

Explainability for Graph Networks

What is Explainability?

Explainability refers to methods and techniques in the application of machine learning or deep learning models such that the results of the solution can be understood by human experts.

Problem Statement

Many people say machine learning models are "black boxes", in the sense that they can make good predictions but you can't understand the logic behind those predictions. This statement is true in the sense that most data scientists don't know how to extract insights from models yet. However, explainability techniques help to extract the following insights from sophisticated machine learning models.

1. What features in the data did the model think are most important?
2. For any single prediction from a model, how did each feature in the data affect that particular prediction?
3. How does each feature affect the model's predictions in a big-picture sense (what is its typical effect when considered over a large number of possible predictions)?
4. Which neurons get activated towards this prediction?

Different Methods of Explainability

1. Contrastive gradient-based saliency maps is perhaps the most straight-forward and well-established approach. In this approach, one simply differentiates the output of the model with respect to the model input. The negative values in the gradient are discarded to only retain the parts of input that positively contribute to the solution.

2. Class Activation Mapping provides an improvement over saliency maps for convolutional neural networks, including GCNNs, by identifying important, class-specific features at the last convolutional layer as opposed to the input space. It is well-known that such features tend to be more semantically meaningful (e.g., faces instead of edges).

3. Layer-wise Relevance Propagation (LRP) is a method that identifies important pixels by running a backward pass in the neural network. The backward pass is a conservative relevance redistribution procedure, where neurons that contribute the most to the higher-layer receive most relevance from it.

Need for Explainability for Graph Networks

Explainability can be particularly helpful for graphs, even more than for images, because non-expert humans cannot intuitively determine the relevant context within a graph, for example, when identifying groups of atoms (a sub-graph structure on a molecular graph) that contribute to a particular property of a molecule.

Though, there are not many methods to explain a graph network, we have tried to modify GradCAM method for images and use it in our graph networks.

Methodology

1. Get the gradient value before the activation function of the required layer and the respective weight param.
2. Do dot product of the gradient value matrix and respective weight matrix to get a node value tensor containing the logits of each atom in the molecule.
3. Map the node value to visualize which atom relatively contributes how much towards the prediction.

Visualisation of effects of each atom in a molecule towards Solubility(property)

COLOR INDEX:

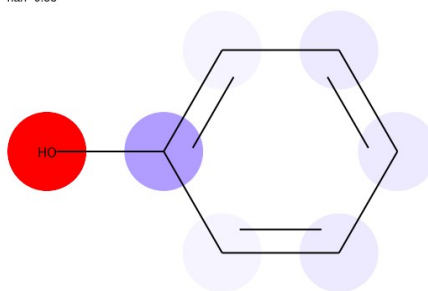
Range-

Dark Red(Highly positive)

to

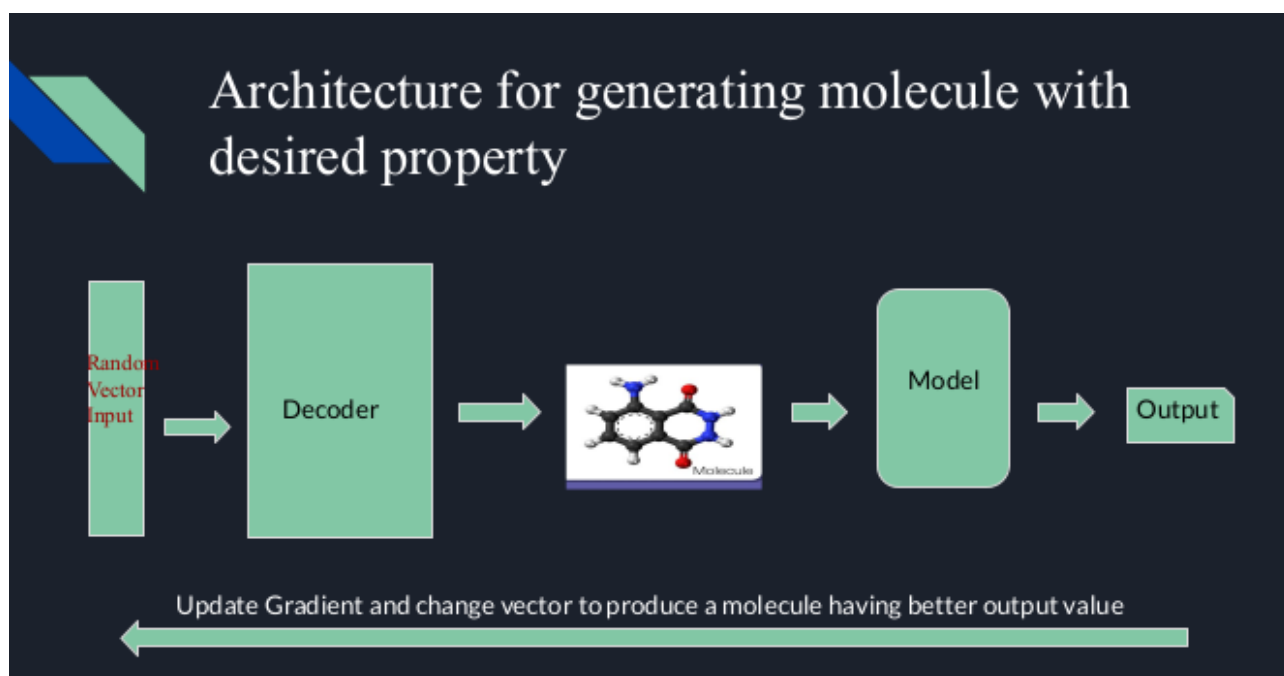
Dark Blue(Highly negative)

