

# A Comprehensive Study of CayleyNets

Alice Valena De Lorenci  
Institut Polytechnique de Paris  
Palaiseau, France  
alice.valenca@telecom-paris.fr

Julian Alvarez de Giorgi  
Institut Polytechnique de Paris  
Palaiseau, France  
julian.alvarez@telecom-paris.fr

## Abstract

The objective of this report was to study the Graph Convolutional Neural Network with complex rational spectral filters proposed by Levie et al. [6], called CayleyNets. With that intent, we start by overviewing the underlying spectral graph theory, leading to the construction of general Spectral Graph Convolutional Neural Network models and of CayleyNets. Our objective was to achieve a clear and intuitive description of this model and its properties. Furthermore we evaluate the CayleyNets through several experiments on the CORA dataset, using our own Python implementation of the model.

**CCS Concepts:** • Computing methodologies → Neural networks; Spectral methods.

**Keywords:** Convolutional Graph Neural Networks; Geometrical Deep Learning; Cayley Filters; Spectral Graph Theory

## 1 Introduction

In recent years, Deep Learning models have enjoyed increasing success in dealing with high dimensional Euclidian data. Models such as Convolutional Neural Networks (CNNs) are able to efficiently handle images, which are inherently high dimensional, by exploring underlying structural characteristics of the data, in this case, the translational invariance of signal classes.

Going beyond Euclidean data, graphs arise as a natural representation for a variety of data, such as social networks, molecular structures and anatomical surfaces. The power of graphs can be geometrically understood by interpreting them as metric spaces that are defined locally. This geometric interpretation promptly raises the question of whether it is possible to generalize CNNs to signals defined over graphs; or, more generally, whether it is possible to generalize Deep Learning to non-Euclidean data. Indeed this is the question behind the growing field of Geometrical Deep Learning [2]. One of the current approaches of this field to dealing with signals defined over graphs is that of generalizing convolutions by means of the spectrum of the graph Laplacian [3].

In this report, we will analyze a type of Spectral Graph Convolutional Neural Networks introduced by [6], called CayleyNets. First, in Section 2, we overview the main mathematical aspects of spectral graph theory; in Section 3 we

introduce CayleyNets; finally, in Section 4 we propose a series of experiments to evaluate the performance of the model and illustrate its properties.

## 2 Background

### 2.1 Spectral graph theory

Consider an undirected graph  $\mathcal{G}(V, E, \mathbf{W})$ , where  $V = \{1, \dots, n\}$  is the set of  $n$  nodes,  $E$  is the set of edges and  $\mathbf{W} = (w_{ij})_{ij} \in \mathbb{R}^{n \times n}$  is the adjacency matrix, defined as:

$$\begin{aligned} w_{ij} &= 0, & \text{if } (i, j) \in E \\ w_{ij} &> 0, & \text{if } (i, j) \notin E \end{aligned} \quad (1)$$

Since the graph is undirected, the adjacency matrix is symmetric. From the adjacency matrix, the degree matrix  $\mathbf{D} = \text{diag}(\sum_{j \neq i} w_{ij})_{i=1}^n$  is defined, it is a diagonal matrix whose entries correspond to the degrees of the nodes.

The Laplacian matrix of the graph is defined as  $\Delta = \mathbf{D} - \mathbf{W} \in \mathbb{R}^{n \times n}$ , and it can be written as  $\Delta = \nabla^T \nabla$ ,  $\nabla$  being the incidence matrix of the graph, for some fixed arbitrary direction of the edges. From this construction of the Laplacian it is easily seen that it is a symmetric positive semi-definite matrix,  $\Delta \geq 0$ . Furthermore, one can define the normalized Laplacian as  $\tilde{\Delta} = \mathbf{D}^{-1/2} \Delta \mathbf{D}^{-1/2}$ .

Since the Laplacian matrix is real and symmetric, it has a spectral decomposition  $\Delta = \Phi \Lambda \Phi^T$ .  $\Phi = (\phi_1, \dots, \phi_n)$  with  $\phi_i \in \mathbb{R}^n$ ,  $i = 1, \dots, n$ , is the matrix whose columns are the right eigenvectors of  $\Delta$ , and  $\Lambda = \text{diag}(\lambda_i)_{i=1}^n$  is the diagonal matrix of eigenvalues, such that,  $0 = \lambda_1 \leq \dots \leq \lambda_n$  [1].

The spectral decomposition of the Laplacian matrix is used to define the graph Fourier transform, in analogous fashion to the classical Fourier transform [4]. If we let  $\mathbf{f} = (f_1, \dots, f_n)^T \in \mathbb{R}^n$  be a signal defined on the nodes of the graph, its graph Fourier transform is  $\hat{\mathbf{f}} = \Phi^T \mathbf{f} = (\hat{f}_1, \dots, \hat{f}_n)^T$ . Correspondingly, the inverse graph Fourier transform is  $\mathbf{f} = \Phi \hat{\mathbf{f}}$ , i.e., any signal on the graph can be expressed as:

