

# **Representation in Supervised Machine Learning** Application to Biological Problems

# **Robert Langlois**

Frank Lab

Howard Hughes Medical Institute

&

**Columbia University** 

August 9, 2011



# What is Machine Learning?

- Computer Algorithms that Learn
  - Think of it as *Programming by Example*
  - Devise a strategy to complete a task
- Example Handwriting
  - Given examples of written letters
  - Optical character recognition
  - U.S. Post Office

7	7	3	1		3	3	7	2	2
3	6	2	S	6	8	4	l	8	4
4	6	9	3	3	3.	1	5	3	6
5	1	5	3	9	3	7	7	ч	0
8	8	9	5	8	4	4	3	3	7
9	8	6	6		9	3	1	7	Y
8	3	5	9	5	1	2	1	3	4
6	133	9	100	1944	2	0	H	2	3
1	0	1	4	0	1	6	1	8	60
0	2	6	5	43	2	9	8	8	9



#### **Short Biology Review**









# **Application: What function does a protein perform?**

- Protein  $\rightarrow$  Molecular machines
  - Catalyze reactions
  - Mechanical work
  - Structural components
- Proteins that bind DNA
  - Transcription Regulation
  - Epigenetics (Encode behavior)
  - DNA Utility: Repair





# **Motivation: Protein Function Prediction**

- Large amount of sequence data
  - Number of solved genomes
  - Personal genomics ( $10^6$  bases < \$0.50)
- Large amount of structural data
  - Protein databank: Exponential Growth
  - Cryo-EM Databank: Protein complexes
- Training data (Experimental Function)
  - Filter binding assays
  - Chromatin Immuno-precipitation











# **Learning Process (Binary Classification)**





# Large Margin Linear Classifier: Support Vector Machine

• Formulation:





# **Non-linear Support Vector Machines**



B. Schölkopf, Canberra, February 2002

COLUMBIA UNIVERSITY IN THE CITY OF NEW YORK

Robert Langlois



# Representation

- Attribute: Numerical description of an object
  Where *j*=1...*m*
- Feature Vector : Fixed length set of attributes
  - Where i=1...n
- Special Attributes
  - Class: Target or unknown attribute

 $(\vec{x}_i, y_i) \hat{I} S$ 

	000		Incanter Dataset				
	Sepal Length	Sepal.Width	Petal.Length	Petal.Width	Species		
	5.1	3.5	1.4	0.2	setosa	0	
	4.9	3.0	1.4	0.2	setosa		
	4.7	3.2	1.3	0.2	setosa	U	
	4.6	3.1	1.5	0.2	setosa		
	5.0	3.6	1.4	0.2	setosa		
7	5.4	3.9	1.7	0.4	setosa		
' i—	4.6	3.4	1.4	0.3	setosa		
v	5.0	3.4	1.5	0.2	setosa		
	4.4	2.9	1.4	0.2	setosa		
	4.9	3.1	1.5	0.1	setosa		
	5.4	3.7	1.5	0.2	setosa		
	4.8	3.4	1.6	0.2	setosa		
	4.8	3.0	1.4	0.1	setosa		
	4.3	3.0	1.1	0.1	setosa		
	5.8	4.0	1.2	0.2	setosa		
	5.7	4.4	1.5	0.4	setosa		
	5.4	3.9	1.3	0.4	setosa		
	5.1	3.5	1.4	0.3	setosa		
	5.7	3.8	1.7	0.3	setosa		
	5.1	3.8	1.5	0.3	setosa		
	5.4	3.4	1.7	0.2	setosa		
	5.1	3.7	1.5	0.4	setosa		

#### http://data-sorcery.org/category/pca

 $\vec{x}_i$ 





# Outline

- Can a positive patch solve the problem?
  - Binary Classification
  - Represent entire sequence or structure
- Is there a more generic representation?
  - Binary Classification
  - Represent by residue content
- Most generic representation
  - Multiple-instance Learning
  - Represent each residue



http://ce.sharif.ir/courses/87-88/1/ce717/syllabus/Logo



# Outline

- Can a positive patch solve the problem?
  - Binary Classification
  - Represent entire sequence or structure
- Is there a more generic representation?
  - Binary Classification
  - Represent by residue content
- Most generic representation
  - Multiple-instance Learning
  - Represent each residue



http://ce.sharif.ir/courses/87-88/1/ce717/syllabus/Logo



# **Protein Representation**

- Structure
  - Size of largest positive patch (1)
  - Surface amino acid composition (20)
- Sequence
  - Charge (1)
  - Amino acid composition (20)
- Represent each protein with
  - 42 numeric features
  - -10, 1, 0, 0, 11, 5 ... 3





# Composition





# **Performance of Positive Patch + 41 other features**

- Support vector machines
  - Outperforms sequence analysis techniques like Blast
- Boosted Trees
  - Compares well compared to Support Vector Machines (SVM)
  - 10-fold cross-validation

