

A graph is denoted by the pair  $(V, E)$ , where  $V$  denotes the set of nodes and  $E$  the set of edges. The connections between vertices can be structured into an "adjacency" matrix  $A$  s.t.  $A_{ij} = \begin{cases} 1, & \text{if } (i, j) \in E \\ 0, & \text{if } (i, j) \notin E \end{cases}$

Often, we are interested in functions  $f: V \rightarrow \mathbb{R}$ .

In general, when we have a finite set  $V$ , the set of functions  $f: V \rightarrow \mathbb{R}$  forms a  $\dim(V)$  dimensional vector space with basis the functions  $\delta_e(x) = \begin{cases} 1, & \text{if } x = e \\ 0, & \text{otherwise} \end{cases}$

Hence, a function is determined by how it acts on every node  $x \in V$ .

$$f \leftrightarrow \left( \begin{array}{c} f(x_1) \\ f(x_2) \\ \vdots \\ f(x_m) \end{array} \right)$$

$x_1, \dots, x_m$  are all nodes in  $V$  and  $m = \dim(V)$

$f(x_m)$

This is an isomorphism between  $(\mathbb{R}^n)$  and the vector space of functions  $f: V \rightarrow \mathbb{R}$ .

Therefore, we can think of  $f$  as a vector in  $(\mathbb{R}^n)$ .

There is an important graph operator that acts on  $f$  called the Laplacian. It is essentially the discrete or finite spaced version of the infinitesimal Laplacian  $\nabla^2$ .

$$\text{so, } (\mathcal{L}f)(x) = \sum_{\substack{x,y \\ (x,y) \in E}} A_{x,y} (f(x) - f(y)) \quad (1)$$

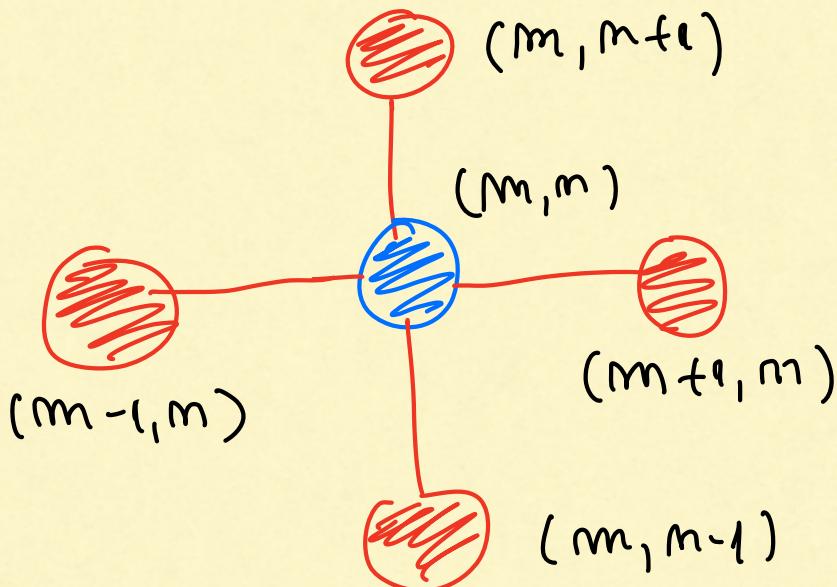
is the reversible operator that could be defined and whose generalization to infinitesimal spacing is the usual Laplacian. We can see this by analysing a grid of points, each having four neighbours with  $a_{x,y} = \frac{1}{(\delta x)^2}$  as a weight.  $\delta x$  being the distance between points on the evenly spaced grid.

$$(1) \quad - (l_1 l_{m-1, m} + l_{m+1, m} - l_{m-1, m})$$

$$(\Delta f)_{m,n} = (f_{m,n} - f_{m+1,n} - f_{m,n-1} + f_{m,n+1}) / (\Delta x)^2$$

which is the approximation of  $\nabla^2 f$ .

For getting (1) we see  $L = D - A$  is the operator, where  $D$  is a diagonal matrix with  $D_{ii} = \sum_{j \in V} A_{ij}$



On  $\mathbb{R}$ ,  $e^{i\omega x}$  are eigenfunctions of the 1D Laplacian operator  $\frac{d}{dx^2}$ . So,  $f(x)$  is written as a sum of these eigenfunctions or linear combination in the finite case. Hence, in analogy it is justified to write the Fourier transform of  $f$  (now back to graphs;  $f$  above was on  $\mathbb{R}$ ) as a linear combination of the

eigenvalues of  $\hat{f}$ .

$$\hat{f} X_i = \lambda_i X_i \quad i \in \overline{1, m}$$

As  $\hat{f}$  is symmetric  $\Rightarrow \lambda_m$  is real.

$$\text{So, } \hat{f}(m) = \langle X_m, f \rangle = \sum_{n=1}^N X_m^* x_n^{(n)} f(n)$$

Fourier transform

From now I will write  $N$  as  $\dim(V)$  instead of  $m$ , so that I can use  $m$  in the sum.

$$\text{The inverse is thus } f(m) = \sum_{n \in V} \hat{f}(l) x_l(m).$$

We would like to have the same interpretation as in CNN by considering  $K * f$  (the convolution between a kernel  $K$  and a function  $f$  in  $\mathbb{R}^N$ ). What should be the convolution in the current context?