$$\mathbf{f} = \hat{f}_1 \phi_1 + \dots + \hat{f}_n \phi_n \quad (2)$$

In fact, the eigenvalues  $\lambda_i$  are a measure of the smoothness of the associated eigenvector  $\phi_i$ , interpreted as a signal over the graph. The smaller the eigenvalue, the smoother the corresponding eigenvector. By extension, the Fourier coefficients  $(\hat{f}_1, \dots, \hat{f}_n)^T$  of the signal  $\mathbf{f}$  on the  $\{\phi_1, \dots, \phi_n\}$  basis reflect the smoothness of the signal. Those considerations give rise to the interpretation of the eigenvalues of the Laplacian as graph frequencies and of  $\{\phi_1, \dots, \phi_n\}$  as

the graph Fourier basis. It is important to highlight that the graph Fourier transform is defined relative to the graph  $\mathcal{G}$ . For a detailed description of the connection between the spectral decomposition of the Laplacian and the smoothness of signals defined over graphs the reader may refer to [3].

Having defined the graph Fourier transform, it is possible to define the spectral convolution of two given signals  $\mathbf{f}, \mathbf{g}$  on the graph:

$$\mathbf{f} * \mathbf{g} := \Phi \left( \left( \Phi^T \mathbf{g} \right) \cdot \left( \Phi^T \mathbf{f} \right) \right) \quad (3)$$

where  $\cdot$  denotes the element-wise multiplication. The spectral convolution can be equivalently represented as:

$$\mathbf{f} * \mathbf{g} = \Phi \hat{\mathbf{G}} \Phi^T \mathbf{f} \quad (4)$$

$$\mathbf{f} * \mathbf{g} = \mathbf{G} \mathbf{f} \quad (5)$$

where  $\hat{\mathbf{G}} := \text{diag}(\hat{g}_1, \dots, \hat{g}_n)$  and  $\mathbf{G} := \Phi \hat{\mathbf{G}} \Phi^T$ . Those equivalent representations will be useful in the sequence. We say that the signal  $\mathbf{f}$  is filtered by  $\mathbf{g}$ .

In order to build intuition about spectral operations on graphs, it is interesting to analyze Eq. 4. Consider the following manipulations of Eq. 4:

$$\begin{aligned} \mathbf{f} * \mathbf{g} &= \Phi \hat{\mathbf{G}} \left( \Phi^T \mathbf{f} \right) \\ &= \Phi \hat{\mathbf{G}} \left( \hat{f}_1, \dots, \hat{f}_n \right)^T \\ &= \Phi \left( \hat{g}_1 \hat{f}_1, \dots, \hat{g}_n \hat{f}_n \right)^T \\ &= \hat{g}_1 \hat{f}_1 \phi_1 + \dots + \hat{g}_n \hat{f}_n \phi_n \end{aligned} \quad (6)$$

From Eq. 6 it can be seen that a diagonal operator on the Fourier space modulates the smoothness of the signal.

## 2.2 Spectral Graph Convolutional Neural Networks

Graph spectral convolutions were first introduced as a means to define CNNs on graphs by [3].

Consider a feature matrix  $\mathbf{F}^{\text{in}} = \left( \mathbf{f}_1^{\text{in}}, \dots, \mathbf{f}_p^{\text{in}} \right) \in \mathbb{R}^{n \times p}$ , where  $\mathbf{f}_j \in \mathbb{R}^n$ ,  $j = 1, \dots, p$ , denotes a signal on the graph, and the  $i^{\text{th}}$  row of  $\mathbf{F}^{\text{in}}$  represents the  $p$ -dimensional feature vector of node  $i$ ,  $i = 1, \dots, n$ .

Then, a convolutional layer transforms an input feature matrix  $\mathbf{F}^{\text{in}} \in \mathbb{R}^{n \times p}$  into an output feature matrix  $\mathbf{F}^{\text{out}} = \left( \mathbf{f}_1^{\text{out}}, \dots, \mathbf{f}_{p'}^{\text{out}} \right) \in \mathbb{R}^{n \times p'}$  according to:

$$\mathbf{f}_{j'}^{\text{out}} = \sigma \left( \sum_{j=1}^p \Phi_k \hat{\mathbf{G}}^{(j',j)} \Phi_k^T \mathbf{f}_j^{\text{in}} \right), \quad j' = 1, \dots, p' \quad (7)$$

On Eq. 7,  $\sigma$  is a non-linearity,  $\Phi_k = (\phi_1, \dots, \phi_k) \in \mathbb{R}^{n \times k}$  is the matrix with the first  $k$  eigenvectors of the Laplacian, and  $\hat{\mathbf{G}}^{(j',j)} \in \mathbb{R}^{k \times k}$ ,  $j = 1, \dots, p$ ,  $j' = 1, \dots, p'$ , are diagonal matrices of learnable parameters. Therefore, such a spectral convolutional layer has  $pp'k$  trainable parameters.

On Eq. 7,  $\hat{\mathbf{G}}^{(j',j)}$  can be interpreted as a learnable filter in the frequency domain. Notice that this definition of a

convolutional layer is based on the representation of the spectral convolution given by Eq. 4.

One drawback of the layer described by Eq. 7 is that it relies on dense matrix multiplications by  $\Phi_k$  and  $\Phi_k^T$ , incurring a cost of  $\mathcal{O}(n^2)$ , assuming that  $k = \mathcal{O}(n)$ .

