



Convolutional neural networks provide a promising approach to scoring protein-ligand interactions. Neural networks are inherently difficult to analyze and understand, often being called 'black boxes'. Running a given protein-ligand complex through a network simply produces a value from 0 to 1 to represent the probability that the pose is correct, without providing any insight as to how that number was generated. Visualizations of the neural network's decision-making process allow for the analysis of its understanding of chemical interactions. We describe multiple visualization workflows and provide examples of how the resulting visualizations can aid in biological and chemical understanding of the interaction while also discussing current limitations of the method.





The trained network substantially outperforms AutoDock Vina at selecting and ranking poses. Visualizations of the network provide insights into what the network has learned and how it processes specific features.



References

Ragoza, M., Hochuli, J., Idrobo, E., Sunseri, J. and Koes, D.R., 2016. Protein-Ligand Scoring with Convolutional Neural Networks. arXiv preprint arXiv:1612.02751

Samek, W., Binder, A., Montavon, G., Lapuschkin, S. and Müller, K.R., 2016. Evaluating the visualization of what a deep neural network has *learned*. IEEE Transactions on Neural Networks and Learning Systems. Bach, S., Binder, A., Montavon, G., Klauschen, F., Müller, K.R. and Samek, W., 2015. On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. PloS one, 10(7), p.e0130140.

Visualization of Convolutional Neural Network **Scoring of Protein-Ligand Binding**

Joshua Hochuli^{1,2}, Matthew Ragoza³, and David Ryan Koes³ ¹Department of Computer Science, ²Department of Biological Sciences, ³Department of Computational and Systems Biology, University of Pittsburgh, Pittsburgh, PA



Pros:

- Easy to implement no need to modify CNN
- Can evaluate large modifications (fragments)

Cons:

- Evaluates non-physical structures
- Inefficient requires many CNN evaluations
- No volumetric output





How the CNN changes the structure to score better

Backpropagation propagates the gradient of the loss throughout the network and can ultimately calculate the gradient with respect to the input values.

$$\delta^N =
abla_a L \odot \sigma'(z^N)$$

Compute output error

 $\delta^{\iota} = ((w^{\iota+1})^T \delta^{\iota+1}) \odot \sigma'(z^{\iota})$ $rac{\partial L}{\partial x_j} = \delta_j^0$

Calculate gradient with respect to input

Backpropagate

$$rac{\partial L}{\partial A_{\{x,y,z\}}} = \sum_j rac{\partial L}{\partial x_j} rac{\partial x_j}{\partial dist_{A,x_j}} rac{\partial dist_{A,x_j}}{\partial A_{\{x,y,z\}}}$$

Pros:

- Implemented in a single backwards test
- Indicates how to improve structure
- Volumetric visualization
- Provides directionality

Cons:

 Intermixes positive and negative effects • Does not explain current score



Relevance propagation propagates the classifier output throughout the network as a relevance quantity. At each network level the total relevance is conserved.

The negative and positive components of the relevance are treated separately.

 R'_{i}



Relevance

How the CNN scores the unchanged structure

$$f(x) = \ldots = \sum_{d \in l+1} R_d^{(l+1)} = \sum_{d \in l} R_d^{(l)} = \ldots = \sum_d R_d^{(1)}$$

Pros:

- Implemented in a single backwards test
- May better delineate classifier decision boundary • Separates positive and negative effects
- Volumetric visualization
- Explains current score

Cons:

- Custom implementation for each network layer • Does not extrapolate beyond input
- No directionality

