# MACHINE LEARNING TO IDENTIFY VARIABLES IN THERMODYNAMICALLY SMALL SYSTEMS

David M. Ford[*], Aditya Dendukuri, Gulce Kalyoncu, Khoa Luu, and Matthew J. Patitz
University of Arkansas - Fayetteville
Fayetteville, AR 72701

*Abstract Overview*

Thermodynamically small systems, with a number $N$ of interacting particles in the range of 1-1000, are increasingly of interest in science and engineering. While the thermodynamic formalism for bulk systems, where $N$ approaches infinity, was established long ago, the thermodynamics of small systems is currently approached by adding new variables in a somewhat *ad hoc* fashion. We propose a more rigorous approach based on machine learning (ML), which we demonstrate by applying both supervised (neural network) and unsupervised (diffusion map) ML methods to large data sets from Monte Carlo simulations of systems comprising $N$=3 Lennard-Jones particles in a three-dimensional periodic box. The ML methods clearly identify structural and energetic changes that occur in this model system and effectively reduce the dimensionality from nine to either two or one. Work is ongoing to correlate the reduced variables with geometric properties of the original system and to study systems with a larger $N$.

*Keywords*

Artificial neural networks, diffusion maps, molecular simulation, clusters.

## Introduction

Systems with a finite number of interacting particles are increasingly of interest due to applications in nanotechnology (Hill, 2001). Such systems require an adjustment to the fundamental thermodynamic equations developed for infinite (bulk) systems. One common example is the addition of a surface energy term to describe bubbles or droplets. However, for very small systems, *e.g.* with a few to a few thousand particles, the usual definitions of surface area and even volume may become ambiguous or non-applicable, so other variables must be identified. Here we demonstrate how machine learning (ML) can be systematically applied to this task.

## Model System and Data Sets

The model system was a set of three particles interacting through the well-known Lennard-Jones (LJ) pairwise potential in a three-dimensional box with periodic boundary conditions (Allen & Tildesley, 1987). With three particles and three dimensions, this is effectively a nine-dimensional system. The box was $10\sigma$ on a side, where $\sigma$ is the LJ particle diameter. The system was simulated using the conventional canonical Monte Carlo algorithm at a dimensionless temperature $T^* = kT/\varepsilon = 0.18$, where $k$ is the Boltzmann constant and $\varepsilon$ is the LJ energy parameter characterizing the well depth or 'bond' strength. The potential energy was tracked throughout the simulation trajectory and configuration 'snapshots' were saved at regular intervals for later analysis by the ML techniques. In the following, the potential energy is reported in a dimensionless, per particle basis as $u^* = U/3\varepsilon$, where $U$ is the total potential energy arising from the sum of the three pairwise contributions.

---

[*] To whom all correspondence should be addressed, daveford@uark.edu

Figure 1(a) shows the different structural motifs that were observed in snapshots of the system, ranging from tightly clustered to completely dissociated. Figure 1(b) shows a histogram of the potential energies that were sampled in the simulation trajectory. Clearly the system sampled the range of possible values, since a tightly clustered state will have three pairwise 'bonds,' each contributing $-\varepsilon$, for a value of $u^* = -1$, while a completely dissociated state will have no 'bonds' for a value of $u^* = 0$. In fact, we chose the temperature value $T^* = 0.18$ because all of these different structural motifs and potential energies were observed in a single trajectory. At this point one might attempt to assign peaks in the histogram to specific structural motifs, but as we will see from the ML results below, such assignment is not straightforward.



*Figure 1. Results from the Monte Carlo simulation. (a) Representative snapshots showing the four major observed structural motifs: (i) tightly clustered, (ii) linear, (iii) partially dissociated, and (iv) completely dissociated. (b) Histogram of dimensionless potential energy u\*.*

**Machine Learning Methods**

We employed two ML methods, namely neural networks and diffusion maps, to see if the configurational data sets from the Monte Carlo simulations could be used to reduce the full nine dimensions of the model (*xyz* coordinates of the three particles) to a lower-dimensional space.

Neural networks (Russel *et al*., 2009) are designed to model complex non-linear transformations. We designed a neural network to be a dimensionality reduction machine. The neural network had four layers, with the first layer being the input layer (***H***) which we defined to be the high dimensional space of the three *xyz* coordinates in a 'snapshot' (dim = 9). The second and third layers (dim = 11 and 6 respectively) were meant for pattern recognition and the fourth layer (dim = 2) was defined as our reduced dimension (***L***). The output layer (dim = 1) is defined to be $u^*$. The network was trained to accurately predict $u^*$ using the gradient-descent optimization technique. The result was the following mapping:

*For every point x in **H**, there exists a point x\* in **L** such that:*
$$u^*(x) = u^*(x^*); \quad dim(L) << dim(H).$$

Diffusion maps (Coifman et al., 2005) was the other ML method employed in this work. In DMap, a kernel matrix is constructed for the data set based on the distance $d_{ij}$ between the $(i,j)$ pairs of data points, as

$$K_{ij} = e^{d_{ij}^2/2\delta^2} \tag{1}$$

where $\delta$ is the kernel bandwidth that sets the scale of connectivity for the data set. A Markov matrix $M_{ij}$ can be created from a proper normalization of the kernel matrix, and a spectral analysis of $M_{ij}$ will indicate whether the data support the existence of a lower-dimensional manifold, and if so, provide a representation of the data in this low-dimensional space.