An alternative approach is to express the filtering of the signal  $\mathbf{f}$  as in Eq. 5, i.e., as  $\mathbf{G} \mathbf{f}$ , where  $\mathbf{G}$  is such that it can be written as:

$$\mathbf{G} = g(\Delta) = \Phi g(\Lambda) \Phi^T = \Phi \text{diag}(g(\lambda_1), \dots, g(\lambda_n)) \Phi^T \quad (8)$$

Depending on the choice of  $g$ , the matrix  $\mathbf{G}$  may be sparse. This characteristic can be exploited to reduce the computational complexity of the layer. Furthermore, this approach does not require the explicit computation of the Laplacian eigenvectors.

A convolutional layer is then defined as:

$$\mathbf{f}_{j'}^{\text{out}} = \sigma \left( \sum_{j=1}^p \mathbf{G}^{(j',j)} \mathbf{f}_j^{\text{in}} \right), \quad j' = 1, \dots, p' \quad (9)$$

## 3 CayleyNets

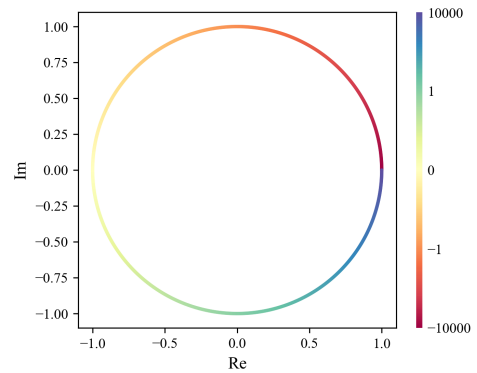
CayleyNets are based on the filtering approach described by Eq. 8 and consist on convolutional layers as described by Eq. 9, with  $\mathbf{G}$  a Cayley filter.

In order to define such a filter, one must first consider the Cayley transform:

$$C : \mathbb{R} \rightarrow e^{i\mathbb{R}} \setminus \{1\}, \quad C(x) = \frac{x - i}{x + i} \quad (10)$$

The Cayley transform defines a smooth bijection between  $\mathbb{R}$  and the complex unit circle, as illustrated on Fig. 1. This definition is extended to nonsingular real square matrices as:

$$C(\mathbf{X}) = (\mathbf{X} - i\mathbf{I})(\mathbf{X} + i\mathbf{I})^{-1} \quad (11)$$



**Figure 1.** Plot of the Cayley transform  $C(x)$ ,  $x \in \mathbb{R}$ , of the real line on the complex plane. The value of  $x$  can be identified via the color bar.

Applying the Cayley transform to the Laplacian  $\Delta$  scaled by a spectral zoom parameter  $h > 0$  and using the spectral

decomposition of  $\Delta$ , the following equality is obtained:

$$\begin{aligned} C(h\Delta) &= (h\Delta - i\mathbf{I})(h\Delta + i\mathbf{I})^{-1} \\ &= \Phi(h\Delta - i\mathbf{I})(h\Delta + i\mathbf{I})^{-1}\Phi^T = \Phi C(h\Delta)\Phi^T \end{aligned} \quad (12)$$

Eq. 12 will guarantee that the Cayley filter has the format of Eq. 8.

Since  $\Lambda$  is a diagonal matrix,  $C(h\Lambda) = \text{diag}(C(h\lambda_i))_{i=1}^n$ . Therefore, it follows from Eq. 12 that the complex matrix  $C(h\Delta)$  has its spectrum on  $e^{i\mathbb{R}}$ , therefore it is unitary. Furthermore, since  $\Delta \geq 0$  its eigenvalues are non-negative, thus the spectrum of  $C(h\Delta)$  lies on the lower-half of the complex unitary circle, as can be seen from Fig. 1.

Having defined the Cayley transform, the Cayley polynomial of order  $r$  is defined as:

$$g_{c,h}(\lambda) = c_0 + 2 \text{Re} \left\{ \sum_{k=1}^r c_k C(h\lambda)^k \right\} \quad (13)$$

where  $\mathbf{c} = (c_0, \dots, c_r) \in \mathbb{R} \times \mathbb{C}^r$  and  $h > 0$ . This definition is naturally extended to nonsingular real square matrices. In particular, the Cayley filter is the Cayley polynomial applied to the Laplacian:

$$\mathbf{G}_{c,h} = g_{c,h}(\Delta) = c_0 + 2 \text{Re} \left\{ \sum_{k=1}^r c_k C(h\Delta)^k \right\} \quad (14)$$

Eq. 12 guarantees that the Cayley filter is of the form of Eq. 8. A set  $(\mathbf{G}^{(j,j')})_{j,j'=1}^{p,p'}$  of Cayley filters then defines a convolutional layer as in Eq. 9, where the parameters  $\mathbf{c}^{(j,j')}$ , associated to each filter  $\mathbf{G}^{(j,j')}$ , and the parameter  $h$ , associated to the layer, are optimized during training. Such a layer has  $pp'(2r+1)+1$  trainable real parameters, considering that each complex parameter corresponds to two real parameters.

