# GRAPH CONVOLUTIONAL NETWORKS WITH AUTOENCODER-BASED COMPRESSION AND MULTI-LAYER GRAPH LEARNING

L. Giusti[1] C. Battiloro[2] P. Di Lorenzo[2] S. Barbarossa[2]

[1] DIAG Department, Sapienza University of Rome, Via Ariosto 25, 00185, Rome, Italy
[2] CDIET Department, Sapienza University of Rome, Via Eudossiana 18, 00184, Rome, Italy

## ABSTRACT

The aim of this work is to propose a novel architecture and training strategy for compress the convolutional features at multiple hidden layers, hinging on a novel end-to-end training procedure that learns different graph representations per each layer.

Contribution: We exploit autoencoders in each layer, before applying the pointwise non-linearity, so that the convolutional features can be tunably compressed in an information-rich embedding. Then, since compression calls for learning a new graph representation to be used in the following layer, we formulate a novel training strategy that jointly optimizes the GNN weight parameters and the graph representations at different layers.

## SIGNALS ON GRAPHS

- Let $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ be a weighted undirected graph
- The sets $\mathcal{V} = \{1, 2, ..., N\}$ and $\mathcal{E} = \{a_{i,j}\}_{i,j \in \mathcal{V}}$ are the sets of vertices and edges, respectively
- The weights $a_{i,j} \geq 0$ if there is a relationship from vertex i to vertex j, or $a_{i,j} = 0$ otherwise.
- The adjacency matrix $\mathbf{A} \in \mathbb{R}^{N \times N}$: a collection of all weights, i.e., $\mathbf{A} = \{a_{i,j}\}$, $i, j = 1, ..., N$
- The Laplacian matrix: $\mathbf{L} = \text{diag}(\mathbf{1}^T \mathbf{A}) - \mathbf{A}$, where $\text{diag}(\mathbf{x})$ is a matrix having $\mathbf{x}$ as main diagonal, and zeros elsewhere
- A graph signal (or data) is defined as a one-to-one mapping from the set $\mathcal{V}$ of vertices to the set of real numbers:

$$\mathbf{x} = \mathcal{V} \to \mathbb{R} \tag{1}$$

- An order $K$ linear shift invariant graph filter (LSIGF) can be written as a $K$-degree polynomial of the shift operator $\mathbf{S}$, with coefficients $\mathbf{h} = [h_0, ..., h_{K-1}]^T$.
- Let $\mathbf{u}$ and $\mathbf{y}$ be the input and the filtered signals, respectively, we have:

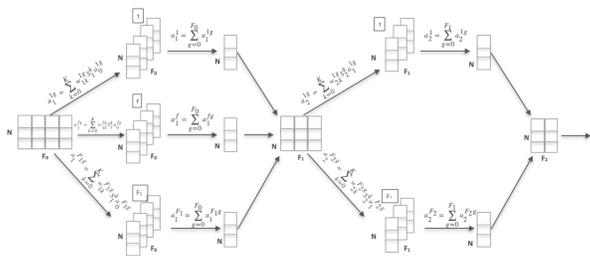$$\mathbf{y} = \sum_{k=0}^{K-1} [\mathbf{h}]_k \mathbf{S}^k \mathbf{u}. \tag{2}$$

- The LSIGF are able to account for the local structure of the graph, requiring information only from the $K$-neighborhood of each node.
- Linear and shift invariant graph filters represent a legit generalization of the convolution operation for signals supported on graphs [1], and are the basic building block of GCNs.

## GRAPH CONVOLUTIONAL NETWORKS

- The $l$-th layer of a GCN, taking as input $\widetilde{\mathbf{Z}}_{l-1} = \{\widetilde{\mathbf{z}}_{l-1}^f\}_{f=1}^{F_{l-1}}$ and yielding as output $\widetilde{\mathbf{Z}}_l = \{\widetilde{\mathbf{z}}_l^g\}_{g=1}^{F_l}$, with pointwise non-linearity $\sigma_l(\cdot)$, reads as: [2]:

$$\widetilde{\mathbf{z}}_l^g := \sigma_l\left(\sum_{f=1}^{F_{l-1}} \sum_{k=0}^{K_l-1} [\mathbf{h}_l^{fg}]_k \mathbf{S}^k \widetilde{\mathbf{z}}_{l-1}^f\right), \quad g = 1, ..., F_l. \tag{3}$$

