



PhD thesis evaluation

The thesis "*Machine learning through geometric mechanics and thermodynamics*" by Martin Šípka presents the work published in three enclosed publications on several related topics of the use machine learning.

The introduction to the thesis is on the shorter side but is written well. It gives a solid overview of the broader context of the work. The introductory paragraphs of the individual chapters are also clear and open the given topics appropriately. The level of detail is a little uneven throughout the thesis. Some topics are explained quite thoroughly and systematically, while other topics of similar relevance are mentioned just in passing, if at all. I would have liked a summary of the contributions of the student to the individual enclosed papers. They are all first-author publications, so assumptions of major contributions seem safe, but given that there are multiple co-authors, it would still be nice to have this stated explicitly and specifically.

The chapter on machine learning in science introduces a basic perceptron as well as other ML concepts clearly and at an appropriate level of detail. One thing that is potentially missing is a mention of the importance of the training set and the need to obtain suitable training data.

In the chapter on Hamiltonian systems, the text becomes quite terse and non-systematic rather quickly. A more complete explanation of the connection of Hamiltonian dynamics and machine learning would be beneficial, including a summary of how training is performed. The terms that are modeled using neural networks are mentioned basically by the way – for example in the sentence starting with "*Since our L is represented by a neural network...*" on page 14. I liked the analysis of dissipative dynamics in section 2.4.

In the chapter on chemical reactions, the transition from PCA to VAEs is a little sudden. Given the importance of the problem of dimensionality reduction, perhaps a brief mention of other existing approaches with some references would give the reader a broader perspective. The explanation of VAEs is done well and is one of the parts that is more detailed than the rest of the thesis. The important concept of an ML interaction potential is introduced basically by the way, without much detail or emphasis. Given how important it is, also as a starting point for this work, perhaps it could have received more attention.

The Conclusion chapter summarizes the work well and also identifies open problems and possible future directions.

The text of the thesis is written in clear English and reads well, with very few mistakes, mostly some typos and minor mistakes in formulations or word choices. The overall quality of the document and figures is quite high, with a few minor typesetting or style issues (for example slanted text, em dashes, references to figure or citation numbers, or capitalization in the list of references).

The quality of both the thesis itself and the published work is clearly deserving of granting a PhD degree. The topics are at the forefront of the application of machine learning to problems in physics and chemistry and the published work presents substantial new contributions. The thesis shows the candidate's ability to independently formulate and present his research and to offer a broader perspective.

Fyzikální ústav UK

Finally, I have some questions and starting points for discussion that I would like the candidate to address.

1. For Hamiltonian systems, can you comment on the possibility of training on derivatives but modelling the scalar value of the Hamiltonian, analogous to ML potentials?
2. Have you tried, and would there be any benefit to, assuming some known parts of the Hamiltonian, such as the kinetic energy, and learning the rest?
3. On page 23 of the thesis, you say that a “transition path is not hard to find with tools like the string method or metadynamics”. There are crucial differences between these two methods, though – can you elaborate on these differences, especially when it comes to dimensionality?
4. Representations of atomic environments (whether learned or not) discard some information and do not allow for a full reconstruction of atomic positions. Does this impact the VAEs used in this work in any way?
5. In the field of enhanced sampling simulations, the problem of determining whether a given collective variable is a good reaction coordinate or order parameter is a well-known issue. What are the criteria for a good reaction coordinate, given its purpose? Can you talk about some options that are available (regardless of ML) to assess a given collective coordinate in this respect and whether and how they could be used in conjunction with your learned collective variables?
6. Even with an ideal reaction coordinate and ideal bias (i.e. one that exactly compensates the corresponding free energy) along it, diffusion in this degree of freedom can still be slow in a high dimensional system. Is the situation potentially different in the DiffSim approach?

RNDr. Ondřej Maršálek, Ph.D.

Fyzikální ústav Univerzity Karlovy
Matematicko-fyzikální fakulta Univerzity Karlovy