

Нейронные сети на графах

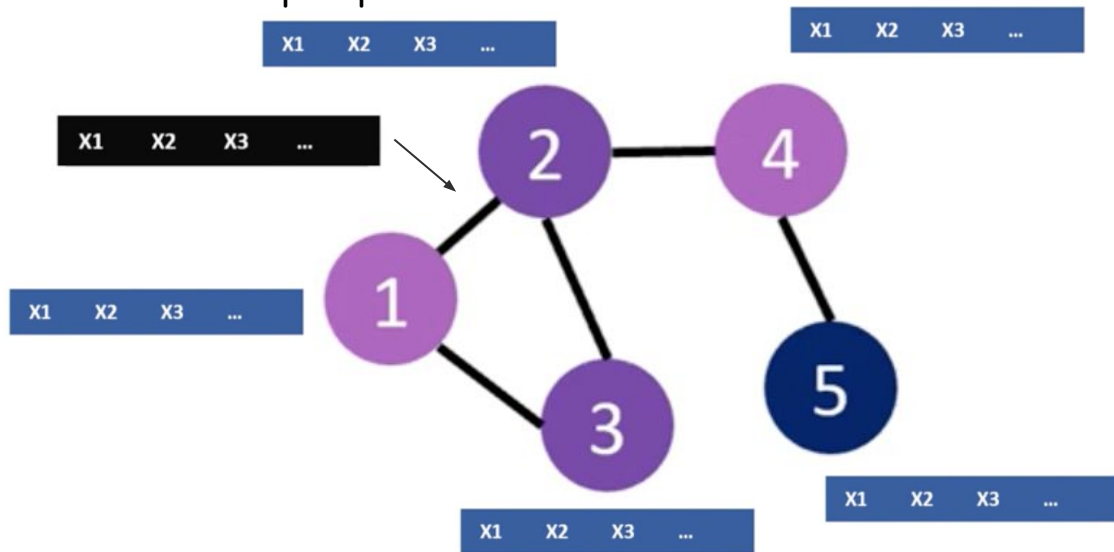
Алекберов Артём, Седашов Данила, 181

Часть 1.

Определения и задачи

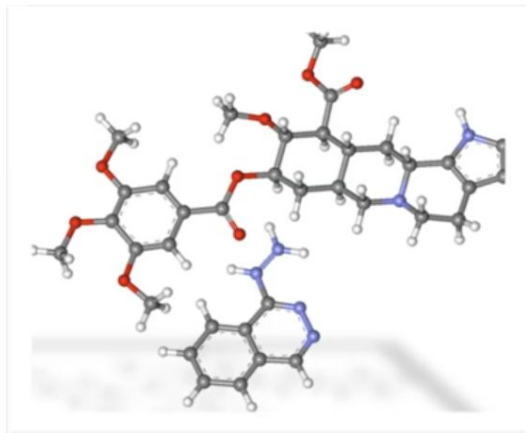
Что такое граф?

- Набор вершин и рёбер $G = (V, E)$
- В памяти представляется матрицей смежности размера $V \times V$
- У каждой вершины и ребра есть набор признаков



	V1	V2	...
V1	0	1	...
V2	1	0	...
V3	1	1	...
...

Примеры графовых данных



Химия / медицина



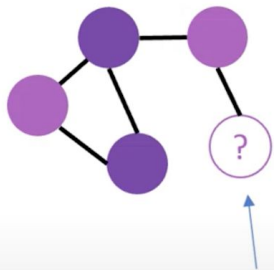
Рекомендательные
системы



Социальные сети

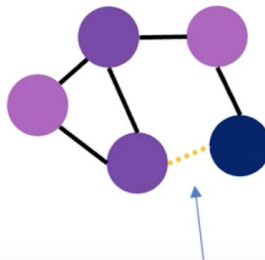
Обзор задач

Node-level predictions



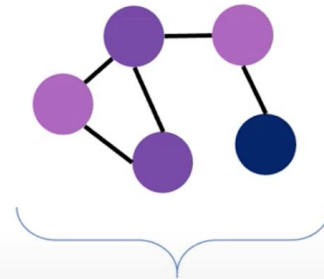
Курит ли этот человек?

Edge-level predictions
(Link prediction)



Следующее видео на YouTube?

Graph-level predictions



Подходит ли эта молекула для лекарства?

Часть 2.

Рекуррентные нейронные сети на графах

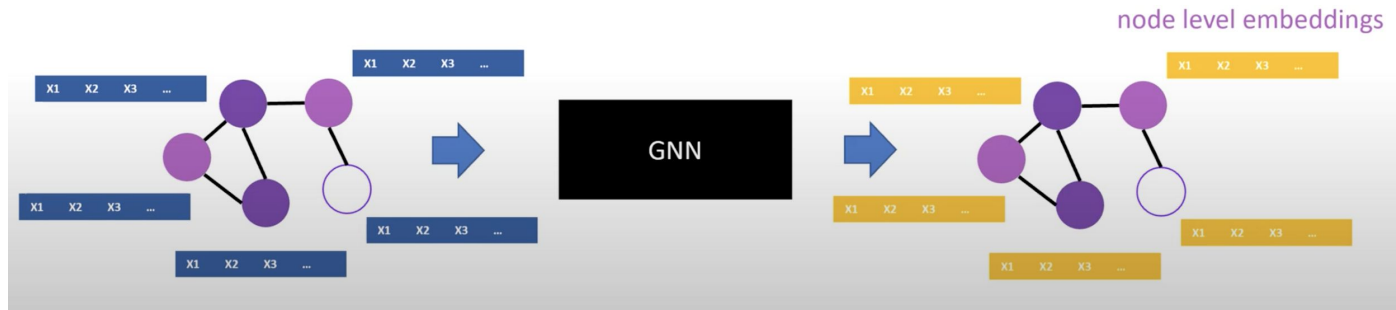
Проблемы графовых данных

- Размерности входов могут меняться (количество вершин)
- Существуют изоморфные графы => модель должна быть инвариантна к ним
- Графы не лежат в евклидовом пространстве

Проблемы графовых данных

- Размерности входов могут меняться (количество вершин)
- Существуют изоморфные графы => модель должна быть инвариантна к ним
- Графы не лежат в евклидовом пространстве

Решение: будем обучать эмбединги для каждой вершины и потом при необходимости агрегировать их в эмбединг графа



Общий алгоритм

Пусть $x_v \in \mathbb{R}^d$ — вектор признаков вершины v ,

$x_{(v,u)}^e \in \mathbb{R}^c$ — вектор признаков ребра, соединяющего вершины u , v .

Также пусть $N(v)$ — множество соседей вершины v .

Общий алгоритм

Пусть $\mathbf{x}_v \in \mathbb{R}^d$ — вектор признаков вершины v ,

$\mathbf{x}_{(v,u)}^e \in \mathbb{R}^c$ — вектор признаков ребра, соединяющего вершины u , v .

