

Iah Sci

Methods for classification and feature extraction

#### Supervised Learning

Multilaver Perceptron(MLP)

- Feedforward artificial neural network with fully connected layers.
- Implementation from sklearn [1] with a stochastic gradient

# Learning Important Features from **Molecular Simulations**

## OLIVER FLEETWOOD, MARINA KASIMOVA, ANNIE WESTERLUND & LUCIE DELEMOTTE

Abstract: Biomolecular simulations are intrinsically high dimensional and generate datasets of ever increasing size. Reducing the number of features in the dataset and gaining insight into the biophysical properties of molecular states is currently a big challenge that many scientists face on a regular basis. Following the recent years' rising interest in machine learning methods there are now many powerful dimensionality reduction tools easily accessible, although such methods are often criticized to resemble black boxes and provide limited human-interpretable insight. In this study we demonstrate how a number of methods from supervised as well as unsupervised machine learning can learn ensemble properties from molecular simulations and provide easily interpretable metrics of what features are actually important. In order to show which methods perform best under different circumstances we first test the performance using a toy model designed to mimic real macromolecular behavior. Finally, we apply the methods to simulations of the β<sub>2</sub> adrenergic receptor to gain insights into its activation mechanism and the effect of ligand binding. The results demonstrate how machine learning methods can produce valuable insights into properties of biomolecular states and we anticipate that our approach can be useful to aid many researchers in demystifying complex simulations.



#### **Unsupervised Learning**

#### **Principal Component Analysis (PCA)**

- Converts the input features to an orthogonal set of linearly uncorrelated variables (principal components) through an orthogonal transformation.
- The first component has the largest variance possible.
- The feature importance is taken to be the components' coefficients of a feature times the variance covered by (i.e. the eigenvalues) the components.

**Restricted Boltzmann Machine** (RBM)

- A generative stochastic neural network trained to maximize the likelihood of the data using a graphical model with a layer of hidden nodes connected to the input nodes.
- Important features are identified using Deep Taylor Decomposition.

Autoencoder (AE)

• A generative neural network trained to reconstruct the features through a set of hidden layers of lower dimensionality. Important features are identified using Deep Taylor Decomposition.



- PCA on the other hand, typically performs well on identifying linear displacements.

Nonlinear displacement + random rotation

### Applications to the $\beta_2$ adrenergic receptor



Hallmarks of  $\beta_2$  activation include the outward movement of transmembrane helix 6 (TM6), the collapse of the cavity around the highly conserved residue Asp79<sup>2.50</sup>, as well as the twist of the conserved NPxxY motif at the bottom of TM7.

**UNSUPERVISED RBM** 

an example of how unsupervised learning

can help demystify molecular simulations.

blue and most important residues in red.





We identify important features using spectral clustering [3] on a trajectory from [4] to derive a number of clusters along the activation path and see which features are important for cluster classification. Similar profiles are obtained for MLP and KL.

### Learnings

- All methods identify that the G-protein binding site undergoes the most significant conformational change, especially TM6.
- The importance profiles change upon increase of the number of clusters, revealing more features important for activation.
- Regarding the effect of ligand binding we see that there are

#### References

- 1. Scikit-learn: Machine Learning in Python, Pedregosa et al., *JMLR* 12, pp. 2825-2830, 2011.
- 2. G. Montavon, S. Lapuschkin, A. Binder, W. Samek, K.-R. Müller, "Explaining NonLinear Classification Decisions with Deep Taylor Decomposition", *Pattern Recognition*, 65:211–222, 2017.
- 3. Westerlund, Annie M., and Lucie Delemotte. "Effect of Ca2+ on the promiscuous target-protein binding of calmodulin." PLoS computational biology 14.4 (2018): e1006072.
- 4. Dror R.O., Arlow D.H., Maragakis P., Mildorf T.J., Pan A.C., Xu H., Borhani D.W., Shaw D.E., "Activation mechanism of the beta2-adrenergic receptor", Proc. Natl. Acad. Sci. USA. 2011;108(46).
- 5. "File:Autoencoder structure.png." Wikimedia Common



Important residues identified by an RBM as We also learn to discriminate between trajectories along the activation path with (holo) and without (apo) an agonist ligand bound. Least important residues for classification are Least important residues are highlighted in highlighted in blue and most important residues in red.

SUPERVISED MLP

important residues in the extracellular region close to the ligand binding site.

- The corresponding input features are more strongly activated when a ligand is present.
- There is also a strong signal in the G-protein binding site, which illustrates how it is allosterically coupled to the ligand binding site.

KTH ROYAL INSTITUTE OF TECHNOLOGY

Oliver Fleetwood	Science for Life Laboratory	Tel: +46 - 708613650
PhD Student	Stockholm	Email: oliver.fleetwood@scilifelab.se
Delemottelab	Sweden	