- The order $K_l$ of the filters, the number $F_l$ of convolutional features of the output, and the non-linearity $\sigma_l(\cdot)$ are hyperparameters to be chosen at each layer
- A GCN of depth $L$ with input data $\mathbf{X}$ is built as the stack of $L$ layers defined as in (3), where $\widetilde{\mathbf{Z}}_0 = \mathbf{X}$
- Based on the learning task, an additional multi layer perceptron (MLP) can be inserted after the last layer
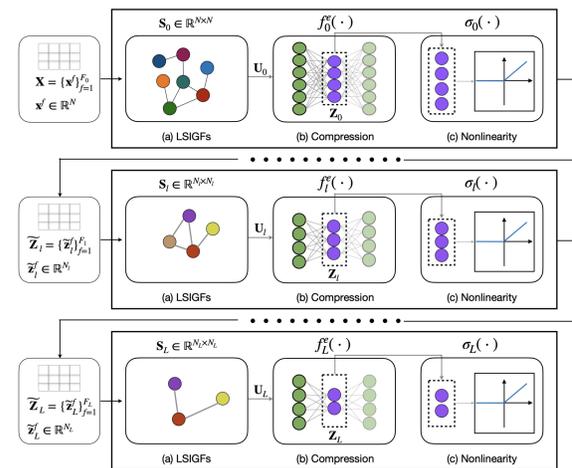


## AUTOENCODER-AIDED GRAPH CONV NETS

- The idea is to exploit autoencoders to perform representation learning and dimensionality reduction in the context of GCNs [3].
- Each layer is composed of three main stages: (i) A linear shift invariant graph filtering stage; (ii) Autoencoder-based compression; (iii) Pointwise non-linearity
- Autoencoders are able to reduce the dimension of the hidden layers' convolutional features, thus learning powerful and task-oriented low-dimensional representations in a data-driven fashion.
- To this aim, the layer in (3) is modified in the following way:

$$\widetilde{\mathbf{z}}_l^g = \sigma_l\big(\underbrace{f_l^e(\mathbf{u}_l^g)}_{\mathbf{z}_l^g}\big), \quad g = 1, ..., F_l, \tag{4}$$

Where $\widetilde{\mathbf{z}}_l^g \in \mathbb{R}^{N_l}$, $N_l \in \mathbb{N}$ is (generally) smaller than $N_{l-1}$, $f_l^e : \mathbb{R}^{N_{l-1}} \to \mathbb{R}^{N_l}$ is the encoder function of an autoencoder $f_l^d \circ f_l^e : \mathbb{R}^{N_{l-1}} \to \mathbb{R}^{N_{l-1}}$ associated with the $l$-th layer of the GCN, and

$$\mathbf{u}_l^g = \sum_{f=1}^{F_{l-1}} \sum_{k=0}^{K_l-1} [\mathbf{h}_l^{fg}]_k \mathbf{S}_l^k \widetilde{\mathbf{z}}_{l-1}^f, \quad g = 1, ..., F_l, \tag{5}$$

where $\mathbf{S}_l$ denotes the shift operator associated with the $l$-th layer. The compressed features at layer $l$ are denoted as $\mathbf{Z}_l = \{\mathbf{z}_l^g\}_{g=1}^{F_l}$, and the final output of the layer is $\widetilde{\mathbf{Z}}_l \in \mathbb{R}^{N_l \times F_l}$. We call the stack of $L$ layers as in (4) an Autoencoder-Aided Graph Convolutional Network (AA-GCN).



## PROBLEM FORMULATION

$$\min_{\{\mathbf{A}_l\}_{l=1}^L, \mathbf{H}, \mathbf{W}} \mathcal{L}(\{\mathbf{A}_l\}_{l=1}^L, \mathbf{H}, \mathbf{W}; \{\mathbf{x}_i, \mathbf{y}_i\}_{i \in \mathcal{T}})$$

$$+ \eta \sum_{l=1}^L \sum_{g=1}^{F_l} \|f_l^d \circ f_l^e(\mathbf{w}_l; \mathbf{u}_l^g) - \mathbf{u}_l^g\|_2^2; \to \text{Autoencoders' Loss}$$