Также пусть $N(v)$ — множество соседей вершины v .

Будем обновлять эмбединги вершин итеративно.

Пусть $\mathbf{h}_v^{(t)}$ — эмбединг вершины v на t -ой итерации.

Общий алгоритм обновления эмбедингов: для каждой из вершин делаем преобразование:

$$\mathbf{h}_v^{(t)} = \sum_{u \in N(v)} f(\mathbf{x}_v, \mathbf{x}_{(v,u)}^e, \mathbf{x}_u, \mathbf{h}_u^{(t-1)})$$

пока не достигнем точки равновесия ($f(\cdot)$ — параметрическая ф-ия)

Теорема Банаха о неподвижной точке

Пусть (X, d) — непустое метрическое пространство,

$T : X \rightarrow X$ — сжимающее отображение

(то есть $\exists 0 \leq \alpha < 1 : d(Tx, Ty) \leq \alpha d(x, y) \forall x, y \in X$)

Теорема Банаха о неподвижной точке

Пусть (X, d) — непустое метрическое пространство,

$T : X \rightarrow X$ — сжимающее отображение

(то есть $\exists 0 \leq \alpha < 1 : d(Tx, Ty) \leq \alpha d(x, y) \forall x, y \in X$)

Тогда:

- 1) существует и единственна неподвижная точка $x^* \in X : Tx^* = x^*$
- 2) итерационная последовательность x, Tx, T^2x, \dots сходится к этой точке с экспоненциальной скоростью

Теорема Банаха о неподвижной точке

Пусть (X, d) — непустое метрическое пространство,

$T : X \rightarrow X$ — сжимающее отображение

(то есть $\exists 0 \leq \alpha < 1 : d(Tx, Ty) \leq \alpha d(x, y) \forall x, y \in X$)

Тогда:

- 1) существует и единственна неподвижная точка $x^* \in X : Tx^* = x^*$
- 2) итерационная последовательность x, Tx, T^2x, \dots сходится к этой точке с экспоненциальной скоростью

Пусть F_w — глобальная функция преобразования (которая переводит множество всех эмбедингов в себя)

Таким образом, при условии, что F_w — сжимающее отображение, алгоритм обновления эмбедингов сойдётся

Общий алгоритм: продолжение

- Введём также параметрическую функцию $g(\cdot)$ которая принимает на вход эмбединги и формирует ответ
- Через эту функцию можем определить функционал ошибки \mathcal{L} (например, MSE для регрессии или кросс-энтропию для классификации)

Общий алгоритм: продолжение

- Введём также параметрическую функцию $g(\cdot)$ которая принимает на вход эмбединги и формирует ответ
- Через эту функцию можем определить функционал ошибки \mathcal{L} (например, MSE для регрессии или кросс-энтропию для классификации)
- Таким образом, алгоритм обучения следующий:
 - итеративно обновляем $h_v^{(t)}$ для каждой вершины, пока не достигнем точки равновесия ($h_v^{(0)}$ инициализируем случайно)
 - считаем градиент функционала по параметрам f и g : $\frac{\partial \mathcal{L}}{\partial w}$
 - делаем шаг градиентного спуска

Реализации f и g

- На g не накладывается никаких ограничений $\Rightarrow g = \text{MLP}$
- F_w должна быть сжимающим отображением

Linear GNN

- Будем считать $f(x_v, x_{(v,u)}^e, x_u, h_u^{(t-1)}) = A_{n,u} h_u^{(t-1)} + b_n$
где $A_{n,u} \in \mathbb{R}^{s \times s}, b_n \in \mathbb{R}^s$ (s — размер эмбединга)

Linear GNN

- Будем считать $f(x_v, x_{(v,u)}^e, x_u, h_u^{(t-1)}) = A_{n,u} h_u^{(t-1)} + b_n$
где $A_{n,u} \in \mathbb{R}^{s \times s}, b_n \in \mathbb{R}^s$ (s — размер эмбединга)
- $A_{n,u}, b_n$ генерируются при помощи двух MLP
- $A_{n,u} = \frac{\mu}{s \cdot |N(u)|} \cdot \text{resize}(MLP_1(x_v, x_{(v,u)}^e), x_u), \mu \in (0, 1)$
- $b_n = MLP_2(x_v)$

Linear GNN

- Будем считать $f(x_v, x_{(v,u)}^e, x_u, h_u^{(t-1)}) = A_{n,u} h_u^{(t-1)} + b_n$
где $A_{n,u} \in \mathbb{R}^{s \times s}, b_n \in \mathbb{R}^s$ (s — размер эмбединга)
- $A_{n,u}, b_n$ генерируются при помощи двух MLP
- $A_{n,u} = \frac{\mu}{s \cdot |N(u)|} \cdot \text{resize}(MLP_1(x_v, x_{(v,u)}^e), x_u), \mu \in (0, 1)$
- $b_n = MLP_2(x_v)$

При условии, что используются подходящие функции активации (например, \tanh), несложно доказываемается, что итоговая F_w — сжимающее отображение

Non-linear GNN

- $f(\cdot)$ — просто MLP
- Необходимо, чтобы выполнялось свойство сжимающего отображения

Non-linear GNN

- $f(\cdot)$ — просто MLP
- Необходимо, чтобы выполнялось свойство сжимающего отображения
- Это решается добавлением к функционалу ошибки регуляризации Якобиана: $\mathcal{L} = \dots + \beta L(\|\frac{\partial F_w}{\partial x}\|)$

$$L(y) = \begin{cases} (y - \mu)^2, & y > \mu \\ 0, & otherwise \end{cases} \quad \mu \in (0, 1)$$

Non-linear GNN

- $f(\cdot)$ — просто MLP
- Необходимо, чтобы выполнялось свойство сжимающего отображения
- Это решается добавлением к функционалу ошибки регуляризации Якобиана: $\mathcal{L} = \dots + \beta L(\|\frac{\partial F_w}{\partial x}\|)$

$$L(y) = \begin{cases} (y - \mu)^2, & y > \mu \\ 0, & otherwise \end{cases} \quad \mu \in (0, 1)$$

В общем случае регуляризатор может быть любым выражением, дифференцируемым по \mathbf{w} и монотонно возрастающим по норме Якобиана

В обеих моделях для оптимизации вычислений используется Almeida-Pineda recurrent backpropagation

Результаты

- Mutagenesis dataset (классификация молекул)
- Молекула = граф, в каждом графе одна вершина является помеченной и по ней формируется ответ

Method	Knowledge	Reference	Accuracy
non-linear GNN	AB+C+PS		90.5%
$1nn(d_m)$	AB	[87]	81%
$1nn(d_m)$	AB+C	[87]	88%
TILDE	AB	[92]	77%
TILDE	AB+C	[92]	82%
RDBC	AB	[88]	83%
RDBC	AB+C	[88]	82%

