## Machine Learning & Application in Biology

### HyunJung (Helen) Shin

shin@tuebingen.mpg.de

College of Medicine, Seoul National University, Seoul, Korea

Friedrich Miescher Laboratory Max-Planck-Society, Tuebingen, Germany

European School of Genetic Medicine, 6<sup>th</sup> Courses in Bioinformatics for Molecular Biologist, Bertinoro di Romagna, GSF-National Research Center for Environment and Health

## Goal of Lecture

# Machine Learning can alleviate the burden of solving many biological problems,

- saving the time and cost required for experiments

- providing predictions that guide new experiments.

## Goal of Lecture

# Machine Learning can alleviate the burden of solving many biological problems,

- saving the time and cost required for experiments

- providing predictions that guide new experiments.

# The goal of this tutorial is to raise awareness and comprehension of machine learning

so that biologists can properly match the task at hand to the corresponding analytical approach

## Abstract

We explore representative models, from traditional statistical models to recent machine learning models,

presenting several up-to-date research projects in bioinfomatics to exemplify how biological questions can benefit from a machine learning approach.

## Content

- 1. Basics
- 2. Tasks
- 3. Learning
- 4. Models with Examples
- 5. Evaluation and Statistical Tests

## Content

- 1. Basics
- 2. Tasks
- 3. Learning
- 4. Models with Examples





Aj attribute, feature, descriptor, input variable, predictor variable, independent variable, exogeneous variable, etc,

			$\wedge$						
1	(					(	or		
	<i>A1</i>	<b>A</b> 2	<b>A</b> 3	•••	A10	у	y		
$\boldsymbol{x}_{I}$	10	5	red	•••	1000	class1	1		
<b>X</b> 2	6	6	blue	•••	3500	class2	20		
<b>X</b> 3	7	7	yellow	•••	400	class1	45		
				•••	•••				
<b>X</b> 18	3	56	red	•••	0	class2	30		
<b>X</b> 19	15	62	red	•••	500	class1	100		
<b>X</b> 20	3	88	blue	•••	700	class2	3		
	• <b>•</b>								

0r

*Xi* input, predictor, etc.

							$\overline{\ }$
	<i>A</i> 1	A2	Аз	•••	A10	y	y
<b>x</b> 1	10	5	red		1000	class1	1
<b>x</b> 2	6	6	blue	•••	3500	class2	20
<b>X</b> 3	7	7	yellow	•••	400	class1	45
••••	•••	•••	•••	•••	•••	•••	•••
<b>X</b> 18	3	56	red	•••	0	class2	30
<b>X</b> 19	15	62	red	•••	500	class1	100
<b>X</b> 20	3	88	blue	•••	700	class2	3

\* Input set:  $X = \{x_1, x_2, ..., x_{20}\}$ 

*Yi* output variable, response, target variable, endogeneous variable, label, etc.

	<b>A</b> 1	<b>A</b> 2	Аз	•••	A10	y	y
<b>x</b> 1	10	5	red	•••	1000	class1	1
<b>X</b> 2	6	6	blue	•••	3500	class2	20
<b>X</b> 3	7	7	yellow	•••	400	class1	45
•••	•••	•••		•••	•••	•••	•••
<b>X</b> 18	3	56	red	•••	0	class2	30
<b>X</b> 19	15	62	red	•••	500	class1	100
<b>X</b> 20	3	88	blue		700	class2	3

\* Output(Target) Set:  $Y = \{y_1, y_2, ..., y_{20}\}$ 

## Content

Basics
Tasks
Learning
Models with Examples



### **Tasks**

### Prediction

- Classification
- Regression

### Description

- Clustering
- Feature Description

### Dimensionality Reduction

- Feature Selection
- Feature Extraction
- Data Reduction (Sample Selection)

### ✤ Data Integration

### Classification

Classification is concerned with the problem of separating distinct sets of data points and allocating new (test or unknown) data points to previously defined group (class)

## Classification

 		In	Target					
	4 -		4 -		4			
	AI	A2	<b>A</b> 3	•••	A10	У	У	
<b>x</b> 1	10	5	red	•••	1000	class1	1	
$\boldsymbol{x}_2$	6	6	blue		3500	class2	20	
<b>X</b> 3	7	7	yellow	•••	400	class1	45	
•••	•••	•••		•••				
<b>X</b> 20	3	88	blue	•••	700	class2	3	

Target variable (y) is categorical

Nominal : (ex) yes or no, 1 or -1 : (ex) blue, red, yellow... Ordinal : (ex) age groups (10-20, 20-30, 30-40, ...)

