

Group meeting - February 14, 2020

Paper review & Research progress

Truong Son Hy *

*Department of Computer Science
The University of Chicago

Ryerson Physical Lab



GNNExplainer: Generating Explanations for Graph Neural Networks



GNNExplainer: Generating Explanations for Graph Neural Networks

Rex Ying, Dylan Bourgeois, Jiaxuan You, Marinka Zitnik, Jure Leskovec
Stanford University, Department of Computer Science



Problem Setup

Understanding GNN's predictions is important and useful for several reasons:
a) Can increase trust in predictions and give insights into important graph structures.
b) Can improve model transparency and interpretability in decision-critical applications.
c) Can identify and correct systematic patterns of mistakes made by GNNs before deploying them in the real world.

- GNNExplainer**, the first approach for explaining predictions made by GNNs:
Any machine learning task on graphs: It is model-agnostic and provides explanations for node classification, link prediction, graph classification.
Any GNN model: It is model-agnostic and can explain predictions of any GNN variant using the message-passing neural architecture.
Rich graph structural information and node features: It incorporates relevant node features and relational information, the essence of graphs.

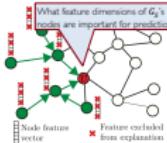
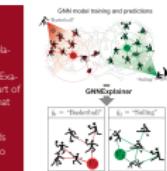
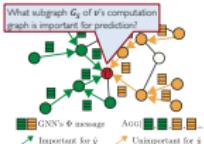
Given a trained GNN and its prediction(s), GNNExplainer returns an explanation in the form of a small subgraph of the input graph together with a small subset of node features that are most influential for the prediction(s).

Key Idea and Explanation Mask

A GNN model Φ is trained on a social interaction graph to predict future sport activities.

Given a prediction $\hat{y}_v = \text{"Basketball"}$ for person v , GNNExplainer finds a small subgraph of the input graph together with a small subset of node features that are most influential for \hat{y}_v . Examining explanation for \hat{y}_v , we see that many friends in one part of v 's social circle enjoy ball games, and so the GNN predicts that v will like basketball.

Similarly, examining explanation for \hat{y}_v , we see that v 's friends and friends of his friends enjoy water and beach sports; and so the GNN predicts $\hat{y}_v = \text{"Sailing"}$.



GNNExplainer: Generating Explanations for Graph Neural Networks. NeurIPS 2019.

GNNExplainer: Design and Results

Objective: Maximize mutual information between the prediction and the explanation:

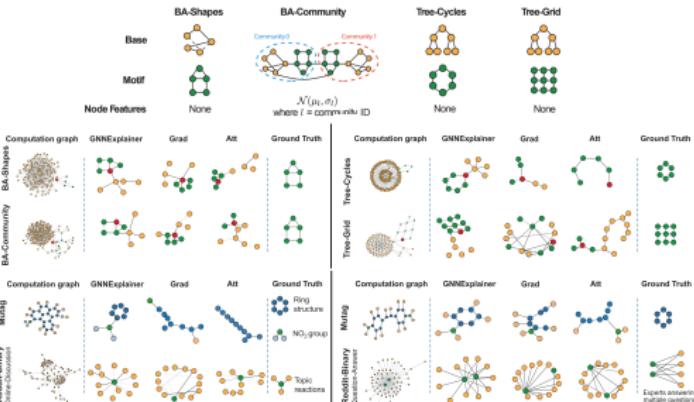
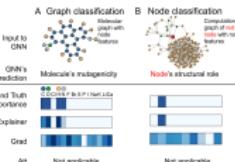
$$\max_{G_S} MI(Y, (G_S, X_S)) = H(Y) - H(Y|G = G_S, X = X_S)$$

Conditional Entropy:

$$\min_{\mathcal{G}} \mathbb{E}_{G_S \sim \mathcal{G}} H(Y|G = G_S, X = X_S)$$

Tractable approximation via masking the graph structure and node features:

$$\min_M - \sum_{c=1}^C \mathbb{I}[y=c] \log P_\Phi(Y=y|G = A_S \odot \sigma(M), X = X_S)$$