Gated graph neural network

- Предыдущие модели могут не очень хорошо работать на больших графах
- Применим GRU

$$\mathbf{h}_v^{(1)} = [\mathbf{x}_v^\top, \mathbf{0}]^\top \quad (1)$$

$$\mathbf{a}_v^{(t)} = \mathbf{A}_{v:}^\top \left[\mathbf{h}_1^{(t-1)\top} \dots \mathbf{h}_{|\mathcal{V}|}^{(t-1)\top} \right]^\top + \mathbf{b} \quad (2)$$

$$\mathbf{z}_v^t = \sigma \left(\mathbf{W}^z \mathbf{a}_v^{(t)} + \mathbf{U}^z \mathbf{h}_v^{(t-1)} \right) \quad (3)$$

$$\mathbf{r}_v^t = \sigma \left(\mathbf{W}^r \mathbf{a}_v^{(t)} + \mathbf{U}^r \mathbf{h}_v^{(t-1)} \right) \quad (4)$$

$$\widetilde{\mathbf{h}}_v^{(t)} = \tanh \left(\mathbf{W} \mathbf{a}_v^{(t)} + \mathbf{U} \left(\mathbf{r}_v^t \odot \mathbf{h}_v^{(t-1)} \right) \right) \quad (5)$$

$$\mathbf{h}_v^{(t)} = (1 - \mathbf{z}_v^t) \odot \mathbf{h}_v^{(t-1)} + \mathbf{z}_v^t \odot \widetilde{\mathbf{h}}_v^{(t)}. \quad (6)$$

Часть 3.

Сверточные нейронные сети на графах

Spectral convolutions

Лапласиан. Определение

Рассмотрим неориентированный граф
его матрицу смежности A и матрицу степеней D

Тогда лапласиан графа L определяется так:

$$L = D - A \in \mathbb{R}^{n \times n}$$

$$L_{i,j} := \begin{cases} \deg(v_i) & \text{if } i = j \\ -1 & \text{if } i \neq j \text{ and } v_i \text{ is adjacent to } v_j \\ 0 & \text{otherwise} \end{cases}$$

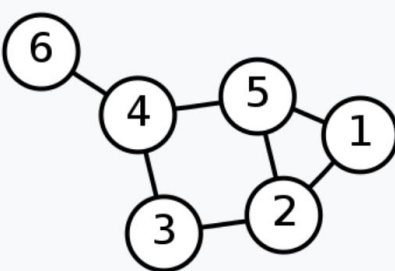
Лапласиан. Определение

Нормализованный Лапласиан

$$L^{\text{sym}} := D^{-\frac{1}{2}} L D^{-\frac{1}{2}} = I - D^{-\frac{1}{2}} A D^{-\frac{1}{2}}$$

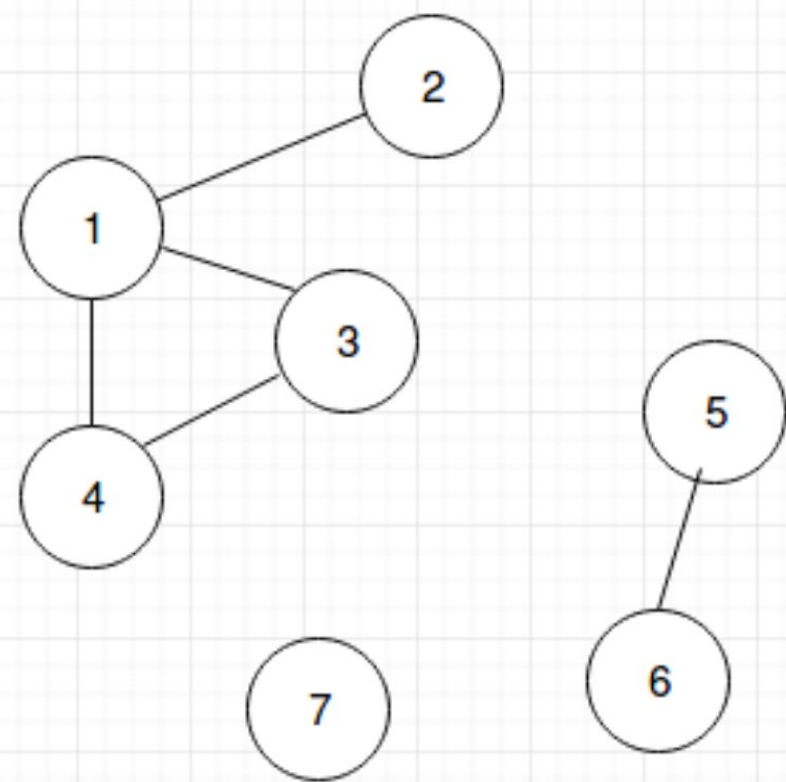
$$L_{i,j}^{\text{sym}} := \begin{cases} 1 & \text{if } i = j \text{ and } \deg(v_i) \neq 0 \\ -\frac{1}{\sqrt{\deg(v_i) \deg(v_j)}} & \text{if } i \neq j \text{ and } v_i \text{ is adjacent to } v_j \\ 0 & \text{otherwise.} \end{cases}$$

Пример

Labelled graph	Degree matrix	Adjacency matrix	Laplacian matrix
	$\begin{pmatrix} 2 & 0 & 0 & 0 & 0 & 0 \\ 0 & 3 & 0 & 0 & 0 & 0 \\ 0 & 0 & 2 & 0 & 0 & 0 \\ 0 & 0 & 0 & 3 & 0 & 0 \\ 0 & 0 & 0 & 0 & 3 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix}$	$\begin{pmatrix} 0 & 1 & 0 & 0 & 1 & 0 \\ 1 & 0 & 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 & 1 & 1 \\ 1 & 1 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \end{pmatrix}$	$\begin{pmatrix} 2 & -1 & 0 & 0 & -1 & 0 \\ -1 & 3 & -1 & 0 & -1 & 0 \\ 0 & -1 & 2 & -1 & 0 & 0 \\ 0 & 0 & -1 & 3 & -1 & -1 \\ -1 & -1 & 0 & -1 & 3 & 0 \\ 0 & 0 & 0 & -1 & 0 & 1 \end{pmatrix}$

Свойства

- Лапласиан - симметричная матрица
- Лапласиан - неотрицательно определенная матрица
- Сумма каждой строки/столбца = 1
- Собственное значение $\lambda_0 = 0$ т.к. $v = (1, 1, \dots, 1)$ - собственный вектор
- число собственных значений, равных 0, = числу компонент графа



[3,	-1,	-1,	-1,	0,	0,	0]
[-1,	1,	0,	0,	0,	0,	0]
[-1,	0,	2,	-1,	0,	0,	0]
[-1,	0,	-1,	2,	0,	0,	0]
[0,	0,	0,	0,	1,	-1,	0]
[0,	0,	0,	0,	-1,	1,	0]
[0,	0,	0,	0,	0,	0,	0]

Преобразование Фурье графа

- Определим сигнал на графе как $f : V \rightarrow \mathbb{R}$
- Преобразование фурье от функции f станет функцией от спектра
- Пусть λ_l - l -ое собственное значение Лапласиана, μ_l - l -ый собственный вектор Лапласиана

Тогда

$$\mathcal{GF}[f](\lambda_l) = \hat{f}(\lambda_l) = \langle f, \mu_l \rangle = \sum_{i=1}^N f(i) \mu_l^*(i),$$

where $\mu_l^* = \mu_l^T$.

