

A Systematic Assessment of Deep Learning Models for Molecule Generation

Davide Rigoni^{1,2}, Nicolò Navarin¹ and Alessandro Sperduti¹ *

¹University of Padua - Department of Mathematics "Tullio Levi-Civita"
via Trieste 63, 35121, Padua - Italy

²Fondazione Bruno Kessler, via Sommarive 18, 38123 Trento, Italy

Abstract. In recent years the scientific community has devoted much effort in the development of deep learning models for the generation of new molecules with desirable properties (i.e. drugs). This has produced many proposals in literature. However, a systematic comparison among the different VAE methods is still missing. For this reason, we propose an extensive testbed for the evaluation of generative models for drug discovery, and we present the results obtained by many of the models proposed in literature.

1 Introduction

The chemical space is so vast that, with the computational resources available nowadays, its complete exploration is impossible. For this reason, in recent years the scientific community has devoted much effort in the study of deep learning models that are capable of generating candidate molecules that are likely to exhibit some pre-specified properties, allowing researchers to focus just on a small part of this chemical space. These methods can be used, for instance, in the process leading to the discovery of molecules that can become new drugs. Thanks to the development of increasingly effective deep learning models and the presence of large data sets, in recent years promising results have been achieved.

Starting from the model in [1], many other works have been proposed [2, 3, 4, 5, 6, 7]. However, often different works use different metrics to evaluate their models, making a fair and objective comparison among models difficult. For this reason we present a systematic approach to evaluate models for drug generation, using a precise set of metrics that aims to detect the various nuances that exist among the models. In this way it is possible to fairly evaluate the models based on the statistics and properties of the molecules they generate. This comparison is not possible just referring to the original works in literature.

In this work we focus mainly on VAE models, reporting the reader to [8] for a comparison considering other types of models.

2 Models for Molecule Generation

State-of-the-art models for drug generation are based on Variational Autoencoders (VAEs) [9] or on Generative Adversarial Networks (GANs) [10]. The

*The authors acknowledge the HPC resources of the Department of Mathematics, University of Padua, made available for conducting the research reported in this paper.

main idea of these approaches is to learn an embedding in a vector space (latent space) of the input data that aims to capture their properties and relations. Samples from the latent space are then used to generate new data that are supposed to exhibit the properties of interest. VAEs use an *encoder* to encode an input molecule in the latent space and a *decoder* to reconstruct the molecule corresponding to a point in the latent space. The GAN architecture is composed of a *discriminator* and a *generator*, both implemented by neural networks. The *generator* performs the same function as the decoder in VAE: starting from a point in the latent space, in this case sampled from a standard normal distribution, it generates the corresponding molecule. The *discriminator* aims to distinguish the input molecules that have been generated by the generator, from the ones in the training set. The generator is trained only via the discriminator model, that provides its loss function. Thanks to an ad-hoc optimization process [1, 2, 3, 4, 7] or by introducing additional terms in the loss function [5, 6], specific properties of the generated molecules can be optimized. In this work, we only focus on comparing the generation capabilities of the different models, leaving the assessment of property optimization techniques as future work.

A brief description of the models we consider follows. The Character VAE [1] exploits, in input and output, SMILES [11] strings describing the structure of the molecule. The encoder uses a convolutional network and the decoder uses a gated recurrent unit (GRU). Since every molecule can be represented as a SMILES string but not vice-versa, this approach generates many strings that do not correspond to actual molecules. Grammar VAE [2] attempts to enforce the syntactic validity of SMILES strings by introducing a context-free grammar to direct the generation process. Syntax Directed VAE [3] improves Grammar VAE using a more expressive grammar (a variant of the attribute grammar) which aims to generate strings that not only are syntactically valid, but also semantically reasonable. Junction Tree VAE [4] represents molecules using graphs, composed of chemical substructures that are extracted from the training set. New molecular graphs are obtained by first generating a tree-structured scaffold formed by substructures (the junction tree), and then combining the substructures together using a graph message passing network. Regularized Graph VAE [6] casts the molecule generation problem as a constrained optimization problem, where chemical constraints are encoded in the VAE loss function. The encoder and decoder are implemented with a convolutional and deconvolutional networks, respectively. Constrained Graph VAE [7] encodes in the latent space single atoms rather than whole molecules. To generate a molecule, first the model samples several nodes in the latent space and assigns them an atom type using a linear classifier; then it connects them using a (constrained) breadth first algorithm. Both the encoder and decoder are implemented by a gated graph sequence neural network (GGNN) [12]. MolGan [5], based on generative adversarial networks, learns via reinforcement learning to directly reconstruct, by a multi-layer perceptron, the molecular graph by predicting directly the atoms type, and the existence of bonds (and their types).

