# Thermodynamic profiling of protein-ligand binding energies

Application of machine learning methods in Bioinformatics

#### **Chaitanya Koppisetty**

Department of Computer Science and Engineering, Chalmers University of Technology

and

**Biognos AB** 

# Overview

- Background
- Challenges
- Tools / methods
- Results / insights
- Summary

Human body is constantly invaded by pathogens.

• "Proteins" on the surface of pathogens are vital in adhesion and proliferation.



Electron micrographs of viruses.

- Drugs / Inhibitors / ligands
  - Small molecules that prevent the adhesion or proliferation of pathogens.
    - "Relenza" is the trade name for infuenza virus inhibitors (ligand) that binds to a surface protein of influenza virus that stops proliferation.

Inhibitors



Two important properties of a drug

- Affinity
- Specificity

#### Affinity

– How strong does a drug bind to the target.

Specificity

- How specific are the drug's interactions to that particular target
  - Is it binding to other proteins in the human body?
  - Main cause of side-effects

- Binding energy
  - Strength of interaction between the protein and ligand (negative value indicates binding)



### Challenges



$$\Delta G_{Estimated} = \Delta H_{Estimated} - T\Delta S_{Estimated} \qquad eq.2$$

where

$$\Delta H_{Estimated} = \sum_{i=1}^{n} H_i f(Descriptor_i) \text{ and } T\Delta S_{Estimated} = \sum_{i=1}^{n} S_i f(Descriptor_i)$$



Accurate estimation of  $\Delta$  H and T  $\Delta$  S is necessary for precise placement of  $\Delta$  G case

Qualitative classification

- Neural networks
- Support vector machines

Quantitative estimation

– Support vector machine regression

#### Neural networks (Multilayered perceptron)



Each node/neuron in hidden layer is a

non-linear activation function

$$\mathcal{Y}i = \frac{1}{1 + e^{-S_i}}$$

Where,

*y<sub>i</sub>* is output of neuron *i*,

S<sub>i</sub> is the weighted sum of all inputs and bias to neuron i

Error back-propagation algorithms

#### Support vector machines classification



#### Support vector machines classification

Linearly difficult to classify



#### SVM regression



### **Feature Generation**



Features / descriptors based on structure based physico chemical properties of protein-ligand structures.

Training dataset of 120 protein-ligand structures.

Validation dataset of 1300 protein-ligand structures independent of training dataset.

220 Features / descriptors.

Descriptors grouped into 8 classes based on method of generation.

# **Feature Selection**

#### Feature selection / elimination

- Feature reduction using Principal Component Analysis (PCA)
  - Reduce the dimensionality of the data by fewer samples but still preserving the variance
- Backward feature elimination (BFE)

#### **Cross validation**

- 2 fold Cross validation
- Leave-one-out Cross validation
- N-fold Stratified sampling cross validation

#### **Results Classification**



#### **Confusion matrix**



Source : Wikipedia

#### **Results Classification**



#### Results – SVR prediction models



Experimental binding energies

#### Results – SVR prediction models

SVR Validation 1300 protein ligand complexes



**Experimental binding energies** 

### Application of SVR to real data

- Kernel choice and parameters
  - Linear
  - Polynomial (degree)
  - Gaussian (width parameter)
- Hyper-parameters
  - •Parameter C
  - •Parameter  $\varepsilon$

#### Parameter optimization by cross-validation and grid search

• Data Normalization (?)

### Summary



Training dataset of 120 protein-ligand structures, 220 Features / descriptors.

Validation dataset of 1300 protein-ligand structures.





SVM regression models



**MLP-NN and SVM Classifiers** 



### Thanks for listening !