

Interpretable Latent Variable Model for Molecular Design

UGP Presentation

Gurpreet Singh

25th April, 2019

Indian Institute of Technology, Kanpur

Background

In Variational Inference, the optimization objective is given by the ELBO

$$\log p(\mathbf{x}) \geq \mathcal{L}(\mathbf{x}, \phi) = \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z} | \mathbf{x}, \phi)} \left[\log p(\mathbf{x}, \mathbf{z}) - \log q(\mathbf{z} | \mathbf{x}, \phi) \right]$$

In Variational Inference, the optimization objective is given by the ELBO

$$\log p(\mathbf{x}) \geq \mathcal{L}(\mathbf{x}, \phi) = \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z} | \mathbf{x}, \phi)} \left[\log p(\mathbf{x}, \mathbf{z}) - \log q(\mathbf{z} | \mathbf{x}, \phi) \right]$$

Black Box Variational Inference is based on stochastic optimization of the variational objective

$$\nabla \mathcal{L}(\mathbf{x}, \phi) = \mathbb{E}_{q(\mathbf{z} | \mathbf{x}, \phi)} \left[\nabla_{\phi} \log q(\mathbf{z} | \mathbf{x}, \phi) \left(\log p(\mathbf{x}, \mathbf{z}) - \log q(\mathbf{z} | \mathbf{x}, \phi) \right) \right]$$

The above term is known as the Score Function.

If we can find a reparameterization $\epsilon = f^{-1}(\mathbf{z}, \phi)$ such that ϵ is independent of ϕ , then

$$\nabla \mathcal{L}(\mathbf{x}, \phi) = \mathbb{E}_{q(\epsilon)} \left[\nabla_{\mathbf{z}} \left[\log p(\mathbf{x}, \mathbf{z}) - \log q(\mathbf{z} \mid \mathbf{x}, \phi) \right]_{\mathbf{z}=f(\epsilon, \phi)} \nabla_{\phi} f(\epsilon, \phi) \right]$$

If we can find a reparametrization $\epsilon = f^{-1}(\mathbf{z}, \phi)$ such that ϵ is independent of ϕ , then

$$\nabla \mathcal{L}(\mathbf{x}, \phi) = \mathbb{E}_{q(\epsilon)} \left[\nabla_{\mathbf{z}} \left[\log p(\mathbf{x}, \mathbf{z}) - \log q(\mathbf{z} | \mathbf{x}, \phi) \right]_{\mathbf{z}=f(\epsilon, \phi)} \nabla_{\phi} f(\epsilon, \phi) \right]$$

If \mathbf{z} is discrete $p(\mathbf{z})$ is a p.m.f. and therefore, the derivatives no longer exist. Hence, reparameterization is **not** possible

Some methods to train discrete latent variable models -

1. Marginalize all discrete latent variables
2. Using local expectation gradients, and reparametrization and marginalization
3. Relaxed computation of discrete densities (straight through approach)
4. Smooth relaxations for discrete variables (e.g. Concrete Distribution)
5. Using the Score Function approach with variance reduction using control variates/non-uniform sampling.

The Concrete Distribution

Sampling from a categorical distribution: Sample U from a $[0, 1]$ uniform distribution, then reparameterize as

$$C = \arg \max_{k \in [K]} \log \alpha_k - \log (-\log (U_k))$$

A softmax version of this can be used as smooth relaxations

$$Z_k = \frac{\exp \left((\log \alpha_k + G_k) / \tau \right)}{\sum_{k'=1}^K \exp \left((\log \alpha_{k'} + G_{k'}) / \tau \right)}$$

where τ is the **temperature** parameter.

The Gumbolt Trick

Gumbel-Softmax trick only works for factorial priors. The Gumbolt trick offer an approach which allows us to work with more expressive priors, the (Restricted) Boltzmann Machine priors.