3 Evaluation Processes and Metrics

We consider two datasets of molecules: QM9 [13], composed by about 134,000 organic molecules with a maximum of 9 atoms, and ZINC [14], composed by 250,000 drug-like molecules with up to 38 atoms. We fix the training and test splits for each dataset, that we release together with the code¹. We choose a set of metrics trying to capture the strengths and weaknesses of the models (and to limit the required computational efforts). These are divided into two main categories: those that evaluate the generated molecules based on chemical properties and those that use only the information about their structure. In the latter we find: *Reconstruction* that, given an input molecule and a set of generated molecules, computes the percentage of generated molecules that are equal to the one in input; *Validity* that, given a set of generated molecules, represents the percentage of them that is *valid*, i.e. that represent actual molecules; *Novelty* that represents the percentage of new generated molecules, i.e. not in the training set; *Uniqueness* that represents (in percentage) the ability of the model to generate different molecules in output, and is computed as the size of the unique set of valid generated molecules divided the total number of valid generated molecules; *Diversity* that measures how much the generated molecules are different from those in the training set. This is a heuristic that uses randomly selected substructures present in the molecules. The metrics used to measure the properties exhibited by the generated molecules are: *Natural Product (NP)* which indicates how much the generated molecules structural space is similar to that covered by natural products [15]; *Solubility (Sol.)* which indicates how much a molecule is soluble in water; *Synthetic Accessibility Score (SAS)* which represents how easy (0) or difficult (100) it is to synthesize a molecule; *Quantitative Estimation Drug-likeness (QED)* which indicates in percentage how likely it is that the molecule is a good candidate to become a drug.

The process used to generate the molecules on which to calculate the *Reconstruction* metric consists of encoding each of the molecules in the test set 20 times (obtaining 20 slightly different representations), and decoding each of these points only once. This process was chosen because both the encoder and the decoder always contain a probabilistic component and in this way we estimate the model’s ability to reconstruct the molecule considering both factors. *GAN* model cannot compute the reconstruction metric since it cannot generate the latent space representation of a molecule. Since we are interested in the generation of new molecules, the other metrics are computed by another process that consists of directly sampling 20,000 points from the standard normal distribution and decoding each point only once. Since the considered models require a high computational load, we adopted the hyperparameters values reported in the original papers, when available. For Graph VAE and Regularized Graph VAE the hyperparameters for the ZINC dataset were not specified. We thus decided to keep the same values provided for QM9 and, after preliminary results, to double the number of epochs.

¹<https://github.com/drigoni/ComparisonsDGM>.