c1cc(O)ccc1
nan -0.33



Making graphs features understandable

We have seen which atoms are responsible for your prediction. Now, we are also trying to get which neurons or channels get activated towards this prediction. For this we came up with a model idea which is discussed below.



Hierarchical Generation of Molecular Graphs using Structural Motifs

Motif Extraction for making Vocabulary

We define a motif $S_i = (V_i, E_i)$ as a subgraph of molecule G induced by atoms in V_i and bonds in E_i . Given a molecule, we extract its motifs S_1, \dots, S_n such that their union covers the entire molecular graph: $V = \bigcup_i V_i$ and $E = \bigcup_i E_i$. To extract motifs, we decompose a molecule G into disconnected fragments by breaking all the bridge bonds that will not violate chemical validity (illustrations in the appendix).

1. Find all the bridge bonds $(u, v) \in E$, where both u and v have degree $\Delta u, \Delta v \geq 2$ and either u or v is part of a ring. Detach all the bridge bonds from its neighbors.

- Now the graph G becomes a set of disconnected subgraphs G_1, \dots, G_N . Select G_i as motif in G if its occurrence in the training set is more than $\Delta = 100$.
- If G_i is not selected as motif, further decompose it into rings and bonds and select them as motif in G .

Hierarchical Graph Generation

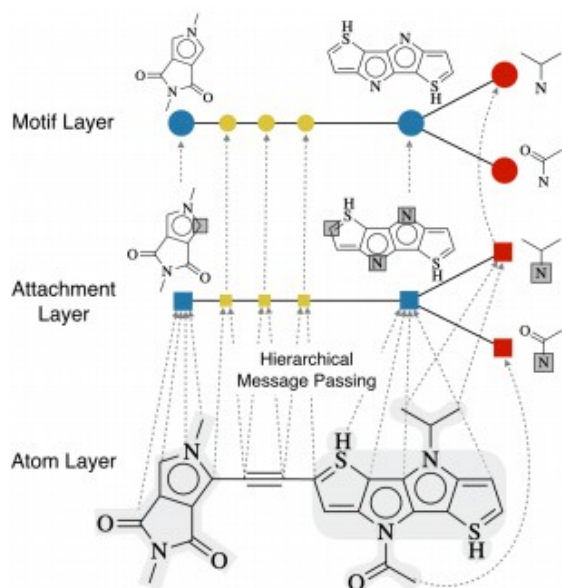


Figure 2. Hierarchical graph encoder. Dashed arrows connect each atom to the motifs it belongs. In the attachment layer, each node \mathcal{A}_i is a particular attachment configuration of motif S_i . The atoms in the intersection between each motif and its neighbors are highlighted in faded block.

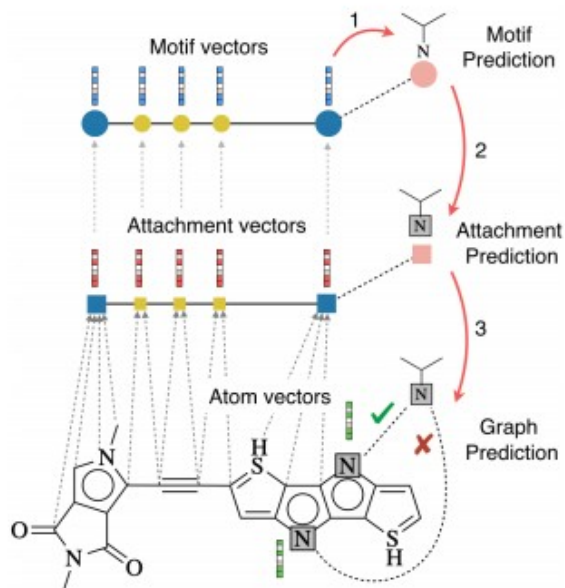


Figure 3. Hierarchical graph decoder. In each step, the decoder first runs hierarchical message passing to compute motif, attachment and atom vectors. Then it performs motif and attachment prediction for the next motif node. Finally, it decides how the new motif should be attached to the current graph via graph prediction.