### \* Predicted Value

$$f \in \frac{model}{function}$$

 $f(x) \begin{cases} predicted value (output), \\ output, \\ score, \\ etc \end{cases}$ 

### Classification





# Regression is concerned with the problem of **predicting** the value of continuous target variable.

### Regression

Input 2												
	<u>A1</u>	<b>A</b> 2	Аз	•••	A10	у	y					
<b>x</b> 1	10	5	red	•••	1000	class1	1					
<b>x</b> 2	6	6	blue	•••	3500	class2	20					
<b>X</b> 3	7	7	yellow	•••	400	class1	45					
•••	•••	•••		•••	•••							
<b>X</b> 20	3	88	blue	•••	700	class2	3					

*Target variable* (y) *is continuous* 

### Regression



HyunJung (Helen) Shin, Max Planck Society, European School of Genetic Medicine, 03. 2006 21

\* Scales of Variable



\*\* Classification can be regarded as a subset of Regression in viewpoint of modeling (not task).

Clustering is concerned with the identification of groups of similar data points based on <u>similarity measures</u>.

Clustering is concerned with the identification of groups of similar data points based on <u>similarity measures</u>.

Clustering is distinct from classification in that

• Classification pertains to a known number of groups and its operational objective is to assign new data points to one of these groups

• Clustering makes no assumption concerning the number of groups

r			In	No Target !!			
		<b>A</b> 1	<b>A</b> 2	Аз	•••	A10	
	<b>x</b> 1	10	5	red	•••	1000	
	<b>x</b> 2	6	6	blue	•••	3500	
	<b>X</b> 3	7	7	yellow	•••	400	
	•••	•••	•••		•••		
	<b>x</b> 20	3	88	blue	•••	700	

*No target variable (y)* 



2 clusters



HyunJung (Helen) Shin, Max Planck Society, European School of Genetic Medicine, 03. 2006

### **Dimensionality Reduction**

Dimensionality reduction is concerned with the process which removes irrelevant or redundant features (attributes) from the original feature set, in order to avoid "curse of dimensionality" complication of learning process, erroneous results, computational burden.

\* note: irrelevant or redundant for learning or modeling

### **Dimensionality Reduction**

### • Feature Selection

A process of finding a subset of relevant features (attributes) from the original set of features.

(Ex) Selected Features: A1, A1000

	Aı	<b>A</b> 2	Аз	•••	A 1000	у			Aı	A 1000	y
<b>x</b> 1	10	5	red		1000	1	$\int f$	<b>x</b> 1	10	1000	1
<b>x</b> 2	6	6	blue		3500	20		<b>x</b> 2	6	3500	20
<b>x</b> 3	7	7	yello w	•••	400	45		<b>X</b> 3	7	400	45
•••		•••		•••		•••			•••		•••
<b>x</b> 20	3	88	blue		700	3		<b>X</b> 20	3	700	3

### **Dimensionality** Reduction

### Feature Extraction

A process of defining new descriptors (features) condensed via transformations of the raw features. The descriptors are represented as the features in the new feature space

(Ex) Extracted Features:  $\begin{cases} P_1 = \beta_1 A_1 + \beta_1 A_2 + \beta_1 A_1 A_{350} \\ P_2 = \varphi(A_1, A_2, \dots, A_{1000}) \end{cases}$ 



HyunJung (Helen) Shin, Max Planck Society, European School of Genetic Medicine, 03. 2006

### Data Reduction (Sample Selection)

Data Reduction is concerned with the process which removes irrelevant or redundant "data points" from the original data set,

in order to avoid complication of learning process or computational burden.

