

## Intro to Graph Neural Networks from a HEP perspective

Joosep Pata (joosep.pata@cern.ch) NICPB/KBFI, Tallinn, Estonia

Lecture at <u>QU Data Science Basics (Hamburg)</u> May 25, 2021



This work is supported by the Mobilitas Pluss Returning Researcher Grant MOBTP187 of the Estonian Research Council.





Moreno, E.A., Cerri, O., Duarte, J.M. *et al.* JEDI-net: a jet identification algorithm based on interaction networks. *Eur. Phys. J. C* **80,** 58 (2020). https://doi.org/10.1140/epjc/s10052-020-7608-4

#### Data representation

#### set of inputs with N constituents, M features

{..., (pT, η, φ, particle ID), ...}

#### feature matrix (N, M)



jet constituents

Set of feature vectors + ordering  $\rightarrow$  feature matrix

## Simple neural network



**Order**: The ordering is important! A feedforward network trained with e.g. pT-descending ordering would not work with pT-ascending. Which ordering is optimal?

**Representation**: What if for each jet you want to classify, the number of constituents N varies? Need to make all feature matrices the same size (e.g. with 0 padding).

**Structure**: All-to-all connectivity. Every constituent in the input layer can affect every other constituent in the next layer.

## Graph structure



#### Where do we get this graph structure?

1. All-to-all connections, in case of small input sets.

- 2. From physics priors: connect "nearby" elements in advance
- 3. Optimize as a part of the learning process (Graph Structure Learning)

## Example graph structures

Jet constituents (all-to-all)

Particle tracking (neighborhood)



event constituents (all-to-all)





Multilayer calorimeter hits (neighborhood)



Graph Neural Networks in Particle Physics, Jonathan Shlomi, Peter Battaglia, Jean-Roch Vlimant, 2007.13681, 10.1088/2632-2153/abbf9a

## Operations on a graph



## Graph problems

Graph-level prediction, Graph generation

> Jet tagging, event tagging



J. Leskovec et al [2021]

## Graph Convolutional Network (GCN)

![](_page_9_Figure_1.jpeg)

## GCN properties

![](_page_10_Figure_1.jpeg)

• A trainable weight matrix W<sub>i</sub> (d<sub>in</sub> x d<sub>out</sub>) in layer *i* shared across all nodes

- The input and output is a graph. The node features are transformed, the graph structure does not change.
- The GCN is permutation-invariant: it does not matter in which order the set of nodes is formatted as a matrix for computations, due to the permutation-invariant aggregation function
- A very nice overview can be found from Kipf & Welling: <u>https://tkipf.github.io/graph-convolutional-networks/</u>

#### Node smoothing

![](_page_11_Figure_1.jpeg)

![](_page_11_Figure_2.jpeg)

Figure 2. Residual learning: a building block.

#### Deep GCN without skip connections $\rightarrow$ oversmoothing, performance drops

Kipf, Thomas N., and Max Welling. "Semi-supervised classification with graph convolutional networks." arXiv preprint arXiv:1609.02907 (2016).

## Message passing

- Different types of graph-related algorithms can be formulated in the message passing language
- Nodes pass messages to their neighbors
- Aggregate the messages and update the node state

$$\begin{array}{c} \begin{array}{c} \mbox{learnable}\\ \mbox{message}\\ \mbox{function} \end{array} \mbox{edge features} \end{array} \\ \mbox{message}\\ \mbox{message} \end{array} \mbox{$mv^{t+1} = \sum_{w \in N(v)} M_t(h_v^t, h_w^t, e_{vw})$} \\ \mbox{node features} \end{array} \mbox{$h_v^{t+1} = U_t(h_v^t, m_v^{t+1})$} \\ \mbox{learnable update}\\ \mbox{function} \end{array}$$

Gilmer, Justin, et al. "Neural message passing for quantum chemistry." International Conference on Machine Learning. PMLR, 2017.

### GCN as message passing

![](_page_13_Figure_1.jpeg)

Gilmer, Justin, et al. "Neural message passing for quantum chemistry." International Conference on Machine Learning. PMLR, 2017.

# Graph Attention (GAT)

Compute an attention coefficient *α<sub>ij</sub>* between two pairs of connected nodes. Trainable attention vector **a**, feature weight vector **W**.

$$\alpha_{ij} = \frac{\exp\left(\text{LeakyReLU}\left(\vec{\mathbf{a}}^T[\mathbf{W}\vec{h}_i \| \mathbf{W}\vec{h}_j]\right)\right)}{\sum_{k \in \mathcal{N}_i} \exp\left(\text{LeakyReLU}\left(\vec{\mathbf{a}}^T[\mathbf{W}\vec{h}_i \| \mathbf{W}\vec{h}_k]\right)\right)}$$

Update the node feature vector based on nearby attention coefficients.

$$\vec{h}_i' = \sigma\left(\sum_{j \in \mathcal{N}_i} \alpha_{ij} \mathbf{W} \vec{h}_j\right)$$

Inputs are graphs: N x d<sub>in</sub> Outputs are graphs: N x d<sub>out</sub> Attention vector **a** can be interpreted as feature-to-feature association.