Обратное преобразование Фурье графа

$$\mathcal{IGF}[\hat{f}](i) = f(i) = \sum_{l=0}^{N-1} \hat{f}(\lambda_l) \mu_l(i)$$

Преобразование Фурье графа

- Т.к. Лапласиан раскладывается в следующем виде

$$\mathcal{L} = U \Lambda U^{-1} = U \Lambda U^T$$

$$\Lambda = \begin{bmatrix} \lambda_1 & & & \\ & \lambda_2 & & \\ & & \ddots & \\ & & & \lambda_n \end{bmatrix} = \text{diag}([\lambda_1, \dots, \lambda_n]) \in \mathbb{R}^{n \times n}$$

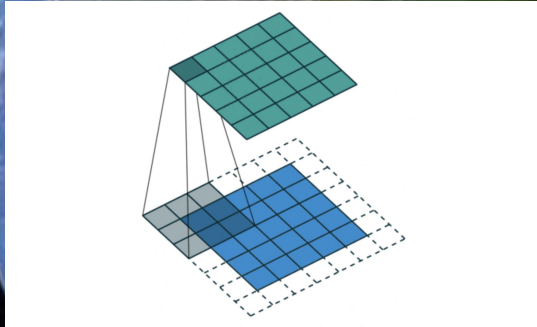
- То преобразование Фурье можно определить как $\hat{f} = U^T f$

Свертка в мат. анализе

$$(f * g)(x) \stackrel{\text{def}}{=} \int_{\mathbb{R}^n} f(y) g(x - y) dy = \int_{\mathbb{R}^n} f(x - y) g(y) dy.$$

Always has been

Is it all a convolution from calculus-2?



Теорема о свертке

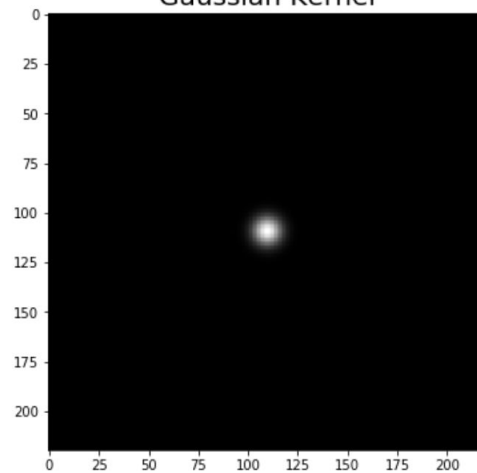
Пусть F - преобразование Фурье, тогда

$$\mathcal{F}\{f * g\} = \mathcal{F}\{f\} \cdot \mathcal{F}\{g\}$$

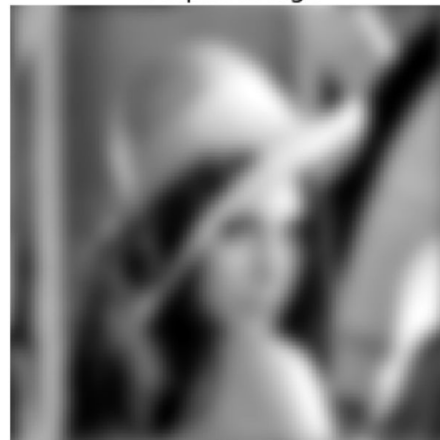
Original Image



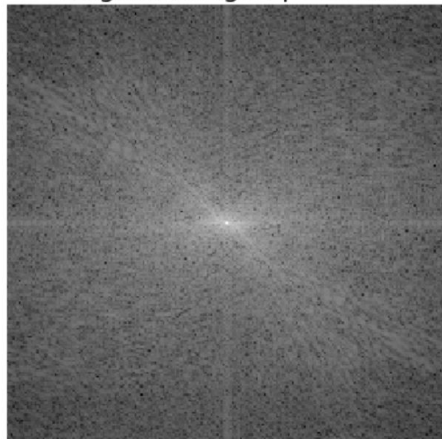
Gaussian Kernel



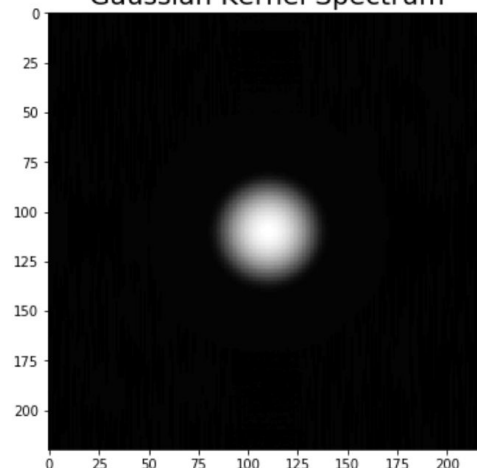
Output Image



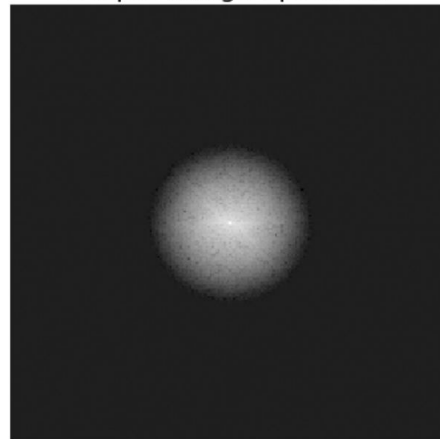
Original Image Spectrum



Gaussian Kernel Spectrum



Output Image Spectrum



Spectral graph convolution

Пользуясь теоремой, если f - сигнал, g - ядро свертки, то определим свертку как

$$\begin{aligned} g(\cdot_G)f &= \mathcal{F}^{-1}(\mathcal{F}(g) \odot \mathcal{F}(f)) = U(U^T g \odot U^T f) \\ &= U g_\theta(\Lambda) U^T f = g_\theta(\mathcal{L})f \end{aligned}$$

где $g_\theta(\Lambda) = \text{diag}(U^T g)$ (делаем из вектора матрицу с этим вектором на диагонали)