<http://cs.stanford.edu/~name> name = {rexy, jiaxuan, marinka, jure}

Graph Transformer Networks

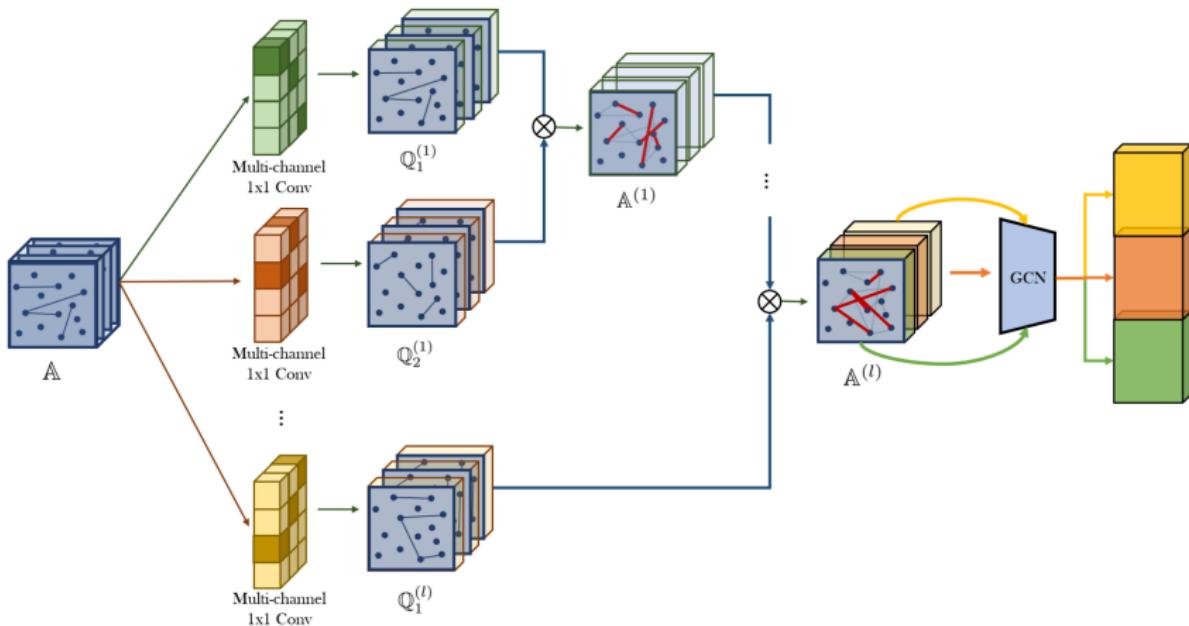


Figure 2: Graph Transformer Networks (GTNs) learn to generate a set of new meta-path adjacency matrices $A^{(l)}$ using GT layers and perform graph convolution as in GCNs on the new graph structures. Multiple node representations from the same GCNs on multiple meta-path graphs are integrated by concatenation and improve the performance of node classification. $Q_1^{(l)} \text{ and } Q_2^{(l)} \in \mathbb{R}^{N \times N \times C}$ are intermediate adjacency tensors to compute meta-paths at the l th layer.



Questions for nCRP

- ① How to initialize the topic-hierarchy tree given a set of documents?
- ② How to adjust the tree given new documents? The algorithm is not clear to me.
- ③ Is each word in a document labeled to a particular topic? Or is it unsupervised in which the topics are **abstract**?
- ④ Maybe this work **Autoencoding Variational Inference For Topic Models** (ICLR 2017) is the right way to go? I did not have time to cover it yet.



Problem

We address the problem of learning topic hierarchies from data. The model selection problem in this domain is daunting - which of the **large collection of possible trees** to use?



Hierarchical Topic Models and the Nested CRP (1)

Problem

We address the problem of learning topic hierarchies from data. The model selection problem in this domain is daunting - which of the **large collection of possible trees** to use?

Proposal

- We take a Bayesian approach, generating an appropriate prior via a distribution on partitions that we refer to as the nested **Chinese restaurant process** (CRP).
- We build a hierarchical topic model by combining this prior with a likelihood that is based on a hierarchical variant of **latent Dirichlet allocation**.



Hierarchical Topic Models and the Nested CRP (2)

Problem: Learning a topic hierarchy from data

Given a collection of **documents**, each of which contains a set of **words** we wish to discover common usage patterns or **topics** in the documents, and to organize these topics into a hierarchy.