\* note: irrelevant or redundant for learning or modeling

### Data Reduction (Sample Selection)

### (Ex) Selected Data Points: x2, x3, x100, x9999

	<i>A1</i>	<b>A</b> 2	Аз	•••	<b>A</b> 10	У								
<b>x</b> 1	10	5	red		1000	class1								
<b>x</b> 2	6	6	blue		3500	class2								
<b>x</b> 3	7	7	yellow		400	class1	C		Aı	<b>A</b> 2	Аз		A10	у
	•••	•••			•••		$\int$	<b>X</b> 2	6	6	blue	•••	3500	class2
<b>X</b> 100	3	88	blue		700	class1		<b>X</b> 3	7	7	yellow	•••	400	class1
						•••		<b>X</b> 100	3	88	blue	•••	700	class1
<b>X</b> 500	60	68	red		1700	class2		<b>X</b> 9999	3	85	green	•••	2500	class2
	•••	•••												
<b>X</b> 9999	3	85	green		2500	class2								
<b>X</b> 10000	3	1	blue	•••	5700	class1								

### Data Reduction (Sample Selection)



### Data Integration

Data Integration is concerned with the **integration of different or heterogeneous data sources** (sets) in order to enhance the total information about the problem at hand.

Each data source contains partly independent and partly complementary pieces of information about the problem.

## Data Integration



HyunJung (Helen) Shin, Max Planck Society, European School of Genetic Medicine, 03. 2006

## Content

Basics
Tasks
Learning
Models with Examples


#### Learning

#### Building a model *f* given dataset {X,Y} is called "Learning" or "Training"



#### Learning

#### Building a model *f* given dataset {X,Y} is called "Learning" or "Training"



Ex) Regression Model  $f(x) = \beta_1 x^2 + \beta_2 x + c$ 

# In other words, given data {X, Y}, finding the values of parameters, $\beta_1$ , $\beta_2$ , and *c* is "Learning"

### Data Set Split

			<b>A</b> 1	<b>A</b> 2	Аз	•••	<b>A</b> 10	у
"Known" data points		$\boldsymbol{x}_{I}$	10	5	red	•••	1000	1
		<b>X</b> 2	6	6	blue	•••	3500	20
		<b>X</b> 3	7	7	yellow	•••	400	45
	$\leq$	• • •	•••	•••	•••	•••	•••	•••
		<b>X</b> 18	3	56	red	•••	0	30
		<b>X</b> 19	15	62	red		500	100
		<b>X</b> 20	3	88	blue		700	3
"Unknown" data points		<b>x</b> 21	5	42	red		560	?
	$\leq$	•••	••••	••••	•••		••••	?
		<b>X</b> 50	25	56	blue		600	?

#### Data Set Split

Training set

*Training (or learning or building) a model f* 

"Known" data points

"Unknown" data points \* Model :  $f(x) = \beta_1 x^2 + \beta_2 x + c$ 

#### Validation set

Model selection (or model parameter selection) \* Best parameters ( $\beta_1$ ,  $\beta_2$ , c)?

Test set

Prediction with a trained model



HyunJung (Helen) Shin, Max Planck Society, European School of Genetic Medicine, 03. 2006 42

Training set

Build a model "f" minimizing the errors min.  $\sum_{i=1}^{n} \varepsilon^{2} = \sum_{i=1}^{n} (f(x_{i}) - y_{i})^{2}$ 



Training set

Build a model "f" minimizing the errors min.  $\sum_{i=1}^{n} \varepsilon^{2} = \sum_{i=1}^{n} (f(x_{i}) - y_{i})^{2}$ 

















fc is the best model?





Why not?



#### **Training** set

Test set

The data points in Test set are assumed to be drawn the same distribution as those in Training set









Training set

Test set

(note that Test set is unknown during Training)

# Then, how can we find a "proper" model with absence of Test set ?



Validation set

# Then, how can we find a "proper" model with absence of Test set ?

Use Validation set (say, a Pseudo Test Set) ! : Temporarily assume that the data set is "Unknown"









If the known data points are <u>large</u> enough for training after separating the validation set off....