Spectral CNN

Теперь введем сигнал как функцию $f : V \rightarrow \mathbb{R}^{C_k}$

И определим сверточный слой нейросети так:

$$f_j^{(k+1)} = \sigma \left(U \sum_{i=1}^{C_k} g_{\theta_{i,j}}^{(k)} U^T f_i^{(k)} \right) = \sigma \left(U \sum_{i=1}^{C_k} g_{\theta_{i,j}}^{(k)} \hat{f}_i^{(k)} \right)$$

$$g_{\theta_{i,j}}^{(k)} = \begin{bmatrix} \theta_1^{(k)} & & & \\ & \theta_2^{(k)} & & \\ & & \ddots & \\ & & & \theta_n^{(k)} \end{bmatrix}$$

Проблемы

- число параметров зависит от числа вершин
- смена индексации вершин приводит к смене Лапласиана, а значит и базиса
- разложение Лапласиана - $O(n^3)$
- фильтр не является локальным
(каждый раз у свертки n параметров)

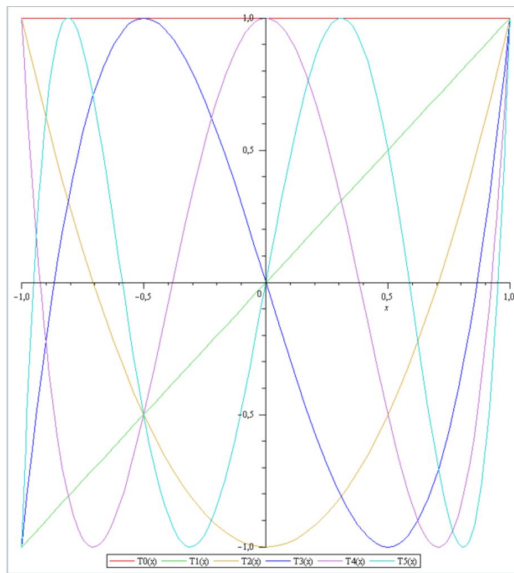
Решение - ChebNet!

Многочлены Чебышева

$$T_0(x) = 1$$

$$T_1(x) = x$$

$$T_{n+1}(x) = 2xT_n(x) - T_{n-1}(x).$$



Свойства многочленов Чебышева

Свойств довольно много, но самое главное, что

- ортогональность, скалярное произведение - $(f, g) = \int_{-1}^1 \frac{1}{\sqrt{1-x^2}} f(x)g(x)dx$
- используются для приближения функций
(перед этим надо перевести функцию в интервал ортогональности многочленов - $[-1, 1]$)

Spectral graph convolution. Наноминание

Если f - сигнал, g - ядро свертки, то определим свертку как

$$\begin{aligned} g(\cdot_G)f &= \mathcal{F}^{-1}(\mathcal{F}(g) \odot \mathcal{F}(f)) = U(U^T g \odot U^T f) \\ &= U g_\theta(\Lambda) U^T f = g_\theta(\mathcal{L})f \end{aligned}$$

где $g_\theta(\Lambda) = \text{diag}(U^T g)$ (делаем из вектора матрицу
с этим вектором на диагонали)

Приближение многочленами Чебышева

Диагональной формы:

$$1. \quad \tilde{\Lambda} = \frac{2\Lambda}{\lambda_{max}} - I_n \in [-1, 1]$$

$$2. \quad g_{\theta}(\Lambda) = \sum_{k=0}^{K-1} \theta_k T_k(\tilde{\Lambda})$$

Лапласиана:

$$1. \quad \tilde{\mathcal{L}} = 2\mathcal{L}/\lambda_{max} - I_n$$

$$2. \quad g_{\theta}(\mathcal{L})f = \sum_{k=0}^{K-1} \theta_k T_k(\tilde{\mathcal{L}})f$$

Приближение многочленами Чебышева

$$g_{\theta}(\Lambda) = \begin{bmatrix} \hat{g}(\lambda_1) & & & \\ & \hat{g}(\lambda_2) & & \\ & & \ddots & \\ & & & \hat{g}(\lambda_n) \end{bmatrix} = \begin{bmatrix} \sum_{k=0}^{K-1} \theta_k T_k(\hat{\lambda}_1) & & & \\ & \sum_{k=0}^{K-1} \theta_k T_k(\hat{\lambda}_2) & & \\ & & \ddots & \\ & & & \sum_{k=0}^{K-1} \theta_k T_k(\hat{\lambda}_n) \end{bmatrix}$$

where θ_k is a vector of Chebyshev coefficients, which is trainable parameter.

Приближение многочленами Чебышева

$$\begin{aligned} f(\cdot_G)g_\theta &= g_\theta(U\Lambda U^T)f = U \sum_{k=0}^{K-1} \theta_k T_k(\tilde{\Lambda}) U^T f = \sum_{k=0}^{K-1} U \theta_k T_k(\tilde{\Lambda}) U^T f \\ &= \sum_{k=0}^{K-1} U \theta_k \left(\sum_{c=0}^k \alpha_{kc} \tilde{\Lambda}^k \right) U^T f = \sum_{k=0}^{K-1} \theta_k \left(\sum_{c=0}^k \alpha_{kc} U \tilde{\Lambda}^k U^T \right) f \\ &= \sum_{k=0}^K \theta_k \left(\sum_{c=0}^k \alpha_{kc} (U \tilde{\Lambda} U^T)^k \right) f = \sum_{k=0}^{K-1} \theta_k T_k(U \tilde{\Lambda} U^T) f = \sum_{k=0}^{K-1} \theta_k T_k(\tilde{\mathcal{L}}) f \end{aligned}$$

Какие преимущества мы получили?

- не надо искать собственные значения Лапласиана
- количество параметров не зависит от размера графа
- фильтры стали K -локализованными

Попробуем сделать еще меньше
параметров...

Graph Convolutional Network

- зафиксируем $K = 1$ в ChebNet
- будем считать, что $\lambda_{max} \sim 2$

Тогда

$$f(\cdot_G)g = \sum_{k=0}^1 \theta_k T_k(\tilde{\mathcal{L}})f = \theta_0 T_0(\tilde{\mathcal{L}})f + \theta_1 T_1(\tilde{\mathcal{L}})f$$

Graph Convolutional Network

- чтобы еще больше уменьшить количество параметров, предположим, что $\theta_0 = -\theta_1 = \theta$
- Тогда

$$f(\cdot_G)g \approx (\theta_0 + \theta_1(\mathcal{L} - I_n))f = \theta_0 - \theta_1 D^{-\frac{1}{2}} A D^{-\frac{1}{2}} = (\theta(D^{-\frac{1}{2}} A D^{-\frac{1}{2}} + I_n))f$$