### 3.1 Properties of Cayley filters

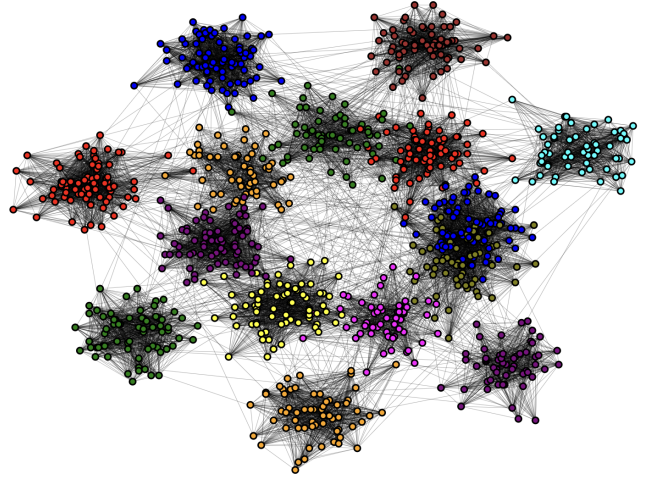
**3.1.1 Spectral zoom.** It is interesting to analyze the role of the spectral zoom parameter  $h$  in the Cayley filter.

A characteristic of the Cayley transform that will be relevant for the analysis is that the mapping of the interval  $[0, 1]$  is spread out on the lower left quadrant of the unit circle. On the other hand,  $[1, +\infty)$  is mapped to the lower right quadrant, resulting in the concentration close to the point  $(1, 0)$  of the mapping of large values. This property can be deduced from Fig. 1.

Recalling Eq. 12, where  $C(h\Lambda) = \text{diag}(C(h\lambda_i))_{i=1}^n$ , it is seen that  $h$  has a direct influence on the spectrum of  $C(h\Delta)$  and amounts to “zooming” into different parts of the spectrum. Small values of  $h$  will better spread apart the larger eigenvalues of  $\Delta$  that are crushed by the Cayley transform, whereas large values of  $h$  will zoom into the smaller eigenvalues.

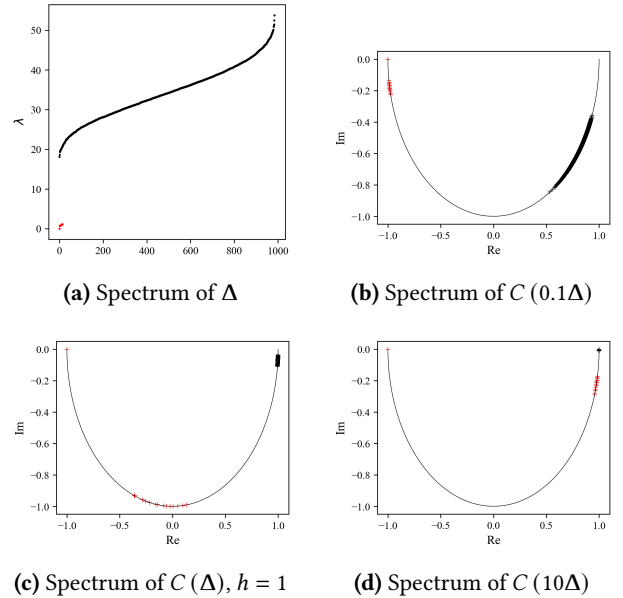
This behavior is illustrated for a synthetic 15-communities graph, represented on Fig. 2. On Fig. 3, the eigenvalues of the graph Laplacian are plotted, with the 15 smallest eigenvalues highlighted in red, the eigenvalues of the Cayley transform

of the scaled graph Laplacian are plotted as well for different values of  $h$ .



**Figure 2.** Synthetic 15-communities graph. The graph was created by assigning each node to a community uniformly at random; then an edge was added for each pair of nodes with probability  $p = 0.5$  for nodes belonging to the same community and probability  $q = 10^{-3}$  for nodes belonging to different communities.

**3.1.2 Localization.** Another interesting property of Cayley filters is that they are not strictly truncated to  $r$ -hop



**Figure 3.** (a) Eigenvalues of the unnormalized Laplacian  $\Delta$  of the 15-communities graph of Fig. 2. (b,c,d) Eigenvalues of  $C(h\Delta)$  for  $h = 0.1, 1$  and  $10$ , respectively. The 15 smallest eigenvalues are marked in red.

neighborhoods, unlike polynomial filters. This gives rise to an alternative notion of localization related to the exponential decay of  $G\delta_m$ , where  $\delta_m$  denotes a delta-function on the graph. In particular, the following result was proven by the authors in the article [6]:

$$\|G\delta_m|_{\mathcal{N}_{k,m}^c}\|_2 \leq 4M\mathcal{K}^{k/r} \quad (15)$$

where:

- $\mathcal{N}_{k,m}$  is the  $k$ -hop neighborhood of  $m$ .
- $M = \sqrt{n} \sum_{j=1}^r |c_j|$
- $\mathcal{K} = \frac{hd}{\sqrt{h^2d^2+1}}$  and  $d = \max_j \{d_{i,j}\}$ , i.e. maximum node degree.