If we think of  $K : \mathbb{R}^+ \rightarrow \mathbb{R}^+$  as a positive

function of only one variable, then we might propose that  $\kappa$  should act on the eigenvalues of  $f$ , i.e.  $K(\lambda_\ell)$  because  $K$  should somehow behave as a filter by which a strength or weights of some of the basis elements in the Fourier transform change. I want to emphasize that this convolution looks different than in CNN's, where we translated inside integrals, here translation is meaningless, since we have only coefficients in  $\mathbb{N}$ , but not in  $\mathbb{Z}$ . Therefore, we are forced to disregard some of the properties since the current setting is different. Furthermore, we want  $K * f$  to be as in CNN's a new function, hence

$$(K * f)(m) = \sum_{l=1}^N K(\lambda_\ell) \hat{f}(l) \chi_\ell(m)$$

is a good candidate for the operation if we combine all the observations above. As you can see  $K * f$  is the function where

Fourier coefficients are  $K(\lambda_l) \hat{f}(l)$ , a scaled version of  $\hat{f}(l)$ . Since  $K$  has a continuous domain we can require  $\hat{f}$  to be a continuous function and also restrict its domain to  $[0, \lambda_N]$  since  $L$  is positive semi-definite with positive eigenvalues and we can relabel n.t.  $\lambda_N$  is the biggest. Using Chebyshev polynomials to approximate  $K$  as a minimax approximation

$$\min_{p \in P_m} \sup_{x \in [0, \lambda_N]} |K(x) - p(x)|.$$

$P_m$  - the space of polynomials of degree  $\leq m$ .

These approximations are indeed essential because it is computationally expensive to compute eigenvectors and eigenvalues. Hence,

$$(K * f)(m) = \sum_{l=1}^N \sum_{i=1}^L \theta_i T_i(\lambda_l) \hat{f}(l) \chi_l^{(m)}$$

There is because  $T_i$  form an orthonormal basis

for the Hilbert space  $L^2([0, \lambda_N], d\mu)$ , where  $d\mu$  is a measure scaled from  $L^2([-1, 1], \frac{dx}{\sqrt{1-x^2}})$  on which  $T_i$  are usually defined.

$\Psi_m(z) \in$  is the maximum degree of  $T_i$  ( $\deg T_i = i$ ). If we sum with  $l \rightarrow \infty$  then we get exactly  $K$ , but instead we will approximate by truncating. Also,  $\theta_i$  are the coefficients corresponding to this linear combination ( $T_i$  forms a basis as I said above).

Now, we have neural networks that are universal function approximators, so we are not limited to only Chebyshev poly.

The order of polynomials increases if we stack multiple operations  $K * (K * \dots * f)$  with  $l = 1$  (linear model). Therefore,

$$(K * f)(m) \approx \sum_{l=1}^N (\theta_0 T_0(\lambda_l) + \theta_1 T_1(\lambda_l)) f(x_{\theta(l)})$$

$$= \sum_{l=1} (\theta_0 + \theta_1 \lambda_e) f(l) \lambda_e^{(1..l)}$$

Which in matrix notation can be written

$$\text{as } \theta (\mathbb{I}_m + \mathcal{L}) f \quad (3)$$

If we replace  $\mathcal{L} \rightarrow \mathcal{L} - \mathbb{I}_m := \tilde{\mathcal{L}}$

$\Rightarrow \theta \tilde{\mathcal{L}} f$ , where  $\theta$  is a vector now.

For our final architecture with input  $x \in \mathbb{R}^{m \times c}$ , where  $c$  is the dimension of the feature space (for example in a graph with papers we can embed the key words in vectors, so every node is associated an embedding or feature), we obtain

$$H^{(l+1)} = \phi(\tilde{\mathcal{L}} H^{(l)} w) \text{ as the}$$

activation in layer  $l+1$  with nonlinearity  $\phi$  and weight matrix  $w$ . One consequence of this construction is the permutation equivariance that our graph should have.

Under a permutation of the nodes, either  
locally or on a whole graph:

$$\hat{L} \rightarrow P \hat{L} P^T$$

$$H^{(l)} \rightarrow P H^{(l)}$$

, where  $P$  is a permutation  
matrix.

We see that  $H^{(l+1)} \rightarrow P H^{(l+1)}$  as desired.

To motivate in (3) why I replaced the  
eigenvalue part with  $\hat{L}$  is because the transformation  
can be written as  $U \Delta U^T$  (the eigendecomposition  
of  $\hat{L}$ ), which is just  $\hat{L}$ .