If the known data points are <u>insufficient</u> for training after taking the validation set out ?

#### (Ex) 5 Cross-Validation (5CV)





HyunJung (Helen) Shin, Max Planck Society, European School of Genetic Medicine, 03. 2006

65

# Learning Schemes

- Supervised
- Unsupervised
- Semi-Supervised

#### Supervised Learning



#### Supervised Learning



# Supervised Learning



HyunJung (Helen) Shin, Max Planck Society, European School of Genetic Medicine, 03. 2006 69

# Unsupervised Learning



HyunJung (Helen) Shin, Max Planck Society, European School of Genetic Medicine, 03. 2006

70

# Models

#### Models

- Traditional Statistical Models
- Neural Networks
- Decision Trees
- Kernel Methods
- Semi-Supervised Learning (SSL) and Transductive Inference Methods
- Ensemble Methods
- Generative (Probabilistic) Methods
### Models: Traditional Statistical Models



### Models: Neural Networks



#### Models: Decision Trees (or Rule-base)



### Models: Kernel Methods



# Models: Semi-Supervised Learning



• *etc, etc, etc,...* 

# Models: Generative (Probabilistic) Methods



#### Models: Ensemble Methods



### The Most Up-To-Date Models

Kernel Methods

Support Vector Machines (SVM), kPCA, kCCA, kICA, etc

Semi-Supervised Learning Methods

Graph-based SSL, Transductive Inference Methods

Kernel Methods: Support Vector Machines (SVM)



• *Kernel methods can operate on very <u>general types of data</u> and can detect very <u>general types of relations</u>* 

• Various tasks-
$$\begin{cases} PCA, \\ FA, \\ DA, \\ Clustering \end{cases}$$
 can be performed on diverse data 
$$\begin{cases} vectors, \\ sequences, \\ text, \\ images, \\ graphs \end{cases}$$

• Integration of different types of data is easy and natural



 $K_{ij} = K(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i) \cdot \phi(\mathbf{x}_j)$  where  $\Phi(.)$  is a mapping function



**Feature Space** The Mapping from Input to Feature space is...

- Highly <u>Nonlinear</u>
- <u>Dimension Expanding</u> (up to infinite dim)
- *Not unique to a Feature Space, Probably Unknown*

*Finding the mapping function* has been <u>the most difficult barrier</u> in the traditional statistics and early machine learning algorithms

*In KM, those difficulties could be circumvented by means of "<u>Kernel</u> <u>Trick</u>" which replaces the <u>dot product</u> between mapping functions* 

 $\boldsymbol{\Phi}(\mathbf{x}) \cdot \boldsymbol{\Phi}(\mathbf{y}) = K(\mathbf{x}, \mathbf{y})$ 

4

Kernel Function:  $K(\mathbf{x}, \mathbf{y}) = \phi(\mathbf{x}) \cdot \phi(\mathbf{y})$ 

**Functions Satisfying Mercers's Theorem** 

Polynomial kernels

Radial Basis kernels

Sigmoid Kernels (3-MLP NN)

Kernel Function

$$K(x,y) = (x \cdot y)^{P}$$

$$K(x,y) = \exp\left(\frac{-\|x - y\|^{2}}{2\sigma^{2}}\right)$$

$$K(x,y) = \tanh\{\kappa(x \cdot y) + \Theta\}$$

#### A Single Kernel Produces Multiple Mappings



The flexible combination of appropriate kernel design and relevant kernel algorithms has given rise to a powerful class of methods, whose computational and statistical properties are well understood

Particularly, KM has increasingly been used in in Bioinformatics as diverse as biosequences and microarray data analysis, etc.

# SVM Classification

Basic Idea of SVM

Properties of SVM ... Optional

8

- Margin
  Convexity
  Duality
  Kernels
  Sparseness

SVM looks for the <u>Separating Hyperplane</u> with the Largest Margin.