Boltzmann Machine Prior

$$q(\mathbf{z}) = \frac{e^{-E(\mathbf{z})}}{Z_{\theta}}, \quad \text{where}$$
$$E(\mathbf{z}) = \mathbf{z}^T \mathbf{b} + \mathbf{z}^T \mathbf{W} \mathbf{z}$$

Restricted Boltzmann Machine Prior

$$E(\mathbf{z}) = \mathbf{z}_1^T \mathbf{b}_1 + \mathbf{z}_2^T \mathbf{b}_2 + \mathbf{z}_1^T \mathbf{W} \mathbf{z}_2$$

The posterior and the prior for the smoothing variables are given as

$$p(\zeta | \theta) = \frac{e^{-E_{\theta}(\zeta)}}{\tilde{Z}_{\theta}}, \text{ where } \tilde{Z}_{\theta} = \int_{\zeta} e^{E_{\theta}(\zeta)} d\zeta$$
$$q(\zeta | \mathbf{x}, \phi) = \prod_i \zeta_i q_i + (1 - \zeta_i)(1 - q_i)$$

The posterior and the prior for the smoothing variables are given as

$$p(\zeta | \theta) = \frac{e^{-E_{\theta}(\zeta)}}{\tilde{Z}_{\theta}}, \text{ where } \tilde{Z}_{\theta} = \int_{\zeta} e^{E_{\theta}(\zeta)} d\zeta$$
$$q(\zeta | \mathbf{x}, \phi) = \prod_i \zeta_i q_i + (1 - \zeta_i)(1 - q_i)$$

The alternate ELBO, therefore, is given as

to
$$\tilde{\mathcal{L}}(\mathbf{x}, \phi) = \mathbb{E}_{q(\zeta | \mathbf{x}, \phi)} \left[\log p(\mathbf{x} | \zeta, \theta) - \log q(\zeta | \mathbf{x}, \phi) - E_{\theta}(\zeta) \right] - \log \tilde{Z}_{\theta}$$

Relaxed ELBO is given as

to

$$\check{\mathcal{L}}(\mathbf{x}, \phi) = \mathbb{E}_{q(\zeta | \mathbf{x}, \phi)} \left[\log p(\mathbf{x} | \zeta, \theta) - \log q(\zeta | \mathbf{x}, \phi) - E_{\theta}(\zeta) \right] - \log Z_{\theta}$$

Relaxed ELBO is given as

to

$$\check{\mathcal{L}}(\mathbf{x}, \phi) = \mathbb{E}_{q(\zeta | \mathbf{x}, \phi)} \left[\log p(\mathbf{x} | \zeta, \theta) - \log q(\zeta | \mathbf{x}, \phi) - E_{\theta}(\zeta) \right] - \log Z_{\theta}$$

We have

$$\log p(\mathbf{x} | \theta) \geq \tilde{\mathcal{L}}(\mathbf{x}, \phi) \geq \check{\mathcal{L}}(\mathbf{x}, \phi)$$

Relaxed ELBO is given as

to

$$\check{\mathcal{L}}(\mathbf{x}, \phi) = \mathbb{E}_{q(\zeta | \mathbf{x}, \phi)} \left[\log p(\mathbf{x} | \zeta, \theta) - \log q(\zeta | \mathbf{x}, \phi) - E_{\theta}(\zeta) \right] - \log Z_{\theta}$$

We have

$$\log p(\mathbf{x} | \theta) \geq \tilde{\mathcal{L}}(\mathbf{x}, \phi) \geq \check{\mathcal{L}}(\mathbf{x}, \phi)$$

Also, as $\tau \rightarrow 0$, $\check{\mathcal{L}} \rightarrow \tilde{\mathcal{L}}$.

