Ablation study: same architecture, different numerical solvers

Model (NFE)	Cora	Citeseer	Pubmed
GCN GCN*	$\begin{array}{c} 81.4\pm 0.5\%\\ 82.8\pm 0.3\%\end{array}$	$\begin{array}{c} 70.9\pm0.5\%\\ 71.2\pm0.4\%\end{array}$	$\begin{array}{c} 79.0 \pm 0.3\% \\ 79.5 \pm 0.4\% \end{array}$
GCDE-rk2 (2) GCDE-rk4 (4) GCDE-dpr5 <b>(158)</b>	$\begin{array}{c} 83.0\pm 0.6\%\\ \textbf{83.8}\pm 0.5\%\\ 81.8\pm 1.2\%\end{array}$	$\begin{array}{c} 72.3\pm0.5\%\\ \textbf{72.5}\pm0.5\%\\ 68.3\pm1.2\%\end{array}$	$\begin{array}{c} \textbf{79.9} \pm 0.3\% \\ \textbf{79.5} \pm 0.4\% \\ \textbf{78.5} \pm 0.7\% \end{array}$

Results across 100 runs.

### Semi-supervised Node Classification



Higher order solvers: generally more performant, provided graph is dense enough to benefit from additional computation.

No direct advantage of solving the ODE accurately with adaptive solvers.

Michael Poli (KAIST)

Graph Neural Ordinary Differential Equation

# Visualizing Node Feature Trajectories



Trajectories defined by a **forward pass** of GCDE on Cora, Citeseer and Pubmed. Color differentiates between node classes.

The trajectories are divergent, suggesting a **non-decreasing** classification **performance** for GCDE models trained with **longer integration intervals**.

### Spatio-Temporal GDEs as Hybrid Systems

In dynamic settings, i.e. direct modeling of dynamical systems, the **depth** variable s assumes the meaning of time: s := t and can be modified depending on the requirements.

For example, given a time window  $\Delta t$ , the **prediction** performed by a **GDE** assumes the form:

$$\mathbf{H}(t + \Delta t) = \mathbf{H}(t) + \int_{t}^{t + \Delta t} \mathbf{F}(\tau, \mathbf{H}(\tau), \Theta) d\tau,$$

We can extend the GDE framework to settings with sequences of graphs:  $\{\mathcal{G}\}$  by leveraging **hybrid dynamical system** machinery.

# Spatio-Temporal GDEs as Hybrid Systems



- $\mathcal{T} := \{t_k\}_{k \in \mathcal{K}}, \ \mathcal{K} \subset \mathbb{N} \setminus \{0\}$ : set of time instants
- $\{(\mathbf{X}_t, \mathcal{G}_t)\}_{t \in \mathcal{T}}$ : state-graph data stream:

GDEs models vector fields defined on graphs. Autoregressive GDEs can handle dynamic topologies (jumps).

Michael Poli (KAIST)

Graph Neural Ordinary Differential Equations

DLGMA'20 2020-02-08 11 / 18

### General Autoregressive GDE

The solution of a general autoregressive GDE model (one timestamp):

$$\begin{cases} \dot{\mathsf{H}}(s) &= \mathsf{F}_{\mathcal{G}_{t_k}}(\mathsf{H}(s), \Theta) \quad s \in [t_{k-1}, t_k] \\ \mathsf{H}^+(s) &= \mathsf{G}_{\mathcal{G}_{t_k}}(\mathsf{H}(s), \mathsf{X}_{t_k}) \quad s = t_k \quad k \in \mathcal{K}, \\ \mathsf{Y} &= \mathsf{K}(\mathsf{H}(s)) \quad s = t_k \end{cases}$$

(2)

• F, G, K are GNN-like operators or general neural network layers

• H<sup>+</sup> represent the value of H after the discrete transition.

#### Idea

GDEs **smoothly steer** latent node features between two time instants and then apply some discrete operator, resulting in a "jump" of **H** which is then processed by an output layer.



Continuous-depth version of GCGRUs, GCDE-GRU:

$$\begin{vmatrix} \dot{\mathsf{H}}(s) &= \mathsf{F}_{\mathsf{GCN}}(\mathsf{H}(s), \Theta) & s \in [t_{k-1}, t_k] \\ \mathsf{H}^+(s) &= \mathsf{GCGRU}(\mathsf{H}(s), \mathsf{X}_{t_k}) & s = t_k \\ \mathsf{Y} &= \mathsf{K}(\mathsf{H}(s)) & s = t_k \end{vmatrix}$$

GCDE-RNNs or GCDE-LSTMs can be obtained in a similar fashion.

# Traffic Forecasting

Evaluation on an undersampled PeMS(M) traffic dataset: 228 sensor stations

To measure robustness to unevenly sampled datasets we turn **regular** observations (5 minute intervals) into **irregular**:

- 70% probability of removal per point
- offline undersampling of training and test data

Model (depth)	MAPE	NRMSE
GRU	$27.52 \pm 0.00$	$1.47\pm0.00$
GCGRU	$24.80 \pm 0.12$	$1.44\pm0.00$
GCDE–GRU	$23.08 \pm 0.11$	$1.40\pm0.01$

Table: Forecasting test results across 5 runs (mean and standard deviation).

#### Future extensions

#### Unknown topology:

- Compatibility with GDEs due to the **algebraic nature** of the relation between the **attention** operator and the **node features**
- If an optimal adaptive graph representation **S**(*s*, **H**) is obtained via some attentive mechanism:

 $\dot{\mathbf{H}} = \sigma \left( \mathbf{SH} \Theta \right).$ 

#### **Control terms:**

$$\begin{cases} \dot{\mathsf{H}}(s) = \mathsf{F}_{\mathcal{G}}\left(s, \mathsf{H}(s), \Theta\right) + \mathsf{U}\left(s\right) \\ \mathsf{H}(0) = \mathsf{X}_{e} \end{cases}, \quad s \in \mathcal{S} \,.$$

This approach encompasses a variety of previously proposed approaches, e.g. **special residual connections**. In particular, a choice is  $\mathbf{U}(s) := \mathbf{U}_{\mathcal{G}}(s, \mathbf{X})$ .

And naturally, other classes of differential equations.

Michael Poli (KAIST)

Graph Neural Ordinary Differential Equations

Thank you Q & A Appendix



Cora accuracy of GCDE models with **different integration times** s. Higher values of S do not affect performance negatively but require a higher number of DLGMA'20

Michael Poli (KAIST)

2020-02-08 17 / 18 Appendix

