

## Graph nets and transformers: intro for HEP folks

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

Lecture at <u>3rd Terascale School of Machine Learning (Hamburg)</u> October 11, 2022



source: L. Gray

#### Data modalities

- Tabular: no particular relation between features of a sample
  - x = [feat<sub>1</sub>, feat<sub>2</sub>, ...]
- Image: features naturally embed into a 2D/3D Cartesian grid
  - x = Matrix(320,320,3)
- Set: each sample consists data points with features, no particular order
  - x = {p<sub>1</sub>, p<sub>2</sub>, p<sub>3</sub>, ...}, p<sub>n</sub> = [feat<sub>1</sub>, feat<sub>2</sub>, ...]
  - Set embaddable in 2D/3D: some features can be easily interpreted as x, y, z in a space

#### Data modalities

• Sequence: data points in a sample are naturally ordered (e.g in time)

- Heterogeneous set: data points in a sample have different features
  - x = {p<sub>1</sub>, p<sub>2</sub>, k<sub>1</sub>, ...}, p<sub>n</sub> = [feat<sub>1</sub>, feat<sub>2</sub>], k<sub>n</sub> = [feat<sub>3</sub>, feat<sub>4</sub>, feat<sub>5</sub>]
- Graph: data points in a sample have meaningful, quantifiable, observable relations

• 
$$x = \{p_1, p_2, p_3, ...\}, A = \{A_{12}, A_{13}, A_{23}, ...\}$$





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 necessarily 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)

#### Sparse tensors



- Not always possible to store NxN adjacency matrix (e.g. if N > few thousands)
- Sparse graph adjacency matrices are typically represented in COO sparse format in DL libraries
- Backprop only on nonzero data values, not on row/column indices

# Ragged/jagged tensor



- Memory-efficient way to represent sequences of different length in a single batch
- Used to avoid zero-padding & masking when doing batch processing
- GPUs generally like inputs to be of the same length: memory / runtime tradeoff

#### 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

## Graph problems

Graph-level prediction, Graph generation

> Jet tagging, event tagging



J. Leskovec et al [2021]

## Operations on a graph



## Graph Convolutional Network (GCN)



trainable weight matrix W: din x dout

update rule:  $X'_i \rightarrow ReLU[A_{ji} \cdot (X_j \cdot W)]$ 

e.g.  $X'_{5} = \text{ReLU}[A_{65}(X_{6} \cdot W) + A_{25}(X_{2} \cdot W)]$ 

## **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.



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/

Adjacency is typically sparse.

Suitable for large inputs (N > 1000).

Typically requires on-the-fly computation (e.g. pytorch).

## GCN properties



• 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





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

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

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

#### GCN as message passing



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

## Graph Attention (GAT)

Compute an attention coefficient  $\alpha_{ij}$  between 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).



## Multi-head GAT

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



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.



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

## IN for jet tagging



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 × 3) R<sub>a</sub>: edge features (edges × 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)





Hitgraph View DeZoort et al



DeZoort, Gage, et al. "Charged particle tracking via edge-classifying interaction networks." arXiv preprint arXiv:2103.16701 (2021).

## Dynamic graph with kNN

- In the previous examples with GCN, GAT and IN, the graph was static and defined/known in advance
- Often, the graph 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  $\{e_{ij}\}$



## Dynamic graph CNN (DGCNN)

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



construct an edge feature using a learnable function

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

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: EdgeConv block

input coordinates (B, N, C)



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



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



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.

#### 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.



Two overlapping showers generated



#### Particle Flow reconstruction



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



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**



#### Performance in simulation

Predicts particle multiplicitly better than the baseline rule-based PF



Runtime (and memory) scale linearly with event size.

#### Performance in CMS



https://arxiv.org/pdf/2203.00330.pdf
#### Recap

# Graph problems

Graph-level prediction, Graph generation

> Jet tagging, event tagging



J. Leskovec et al [2021]

#### Graph operations



Invariance with respect to permutations!

### Graph structure

Defined by the process (but not necessarily observable)

Assumed (static)

Learned (dynamic)











# Advantages/disadvantages

- + Encode physics priors in the graph structure
- + Insensitive to the ordering of inputs
- + Sparse and irregular problem geometries
- + Efficient computation and memory representation

- Support for sparse data structures on GPUs and in DL libraries is not always great (but it's growing)
- It's not always obvious how to best cast the physics problem as a graph problem, i.e. what defines the graph
  - Graph structure may be mismeasured or not known accurately
- Deeper nets often do not perform better

# Useful references

- HEPML Living Review: <u>https://iml-wg.github.io/HEPML-</u> LivingReview/
- 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>
- Graph Neural Networks in Particle Physics: <u>https://arxiv.org/</u> <u>pdf/2203.12852.pdf</u>

#### Practical exercise



Jupyter notebook: github, colab

#### **GNNs to Transformers**



The GNN updates (encodes) each node, given information from nearby nodes, in a learnable way, to minimize an overall objective function.