Since as $\tau \rightarrow 0$, we have $\zeta \rightarrow \mathbf{z}$, therefore $\check{\mathcal{L}} \rightarrow \tilde{\mathcal{L}} \rightarrow \mathcal{L}$



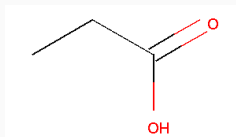
Figure 1: Regeneration Plot for Gumbolt-Vae



Figure 2: Sampling Plot for Gumbolt-Vae

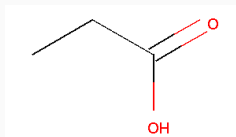
Molecular Design

1. SMILES - offers a string representation of the molecule's skeletal structure



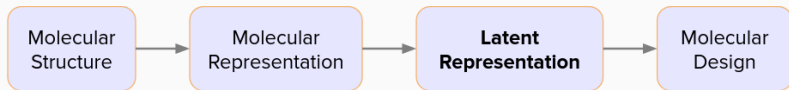
The SMILES string of the above molecule is given as “CCC(=O)O”

1. SMILES - offers a string representation of the molecule's skeletal structure



The SMILES string of the above molecule is given as “CCC(=O)O”

2. Molecular Graphs - assumes each atom to be a node and each bond to be an edge.



Many deep learning based generative models have been proposed lately for automatic molecule design. Two common approaches are GANs and VAEs (our focus). Two categorizations -

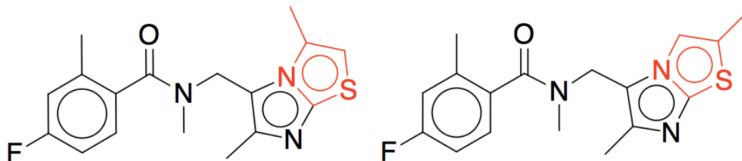
1. SMILES based models
2. Graph based model

1. **CVAE** – Modifies vanilla VAEs to use convolutional layers in the encoder and layers of GRU networks for the decoder to predict sequences of strings (SMILES)

1. **CVAE** – Modifies vanilla VAEs to use convolutional layers in the encoder and layers of GRU networks for the decoder to predict sequences of strings (SMILES)
2. **GVAE** – Improves upon CVAE by utilizing the parse trees (CFG) generated from SMILES strings as the molecular representation and using Grammar Variational Autoencoders as their architecture. Generates syntactically correct SMILES strings.

1. **CVAE** – Modifies vanilla VAEs to use convolutional layers in the encoder and layers of GRU networks for the decoder to predict sequences of strings (SMILES)
2. **GVAE** – Improves upon CVAE by utilizing the parse trees (CFG) generated from SMILES strings as the molecular representation and using Grammar Variational Autoencoders as their architecture. Generates syntactically correct SMILES strings.
3. **SD-VAE** – Incorporates semantics into the parse trees (Attribute Grammars) to generate syntactically as well as semantically coherent SMILES strings.

Problem with SMILES



Cc1cn2c(CN(C)C(=O)c3ccc(F)cc3C)c(C)nc2s1

Cc1cc(F)ccc1C(=O)N(C)Cc1c(C)nc2scc(C)n12

1. **VAE with Regularization Framework** – Uses Variational Graph Autoencoders with additional constraints on graph connectivity and edge types

1. **VAE with Regularization Framework** – Uses Variational Graph Autoencoders with additional constraints on graph connectivity and edge types
2. **CGVAE** – Generates molecular graphs using sequential graph generation utilizing Gated Graph Neural Networks in both the encoder and the decoder.

1. **VAE with Regularization Framework** – Uses Variational Graph Autoencoders with additional constraints on graph connectivity and edge types
2. **CGVAE** – Generates molecular graphs using sequential graph generation utilizing Gated Graph Neural Networks in both the encoder and the decoder.

Problem: Invalid intermediate structures while sequential generation

Junction Tree Variational Autoencoder (JT-VAE)

JT-VAE generates molecular graphs in two phases

1. It generates a tree-structured object (a junction tree) that represents the scaffold of subgraph components and their coarse relative arrangements

Junction Tree Variational Autoencoder (JT-VAE)

JT-VAE generates molecular graphs in two phases

1. It generates a tree-structured object (a junction tree) that represents the scaffold of subgraph components and their coarse relative arrangements
2. The subgraphs (nodes in the junction tree) are assembled into a coherent molecule graph using a graph message passing network. This approach incrementally generates the molecular graph while maintaining molecular validity at every step.