	Accuracy	Sensitivity	Specificity	AUC
SVM	83.0	68.5	88.0	83.5
Boosted Tree	88.2	64.4	96.3	89.8
ADTree	85.1	66.7	91.3	85.3



# **Related Work: Classification with simple features**

- Structural motifs (Shanahan 2005)
- Structural motifs and HMM (Pellegrini-Calace 2005)
- Structure attributes and NN (Ahmad 2002, 2004)
- Low resolution homology and LR (Szilagyi 2006)
- Structure attributes and SVM (Bhardwaj, Langlois, 2005)
   Natures Highlights
- Structure attributes and AB, Tree, SVM (Langlois, 2006)



#### What Rules are Used?

- Unexpected Results
  - Mostly sequence features
  - Exclusionary rules
  - Patch does not dominate
- Caveats
  - Only an approximation
  - Still performs fairly well
- Back to the drawing board



Langlois, R., M. Carson, N. Bhardwaj and H. Lu. Learning to translate sequence and structure to function: Identifying DNA binding and membrane binding proteins. Annals of Biomedical Engineering. 35:1043-1052, 2007.

# Outline

- Can a positive patch solve the problem?
  - Binary Classification
  - Represent entire sequence or structure
- Is there a more generic representation?
  - Binary Classification
  - Represent by residue content
- Most generic representation
  - Multiple-instance Learning
  - Represent each residue



http://ce.sharif.ir/courses/87-88/1/ce717/syllabus/Logo

3 residue window



# **Conditional Amino Acid Composition**

- Composition of both
  - Туре
  - Environment

# • Composition of Proline given

- Helix environment
- Sheet environment
- Neither
- Both

LLL TGFAYWL I PMNGSLGAAOAYMATY IV Both Helical Sheet None Α Α Α Α С С C C D D D D Y Y Y Y

Robert E Langlois, Hui Lu (2010) Boosting the prediction and understanding of DNAbinding domains from sequence, 3149-3158. In Nucleic Acids Research 38 (10). 400



# **Performance of generic sequence features**

• Performance with Boosted Trees



- Structure features should encode more about function
- However, less information with sequence competitive

Robert E Langlois, Hui Lu (2010) Boosting the prediction and understanding of DNA-binding domains from sequence, 3149-3158. In Nucleic Acids Research 38 (10).



#### What rules do the machine learning algorithms use?





# **Genome-wide Results: 2000 Sequence < 25% Similarity**



# Outline

- Can a positive patch solve the problem?
  - Binary Classification
  - Represent entire sequence or structure
- Is there a more generic representation?
  - Binary Classification
  - Represent by residue content
- Most generic representation
  - Multiple-instance Learning
  - Represent each residue



http://ce.sharif.ir/courses/87-88/1/ce717/syllabus/Logo



# **Multiple-instance Learning**

- Examples grouped into bags
- Label on group (bag) not example
- Bag is positive if at least one example is positive
- Otherwise bag is negative
- Only weak information about positive instances!



# ©1998 by Oded Maron



# **Intuitive Example**

#### • Illustration

- Each letter is an example
- Each line indicates a bag (or group)
- Blue: Examples in a positive bag
- Red: Examples in a negative bag
- Classify bags by finding positive instances
  - Close to other instances in the positive bags
  - Far from instances in negative bags





# **Multiple Instance Learning Example**







# **Amino Acid Representation**

- Residue Identity (20)
- Secondary Structure (3)
- Structure Neighbors (20)
- PSSM for residue at position (20)
- Blosum for positions -3 ... 3 (140)
- Properties: Charge, Surface Area (2)

# **Experimental Results**

- Similar Structure based features
- Achieves over 98 AUR!
- Change in representation effective in improving accuracy.

JMB06	Blast	82.5	
	AdaBoost	93.9	
	Szilagyi (2006)	93.0	

Functional Site Discovery from Incomplete Training Data (2011) Robert Langlois Marina Langlois, and Hui Lu. *In preparation*.



# **Instance-level (Binding Residue) Results**

Reciever Operating Characteristic: Residue-DNA

False Positive Rate



Functional Site Discovery from Incomplete Training Data (2011) Robert Langlois Marina Langlois, and Hui Lu. *In preparation*.



# **Open Problems**

- How to incorporate prediction dependencies in MIL?
  - A protein (bag) is DNA-binding if **a set (not a single)** residue binds DNA
  - The prediction of one residues increases likelihood of one nearby!
  - However, both classification and MIL assume IID assumption not valid

- Can we improve learning by knowing the label/function of certain proteins?
  - Multiple-instance Active Learning
  - Expensive to label all proteins possible to label a subset
  - Learning algorithm to find the *optimal* subset



#### CURRENT STATE OF AUTOMATED PARTICLE PICKING

# Acknowledgements

- Joachim Frank (PI) Current
  - Everyone in the Frank Lab

- Hui Lu (PI) Previous
  - Everyone in the Lu Lab



http://www.mendeley.com/profiles/robert-langlois/

😼 Columbia University

**Robert Langlois**