$$+ \beta \sum_{l=1}^L \text{Tr}\{\widetilde{\mathbf{Z}}_l^T \mathbf{L}_l \widetilde{\mathbf{Z}}_l\}; \to \text{Promote signal smoothness}$$

$$- \gamma \sum_{l=1}^L \mathbf{1}^T \log(\mathbf{A}_l \mathbf{1}); \to \text{Penalize disconnected components}$$

$$+ \lambda \sum_{l=1}^L \|\mathbf{A}_l\|_F^2; \to \text{Weights' regularizer}$$

subject to

$$[\mathbf{A}_l]_{i,i} = 0 \to \text{no self-loops in the learnt graphs} \tag{6}$$
$$[\mathbf{A}_l]_{i,j} = [\mathbf{A}_l]_{j,i} \geq 0, \quad \forall i, j, l \to \text{edge weights must be positive and symmetric}$$
$$\text{Tr}\{\mathbf{L}_l\} = d_l, \quad \forall l \to \text{avoid null solutions}$$

where $\lambda, \beta, \gamma, \eta$, and $d_l$ are non-negative parameters to be tuned.

## HALF-VECTOR

- Since the adjacency matrices are symmetric, the number of variables of the optimization problem can be greatly reduced (approximately by a factor of two) solving for the lower triangular parts of $\mathbf{A}_l$ for all $l = 1, ..., L$
- Let $\boldsymbol{\alpha}_l := \text{vech}(\mathbf{A}_l) \in \mathbb{R}^{\frac{N(N+1)}{2}}$ be the half-vectorization of $\mathbf{A}_l$, obtained by vectorizing only the lower triangular part of $\{\mathbf{A}_l\}_l$
- Then, the following relations hold:

$$\text{vec}(\mathbf{A}_l) = \mathbf{M}_d \boldsymbol{\alpha}_l \iff \mathbf{A}_l = \text{vec}^{-1}(\mathbf{M}_d \boldsymbol{\alpha}_l), \tag{7}$$

- $\text{vec}(\cdot)$ and $\text{vec}^{-1}(\cdot)$ are the vectorization and the inverse vectorization operators, respectively
- $\mathbf{M}_d \in \mathbb{R}^{N^2 \times \frac{N(N+1)}{2}}$ is the (highly sparse) duplication matrix

All the objective terms and the constraints can be easily recast in terms of the variables $\boldsymbol{\alpha}_l$, for $l = 1, ..., L$

## ALGORITHMIC SOLUTION

**Algorithm 1 : AA-GCN TRAINING**

Inputs:
  $\mu \in \mathbb{R}$: Learning rate.
  $\Delta_\mu(\cdot)$: Optimizer-dependent backpropagation step.
  $\Pi(\cdot)$: Projection operator on the feasible set.
  $E \in \mathbb{N}_+$: Maximum number of training iterations.
  $\{\mathcal{B}_t\}_{t=1}^E$: Training dataset batches
  Estimates initializations $\widehat{\mathbf{H}}_0$, $\widehat{\mathbf{W}}_0$ and $\{\widehat{\boldsymbol{\alpha}}_{l,0}\}_l$.
  Loss $\mathcal{L}(\cdot)$.
Outputs:
  $\{\widehat{\boldsymbol{\alpha}}_l\}_l$: Learned graph encodings.
  $\widehat{\mathbf{H}}$: Learned graph filters weights.
  $\widehat{\mathbf{W}}$: Learned autoencoders weights.
1: **function** AA-GCN TRAINING(**Inputs**)
2:    **for** $t \in [1, E]$ **do**
3:      $\widehat{\mathbf{H}}_{t+1} = \Delta_\mu\big(\nabla_\mathbf{H} \mathcal{L}(\widehat{\mathbf{H}}_t; \mathcal{B}_t, \{\widehat{\boldsymbol{\alpha}}_{l,t}\}_l, \widehat{\mathbf{W}}_t)\big)$
4:      $\widehat{\mathbf{W}}_{t+1} = \Delta_\mu\big(\nabla_\mathbf{W} \mathcal{L}(\widehat{\mathbf{W}}_t; \mathcal{B}_t, \{\widehat{\boldsymbol{\alpha}}_{l,t}\}_l, \widehat{\mathbf{H}}_t)\big)$
5:      $\widehat{\boldsymbol{\alpha}}_{l,t+1} = \Pi\big(\Delta_\mu\big(\nabla_{\boldsymbol{\alpha}_l} \mathcal{L}(\{\widehat{\boldsymbol{\alpha}}_{l,t}\}_l, \mathcal{B}_t, \widehat{\mathbf{W}}_t, \widehat{\mathbf{H}}_t)\big)\big), \forall l$
6:    **return** $\{\widehat{\boldsymbol{\alpha}}_l\}_l = \{\widehat{\boldsymbol{\alpha}}_{l,E}\}_l$, $\widehat{\mathbf{W}} = \widehat{\mathbf{W}}_E$, $\widehat{\mathbf{H}} = \widehat{\mathbf{H}}_E$