Junction Tree Variational Autoencoder (JT-VAE)

JT-VAE generates molecular graphs in two phases

1. It generates a tree-structured object (a junction tree) that represents the scaffold of subgraph components and their coarse relative arrangements
2. The subgraphs (nodes in the junction tree) are assembled into a coherent molecule graph using a graph message passing network. This approach incrementally generates the molecular graph while maintaining molecular validity at every step.

Thus, the JT-VAE encoder has two parts: graph encoder and tree encoder, so has the decoder: tree decoder and graph decoder. The graph and tree encoders are closely related to message passing neural networks (MPNNs)

- The molecular structure is encoded into a two-part continuous latent representation $\mathbf{z} = (\mathbf{z}_T, \mathbf{z}_G)$ where \mathbf{z}_T encodes the tree structure and what the subgraph components are in the tree. \mathbf{z}_G encodes the graph to capture the fine-grained connectivity

- The molecular structure is encoded into a two-part continuous latent representation $\mathbf{z} = (\mathbf{z}_T, \mathbf{z}_G)$ where \mathbf{z}_T encodes the tree structure and what the subgraph components are in the tree. \mathbf{z}_G encodes the graph to capture the fine-grained connectivity
- Both the molecular graph and junction tree are encoded via message passing using adapted gated recurrent units (GRUs).

- The molecular structure is encoded into a two-part continuous latent representation $\mathbf{z} = (\mathbf{z}_T, \mathbf{z}_G)$ where \mathbf{z}_T encodes the tree structure and what the subgraph components are in the tree. \mathbf{z}_G encodes the graph to capture the fine-grained connectivity
- Both the molecular graph and junction tree are encoded via message passing using adapted gated recurrent units (GRUs).
- The junction tree is sequentially generated with message passing in a top-down as well as bottom-up (Contains topological and label predictions)

- The molecular structure is encoded into a two-part continuous latent representation $\mathbf{z} = (\mathbf{z}_T, \mathbf{z}_G)$ where \mathbf{z}_T encodes the tree structure and what the subgraph components are in the tree. \mathbf{z}_G encodes the graph to capture the fine-grained connectivity
- Both the molecular graph and junction tree are encoded via message passing using adapted gated recurrent units (GRUs).
- The junction tree is sequentially generated with message passing in a top-down as well as bottom-up (Contains topological and label predictions)
- Graph decoding is performed in an alternating manner where a subgraph (for each node of the tree) is updated using messages generated by the tree encode

Gumbolt JT-VAE

Binary latent vectors (\mathbf{b}) are used in company with real latent vectors (\mathbf{r}), and the Gumbolt trick is employed to reparameterize the binary latent vectors (with RBM prior). Final latent representation \mathbf{z} is given as $\mathbf{z} = \mathbf{b} \odot \mathbf{r}$

The Decoder

$$[\mathbf{b}_1, \mathbf{b}_2] \sim \text{RBM}(\mathbf{a}_1, \mathbf{a}_2, \mathbf{W}), \quad \mathbf{r} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

The Encoder

$$q(\mathbf{b}_n, \mathbf{r}_n \mid \mathbf{x}, \mathbf{A}, \phi) = \prod_{k=1}^K q(b_{n,k} \mid \mathbf{x}_n, \mathbf{A}, \phi) q(r_{n,k} \mid \mathbf{x}_n, \mathbf{b}_n, \mathbf{A}, \phi)$$

and

$$q(b_{n,k} \mid \mathbf{x}, \mathbf{A}, \phi) = \text{Bernoulli}(\pi_{n,k})$$

$$q(r_{n,k} \mid \mathbf{x}, \mathbf{A}, \phi) = \mathcal{N}(\mu_{n,k}, \sigma_{n,k}^2)$$

where $\{\pi_{n,k}\}_{k=1}^K = \text{FC}(\mathbf{h}_n)$, $\{\mu_{n,k}, \sigma_{n,k}\}_{k=1}^K = \text{FC}(\mathbf{h}_n, \mathbf{b}_n)$ and $\mathbf{h}_n = \text{MPNN}(\mathbf{x}_n, \mathbf{A})$