Model trained on QM9										
	↑%Rec.	↑%Val.	↑%Nov.	↑%Uniq.	↑%Div.	↑%NP	↑%Sol.	↓%SAS	↑%QED	
Character VAE	2.99 ±17.02	6.41 ±24.48	99.38 ±7.88	92.27	98.03 ±9.25	81.92 ±11.40	32.14 ±25.13	43.48 ±30.32	30.30 ±15.95	
Grammar VAE	58.54 ±49.27	4.45 ±20.73	94.22 ±23.33	83.22	98.88 ±7.71	79.91 ±14.60	28.40 ±20.66	33.75 ±31.15	31.91 ±12.05	
Syntax Directed VAE	52.54 ±49.94	15.00 ±35.71	100.00 ±0	100.00	97.66 ±4.90	82.60 ±14.67	26.99 ±22.11	22.92 ±35.15	35.03 ±11.18	
Graph VAE*	0.60 ±7.72	89.06 ±31.22	42.75 ±49.47	85.74	66.94 ±28.94	94.96 ±10.61	37.28 ±13.54	32.11 ±23.51	48.34 ±7.67	
Regularized GVAE*	0.66 ±8.09	87.71 ±32.83	41.26 ±49.23	83.13	63.00 ±27.91	96.32 ±9.00	37.85 ±13.24	28.89 ±23.40	48.81 ±7.14	
Junction Tree VAE	53.88 ±49.85	91.14 ±2.24	99.95 ±28.36	90.27	57.60 ±30.75	91.46 ±15.23	27.05 ±13.59	18.97 ±20.70	46.15 ±7.88	
Constrained GVAE*	33.86 ±47.32	100.00 ±0	92.82 ±25.82	98.86	79.13 ±21.61	93.05 ±12.26	27.96 ±13.36	13.81 ±19.20	46.78 ±21.30	
MolGAN*	NA	76.74 ±42.25	56.22 ±49.61	20.00	61.11 ±35.94	96.27 ±41.30	31.62 ±17.13	31.09 ±21.28	48.39 ±14.58	
Properties' Scores for Dataset QM9						88.52 ±17.75	27.91 ±13.76	21.86 ±22.88	46.12 ±7.76	
Model trained on ZINC										
	↑%Rec.	↑%Val.	↑%Nov.	↑%Uniq.	↑%Div.	↑%NP	↑%Sol.	↓%SAS	↑%QED	
Character VAE*	25.28 ±43.46	0.93 ±9.60	100.00 ±0	91.40	98.19 ±7.02	80.82 ±12.83	29.60 ±17.60	31.11 ±30.14	38.70 ±10.63	
Grammar VAE*	55.82 ±49.66	5.06 ±22.99	100.00 ±0	94.64	99.21 ±4.47	80.99 ±11.40	50.24 ±33.65	26.75 ±33.14	25.42 ±14.91	
Syntax Directed VAE*	77.38 ±41.84	19.00 ±39.23	100.00 ±0	100.00	93.56 ±18.50	77.84 ±19.76	55.94 ±27.51	14.46 ±24.14	39.45 ±20.98	
Graph VAE	0.27 ±4.58	62.63 ±48.38	100.00 ±0	99.99	71.49 ±25.36	90.68 ±11.71	80.79 ±17.33	28.07 ±20.14	45.96 ±18.69	
Regularized GVAE	0.01 ±0.77	86.47 ±34.21	100.00 ±0	90.33	97.88 ±6.96	95.88 ±6.84	94.42 ±9.61	44.64 ±25.14	34.41 ±13.26	
Junction Tree VAE*	50.23 ±50.00	99.59 ±6.35	99.98 ±1.23	99.75	32.96 ±21.78	52.20 ±17.12	48.06 ±18.48	44.74 ±24.39	75.05 ±13.40	
Constrained GVAE*	0.35 ±5.91	100.00 ±0	100.00 ±0	99.92	65.98 ±22.78	81.38 ±15.98	57.76 ±20.04	16.25 ±21.63	65.14 ±16.39	
Properties' Scores for Dataset ZINC						42.08 ±18.37	56.11 ±17.44	55.95 ±22.90	73.18 ±13.86	

Table 1: Average and standard deviation of different metrics computed on the QM9 and ZINC dataset. The symbol '*' denotes models where we used values for the parameters tuned by the authors; entries with blue background highlight the best score; up and down arrows denote whether the metric should be maximized (↑) or minimized (↓).

4 Results

Tables 1 reports the average and standard deviation of the experimental assessments (using the procedure in Section 3) on the models presented in Section 2. We have shown the strengths and weaknesses of existing models and for the first time we have also reported the standard deviation. The last line of each table reports the properties' scores obtained from the molecules in the datasets. The models should learn the distribution of the input data and for this reason it is expected that each model scores reflect those in the datasets. Note that some of the results we report are slightly different from the ones reported in the original papers. Depending on the cases, this is due to the different evaluation procedure, high variance in the results, or bugs in the original code that we fixed. Character VAE shows low *validity* and *reconstruction* in both datasets. On the contrary, since a small modification to the SMILES string can correspond to a large modification of the molecular structure, it presents a high value of *uniqueness* and *novelty*. Grammar VAE, compared to Character VAE, increases the *reconstruction* and the *validity* scores while maintaining a high *novelty* and a high *uniqueness* in both datasets. However, in QM9 this model has a similar *validity* value as Character VAE, probably because of the observed very high variance. Syntax Directed VAE improves Grammar VAE results in both datasets, even though in the QM9 the *reconstruction* value is a bit lower than the one of Grammar VAE. This is probably due to the fact that the model parameters are not tuned on the QM9 dataset, and again to the presence of high variance. Regularized Graph VAE reaches similar results to Graph VAE in the QM9 dataset, presenting good *validity* and *uniqueness* values, but very low *reconstruction* and only about 40% of *novelty*. In the ZINC dataset, the differences w.r.t. Graph VAE model are more evident. In fact Regularized Graph VAE presents higher *validity* value and a lower *uniqueness* value. Junction Tree VAE presents high *validity*, *novelty* and *uniqueness* values, in both datasets, even if it is optimized only on ZINC. However, since it generates molecules using substructures extracted from the training set, it tends to present a lower value on the *diversity* score. Constrained Graph VAE presents high *validity*, *novelty* and *uniqueness* values, in both datasets. Considering the *reconstruction* values, this model has troubles reconstructing the complex molecules in ZINC. MolGAN is trained only on the QM9 dataset because, as reported from the authors, this model doesn't scale well with larger molecules. It presents good *validity* and *novelty* values, but low *uniqueness* due to the problem of the collapse of the model that is often present in GAN models. Since the goal of the generation process is to find new molecules with certain properties for high-throughput screening, we argue that this model is not particularly suited for this task.