θ это просто одно число

Улучшаем стабильность

- добавляем петли: $\tilde{A} = A + I_n$
- пересчитываем матрицу степеней: $\tilde{D}_{i,i} = \sum_{j=1}^n \tilde{A}_{i,j}$
- в итоге:

$$f(\cdot_G)g = \theta \tilde{D}^{-\frac{1}{2}} \tilde{A} \tilde{D}^{-\frac{1}{2}} f$$

Пример

- semi-supervised классификация вершин с помощью двух слоев GCN:

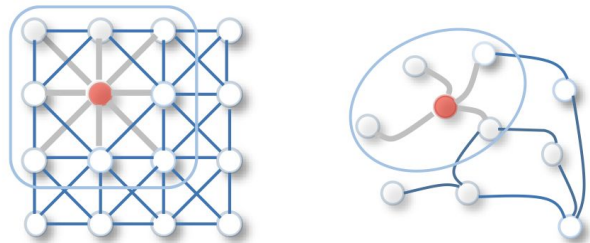
$$Z = \text{softmax}(\hat{A} \text{ReLU}(\hat{A} f W^{<0>}) W^{<1>}) \quad (48)$$

$W^{<0>} \in \mathbb{R}^{C \times H}$ is an input-to-hidden weight matrix for a hidden layer with H feature maps. $W^{<1>} \in \mathbb{R}^{H \times F}$ is a hidden-to-output weight matrix, F is the dimension of feature maps in the output layer.

Spatial convolutions

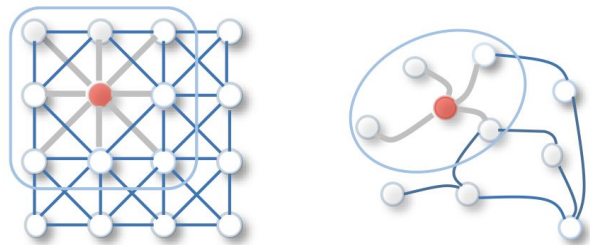
NN4G

- Основная идея везде одна: определим свёрточный слой по аналогии с изображениями



NN4G

- Основная идея везде одна: определим свёрточный слой по аналогии с изображениями
- Выполним свёртку, просуммировав информацию со всех соседей и добавив skip connections
- Получаем формулу обновления эмбединга на k -ом шаге:



$$\mathbf{h}_v^{(k)} = f(\mathbf{W}^{(k)T} \mathbf{x}_v + \sum_{i=1}^{k-1} \sum_{u \in N(v)} \Theta^{(k)T} \mathbf{h}_u^{(k-1)}),$$

Здесь f — функция активации, эмбединг на первой итерации нулевой

NN4G

- Выражение с предыдущего слайда можно переписать в матричной форме:

$$\mathbf{H}^{(k)} = f(\mathbf{X}\mathbf{W}^{(k)} + \sum_{i=1}^{k-1} \mathbf{A}\mathbf{H}^{(i)} \mathbf{\Theta}^{(i)})$$

Матрица смежности не нормирована, поэтому могут возникнуть проблемы с разными порядками эмбедингов разных вершин

Diffusion convolutional neural network

- Предполагается, что информация передаётся от одной вершины к другой с некоторой вероятностью, так что за какое-то количество итераций система достигнет равновесия:

$$\mathbf{H}^{(k)} = f(\mathbf{W}^{(k)} \odot \mathbf{P}^k \mathbf{X}).$$

Здесь $\mathbf{P} = \mathbf{D}^{-1} \mathbf{A} \in \mathbb{R}^{n \times n}$ — матрица вероятностей передачи информации

Diffusion convolutional neural network

- Предполагается, что информация передаётся от одной вершины к другой с некоторой вероятностью, так что за какое-то количество итераций система достигнет равновесия:

$$\mathbf{H}^{(k)} = f(\mathbf{W}^{(k)} \odot \mathbf{P}^k \mathbf{X}).$$

Здесь \mathbf{P} — матрица вероятностей передачи информации, $\mathbf{A} \in \mathbb{R}^{n \times n}$ — матрица смежности графа.

Отметим, что на каждом шаге эмбединг того же размера, что и вектор признаков, и не зависит от предыдущего эмбединга. Поэтому в конце матрицы эмбедингов со всех слоёв стакаются

GraphSage

Количество соседей вершины может быть огромным, поэтому не эффективно рассматривать их всех. Предлагается выбрать случайное подмножество фиксированного размера:

$$\mathbf{h}_v^{(k)} = \sigma(\mathbf{W}^{(k)} \cdot f_k(\mathbf{h}_v^{(k-1)}, \{\mathbf{h}_u^{(k-1)}, \forall u \in S_{\mathcal{N}(v)}\})), \quad (24)$$

where $\mathbf{h}_v^{(0)} = \mathbf{x}_v$, $f_k(\cdot)$ is an aggregation function, $S_{\mathcal{N}(v)}$ is a random sample of the node v 's neighbors. The aggregation function should be invariant to the permutations of node orderings such as a mean, sum or max function.

Graph attention network

До этого мы полагали, что у всех соседей одинаковый вес, хотя это может быть не так. Попробуем применить attention:

$$\mathbf{h}_v^{(k)} = \sigma\left(\sum_{u \in \mathcal{N}(v) \cup v} \alpha_{vu}^{(k)} \mathbf{W}^{(k)} \mathbf{h}_u^{(k-1)}\right),$$

Веса вычисляются как:

$$\alpha_{vu}^{(k)} = \text{softmax}(g(\mathbf{a}^T [\mathbf{W}^{(k)} \mathbf{h}_v^{(k-1)} \parallel \mathbf{W}^{(k)} \mathbf{h}_u^{(k-1)}]))$$

где $g(\cdot)$ — *LeakyReLU*, \mathbf{a}^T — обучаемый вектор параметров,
 \parallel обозначает конкатенацию векторов

Spectral vs Spatial

- Спектральные модели имеют хорошее теоретическое обоснование
- Спектральные модели менее эффективны: им нужно либо вычислять собственные векторы, либо обрабатывать весь граф одновременно
- Пространственные модели более применимы к большим графам, так как вычисления могут производиться батчами
- Спектральные модели имеют плохую обобщающую способность на новые графы
- Спектральным моделям нужны неориентированные графы, пространственные могут работать с любыми

Часть 3.

Graph Pooling

Что мы умеем?

Что мы умеем?

Получать векторные представления вершин

Что мы умеем?

Получать векторные представления вершин

Но как тогда делать классификацию графов?...