Veličković, Petar, et al. "Graph attention networks." arXiv preprint arXiv:1710.10903 (2017).

![](_page_14_Picture_7.jpeg)

## Multi-head GAT

Instead of a single attention coefficient  $\alpha_{ij}$  per a node pair, compute K independent values  $\alpha_{ij}^k$ .

![](_page_15_Figure_2.jpeg)

Veličković, Petar, et al. "Graph attention networks." arXiv preprint arXiv:1710.10903 (2017).

## Interaction network (IN)

In the Interaction Network (2016), the message function M<sub>t</sub> and the node update function U<sub>t</sub> are given as generic neural networks operating on concatenated node and edge inputs.

![](_page_16_Figure_2.jpeg)

Battaglia, Peter W., et al. "Interaction networks for learning about objects, relations and physics." arXiv preprint arXiv:1612.00222 (2016).

## IN for jet tagging

![](_page_17_Figure_1.jpeg)

Moreno, E.A., Cerri, O., Duarte, J.M. *et al.* JEDI-net: a jet identification algorithm based on interaction networks. *Eur. Phys. J. C* **80**, 58 (2020). https://doi.org/10.1140/epjc/s10052-020-7608-4

## IN for particle tracking

Fully connected: 1000 nodes -> 500k edges, not feasible!

Set up an initial sparse hit graph based on node proximity.

Classify possible edges as true/false based on actual track information, predict edge weight

X: node features (nodes  $\times$  3)  $R_a$ : edge features (edges  $\times$  4) R<sub>i</sub>, R<sub>o</sub>: incoming/outgoing edge matrix R<sub>i,o</sub>X: incoming/outgoing nodes (edges x 3)

![](_page_18_Figure_5.jpeg)

![](_page_18_Figure_6.jpeg)

Hitgraph View DeZoort et al

![](_page_18_Figure_8.jpeg)

![](_page_18_Figure_9.jpeg)

Edge Weight Prediction

## Dynamic graph with kNN

- In the previous examples with GCN, GAT and IN, the graph was static and defined/known in advance
- The structure may not be known in advance, or may be inaccurate
- Construct dynamically: point cloud  $\{x_i\} \rightarrow$  for each point  $x_i$ , find k closest neighbors  $\{x_j\}$ , edges  $\{e_{ij}\}$

![](_page_19_Figure_4.jpeg)

#### Detector reconstruction

- kNN + sparse graph adjacency matrix: GravNet
- Cluster energy deposits from overlapping showers in a highly granular, layered tungsten detector simulation
- Predict the energy fraction of each sensor (I) belonging to each shower
  (K): p<sub>ik</sub> vs t<sub>ik</sub>

Qasim, Shah Rukh, et al. "Learning representations of irregular particle-detector geometry with distance-weighted graph networks." *The European Physical Journal C* 79.7 (2019): 1-11.

![](_page_20_Figure_5.jpeg)

Two overlapping showers generated

![](_page_20_Figure_7.jpeg)

### Particle Flow reconstruction

![](_page_21_Figure_1.jpeg)

The Particle Flow algorithm combines elements across different detectors to a global particle-level representation of the collision.

![](_page_22_Figure_0.jpeg)

Pata, J., Duarte, J., Vlimant, JR. et al. MLPF: efficient machine-learned particle-flow reconstruction using graph neural networks. *Eur. Phys. J. C* **81,** 381 (2021)

### **GNNs for Particle Flow**

![](_page_23_Figure_1.jpeg)

#### Performance in simulation

![](_page_24_Figure_1.jpeg)

## **Computational modes**

Suppose we have a dataset of jets we want to classify, each jet having N<sub>i</sub> constituents. NN training often requires batching the data to average gradient updates.

![](_page_25_Figure_2.jpeg)

Adjacency is typically sparse. Suitable for large inputs (N > 1000). Typically requires on-the-fly computation (e.g. pytorch). Adjacency is typically dense. Graphs may be zero-padded / masked to size N. Suitable for small inputs (<1000) and static computational graphs (e.g. tensorflow).

https://graphneural.network/data-modes/

#### Recap

## Graph problems

Graph-level prediction, Graph generation

> Jet tagging, event tagging

![](_page_27_Figure_3.jpeg)

J. Leskovec et al [2021]

### Graph operations

![](_page_28_Figure_1.jpeg)

Invariance with respect to permutations!

