Yash Gandhi Lenny Fobe

Final Report: Boltzmann Generators

Final State of Project:

In the final code sprint for our implementation of Boltzmann generators we were able to complete a sizable part of the initial goals of the project. In our code's final state, it is able to simulate 3 toy molecular dynamics (MD) simulations and train Boltzmann generators using a variety of loss functions on short MD simulations. The Boltzmann generator training process can be monitored using TensorBoard and final trained models can be logged and hosted online using the Streamlit app builder. Informative output can be generated, such as loss graphs, generated sample distributions and free energy profiles as the Botlzman generator is trained. Most behaviors of this program are manifested as interactions between objects.

While getting the molecular simulation and machine learning parts of the program working together in this final code sprint, the construction of the various loss functions used in the paper required a closer relationship between the SimulationDate object and the Model objects. Due to the need to calculate energy of generated configurations on the fly while training, we were only able to train Boltzman generators for systems we simulated in our simulation_library. If we wanted to evaluate configurations from a real MD engine, like Gromacs, we would have to build an adapter which would wrap Gromacs and return values for the common interfaces needed to interact with our program. In the time left, we were not able to complete this task, and focused more energy on completing goals on the other toy systems. Also, because of time constraints, we were not able to fully train a RealNVP model on the dimer simulation. This simulation requires many more samples, has many higher dimensions (72 dimensions versus 2 dimensions for other simulations), and requires a significant search on the hyperparameter space, but the code does allow for this type of model to be trained.

Final Class Diagram:

Our final class diagram is attached below, along with previous iterations of our class diagram. We also added a class diagram for the facade pattern we used for the simulation objects used in this project.

We used a number of object oriented design principles and patterns in our project. Some of these include:

- 1.) We used inheritance in a number of ways. There were multiple levels of inheritance when building the neural network because the RealNVP network consists of multiple layers each with their own set of neural network blocks. Each section - the neural network block, the single RealNVP layer, and the RealNVP model - inherit from base tensorflow classes and then uses references to build a working network.
- 2.) When designing the components for tensorboard, it became clear that not all types of logging were necessary for all types of training. Therefore, we built the tensorboard logging as a decorator pattern. This way the user can decide the different types of logging they would like to use to keep track of training the model.
- 3.) As a complement to the logging, we implemented a subject/observer pattern to update the tensorboard at every training iteration. Once a single training iteration happened, the network

- called update on the tensorboard and the tensorboard got a dictionary of values to plot. With an internal processing of each value, it appropriately plots the results from that iteration.
- 4.) For our model, we are often changing the loss and optimizer functions to test new methods for different data sets. Therefore, we implemented a strategy pattern for both the loss and optimizers. This allows us to easily add new loss functions and optimizers without worrying about the repercussions on the network because they are loosely coupled. We have 4 different loss functions, two of which inherit from one of the loss functions and 3 optimization functions (although we did find that RMSProp was often the only useful optimization function, we did not do a full search on the hyperparameters, so we left the other optimization functions).
- 5.) The Integrator, Thermostat, BoundaryCondition and Potential objects use the strategy pattern. Both have common interfaces which are used by the System and Integrator objects to easily swap between integrator and thermostat methods. In this project we have 4 integrator methods implemented and 1 thermostat implemented. With the use of the strategy pattern, more integrators and thermostats can easily be added to the program without having to rewrite other objects.
- 6.) The Thermostat, Integrator and System objects can all be built from simple factory patterns, where their construction is put into a dedicated class. This simplifies the instantiation of these classes to a simple passing of strings and keyword arguments.
- 7.) The Logger objects used for recording data from the System class were modeled after a Subject/Observer pattern. The concrete observers (CoordinateLogger, EnergyLogger, DistanceLogger, etc.) all observe the System and request pertinent data from the system after a certain number of simulation steps.

Class Diagram Comparison:

At the beginning of the project, we already had an idea about what types of classes and the overall structure we would need. Our early class diagram, while simplistic, did provide a good general overview of the project. The objects responsible for simulations, model parameters, and data output were all clearly segregated into various objects that interacted. As the development progressed, the intricacies of the project revealed themselves and we slowly incorporated more and more detailed class structure into our diagram. This included all of the sub-modules that the simulation objects used and the number of different loss functions and possible logging techniques we would use to track the model during training. Now, we have included a number of different, and more detailed, classes into our diagram and the connections between different classes and sections of the program have become more defined, as we implemented them.

Third-Party Code Statement:

Python has a host of libraries and API to extend different programs. Our project uses a number of these packages especially for scientific calculations and neural network training. We use the following pacakges throughout our code:

- 1.) tensorflow == 2.1.0
- 2.) tensorflow-probability == 0.8.0
- 3.) numpy == 1.17.0
- 4.) streamlit == 0.57.3

- 5.) pyyaml == 5.1.2
- 6.) scipy == 1.4.1
- 7.) matplotlib == 3.1.1
- 8.) Pillow == 6.2.0

In addition to these packages, we use examples from RealNVP tutorials to build our architecture and get the functions for generating the double moon distribution. The link to the specific tutorial is within the network_base.py file. The simulation library we used in this program was in part developed for Computational Statistical Physics (PHYS7810) and modified and expanded for the purposes of this project. All other code is original code built for this project.

OOAD Process:

- 1. Before we became fully invested in this project, we spent time researching RealNVP networks, neural network normalizing flows, molecular dynamics, and other related topics to better understand the underlying equations and algorithms that we were going to be implementing. This analysis of the field gave us guidelines for how the code should behave if it had been designed correctly. By identifying the relevant components that drive the research, we were able to conceptualize the problems within the research and formalize the methods that we would need to develop to build a software that could address those problems.
- 2. We used Trello, a Kanban-like workflow manager, to delegate tasks and notify each other when coding tasks were completed. This development tool ensured we were both on the same page when meeting to discuss code.
- 3. The required UML diagrams in project 4 greatly helped us organize the class structure and data responsibility of our program early on in the development process. Building the sequence diagrams had us think about the order of interactions between high level objects. The activity and use case diagrams gave us general tasks we wanted our system to be able to accomplish and how users would interact with our code, well before any coding was done. The architecture diagram gave us an initial view as to how data would be segregated between various levels of abstraction.
- 4. The multiple planning and coding sprints we did throughout projects 4, 5 and 6, resulted in a project on which we had iterated on multiple ideas. This iterative planning resulted in some improvements from coding spring to coding sprint, as we encountered problems along the way. Notably, the initial independence of simulation logic and neural network logic was scrapped as we learned the loss functions in the paper required energy calculation on the fly of generated configurations.