#### GNNs need structure

#### explicit / predefined









k-nearest neighbors, k = 5

0

k nearest neighbors graph (k = 3)

implicit / learnable



#### Self-attention

Project each element from input space X to Q,K,V with learnable weights.



https://arxiv.org/pdf/2108.04253.pdf

https://jalammar.github.io/illustrated-transformer/

### Self-attention outputs

 $N_{elem} \ x \ N_{elem}$  attention matrix



- Compute the matrix product of per-element queries with per-element keys
- Retrieve the corresponding values according to softmax-normalized attention matrix



#### Multi-head attention



#### Multi-head outputs



# Input to a subsequent task

#### 1) Concatenate all the attention heads Z<sub>0</sub> Z<sub>1</sub> Z<sub>2</sub> Z<sub>3</sub> Z<sub>4</sub> Z<sub>5</sub> Z<sub>6</sub> Z<sub>7</sub>

2) Multiply with a weight matrix W<sup>o</sup> that was trained jointly with the model

Х

WO

3) The result would be the Z matrix that captures information from all the attention heads. We can send this forward to the FFNN



=

#### Self-attention to encoder



#### Permutation invariance



https://www.scipost.org/SciPostPhys.12.6.188/pdf

### Point cloud transformer



https://arxiv.org/abs/2102.05073

### Jet tagging via transformers



#### Add explicit particle-particle interaction terms.

task. Specifically, for a pair of particles a, b with 4-vectors  $p_a$ ,  $p_b$ , we calculate the following 4 features:

$$\Delta = \sqrt{(y_a - y_b)^2 + (\phi_a - \phi_b)^2},$$
  

$$k_{\rm T} = \min(p_{{\rm T},a}, p_{{\rm T},b})\Delta,$$
  

$$z = \min(p_{{\rm T},a}, p_{{\rm T},b})/(p_{{\rm T},a} + p_{{\rm T},b}),$$
  

$$m^2 = (E_a + E_b)^2 - \|\mathbf{p}_a + \mathbf{p}_b\|^2,$$
(3)



https://arxiv.org/pdf/2202.03772.pdf

#### Particle attention block



#### Class attention



to the classifier. We put this problem in evidence by showing that inserting CLS later improves performance (*middle*). In the **CaiT** architecture (*right*), we further propose to freeze the patch embeddings when inserting CLS to save compute, so that the last part of the network (typically 2 layers) is fully devoted to summarizing the information to be fed to the linear classifier.

class token

encoded

particle data

#### Attention as a graph



#### Full self attention ~ all-to-all dense graph. Naively N<sup>2</sup> time/memory complexity!

https://ai.googleblog.com/2021/03/constructing-transformers-for-longer.html

## Global vs. local attention

In case the input is an ordered sequence (e.g. a long text, big image), attention can be constrained to local regions.



# GNNs vs. 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.



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 / attention matrix construction (e.g.  $tf.nn.top_k$ ) scales as O(N<sup>2</sup>) with the number of input nodes N.



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

#### Approximating the attention



- The hyperparameter m << L of random projections to make for each Q, K tunes how closely the attention mechanism is approximated
- Does not rely on sparsity / structure in the attention matrix

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

#### Training on large sequences



# Domain-agnostic models

To what extent can one use a single model architecture on different data domains?



Figure 2: The Perceiver IO architecture. Perceiver IO maps arbitrary input arrays to arbitrary output arrays in a domain agnostic process. The bulk of the computation happens in a latent space whose size is typically smaller than the inputs and outputs, which makes the process computationally tractable even for very large inputs & outputs.

PerceiverIO https://arxiv.org/abs/2107.14795v3

#### Summary

dynamic graph (e.g. EdgeConv, GarNet) Particle Transformer

GNN with a fixed neighborhood graph

Point Cloud Transformer

Generic models (e.g. Perceiver-IO)

4-----

- more assumed structure
- fewer parameters
- more scalable on large inputs

- less assumed structure
- more parameters
- less scalable on large inputs





#### Transformers exercise



Jupyter notebook: github, colab

#### Backup

As an example (batches, elements, features) = (2, 6400, 25)





Requires batch-mode graphs. No N<sup>2</sup> allocation or computation needed.

One scalable combined graph layer. The input elements are projected into a learnable embedding space. Nearby elements in the embedding space are binned to fixed-size bins. A fully-connected graph is built in each bin, which is used for one or multiple graph convolutions that are used to transform the input elements. Finally, the transformed elements are unbinned.



The learned graph structure in the first layer of the model. Each plot corresponds to a bin where an all-to-all graph adjacency matrix is built between the PFElements in the bin. Emtpy bins are present since for efficiency reasons in GPU training and inference, the model operates on a fixed-size zero-padded input.


The learned binning structure in the first two layers of the model. We show one simulated ttbar event, with each point corresponding to a PFElement in the event. The colors correspond to the assignment of the PFElements into the bins in each layer.