Looking at the metrics measuring the chemical properties (NP, Sol.,SAS and QED), it seems that Junction Tree VAE and Constrained Graph VAE are the models able to best capture the characteristics of both datasets, even if the latter model tends to generate molecules that are simple to synthesize. Although in the QM9 dataset the considered models do not differ too much in the measures

of chemical properties, in ZINC these differences are more evident.

5 Conclusions

Deep learning models for the generation of molecules are still in their infancy and they are not properly compared to each other. In fact, different models are tested by the authors on different datasets, with different evaluation processes and metrics, making it difficult to objectively compare them. In this paper, we have proposed a set of processes for the evaluation of existing models according to relevant metrics, for which we have reported experimental results on two commonly adopted datasets. To ease future comparisons, we publicly released the code used for the reported assessment.

References

- [1] R. Gómez-Bombarelli et al. Automatic chemical design using a data-driven continuous representation of molecules. *ACS central science*, 4(2):268–276, 2018.
- [2] M. J. Kusner et al. Grammar variational autoencoder. In *Proceedings of International Conference on Machine Learning*, pages 1945–1954, 2017.
- [3] H. Dai et al. Syntax-directed variational autoencoder for structured data. In *International Conference on Learning Representations*, 2018.
- [4] W. Jin et al. Junction tree variational autoencoder for molecular graph generation. In *Proceedings of the International Conference on Machine Learning*, pages 2328–2337, 2018.
- [5] N. De Cao and T. Kipf. MolGAN: An implicit generative model for small molecular graphs. *ICML 2018 workshop on Theoretical Foundations and Applications of Deep Generative Models*, 2018.
- [6] T. Ma et al. Constrained generation of semantically valid graphs via regularizing variational autoencoders. In *Advances in NeurIPS*, pages 7113–7124, 2018.
- [7] Q. Liu et al. Constrained graph variational autoencoders for molecule design. In *Advances in NeurIPS*, pages 7806–7815, 2018.
- [8] N. Brown et al. Guacamol: benchmarking models for de novo molecular design. *J. Chem. Inf. Model*, 59(3):1096–1108, 2019.
- [9] P. K. Diederik et al. Auto-encoding variational bayes. In *Proceedings of the International Conference on Learning Representations*, 2014.
- [10] M. Arjovsky et al. Wasserstein generative adversarial networks. In *International conference on machine learning*, pages 214–223, 2017.
- [11] D. Weininger. Smiles, a chemical language and information system. 1. introduction to methodology and encoding rules. *J Chem Inf Comput Sci*, 28(1):31–36, 1988.
- [12] Y. Li et al. Gated graph sequence neural networks. In *International Conference on Learning Representations*, 2016.
- [13] R. Ramakrishnan et al. Quantum chemistry structures and properties of 134 kilo molecules. *Scientific Data*, 1, 2014.
- [14] J. J. Irwin and B. K. Shoichet. Zinc- a free database of commercially available compounds for virtual screening. *J. Chem. Inf. Model*, 45(1):177–182, 2005.
- [15] P. Ertl et al. Natural product-likeness score and its application for prioritization of compound libraries. *J. Chem. Inf. Model*, 48(1):68–74, 2008.