Our approach extends the variational autoencoder (Kingma & Welling, 2013) to molecular graphs by introducing a hierarchical decoder and a matching encoder. In our framework, the probability of a graph G is modeled as a joint distribution over structural motifs S_1, \dots, S_n constituting G , together with their attachments A_1, \dots, A_n . Each attachment $A_i = \{v_j \mid v_j \in S_k \text{ } S_i \cap S_k\}$ indicates the intersecting atoms between S_i and its neighbor motifs. To capture complex dependencies involved in the joint distribution of motifs and their attachments, we propose an auto-regressive factorization of $P(G)$:

$$P(G) = \int_z P(z) \prod_k P(S_k, A_k \mid S_{<k}, A_{<k}, z) dz \quad (1)$$

As illustrated in Figure 3, in each generation step, our decoder adds a new motif S_k (motif prediction) and its attachment configuration A_k (attachment prediction). Then it decides how the new motif should be attached to the current graph (graph prediction).

To support the above hierarchical generation, we need to design a matching encoder representing molecules at multiple resolutions in order to provide necessary information for each decoding step. Therefore, we propose to represent a molecule G by a hierarchical graph HG with three layers (see Figure 2):

1. Motif layer: This layer represents how the motifs are coarsely connected in the graph. This layer provides essential information for the motif prediction in the decoding process. Specifically, this layer contains n nodes S_1, \dots, S_n and m edges $\{(S_i, S_j) \mid S_i \cap S_j \neq \emptyset\}$ for all intersecting motifs S_i, S_j . This layer is tree-structured due to our way of constructing motifs.

2. Attachment layer: This layer encodes the connectivity between motifs at a fine-grained level. Each node $A_i = (S_i, \{v_j\})$ in this layer represents a particular attachment configuration of motif S_i ,

where $\{v_j\}$ are atoms in the intersection between S_i and one of its neighbor motifs (see Figure 2). This layer provides crucial information for the attachment prediction step during decoding, which helps reducing the space of candidate attachments between S_i and its neighbor motifs. Just like the motif vocabulary V_S , all the attachment configurations of S_i form a motif-specific vocabulary $V_A(S_i)$, which is computed from the training set.

3. Atom layer: The atom layer is the molecular graph G representing how its atoms are connected. Each atom node v is associated with a label a_v indicating its atom type and charge. Each edge (u, v) in the atom layer is labeled with b_{uv} indicating its bond type. This layer provides necessary information for the graph prediction step during decoding.

We further introduce edges that connect the atoms and motifs between different layers in order to propagate information in between. In particular, we draw a directed edge from atom v in the atom layer to node A_i in the attachment layer if $v \in S_i$. We also draw edges from A_i to S_i in the motif layer. This gives us the hierarchical graph H_G for molecule G , which will be encoded by a hierarchical message passing network (MPN). During encoding, each node S_i is represented as a one-hot encoding in the motif vocabulary V_S . Likewise, each node A_i is represented as a one-hot encoding in the attachment vocabulary $V_A(S_i)$.

Hierarchical Graph Encoder

Our encoder contains three MPNs that encode each of the three layers in the hierarchical graph. For simplicity, we denote the MPN encoding process as $\text{MPN}_\psi(\cdot)$ with parameter ψ , and denote $\text{MLP}(x, y)$ as a multi-layer neural network whose input is the concatenation of x and y .

Atom Layer MPN We first encode the atom layer of H_G (denoted as H_G^a). The inputs to this MPN are the embedding vectors $\{e(a_u)\}$, $\{e(b_{uv})\}$ of all the atoms and bonds in G . During encoding, the network propagates the message vectors between different atoms for T iterations and then outputs the atom representation h_v for each atom v :

$$c_G^a = \{h_v\} = \text{MPN}_{\psi_1}(\mathcal{H}_G^a, \{e(a_u)\}, \{e(b_{uv})\}) \quad (2)$$