This way, while the support is the entire graph we can still define a notion of localization, which is particularly pertinent considering that having well-localized filters on the space domain is an important aspect of convolutional kernels.

**3.1.3 Cayley filters are complete on the space of spectral filters.** Another of the motivating factors for the use of Cayley filters is that any spectral filter can be formulated as one. This can be directly seen through a reformulation of Eq. 14.

We will use the fact that  $2\operatorname{Re}\{C(h\Delta)\} = C(h\Delta) + \overline{C(h\Delta)}$ , indeed this is true for any complex valued matrix. However, the spectrum of  $C(h\Delta)$  lies in  $e^{i\mathbb{R}}$ , therefore  $\overline{C(h\Delta)} = C(h\Delta)^{-1}$ . Thus Eq. 14 can be reformulated as the following Laurent polynomial:

$$\begin{aligned} G_{c,h}(\Delta) &= g_{c,h}(\Delta) = c_0 + \sum_{k=1}^r c_k C(h\Delta)^k + \overline{c_k} \overline{C(h\Delta)}^k \\ &= c_0 + \sum_{k=1}^r c_k C(h\Delta)^k + \overline{c_k} C(h\Delta)^{-k} \end{aligned} \quad (16)$$

Let us analyze each term in the summation of Eq. 16 in terms of the spectral decomposition of  $\Delta$ :

$$\begin{aligned} p_k(h\Delta) &:= c_k C(h\Delta)^k + \overline{c_k} C(h\Delta)^{-k} \\ &= \Phi \left( c_k C(h\Delta)^k + \overline{c_k} C(h\Delta)^{-k} \right) \Phi^T \end{aligned} \quad (17)$$

In Eq. 17,  $C(h\Delta)^k = \operatorname{diag}(C(h\lambda_1)^k, \dots, C(h\lambda_n)^k)$  and  $C(h\lambda_p)^k \in e^{i\mathbb{R}}$ ,  $p = 1, \dots, n$ . Therefore, the power operator can be interpreted as a multiplication by a pure harmonic in the frequency domain.

With that interpretation in mind, let  $\omega_p \in [0, 2\pi)$  be the angular frequency associated to the eigenvalue  $\lambda_p$  by the Cayley transform, i.e.  $C(h\lambda_p) = e^{i\omega_p}$ . Then, from Eq. 17, we have:

$$\begin{aligned} p_k(h\lambda_p) &= c_k e^{i\omega_p k} + \overline{c_k} e^{-i\omega_p k} \\ &= 2\operatorname{Re}\{c_k\} \cos(\omega_p k) - 2\operatorname{Im}\{c_k\} \sin(\omega_p k) \end{aligned} \quad (18)$$

Therefore we conclude that  $G_{c,h}$  is a real-valued trigonometric polynomial. This implies that any spectral filter  $g(\Delta)$  can be written as a Cayley filter, since the values  $g(\lambda_1), \dots, g(\lambda_n)$

that define it can be interpolated by a trigonometric polynomial. The implications of this property of Cayley filters will be analysed on Section 4.

### 3.2 Computational method

CayleyNets consist in convolutional layers that transform an input feature matrix  $F^{\text{in}} \in \mathbb{R}^{n \times p}$  into an output feature matrix  $F^{\text{out}} \in \mathbb{R}^{n \times p'}$  according to Eq. 9, where  $G^{(j',j)}$  is a Cayley filter, defined in Eq. 14.

Therefore, the forward pass of the network relies on computing  $C(h\Delta)^k f$ ,  $k = 1, \dots, r$ . This can be done in an iterative fashion:

$$\begin{aligned} y_0 &:= f \\ y_k &:= C(h\Delta)^k f = C(h\Delta) y_{k-1}, \quad k = 1, \dots, r \end{aligned} \quad (19)$$

Recalling the definition of the Cayley transform,  $C(h\Delta) = (h\Delta - iI)(h\Delta + iI)^{-1}$ , therefore in order to avoid computing inverses, which has a computational complexity of  $O(n^3)$ ,  $y_0, \dots, y_r$  can be expressed instead as solutions to the following systems of linear equations:

$$\begin{aligned} y_0 &= f \\ (h\Delta + iI) y_k &= (h\Delta - iI) y_{k-1}, \quad k = 1, \dots, r \end{aligned} \quad (20)$$

Jacobi's iterative method can then be used to provide approximate solutions to the systems of linear equations of Eq. 20. Running the method with  $K$  iterations, the computational cost of computing each  $y_j$  is  $O(Kn^2)$ . However, if the sparsity of the Laplacian is exploited, the computational cost of the method is reduced to  $O(Kn)$ , assuming that each row of the Laplacian has  $O(1)$  non zero entries. Therefore the overall computational cost of a layer is  $O(Krn)$ .

In practice, an approximate solution is sufficient because the approximation error of Jacobi's method can be compensated by the learned parameters.

## 4 Experiments

The code associated to the original CayleyNets article [6] is incompatible with current versions of Python frameworks for Machine Learning and was proven difficult to adapt. Therefore, we proposed our own implementation of CayleyNets using Pytorch Geometric<sup>1</sup>. Even though our implementation strictly follows the article's formulation of CayleyNets, our results do not exactly match those of [6]. Due to this discrepancy, considerable effort was dedicated to testing the different components of the code and implementing possible variations, however no implementation errors were identified. This difficulty highlights the importance of making available a well documented source code, in order to guarantee the reproducibility of a scientific publication.