Problem: Learning a topic hierarchy from data

Given a collection of **documents**, each of which contains a set of **words** we wish to discover common usage patterns or **topics** in the documents, and to organize these topics into a hierarchy.

Approach

- Specify a **generative probabilistic model for hierarchical structure.**
- **Our hierarchies are random variables;** moreover, these random variables are specified procedurally, according to an algorithm that constructs the hierarchy as data are made available.



Approach

- A distribution on partitions of integers known as the **Chinese restaurant process**.
- Extend the Chinese restaurant process to a hierarchy of partitions.
- Each node in the hierarchy is associated with a topic, where a topic is a distribution across words.
- A document is generated by choosing a path from the root to a leaf, repeatedly sampling topics along that path, and sampling the words from the selected topics.
- This approach differs from models of topic hierarchies which are built on the premise that the distributions associated with parents and children are similar → No such constraint.



Hierarchical Topic Models and the Nested CRP (4)

The Chinese restaurant process (CRP) is a distribution on partitions of integers obtained by imagining a process by which M customers sit down in a Chinese restaurant with an infinite number of tables. The m -th subsequent customer sits at a table drawn from the following distribution:

$$p(\text{occupied table } i \mid \text{previous customers}) = \frac{m_i}{\gamma + m - 1}$$
$$p(\text{next unoccupied table} \mid \text{previous customers}) = \frac{\gamma}{\gamma + m - 1}$$

where m_i is the number of previous customers at table i , and γ is a parameter. After M customers sit down, the seating plan gives a partition of M items.



Hierarchical Topic Models and the Nested CRP (4)

The Chinese restaurant process (CRP) is a distribution on partitions of integers obtained by imagining a process by which M customers sit down in a Chinese restaurant with an infinite number of tables. The m -th subsequent customer sits at a table drawn from the following distribution:

$$p(\text{occupied table } i \mid \text{previous customers}) = \frac{m_i}{\gamma + m - 1}$$
$$p(\text{next unoccupied table} \mid \text{previous customers}) = \frac{\gamma}{\gamma + m - 1}$$

where m_i is the number of previous customers at table i , and γ is a parameter. After M customers sit down, the seating plan gives a partition of M items.

Unrelated

This makes me remember **Hilbert's paradox of the Grand Hotel!**



Hierarchical Topic Models and the Nested CRP (5)

A nested Chinese restaurant process:

- There are an infinite number of infinite-table Chinese restaurants in a city.
- One restaurant is determined to be the **root restaurant** and on each of its infinite tables is a card with the name of another restaurant.
And so on.
- On the first evening, a customer enters the root Chinese restaurant and selects a table. On the second evening, he goes to the restaurant identified on the first night's table and chooses another table.
- After M tourists take L -day vacations, the collection of paths describe a particular L -level subtree of the infinite tree.



Hierarchical Topic Models and the Nested CRP (6)

Consider a data set composed of a **corpus** of documents. Each document is a collection of **words**, where a **word** is an item in a **vocabulary**.



Consider a data set composed of a **corpus** of documents. Each document is a collection of **words**, where a **word** is an item in a **vocabulary**.

Assumption

Our basic assumption is that the words in a document are generated according to a mixture model where the mixing proportions are random and document-specific.



Hierarchical Topic Models and the Nested CRP (7)

Document-specific mixture distribution:

$$p(w|\theta) = \sum_{i=1}^K \theta_i p(w|z=i, \beta_i)$$

where:

- θ : document-specific mixing proportions.
- Assume K different topics in the corpus.
- z ranges over K possible values.
- β is a parameter.
- $p(w|z, \beta)$ is a distribution over words.



Hierarchical Topic Models and the Nested CRP (8)

Two-level generative probabilistic process for generating a document:

- ① Choose a K -vector θ of topic proportions from a distribution $p(\theta|\alpha)$, where α is a corpus-level parameter.
- ② Repeatedly sample words from the mixture distribution $p(w|\theta)$ for the chosen value of θ .

When the distribution $p(\theta|\alpha)$ is chosen to be a Dirichlet distribution, we obtain the **latent Dirichlet allocation** model (LDA).