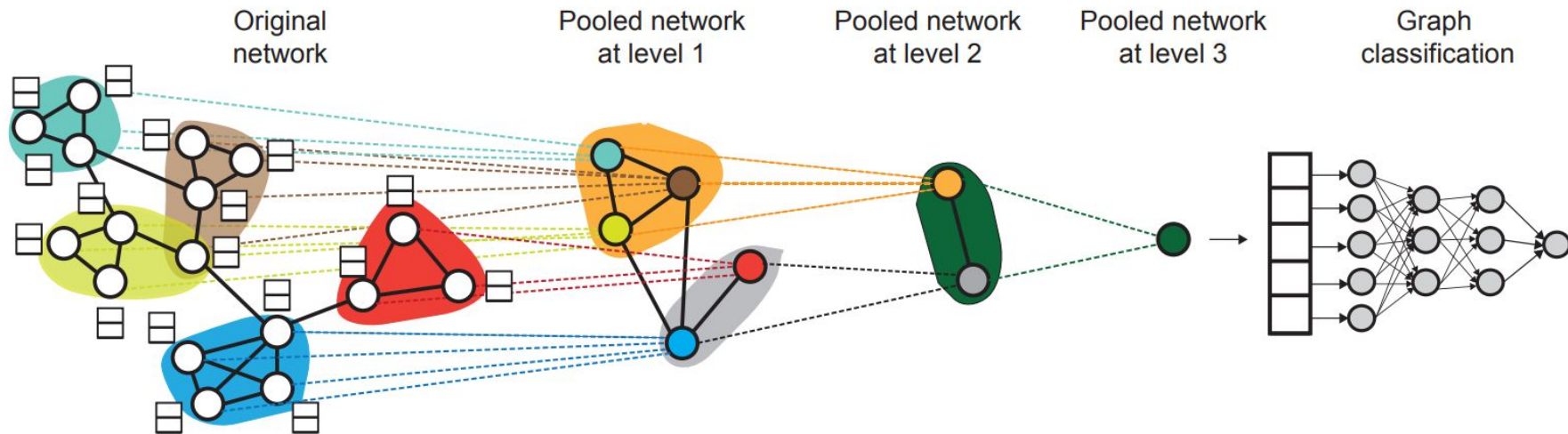
Что мы умеем?

Получать векторные представления вершин

Но как тогда делать классификацию графов?...

Надо научиться агрегировать информацию с вершин

Pipeline



Самое простое

$$\mathbf{h}_G = \text{mean/max/sum}(\mathbf{h}_1^{(K)}, \mathbf{h}_2^{(K)}, \dots, \mathbf{h}_n^{(K)})$$

GRACLUS

- быстрая кластеризация без подсчета собственных значений
- основано на расширении kernel k-means на графы

Pooling with GRACLUS

- с помощью GRACLUS получаем матрицу кластеризации $\mathbf{S}^{(l)} \in \{0, 1\}^{n_{l-1} \times n_l}$ mapping each node to its cluster index in $\{1, \dots, n_l\}$, with $n_l < n_{l-1}$ clusters.
- делаем простой max/average пуллинг
- пересчитываем матрицу смежности

$$\mathbf{A}^{(l)} = \mathbf{S}^{(l)\top} \mathbf{A}^{(l-1)} \mathbf{S}^{(l)}$$

Pooling with GRACLUS

- не может выучить специфичные для данного набора графов структура из-за фиксированного алгоритм
(хочется выучить кластеризацию в end-to-end подходе)
- выдает бинарную матрицу
(хочется понимать влияние каждой вершины на каждую)
- на большим графах все равно долго

DIFFPOOL

- попробуем выдавать кластеризацию (вероятности) с помощью еще одной GNN:

$$\mathbf{S}^{(l)} = \text{softmax} \left(\text{GNN}_1^{(l)}(\mathbf{A}^{(l-1)}, \mathbf{X}^{(l-1)}) \right)$$

- пересчитываем векторные представления

$$\mathbf{X}^{(l)} = \mathbf{S}^{(l)\top} \text{GNN}_2^{(l)}(\mathbf{A}^{(l-1)}, \mathbf{X}^{(l-1)}) \quad \text{and} \quad \mathbf{A}^{(l)} = \mathbf{S}^{(l)\top} \mathbf{A}^{(l-1)} \mathbf{S}^{(l)}.$$

DIFFPOOL

Какие лоссы использовать в качестве мотивации хорошо кластеризовать?

- $\sum_l \| \mathbf{A}^{(l)} - \mathbf{S}^{(l)} \mathbf{S}^{(l)\top} \|$ (вершины рядом должны принадлежать одному кластеру)

DIFFPOOL

Какие лоссы использовать в качестве мотивации хорошо кластеризовать?

- $\sum_l \| \mathbf{A}^{(l)} - \mathbf{S}^{(l)} \mathbf{S}^{(l)\top} \|$ (вершины рядом должны принадлежать одному кластеру)
- $L_E = \frac{1}{n} \sum_{i=1}^n H(S_i)$ (тянем кластеризацию к более вырожденному one-hot вектору)

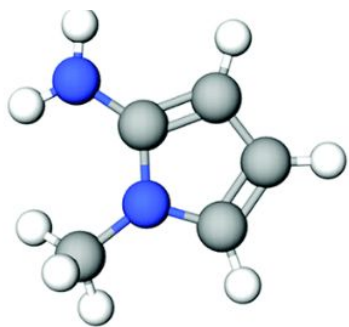
	Method	Data Set					Gain
		ENZYMES	D&D	REDDIT-MULTI-12K	COLLAB	PROTEINS	
Kernel	GRAPHLET	41.03	74.85	21.73	64.66	72.91	
	SHORTEST-PATH	42.32	78.86	36.93	59.10	76.43	
	1-WL	53.43	74.02	39.03	78.61	73.76	
	WL-OA	60.13	79.04	44.38	80.74	75.26	
GNN	PATCHYSAN	—	76.27	41.32	72.60	75.00	4.17
	GRAPHSAGE	54.25	75.42	42.24	68.25	70.48	—
	ECC	53.50	74.10	41.73	67.79	72.65	0.11
	SET2SET	60.15	78.12	43.49	71.75	74.29	3.32
	SORTPOOL	57.12	79.37	41.82	73.76	75.54	3.39
	DIFFPOOL-DET	58.33	75.47	46.18	82.13	75.62	5.42
	DIFFPOOL-NO LP	61.95	79.98	46.65	75.58	76.22	5.95
	DIFFPOOL	62.53	80.64	47.08	75.48	76.25	6.27

Часть 4.