Comparison with JT-VAE (ZINC Dataset)

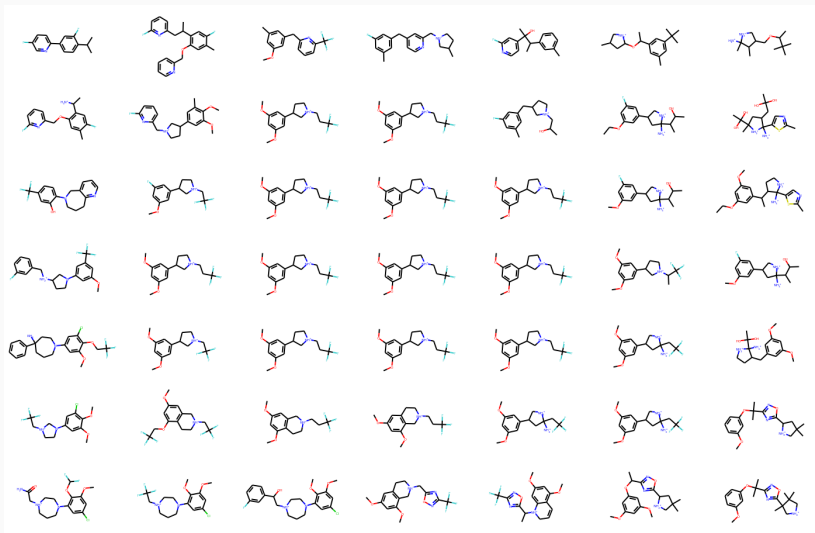
Gumbolt JT-VAE	JT-VAE
3.889	4.019

Table 1: Reconstruction Loss

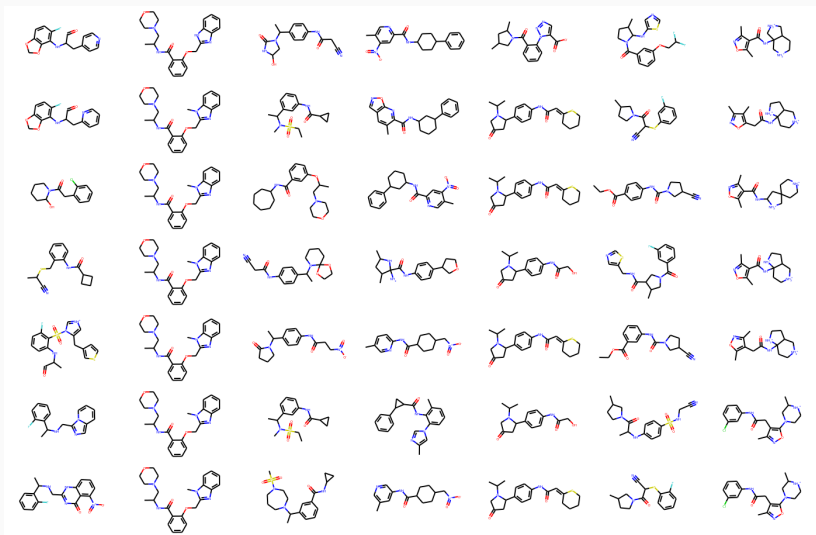
	Gumbolt JT-VAE	JT-VAE
Topological Accuracy	93.38	94.16
Classification Accuracy	99.21	97.57
Assembly Accuracy	95.34	94.31

Table 2: Model Accuracies

Qualitative Analysis - Neighborhood Sampling



Qualitative Analysis - Sampling from the Prior



1. Replace Teacher Forcing with Scheduled Sampling
2. Avoid choosing a node as root in the Junction Tree