## Machine learning techniques for analyzing and interpreting genomics and proteomics data

Pawar Shrikant, Ph.D.



### Sequence Analysis: NGS

• Utilizing neural networks (Restricted Boltzmann machine's) and clustering algorithms to identify certain important, representative HIV-1 PR sequences from a pool of several hundred sequences.

- 1. Analysis of drug resistance in HIV protease, *Shrikant Pawar*, Chris Freas, Robert W. Harrison, and Irene T. Weber, *BMC: Bioinformatics*
- 2. Structural studies of antiviral inhibitor with HIV-1 protease bearing drug resistant substitutions of V32I, I47V and V82I, *Shrikant Pawar*, Yuan-FangWang, Andres Wong-Sam, Johnson Agniswamy, Arun K. Ghosh, Robert W. Harrison, and *Irene T. Weber*, *Elsevier: Biochemical and Biophysical Research Communications*

Next Generation Sequencing Technology (NGS): Align unknown sample contig sequence with known entire human genome to understand known and unknown



### HIV-1 Protease Action



### Drug resistance is a severe problem



~100,000 sequences

#### Major and minor mutations

associated with resistance to all clinical protease inhibitors

Adapted from Weber, Kneller, Wong-Sam. Future Med Chem 2015



Can machine learning help in selecting few drug resistant PR sequences for structure guided drug design?

## Analysis Pipeline



#### **Hierarchical Clusters**



**Divisive Clusters** 

©Pawar/Claflin University

# Representative plot with separation of sequences for inhibitor FPV with Support Vector Machine



X and Y axis are weight vectors

- Three, five- and ten-fold accuracies for SVM and RF ranged from 0.95-0.99. RBM also selected sequences in range 0.91-0.97.
- These accuracies are better than *Yu. Et. al. 2014* where she got SVM accuracies ranging from 0.93-0.96.
- These accuracies are comparable with *Shen. Et. al. 2016* where he got RF accuracies ranging from 0.98-0.99.

Most of the high resistance fold sequences with class 2 were clustered in first 10 clusters for most of the selected inhibitors through both hierarchical and divisive clustering delineating a clean separation between non-resistant and resistant sequences.



**Divisive Clustering** 

# From a pool of 100,000 only 2-35 sequences were selected common through all the 3 approaches

Category	ATV	DRV	FPV	IDV	LPV	NFV	SQV	TPV
H, D and K	0	0	20 (66)	0	35 (61) <i>,</i>	2 (12)	5 (58),	0
					2 (31)		6 (63),	
							3 (73),	
							21 (85)	

Numbers in parenthesis are the cluster from which they were selected.

- 1. The resistance status of the selected sequences should be identified.
- 2. Minimum number of sequences selected for inhibitors, NFV, SQV or LPV would be some of the ideal candidates for testing in laboratory.

### Sequence Analysis: Microarray

• KIFCI, a novel putative prognostic biomarker for ovarian adenocarcinomas.

KIFCI, a novel putative prognostic biomarker for ovarian adenocarcinomas: delineating protein interaction networks and signaling circuitries, *Shrikant Pawar*, Shashikiran Donthamsetty, Vaishali Pannu, Padmashree Rida, Angela Ogden, Nathan Bowen, Remus Osan, Guilherme Cantuaria, and Ritu Aneja, *BMC: Journal of Ovarian Research* A centrosome clustering protein, KIFC1, predicts aggressive disease course in serous ovarian adenocarcinomas, Karuna Mittal, Da Hoon Choi, Sergey Klimov, *Shrikant Pawar*, Ramneet Kaur, Anirban K. Mitra, Meenakshi V. Gupta, Ralph Sams, Guilherme Cantuaria, Padmashree C. G. Rida, Ritu Aneja, *BMC: Journal of Ovarian Research*

## Microarray technology



## Centrosome amplification in ovarian cancer and high KIFC1 expression in ovarian cancer and normal tissue.



Increased KIFC1 expression is associated with poorer overall survival in age-specific ovarian cancer patients and pathways associated with first degree neighbors of KIFC1 protein



## Computer Vision in Biomedical Imaging

### Single shot detector application for image disease localization

- Single shot detector application for image disease localization, Shrikant Pawar, Rushikesh Chopade, Aditya Stanam, bioRxiv
- 2. Cyclical Learning Rates (CLR'S) for Improving Training Accuracies and Lowering Computational Cost, *Shrikant Pawar*, Rushikesh Chopade, Aditya Stanam, *Springer Lecture Notes in Computer Science*

## Bounding boxes for disease localization

- Object localization is a subfield of computer vision that is used to detect the location of object in an image.
- Bounding box algorithms are useful in localization of image patterns. Recently, utilization of convolutional neural networks on X-ray images has proven a promising disease prediction technique.



### Neural Network architecture



The dataset consists of 112,120 chest X-ray images, each image with a 1024\*1024pixel resolution. The images are divided into 15 classes ('No Finding', 'Atelectasis', 'Cardiomegaly', 'Consolidation', 'Effusion'; 'Emphysema', 'Edema', 'Fibrosis', 'Infiltration', 'Mass', 'Nodule', 'Pneumonia', 'Pneumothorax', 'Pleural Thickening' and 'Hernia')





Image of a patient suffering from cardiomegaly showing intersection of original condition and prediction. The prediction accuracy is >95%

Interested in learning and working with above research projects?

- Students interested in working on biomedical imaging projects will receive a \$6000/AY & summer stipend. Supported by 2023-2028 NSF SC EPSCoR ADAPT Track 1 Award.
- The research work can be a part of their thesis or project reports.
- Opportunities are also available to collaborate on microarray and sequencing projects but without stipend.

### Acknowledgment's and Collaborators











Sequence Analysis





Structural Biology



Sequence Analysis





Network Biology





Network Biology



**HPC Resources**