The definition of the distance between data points, $d_{ij}$, is an important step in the DMap process. For our model system this step is nontrivial, as we must define the distance between two 'snapshots' comprising the *xyz* positions of three particles. We considered three different definitions of distance. The first was Hausdorff distance, which is a simple and general metric for the distance between sets of points in Euclidean space. Hausdorff distance has been used in past studies of particle clusters (Bevan *et al*., 2015); as in the previous work, here each configuration was mean-centered and aligned by its principal axes of rotation before analysis. The second distance definition was based on potential energy. Specifically, each configuration was assigned a numerical value $F = \exp[-u^*]-1$, and the distance between two configurations was defined as the difference between their respective $F$ values. The third distance definition was based on the pairwise connectivity of the particles through a spectral graph matching algorithm called IsoRank (Long *et al*., 2014). Each configuration was represented by a 3x3 matrix with binary entries describing the existence (or lack thereof) of a 'bond' between the *ij* particle pair, based on a threshold distance. The distance between two configurations was computed as a metric of the difference between the two matrices. We expect that our DMap results using the IsoRank distance should identify the four distinct states corresponding to the structural motifs shown in Fig. 1(a).

**Results**

Figure 2 shows the results from the neural network analysis. The first reduced variable is strongly correlated with the potential energy of the configuration, which is perhaps not surprising based on the training procedure. The second reduced variable is likely capturing the axial symmetry of the configurations.

*Figure 2. Results from the neural network analysis. (left) Data in the H space as shown by plotting all the xyz particle locations for each of the 5001 configurations in the set. (right) Data in the L space generated by the neural network. Both of the images are colored by u\*.*

Figure 3 shows the DMap results. Although the eigenvalue spectra in all cases suggest that at most two dimensions are needed to represent the system, the data are plotted in the coordinates of the top three eigenvectors for better contextual understanding. Figure 3(a) shows that first nontrivial eigenvector, $v_2$, successfully segregates the configurations that have lower potential energy values (blue) from the ones with moderate (green) and high (yellow) values. However, this eigenvector does not discriminate between the moderate and high energy structures. The next nontrivial eigenvector, $v_3$, does not resolve the data any further and seems to be capturing only a symmetry effect. These results are consistent with the nature of the Hausdorff distance. Figure 3(b) indicates that only one dimension is needed to describe the data and that the corresponding eigenvector is strongly correlated with potential energy, which is perhaps not surprising since the distance was based directly on the potential energy function $F$. Figure 3(c) shows that employing the IsoRank distance metric indeed yields a result in which the data are separated into four distinct clusters, which we have verified as corresponding to the four structural motifs (three, two, one, or no pairwise bonds) shown in Fig. 1(a). Interestingly, the correlation of these structures with the potential energy values is much weaker than we originally anticipated. This can be seen more clearly in the parametric plot of Fig. 3(d), where the substantial overlap of $u\*$ values between neighboring structures is evident.

## Conclusions

ML can play a valuable role in identifying the number and type of variables needed to describe the thermodynamics of small systems, as demonstrated here using large data sets of 'snapshots' from a Monte Carlo simulation of three interacting Lennard-Jones particles in a three-dimensional periodic box. A neural network that enforced the equality of potential energy in the high- and low-dimensional spaces reduced the dimensionality from nine to two. Diffusion maps also identified one- or two-dimensional spaces, although the results differed based on the definition of distance between configurations. The next step in this work is to identify geometric properties of the system that correlate with the reduced-space variables identified by the ML methods.



*Figure 3. (a, b, c) The configuration data points embedded in the space of the top three eigenvectors from DMap analysis based on the following distance metrics: (a) Hausdorff, (b) energy function F, and (c) IsoRank. Coloring is by value of the potential energy. (d) IsoRank data parametrically plotted as potential energy value vs. second eigenvector value.*

## References

Allen, M.P. and Tildesley, D.J. (1987). Computer Simulation of Liquids, Oxford University Press.

Bevan, M.A., Ford, D.M., Grover, M.A., Shapiro, B., Maroudas, D., Yang, Y., Thyagarajan, R., Tang, X., and Sehgal, R.M. (2015). Controlling assembly of colloidal particles into structured objects: Basic strategy and a case study. *Journal of Process Control* 27, pp. 64-75

Coifman, R.R., Lafon, S., Lee, A.B., Maggioni, M., Nadler, B., Warner, F. and Zucker, S.W. (2005). Geometric diffusions as a tool for harmonic analysis and structure definition of data: Multiscale methods. *Proceedings of the National Academy of Sciences*, *102*(21), pp.7432-7437.

Hill, T.L. (2001). Perspective: Nanothermodynamics, *Nano Letters* 1(3), pp. 111-112

Long, A.W. and Ferguson, A.L. (2014). Nonlinear Machine Learning of Patchy Colloid Self-Assembly Pathways and Mechanisms, *Journal of Physical Chemistry B* 118 (15), 4228-4244

Russell, S. and Norvig, P. (2009). Artificial Intelligence: A Modern Approach (3rd ed.), Prentice Hall Press.