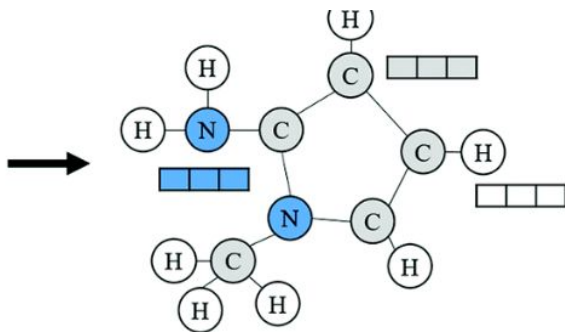
Applications

Химия

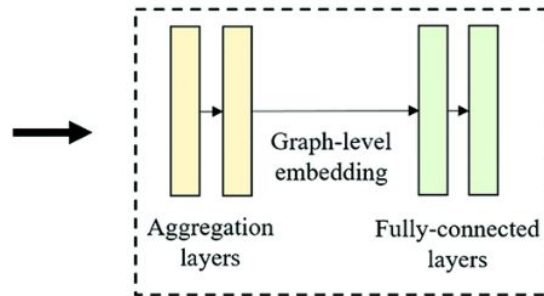
- предсказания какие-то свойства молекул
- генерация новых молекул



Molecular or crystal structure



Molecular graph



Graph neural network

y
Corresponding
target

“A deep learning approach to antibiotic discovery”

2020

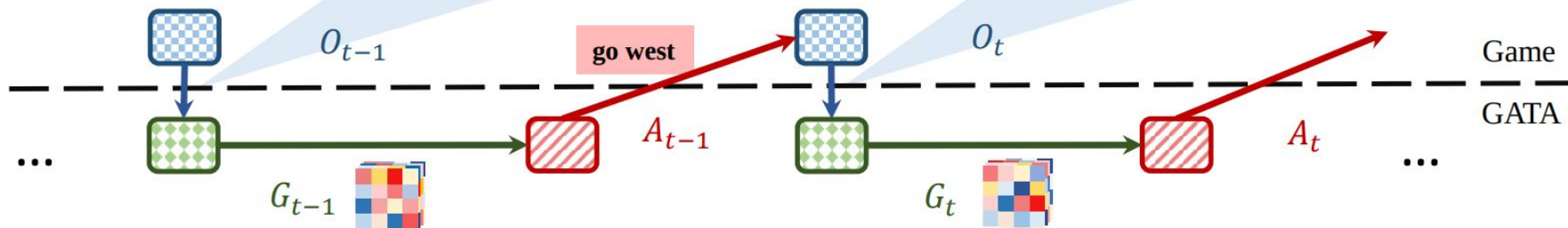
- классификация на то, проявляет ли молекула свойства антибиотика
- найдена новая молекула - Halicin (уже проверили на мышах)

“Learning Dynamic Belief Graphs to Generalize on Text-Based Games”

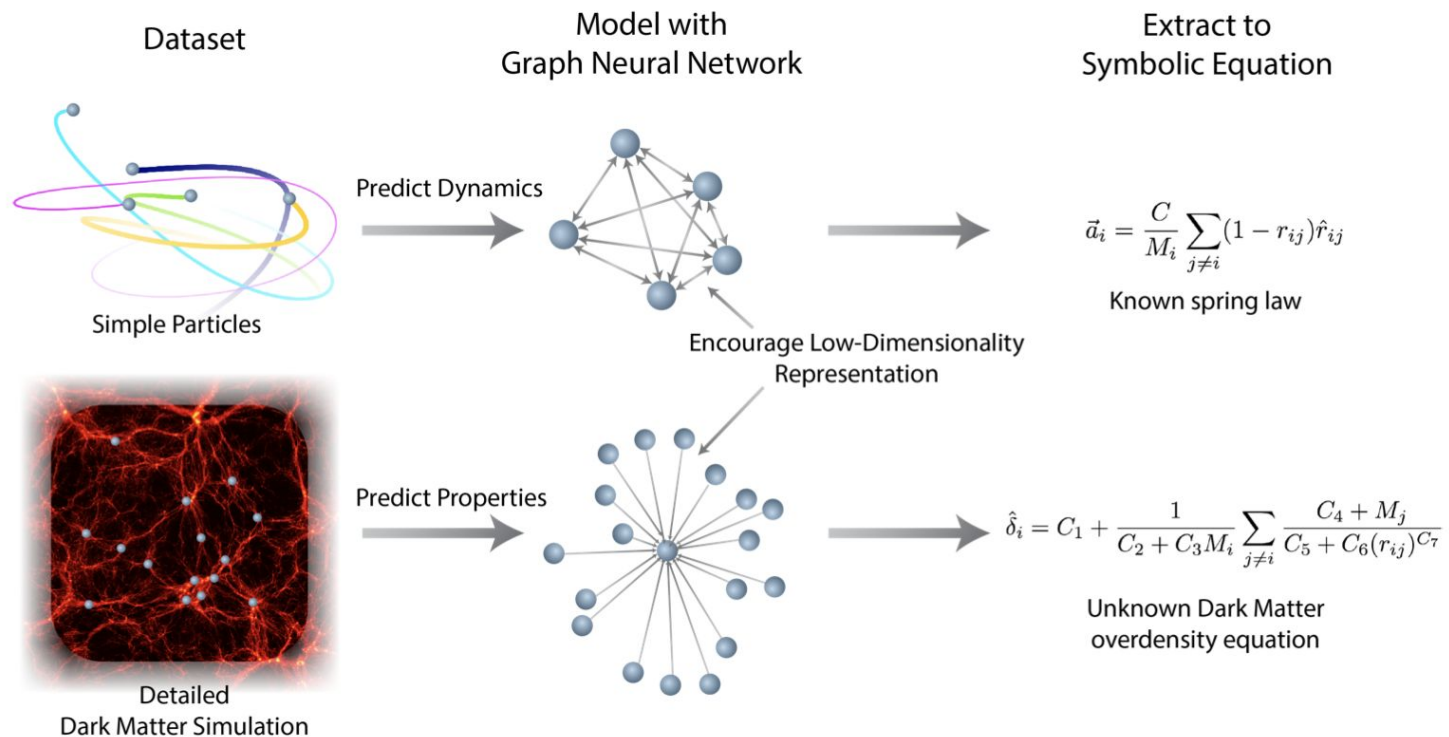
2020

You find yourself in a **backyard**. You make out a **patio table**. But it is empty. You see a **patio chair**. The **patio chair** is stylish. But there isn't a thing on it. You see a gleam over in a corner, where you can see a **BBQ**. There is a **closed screen door** leading **south**. There is an **open wooden door** leading **west**.

Welcome to the **shed**. You can barely contain your excitement. You can make out a **closed toolbox** here. You can see a **workbench**. The **workbench** is wooden. Looks like someone's already been here and taken everything off it, though. You swear loudly. There is an **open wooden door** leading **east**.



Открытие законов Вселенной?



Список источников

- [A Comprehensive Survey On Graph Neural Networks](#)
- [The Graph Neural Network Model](#)
- [Gated Graph Sequence Neural Networks](#)
- [Graph Neural Networks: A Review of Methods and Applications](#)
- [Understanding Graph Neural Networks \(video\)](#)
- [Graph Laplacian and its application in Machine learning](#)
- [Anisotropic, Dynamic, Spectral and Multiscale Filters Defined on Graphs](#)
- [Discovering Symbolic Models from Deep Learning with Inductive Biases](#)
- [Understanding Spectral Graph Neural Network](#)