*Training data*  $\{\mathbf{x}_i, y_i\}, i = 1,...,l, y_i \in \{-1,1\}$ 

<u>Separating Hyperplane</u>  $f(x) = \mathbf{w} \cdot \mathbf{x} + b = 0$ 

 $\operatorname{sign}(\mathbf{f}(\mathbf{x})) = \begin{cases} +1 & \text{if } \mathbf{x} \cdot \mathbf{w} + b \ge 0 \\ -1 & \text{if } \mathbf{x} \cdot \mathbf{w} + b < 0 \end{cases}$ 

SVM looks for the Separating Hyperplane with the Largest Margin.



Supporting Hyperplanes  $H1: \mathbf{x}_i \cdot \mathbf{w} + b \ge +1$  for  $y_i = +1$  $H2: \mathbf{x}_i \cdot \mathbf{w} + b \le -1$  for  $y_i = -1$ 

#### <u>Margin</u>

Distance between H1 and H2

 $\frac{\overline{|\mathbf{l} - \mathbf{b}|}}{\|\mathbf{w}\|} - \frac{|-\mathbf{l} - \mathbf{b}|}{\|\mathbf{w}\|} = \frac{2}{\|\mathbf{w}\|}$ 

### Find the Pair of Hyperplanes (Support Vectors) $\mathbf{x}_i \cdot \mathbf{w} + b \ge +1$ for $\mathbf{y}_i = +1$ $\mathbf{x}_i \cdot \mathbf{w} + b \le -1$ for $\mathbf{y}_i = -1$ under the constraints which gives Maximum Margin $\frac{2}{\|\mathbf{w}\|}$ !

#### Separable Case





#### Non-Separable Case ?



Use Slack Variables !



#### Nonlinear Case ?



Solve (linear) problem in the Feature Space !



#### *Feature Space*

SVMs map the training data nonlinearly into a higher-dimensional feature space via  $\phi$  and construct a separating hyperplane with maximum margin there.

This yields a <u>nonlinear decision boundary</u> in input space.

📕 ---- 🖬 margin 🗖 < Example > Nonlinear & NonSeparable decision boundary support vector of class(1) support vector of class(2) 97



#### Properties of SVM ... optional

#### SVM looks for the Separating Hyperplane with the Largest Margin.



#### SVM looks for the Separating Hyperplane with the Largest Margin.



Supporting Hyperplanes  $H1: \mathbf{x}_i \cdot \mathbf{w} + b \ge +1 \quad \text{for } \mathbf{y}_i = +1$  $H2: \mathbf{x}_i \cdot \mathbf{w} + b \le -1 \quad \text{for } \mathbf{y}_i = -1$ 

#### <u>Margin</u>

Distance between H1 and H2

 $\frac{|\mathbf{1} - \mathbf{b}|}{\|\mathbf{w}\|} - \frac{|-\mathbf{1} - \mathbf{b}|}{\|\mathbf{w}\|} = \frac{2}{\|\mathbf{w}\|}$ 

#### SVM looks for the Separating Hyperplane with the Largest Margin.



Support Vectors  $H1: \mathbf{x}_i \cdot \mathbf{w} + b - 1 = 0 \quad \text{for } \mathbf{y}_i = +1$  $H2: \mathbf{x}_i \cdot \mathbf{w} + b + 1 = 0 \quad \text{for } \mathbf{y}_i = -1$ 

 $x_i$ 's are the Closest Data from Separating Hyperplane,  $\mathbf{w} \cdot \mathbf{x} + b = 0$ 

Find the Pair of Hyperplanes (Support Vectors)  $\mathbf{x}_i \cdot \mathbf{w} + b \ge +1$  for  $\mathbf{y}_i = +1$   $\mathbf{x}_i \cdot \mathbf{w} + b \le -1$  for  $\mathbf{y}_i = -1$ under the constraints which gives Maximum Margin  $\frac{2}{\|\mathbf{w}\|}$  !