Hierarchical Topic Models and the Nested CRP (9)

Extend the process to a hierarchy of topics:

- ① Suppose we are given an L -level tree and each node is associated with a topic.
- ② Choose a path from the root of the tree to a leaf.
- ③ Draw a vector of topic proportions θ from an L -dimensional Dirichlet.
- ④ Generate the words in the document from a mixture of the topics along the path from root to leaf, with mixing proportions θ .



Hierarchical Topic Models and the Nested CRP (10)

Use the nested CRP to relax the assumption of a fixed tree structure. The nested CRP can be used to place a prior on possible trees → **Hierarchical LDA** (hLDA):

- ① Let c_1 be the root topic.
- ② For each level $\ell \in \{2, \dots, L\}$:
 - Draw a topic from $c_{\ell-1} \rightarrow c_\ell$.
- ③ Draw an L -dimensional topic proportion vector θ from $\text{Dir}(\alpha)$.
- ④ For each word $n \in \{1, \dots, N\}$:
 - Draw $z \in \{1, \dots, L\}$ from $\text{Multinomial}(\theta)$.
 - Draw w_n from the topic w_z .



Hierarchical Topic Models and the Nested CRP (11)

Suppose we are given a corpus of M documents $\mathbf{w}_1, \dots, \mathbf{w}_M$:

- $w_{m,n}$ is the n -th word of the m -th document.
- $c_{m,\ell}$ is the restaurant/node corresponding to the ℓ -th topic (at level ℓ -th of the tree) in document m .
- $z_{m,n}$ is the assignment of the n -th word in the m -th document to one of the L available topics.



Hierarchical Topic Models and the Nested CRP (12)

Conceptually, we divide the Gibbs sampler into two parts. First, given the current state of the CRP, we sample the $z_{m,n}$ variables of the underlying LDA model following the algorithm developed in [12], which we do not reproduce here. Second, given the values of the LDA hidden variables, we sample the $c_{m,\ell}$ variables which are associated with the CRP prior. The conditional distribution for \mathbf{c}_m , the L topics associated with document m , is:

$$p(\mathbf{c}_m | \mathbf{w}, \mathbf{c}_{-m}, \mathbf{z}) \propto p(\mathbf{w}_m | \mathbf{c}, \mathbf{w}_{-m}, \mathbf{z}) p(\mathbf{c}_m | \mathbf{c}_{-m}),$$

where \mathbf{w}_{-m} and \mathbf{c}_{-m} denote the \mathbf{w} and \mathbf{c} variables for all documents other than m . This expression is an instance of Bayes' rule with $p(\mathbf{w}_m | \mathbf{c}, \mathbf{w}_{-m}, \mathbf{z})$ as the likelihood of the data given a particular choice of \mathbf{c}_m and $p(\mathbf{c}_m | \mathbf{c}_{-m})$ as the prior on \mathbf{c}_m implied by the nested CRP. The likelihood is obtained by integrating over the parameters β , which gives:

$$p(\mathbf{w}_m | \mathbf{c}, \mathbf{w}_{-m}, \mathbf{z}) = \prod_{\ell=1}^L \left(\frac{\Gamma(n_{c_{m,\ell}, -m}^{(\cdot)} + W\eta)}{\prod_w \Gamma(n_{c_{m,\ell}, -m}^{(w)} + \eta)} \frac{\prod_w \Gamma(n_{c_{m,\ell}, -m}^{(w)} + n_{c_{m,\ell}, m}^{(w)} + \eta)}{\Gamma(n_{c_{m,\ell}, -m}^{(\cdot)} + n_{c_{m,\ell}, m}^{(\cdot)} + W\eta)} \right),$$

where $n_{c_{m,\ell}, -m}^{(w)}$ is the number of instances of word w that have been assigned to the topic indexed by $c_{m,\ell}$, not including those in the current document, W is the total vocabulary size, and $\Gamma(\cdot)$ denotes the standard gamma function. When \mathbf{c} contains a previously unvisited restaurant, $n_{c_{m,\ell}, -m}^{(w)}$ is zero.