In this section, the CayleyNets will be evaluated on a node classification task, using the CORA dataset. All experiments were conducted on a machine with a 3.00GHz Intel Xeon

<sup>1</sup>The code is available at <https://github.com/AliceDeLorenci/CayleyNets>



Gold 6136 CPU, 187GB of RAM, and a NVIDIA Tesla V100 GPU with 16GB of RAM.

#### 4.1 CORA dataset

The CORA dataset [7] consists in a graph of citations between 2708 scientific publications, with 5429 edges. Each node is associated to a 1433 dimensional binary feature vector, that indicates the presence of the corresponding word on the scientific publication; and is classified into one of seven classes, corresponding to the category of the paper. This dataset is a common benchmark in the Graph Machine Learning community, therefore we chose to use it to evaluate the CayleyNets. An alternative formulation of Spectral Graph Neural Network known as ChebNets was used as a baseline, this model is described on Appendix A.

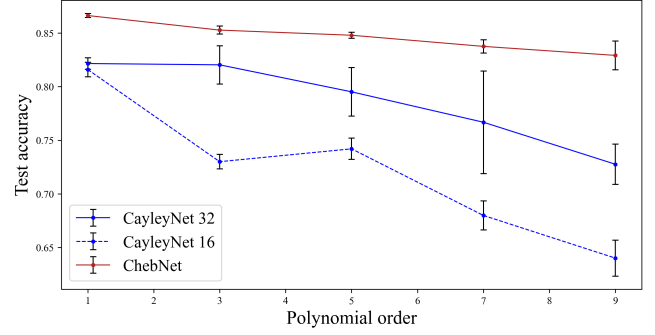
The problem consists in a node classification task in which the labels are given for only a subset of the nodes (training data), and the model must predict the labels for the remaining nodes (test and validation data). For the numerical experiments the CORA dataset was split into 1708 training nodes, 500 validation nodes and 500 test nodes.

#### 4.2 Impact of the polynomial order

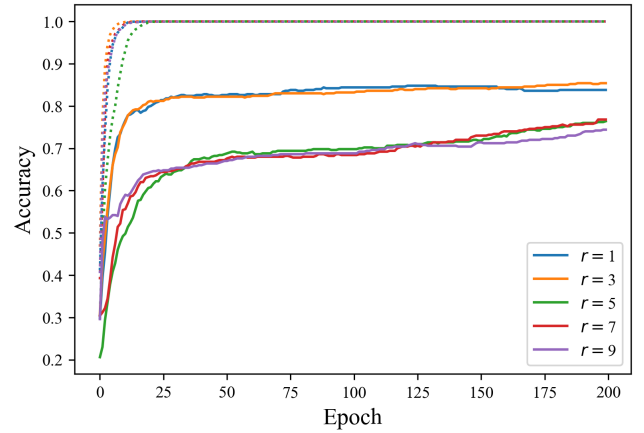
First of all, we analyzed the impact of the polynomial order  $r$  of the spectral filter on the performance of the network. This was done by measuring the test accuracy for  $r = 1, 3, 5, 7, 9$ . For that, the network architecture proposed by [6] was used, i.e., CayleyNets and ChebNets with one hidden layer with 32 features; the hidden layer uses ReLu non-linearity and the output layer, softmax non-linearity. For both networks the normalized Laplacian was used for the spectral filters. The models were trained for 200 epochs using the Adam optimization method, with weight decay  $\lambda = 5 \cdot 10^{-4}$  and learning rate  $\gamma = 5 \cdot 10^{-3}$ , as well as dropout regularization with 0.5 probability. The results are summarized in Fig. 4.

Analyzing Figure 4 we observe that ChebNets consistently outperform CayleyNets, for all tested polynomial degrees. Furthermore, we observe a strong degradation of the test accuracy of CayleyNets for increasing polynomial order  $r$ . This might be an indicator of overfitting. Indeed, the number of parameters for a CayleyNet layer is  $pp'(2r+1) + 1 = O(pp'r)$ , where  $p$  is the number of input features and  $p'$  the number of output features. Therefore the resulting total number of parameters for a CayleyNet with one 32-nodes hidden layer is of the order of 70K, 160K, 250K, 350K and 440K parameters for  $r = 1, 3, 5, 7, 9$ , respectively. In order to verify the overfitting hypothesis, we plotted the train and validation accuracy over the training epochs for the CayleyNet, with different polynomial orders. The results are illustrated on Fig. 5, and confirm that the network has overfitted to the training data.

In an attempt to avoid overfitting while preserving the expressive power of filters of higher polynomial order, we repeated the same experiment, but halving the number of



**Figure 4.** Evaluation of CayleyNets and ChebNets for a node classification task on the CORA dataset for different polynomial orders  $r$  ( $r = 1, 3, 5, 7, 9$ ). Each experiment was repeated 10 times, the mean test accuracy is plotted in red for ChebNets with 32 hidden nodes, in solid blue line for CayleyNets with 32 hidden nodes and in dashed blue line for CayleyNets with 16 hidden nodes. The corresponding standard deviations are displayed as error bars.