*Minimize*  $||w||^2$  *under the constraints !!* 

$$\min \frac{1}{2} \|\mathbf{w}\|^2$$
  
s.t.  $y_i(\mathbf{x}_i \cdot \mathbf{w} + b) - 1 \ge 0 \quad \forall_i$ 

#### Separable Case


Minimize  $||w||^2$  under the constraints !!  $\min \frac{1}{2} ||\mathbf{w}||^2$ s.t.  $y_i(\mathbf{x}_i \cdot \mathbf{w} + b) - 1 \ge 0 \quad \forall_i$ 

*Quadratic Programming* (convex QP : obj ftn is convex, constraints form a convex set)

#### Non-Separable Case ?



#### <u>Use Ślack Variables !</u>





*How to Solve?* 

## Use Lagrange theory ! (Karush-Kuhn-Tucker Condition)

#### Karush-Kuhn-Tucker Condition

Min: f(x)

s.t. h(x) = 0 (m equality constraints)

 $g(x) \leq 0$  (k inequality constraints)

Lagrangian:

$$L(x,a,m) = f(x) + a h(x) + \sum u_i (g_i(x) + s_i)$$

- 1) Gradient of the Lagrangian = 0
- 2) Constraints:  $h(x) = 0 \& g(x) \le 0$
- 3) Complementary Slackness: u.s = 0
- 4) Feasibility for the inequality constraints:  $s \ge 0$
- 5) Sign condition on the inequality multipliers:  $u \ge 0$

33

KKT conditions are satisfied at the solution of <u>any constrained optimization</u> problem

For convex problem, KKT conditions are necessary and sufficient condition for primal, dual solution.

 $\begin{array}{ll} \text{Primal Problem} \\ \text{min } \frac{1}{2} \|\mathbf{w}\|^2 \\ s.t. \ y_i(\mathbf{x}_i \cdot \mathbf{w} + b) \cdot 1 \ge 0 \quad \forall_i \end{array}$ 

Lagrangian

$$L(w,b) \equiv \frac{1}{2} \| \mathbf{w} \|^2 - \sum_{i}^{l} \alpha_i y_i (\mathbf{x}_i \cdot \mathbf{w} + b) + \sum_{i}^{l} \alpha_i$$

$$L(w,b) = \frac{1}{2} ||w||^{2} -\sum_{i}^{l} \alpha_{i} y_{i} (\mathbf{x}_{i} \cdot w + b) + \sum_{i}^{l} \alpha_{i} \qquad \dots Lagrangian$$

$$\frac{\partial}{\partial w_{v}} L_{P} = w_{v} - \sum_{i}^{l} \alpha_{i} y_{i} \mathbf{x}_{i} = 0 \qquad \dots Gradient of the$$

$$Lagrangian = 0$$

$$y_{i} (\mathbf{x}_{i} \cdot \mathbf{w} + b) - \mathbf{1} \ge \mathbf{0} \quad \forall \mathbf{i} \qquad \dots Primal Feasibility$$

$$\alpha_{i} \ge \mathbf{0} \quad \forall \mathbf{i} \qquad \dots Dual Feasibility$$

$$\alpha_{i} (\mathbf{y}_{i} (\mathbf{x}_{i} \cdot \mathbf{w} + b) - \mathbf{1}) = \mathbf{0} \quad \forall \mathbf{i} \qquad \dots Dual Feasibility$$

Solving the SVM problem is equivalent to finding a solution KKT conditions.

Lagrangian L has to be minimized w.r.t. the primal variables w and b and maximized w.r.t. the dual variables  $\alpha_i$ 

• Minimize Lp with respect to w, b :

$$\min \quad L_P \equiv \frac{1}{2} \|w\|^2 - \sum_i^l \alpha_i y_i (x_i \bullet w + b) + \sum_i^l \alpha_i$$
$$\mapsto \quad w = \sum_i^l \alpha_i y_i x_i \quad , \quad \sum_i^l \alpha_i y_i = 0$$

• Maximize  $L_D$  with respect to  $\alpha_i$ :

$$\max \quad L_D \equiv \sum_{i}^{l} \alpha_i - \frac{1}{2} \sum_{i,j}^{l} \alpha_i \alpha_j y_i y_j x_i \bullet x_j$$
  
s.t.  $\alpha_i \ge 0$ ,  $\sum_{i}^{l} \alpha_i y_i = 0$ ,  $\forall i$ 

Why Dual?