Variational Inference for the Nested Chinese Restaurant Process (1)

CRP:

$$p(c_d = k | c_{1:(d-1)}) \propto \begin{cases} m_k & \text{if } k \text{ is previous occupied} \\ \gamma & \text{if } k \text{ is a new table,} \end{cases}$$

nCRP:

1. Draw a path $c_n | c_{1:(n-1)} \sim \text{nCRP}(\gamma, c_{1:(n-1)})$, which contains L nodes from the tree.
2. Draw a latent variable $\mathbf{x}_n \sim p(\mathbf{x}_n | \lambda)$.
3. Draw an observation $t_n \sim p(t_n | W_{c_n}, \mathbf{x}_n, \tau)$.

The parameters λ and τ are associated with the latent variables \mathbf{x} and data generating distribution, respectively. Note that W_{c_n} contains the w_i s selected by the path c_n . Specific applications of the nCRP mixture depend on the particular forms of $p(w)$, $p(x)$ and $p(t|W_c, x)$.



Variational Inference for the Nested Chinese Restaurant Process (2)

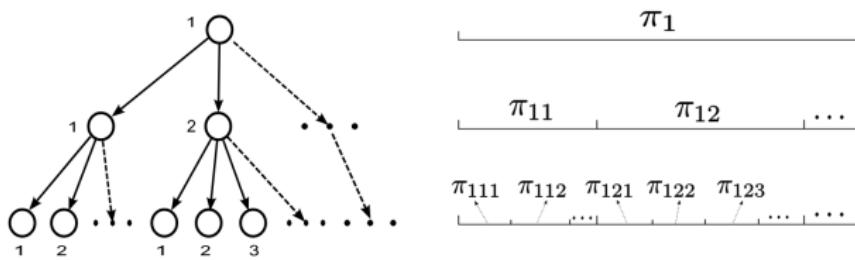


Figure 1: **Left.** A possible tree structure in a 3-level nCRP. **Right.** The tree-based stick-breaking construction of a 3-level nCRP.



Variational Inference for the Nested Chinese Restaurant Process (3)

CRP mixtures can be equivalently formulated using the Dirichlet process (DP) as a distribution over the distribution of each data point's random parameter [21, 4]. An advantage of expressing the CRP mixture with a DP is that the draw from the DP can be explicitly represented using the stick-breaking construction [22]. The DP bundles the scaling parameter γ and base distribution G_0 . A draw from a $\text{DP}(\gamma, G_0)$ is described as

$$v_i \sim \text{Beta}(1, \gamma), \quad \pi_i = v_i \prod_{j=1}^{i-1} (1 - v_j), \quad \mathbf{w}_i \sim G_0, \quad i \in \{1, 2, \dots\}, \quad G = \sum_{i=1}^{\infty} \pi_i \delta_{\mathbf{w}_i},$$

where $\boldsymbol{\pi}$ are the stick lengths, and $\sum_{i=1}^{\infty} \pi_i = 1$ almost surely. This representation also illuminates the discreteness of a distribution drawn from a DP.



Variational Inference for the Nested Chinese Restaurant Process (4)

The tree-based stick-breaking construction lets us calculate the conditional probability of a path given \mathbf{V} . Let the path $\mathbf{c} = [1, c_2, \dots, c_L]$,

$$p(\mathbf{c}|\mathbf{V}) = \prod_{\ell=1}^L \pi_{1,c_2,\dots,c_\ell} = \prod_{\ell=1}^L v_{1,c_2,\dots,c_\ell} \prod_{j=1}^{c_\ell-1} (1 - v_{1,c_2,\dots,j}). \quad (2)$$

By integrating out \mathbf{V} in Equation 2, we recover the nCRP. Given Equation 2, the joint probability of a data set under the nCRP mixture is

$$p(\mathbf{t}_{1:N}, \mathbf{x}_{1:N}, \mathbf{c}_{1:N}, \mathbf{V}, \mathbf{W}) = p(\mathbf{V})p(\mathbf{W}) \prod_{n=1}^N p(\mathbf{c}_n|\mathbf{V})p(\mathbf{x}_n)p(\mathbf{t}_n|\mathbf{W}_{\mathbf{c}_n}, \mathbf{x}_n). \quad (3)$$

This representation is the basis for variational inference.



Autoencoding Variational Inference For Topic Models

Maybe this work is the right way to go? I did not have time to cover this paper yet.



Q & A

Thank you very much for your attention!