### Graph structure

Defined by the process (but not necessarily observable)

Assumed (static)

Learned (dynamic)

![](_page_29_Figure_4.jpeg)

![](_page_29_Figure_5.jpeg)

![](_page_29_Figure_6.jpeg)

![](_page_29_Figure_7.jpeg)

![](_page_29_Figure_8.jpeg)

## Advantages of GNNs

- Encode physics priors in the graph structure
- Invariance to permutations / ordering
- Sparse and irregular problem geometries
- Efficient computation and memory representation

### Useful references

- HEPML Living Review: <u>https://iml-wg.github.io/HEPML-</u> <u>LivingReview/</u>
- ML on Graphs @ Stanford: <u>http://web.stanford.edu/class/</u> <u>cs224w/</u>
- Graph Representation Learning book (WIP): <u>https://</u> <u>www.cs.mcgill.ca/~wlh/grl\_book/</u>

#### Practical exercise

![](_page_32_Figure_1.jpeg)

Jupyter notebook: link

## Acknowledgements

Thanks to Jan Kieseler and Jean-Roch Vlimant for comments & discussion

#### Backup

## Dynamic graph CNN (DGCNN)

Construct neighbor graph: for each point  $x_i$ , find k closest neighbors  $\{x_j\}$ , edges  $\{e_{ij}\}$ 

![](_page_35_Figure_2.jpeg)

construct an edge feature using a learnable function

$$\boldsymbol{h}_{\boldsymbol{\Theta}}(\boldsymbol{x}_i, \boldsymbol{x}_{i_j}) = \bar{\boldsymbol{h}}_{\boldsymbol{\Theta}}(\boldsymbol{x}_i, \boldsymbol{x}_{i_j} - \boldsymbol{x}_i),$$

Compute the new point features  $x_i$  using an aggregation over the edges

$$oldsymbol{x}_i' = igsqcap_{j=1}^k oldsymbol{h}_{oldsymbol{\Theta}}(oldsymbol{x}_i,oldsymbol{x}_{i_j}),$$

#### Wang, Yue, et al. "Dynamic graph CNN for learning on point clouds." Acm Transactions On Graphics (tog) 38.5 (2019): 1-12.

## ParticleNet: DGCNN in HEP

input coordinates (B, N, C)

![](_page_36_Figure_2.jpeg)

Qu, Huilin, and Loukas Gouskos. "Jet tagging via particle clouds." Physical Review D 101.5 (2020): 056019.

## ParticleNet full model

- Up to 100 highest-pT constituents of each jet
- relative η, φ coordinates wrt. the jet axis as coordinates
- Features are derived from 4-momentum (log transforms, ratios)
- Coordinates in subsequent layers are derived from previous layer outputs

![](_page_37_Figure_5.jpeg)

Qu, Huilin, and Loukas Gouskos. "Jet tagging via particle clouds." *Physical Review D* 101.5 (2020): 056019.

## GravNet/GarNet

- Full, dense NxN distance matrix can be too large to store for N>few hundred, kNN can be expensive in a highdimensional input space
- In case low latency, low memory consumption is desirable, optimize by using a sparse adjacency matrix, separating spatial components and feature components

![](_page_38_Figure_3.jpeg)

Qasim, Shah Rukh, et al. "Learning representations of irregular particle-detector geometry with distance-weighted graph networks." *The European Physical Journal C* 79.7 (2019): 1-11.

## **GNNs and Transformers**

Most state-of-the-art language processing models use an attention-based "transformer" architecture: a dense attention matrix with elements  $A_{ij}$  is computed between input elements  $x_{i.}$ . The attention matrix **A** is used to successively transform the input elements.

![](_page_39_Figure_2.jpeg)

In GNNs, the learned graph adjacency is usually sparse, but is similarly used to propagate information between associated input elements to transform them.

https://ai.googleblog.com/2020/10/rethinking-attention-with-performers.html

## Scalable pairwise operations

- Naive kNN graph construction (e.g. *tf.nn.top\_k*) scales as O(N<sup>2</sup>) with the number of input nodes N
- For N>few hundred, this can be prohibitive (in memory and computation) and thus require splitting up the data
- To process long sequences or full events, there has been recent interest in models that scale better than quadratically, e.g. using an approximate bucketing based on Locality Sensitive Hashing

![](_page_40_Figure_4.jpeg)

Kitaev, Nikita, Łukasz Kaiser, and Anselm Levskaya. "Reformer: The efficient transformer." arXiv preprint arXiv:2001.04451 (2020).