**Figure 5.** Train and validation accuracies over the training epochs for CayleyNets with 32 hidden nodes for different polynomial orders  $r$  ( $r = 1, 3, 5, 7, 9$ ). The train and validation accuracies are plotted in dotted and solid lines, respectively, and the polynomial order can be identified via the color of the plots.

hidden nodes of the CayleyNet. The results are also illustrated on Fig. 4. It can be seen that reducing the size of the hidden layer negatively impacted the performance of the network.

#### 4.3 Comparison of different methods

On Section 4.2, CayleyNets were compared to ChebNets, an alternative Spectral Graph Convolutional Neural Network model. However, it is interesting to compare them as well to simpler Graph Neural Network (GNN) models, such as message passing GNNs.

**Table 1.** Test accuracy obtained with different methods on the CORA dataset.

Model	Order	Accuracy	#Parameters
CayleyNet	1	$0.816 \pm 0.006$	138K
ChebyNet	1	$0.872 \pm 0.004$	92K
GNN	-	$0.812 \pm 0.10$	46K

In this context, we evaluated the performance of a standard message passing GNN with one hidden layer with 32 hidden nodes and ReLU activation function, followed by a fully connected layer with 7 output units and softmax activation, mimicking the architecture used for CayleyNets and ChebNets. The training setup is the same as that described on Section 4.2, in particular the same hyperparameters were used.

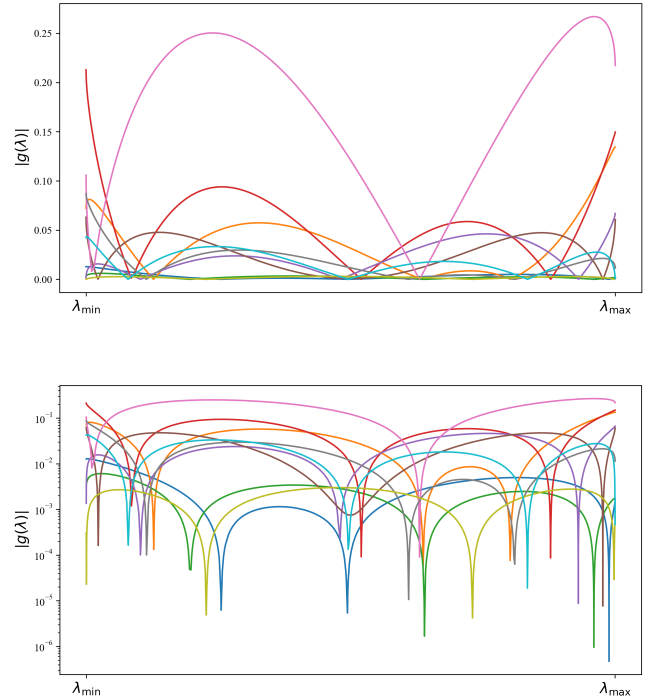
The experiment was repeated 10 times, and the mean test accuracy and standard deviation are displayed in Table 1. In this table we display as well the best mean test accuracy and standard deviation obtained on Section 4.2 for Cayleynets and ChebNets, i.e., using an architecture of 32 hidden nodes and polynomial order  $r = 1$ . It can be seen that the performance of CayleyNets is comparable to that of much simpler message passing networks.

#### 4.4 Cayley filter response

Another aspect of CayleyNets that is interesting to analyze is the spectral response of the learned filters. Recalling the definition of Eq. 9 of a spectral convolutional layer and the expression of Eq. 14 for the Cayley filter, we see that each layer corresponds to  $pp'$  learned filters, where  $p$  denotes the number of input features and  $p'$  the number of output features.

In order to visualize the response of the filters, we retrieved the filter parameters from the trained weights of a CayleyNet trained on the CORA dataset. More specifically, the same experimental setup of Section 4.2 was used and the chosen architecture was a CayleyNet with a hidden layer of 32 nodes and polynomial order  $r = 3$ .

The spectral response of ten filters randomly chosen among all of the learned filters is plotted in Fig. 6. From the figure we observe that the network is able to learn a variety of spectral responses, spanning both narrow band and wide band filters, and focusing on different frequency bands. This variability is a result of the characteristics of Cayley filters discussed on Sections 3.1.1 and 3.1.3, i.e., the ability of focusing on different parts of the spectrum due to the spectral zoom parameter and the property that Cayley filters can approximate any spectral filter.



**Figure 6.** Responses of ten different filters learned by a CayleyNet, in linear (top) and log (bottom). The scaling of the abscissa was obtained by unrolling to the real line the Cayley transform of the Laplacian eigenvalues.

## 5 Conclusion

CayleyNets are a type of Spectral Convolutional Graph Neural Networks based on a sophisticated spectral filter formulation, with sound theoretical motivations. However, in spite of the theoretical properties of Cayley filters, in practice the performance of CayleyNets is inferior to that of simpler architectures, such as ChebNets, or even message passing Graph Neural Networks.