Attachment Layer MPN The input feature of each node A_i in the attachment layer $H_a G$ is an concatenation of the embedding $e(A_i)$ and the sum of its atom vectors $\{h_v \mid v \in S_i\}$:

$$f_{A_i} = \text{MLP}\left(e(A_i), \sum_{v \in S_i} h_v\right) \quad (3)$$

The input feature for each edge (A_i, A_j) in this layer is an embedding vector $e(d_{ij})$, where d_{ij} describes the relative ordering between node A_i and A_j during decoding. Specifically, we set $d_{ij} = k$ if node A_i is the k -th child of node A_j and $d_{ij} = 0$ if A_i is the parent. We then run T iterations of message passing over $H_a G$ to compute the motif representations:

$$c_G^a = \{h_{A_i}\} = \text{MPN}_{\psi_2}(\mathcal{H}_G^a, \{f_{A_i}\}, \{e(d_{ij})\}) \quad (4)$$

Motif Layer MPN Similarly, the input feature of node S_i in this layer is computed as the concatenation of embedding $e(S_i)$ and the node vector h_{A_i} from the previous layer. Finally, we run message passing over the motif layer H^s_G to obtain the motif representation:

$$f_{S_i} = \text{MLP}(e(S_i), h_{A_i}) \quad (5)$$

$$c_G^s = \{h_{S_i}\} = \text{MPN}_{\psi_3}(\mathcal{H}_G^s, \{f_{S_i}\}, \{e(d_{ij})\}) \quad (6)$$

Finally, we represent a molecule G by a latent vector z_G sampled through reparameterization trick with mean $\mu(h_{S_1})$ and log variance $\Sigma(h_{S_1})$:

$$z_G = \mu(h_{S_1}) + \exp(\Sigma(h_{S_1})) \cdot \epsilon; \quad \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad (7)$$

where S_1 is the root motif (i.e., the first motif to be generated during reconstruction)

Hierarchical Graph Decoder

As illustrated in Figure 3, our graph decoder generates a molecule G by incrementally expanding its hierarchical graph. In t th generation step, we first use the same hierarchical MPN architecture to encode all the motifs and atoms in $H^{(t)}_G$, the (partial) hierarchical graph generated till step t . This gives us motif vectors h_{S_k} and atom vectors h_{v_j} for the existing motifs and atoms.

During decoding, the model maintains a set of frontier nodes F where each node $S_k \in F$ is a motif that still has neighbors to be generated. F is implemented as a stack because motifs are generated in their depth-first order. Suppose S_k is at the top of stack F in step t , the model makes the following predictions conditioned on latent representation z_G :

Motif Prediction: The model predicts the next motif S_t to be attached to S_k . This is cast as a classification task over the motif vocabulary V_s :

$$p_{S_t} = \text{softmax}(\text{MLP}(h_{S_k}, z_G)) \quad (8)$$

Attachment Prediction: Now the model needs to predict the attachment configuration A_t of motif S_t (i.e., what atoms $v_j \in S_t$ belong to the intersection of S_t and its neighbor motifs). This is also cast as a classification task over the attachment vocabulary

$$p_{A_t} = \text{softmax}(\text{MLP}(h_{S_k}, z_G)) \quad (9)$$

This prediction step is crucial because it significantly reduces the space of possible attachments between S_t and its neighbor motifs.

Graph Prediction: Finally, the model must decide how S_t should be attached to S_k . The attachment between S_t and S_k is defined as atom pairs $M_{tk} = \{(u_j, v_j) \mid u_j \in A_k, v_j \in A_t\}$ where atom u_j and v_j are attached together. The probability of a candidate attachment M is computed based on the atom vectors h_{u_j} and h_{v_j}

$$p_M = \text{softmax}(h_M \cdot z_G) \quad (10)$$

$$h_M = \sum_j \text{MLP}(h_{u_j}, h_{v_j}) \quad (11)$$

The number of possible attachments are limited because the number of attaching atoms between two motifs is small and the attaching points must be consecutive.

The above three predictions together give an autoregressive factorization of the distribution over the next motif and its attachment. Each of the three decoding steps depends on the outcome of previous step, and predicted attachments will in turn affect the prediction of subsequent motifs.

Training During training, we apply teacher forcing to the above generation process, where the generation order is determined by a depth-first traversal over the ground truth molecule. Given a training set of molecules, we seek to minimize the negative ELBO:

$$- \mathbb{E}_{z \sim Q} [\log P(G|z)] + \lambda_{KL} \mathcal{D}_{KL}[Q(z|G) || P(z)] \quad (12)$$

References

1. <https://github.com/baldassarreFe/graph-network-explainability/tree/master/notebooks>
2. <https://graphreason.github.io/papers/25.pdf>
3. <https://arxiv.org/pdf/2002.03230.pdf>
4. <https://github.com/wengong-jin/hgraph2graph>