## PROJECTION ALGORITHM

**Algorithm 2** Euclidean projection of $\widetilde{\boldsymbol{\alpha}}_l$ onto $\mathcal{X}$ [4]

Input:
  $\{\widetilde{\boldsymbol{\alpha}}_l\}_l$: Updated graphs parameters.
  $\{N_l\}_l$: Number of nodes of $\{\mathcal{G}_l\}_l$.
  $\{d_l\}_l$: Expected degrees of nodes at layer $l$.
Output:
  $\{\widehat{\boldsymbol{\alpha}}_l\}_l$: Projected updated graphs estimates.
Operator $\Pi_\mathcal{X}(\{\widetilde{\boldsymbol{\alpha}}_{l,t+1}\}_l)$:
  **for** $l = 1, ..., L$ **do**:
    Set $[\mathbf{A}_l]_{i,i} = 0$;
    Sort $\widetilde{\boldsymbol{\alpha}}_l$ in increasing order;

$$\rho := \max_{1 \leq j \leq dim(\widetilde{\boldsymbol{\alpha}}_l)} \underbrace{[\widetilde{\boldsymbol{\alpha}}_l]_j + \frac{1}{j}\left(\frac{d_l}{2} - \sum_{i=1}^j [\widetilde{\boldsymbol{\alpha}}_l]_i\right)}_{[\tilde{\boldsymbol{\alpha}}_l^*]_j} \text{ s.t. } [\tilde{\boldsymbol{\alpha}}_l^*]_j > 0$$

$$\lambda := \frac{1}{\rho}\left(\frac{d_l}{2} - \sum_{i=1}^\rho [\widetilde{\boldsymbol{\alpha}}_l]_i\right)$$

    Set $[\widehat{\boldsymbol{\alpha}}_l]_i = \max\left([\widetilde{\boldsymbol{\alpha}}_l]_i + \lambda, 0\right)$

## NUMERICAL RESULTS





- The first results show the accuracy score compared to the compression ratio: $\rho = N_1/N$
- The results show the accuracy score compared to the SNR of the training data: $\text{SNR} = 10\log_{10}(\sigma_T^2/\sigma_\varepsilon^2)$, $\sigma_T^2$ and $\sigma_\varepsilon^2$ are the variance of the data used for training our model and the variance of the AWGN, respectively.
- For a fair comparison, the second hidden layer does not provide a coarser version of the first one

## CONCLUSIONS

- We have enabled tunable compression of the convolutional features, while learning different graph representations jointly with the GNN parameters
- The architecture scales well with the number of nodes of the input graph, extracting higher level representations of the convolutional features
- Experiments illustrate the competitive performance of our architecture with respect to state of the art methods
- Future developments of this research trend include: Topological Neural Networks, Explainability, include additional regularisations to the autoencoders' loss

## REFERENCES

[1] S. Segarra, A. G. Marques, and A. Ribeiro. Optimal graph-filter design and applications to distributed linear network operators. IEEE Transactions on Signal Processing, 65(15):4117–4131, 2017.

[2] F. Gama, A. G. Marques, G. Leus, and A. Ribeiro. Convolutional neural network architectures for signals supported on graphs. IEEE Transactions on Signal Processing, 67(4):1034–1049, 2019.

[3] D. E. Rumelhart, G. E. Hinton, and R. J. Williams. Learning internal representations by error propagation. In Parallel Distributed Processing: Explorations in the Microstructure of Cognition, Volume 1: Foundations, pages 318–362. MIT Press, Cambridge, MA, 1986.

[4] W. Wang and M. A. Carreira-Perpinán. Projection onto the probability simplex: An efficient algorithm with a simple proof and an application. arXiv:1309.1541, 2013.