Therefore the applicability of CayleyNets is put into doubt, specially when one considers their difficult implementation and high computational cost. With respect to the first point, CayleyNets notably rely on complex parameters, which are not natively supported by Python Machine Learning frameworks, demanding ad hoc optimization of the code. Regarding the second point, the computational complexity of a CayleyNet is higher than that of an equivalent ChebNet by a factor of  $K$ , where  $K$  is the chosen number of Jacobi iterations.

However, we emphasize that our considerations regarding the performance of the networks are based on experiments carried out using the CORA dataset. A future directions of work could involve testing CayleyNets on different tasks and benchmarks, for instance, on graph classification tasks

instead of node classification ones, to determine the domain of applicability of those networks.

## References

- [1] Thomas Bonald. 2022–2023. Graph Mining. Course Notes. Course code: SD212.
- [2] Michael M. Bronstein, Joan Bruna, Taco Cohen, and Petar Veličković. 2021. *Geometric Deep Learning: Grids, Groups, Graphs, Geodesics, and Gauges*. <https://doi.org/10.48550/arXiv.2104.13478> arXiv:2104.13478 [cs.LG] Submitted on 27 Apr 2021 (v1), last revised 2 May 2021 (this version, v2).
- [3] Joan Bruna, Wojciech Zaremba, Arthur Szlam, and Yann LeCun. 2014. Spectral Networks and Locally Connected Networks on Graphs. *International Conference on Learning Representations (ICLR)* (April 2014).
- [4] Xinye Chen. 2020. Understanding Spectral Graph Neural Network. *arXiv preprint arXiv:2012.06660* (2020). arXiv:2012.06660
- [5] Michaël Defferrard, Xavier Bresson, and Pierre Vandergheynst. 2016. Convolutional Neural Networks on Graphs with Fast Localized Spectral Filtering. *CoRR* abs/1606.09375 (2016). arXiv:1606.09375 <http://arxiv.org/abs/1606.09375>
- [6] Ron Levie, Federico Monti, Xavier Bresson, and Michael M. Bronstein. 2017. CayleyNets: Graph Convolutional Neural Networks with Complex Rational Spectral Filters. *CoRR* abs/1705.07664 (2017). arXiv:1705.07664 <http://arxiv.org/abs/1705.07664>
- [7] Telecom Paris NETSET. 2023. CORA. <https://netset.telecom-paris.fr/pages/cora.html>

## A ChebNets

ChebNets are a type of Spectral Graph Convolutional Neural Networks introduced by [5]. Analogously to CayleyNets, they adopt the approach of Eq. 9, however, the filter they propose is based on the Chebyshev polynomials instead of Cayley polynomials.

A Chebyshev polynomial has the form:

$$g_{\alpha}(\tilde{\lambda}) = \sum_{k=0}^r \alpha_k T_k(\tilde{\lambda}) \quad (21)$$

On Eq. 21,  $\alpha \in \mathbb{R}^{r+1}$  are the trainable parameters and the functions  $T_k : [-1, 1] \rightarrow \mathbb{R}$  are recursively defined:

$$\begin{aligned} T_0(\tilde{\lambda}) &= 1 \\ T_1(\tilde{\lambda}) &= \tilde{\lambda} \\ T_k(\tilde{\lambda}) &= 2\tilde{\lambda}T_{k-1}(\tilde{\lambda}) - T_{k-2}(\tilde{\lambda}), \quad k > 1 \end{aligned} \quad (22)$$

Since the functions  $T_k$  are defined on the interval  $[-1, 1]$ , the Chebyshev polynomial must be evaluated on the rescaled Laplacian  $\tilde{\Delta} = 2\lambda_{\max}^{-1}\Delta - \mathbf{I}$ , whose eigenvalues lie on  $[-1, 1]$ . The resulting filter is then:

$$\mathbf{Gf} = g_{\alpha}(\tilde{\Delta})\mathbf{f} = \sum_{k=0}^r \alpha_k T_k(\tilde{\Delta})\mathbf{f} \quad (23)$$

From the formulas for  $T_k$ , we see that  $g_{\alpha}(\tilde{\Delta})$  is a polynomial of order  $r$  of the scaled Laplacian, which implies that the filter is an operator that affects only the  $r$ -hops neighborhood of a node.

Leveraging the recurrent formulation of  $T_k$ , the output of the filter can be computed by setting:

$$\begin{aligned} \tilde{\mathbf{f}}_0 &= \mathbf{f} \\ \tilde{\mathbf{f}}_1 &= \tilde{\Delta}\mathbf{f} \\ \tilde{\mathbf{f}}_k &= 2\tilde{\Delta}\tilde{\mathbf{f}}_{k-1} - \tilde{\mathbf{f}}_{k-2}, \quad k > 1 \end{aligned} \quad (24)$$

The computational complexity of the filter is thus  $\mathcal{O}(rn^2)$ . However, exploiting the sparsity of the Laplacian the complexity can be reduced to  $\mathcal{O}(rn)$ , if we assume that the nodes have  $\mathcal{O}(1)$  neighbors on average.

## B Online Resources

The source code is available at <https://github.com/AliceDeLorenci/CayleyNets>.

Received 10 December 2023