$$\max \quad L_D \equiv \sum_{i}^{l} \alpha_i - \frac{1}{2} \sum_{i,j}^{l} \alpha_i \alpha_j y_i y_j x_i \bullet x_j$$
  
s.t.  $\alpha_i \ge 0$ ,  $\sum_{i}^{l} \alpha_i y_i = 0$ ,  $\forall i$ 

Why Dual?

$$\max \quad L_D \equiv \sum_{i}^{l} \alpha_i - \frac{1}{2} \sum_{i,j}^{l} \alpha_i \alpha_j y_i y_j \mathbf{x}_i \bullet \mathbf{x}_j$$
  
s.t.  $\alpha_i \ge 0$ ,  $\sum_{i}^{l} \alpha_i y_i = 0$ ,  $\forall i$ 

#### Dot Product between Training Vectors: We can use Kernel functions !

#### Nonlinear Case ?



Nonlinear Mapping from Input space( $\mathbb{R}^N$ ) to Feature Space(F) (EX)  $\boldsymbol{\Phi} : \mathbb{R}^2 \to \mathbb{R}^3$ ,  $(x_1, x_2) \mapsto (z_1, z_2, z_3)$ 



#### <u>Feature Space</u>

SVMs map the training data nonlinearly into a higher-dimensional feature space via  $\phi$  and construct a separating hyperplane with maximum margin there.

This yields a nonlinear decision boundary in input space.



#### <u>Mapping Function ( $\phi$ ) ?</u>

However,

Mapping function is <u>not unique</u>.

Feature Space could be (possibly) infinite dimensional.



<u>Mapping Function ( $\phi$ )?</u>

## How can we know the mapping function ? How can we to handle the infinite dimensionality?

#### <u>Use Kernel Functions</u> !

SVM depends only on <u>Dot Products</u> between patterns.

$$\underline{Problem} \qquad \max \ L_D \equiv \sum_{i}^{l} \alpha_i - \frac{1}{2} \sum_{i,j}^{l} \alpha_i \alpha_j y_i y_j \phi(\mathbf{x}_i) \cdot \phi(\mathbf{x}_j)$$

$$\underline{Problem} \qquad s.t. \ \alpha_i \ge 0, \quad \sum_{i}^{l} \alpha_i \ y_i = 0, \quad \forall i$$

$$\underline{Decision \ Function} \qquad f(x) = sign(w \cdot \phi(x) + b) = sign(\sum_{i=1}^{l} \alpha_i y_i \phi(x_i) \cdot \phi(x) + b)$$

*By the use of a kernel function, it is possible to compute the dot product in input space <u>without explicitly carrying out the map</u> <u>into the feature space</u>* 

*Kernel Function:*
$$k(\mathbf{x}, \mathbf{y}) = \phi(\mathbf{x}) \cdot \phi(\mathbf{y})$$

#### Functions Satisfying Mercers's Theorem

Polynomial kernels

$$k(\mathbf{x},\mathbf{y}) = (\mathbf{x} \cdot \mathbf{y})^P$$

Radial Basis kernels

$$k(\mathbf{x}, \mathbf{y}) = \exp\left(\frac{-\|\mathbf{x} - \mathbf{y}\|^2}{2\sigma^2}\right)$$

Sigmoid Kernels (3-MLP NN)

$$k(\mathbf{x}, \mathbf{y}) = \operatorname{tanh}(\kappa (\mathbf{x} \bullet \mathbf{y}) + \Theta)$$

# Margin Convexity Duality [Kernel] Sparseness <u>Nonlinear & Nonseparable Case</u>



HyunJung (Helen) Shin, Max Planck Society, European School of Genetic Medicine, 03. 2006

# Only the points nearest to the hyperplane have positive weight !

They are called Support Vectors !

Remind the Complementarity Conditions  $\alpha_i(y_i(\mathbf{x_i} \cdot \mathbf{w} + b) - 1) = 0 \quad \forall i \text{ also note that } \mathbf{w} = \sum_{i=1}^{l} \alpha_i y_i \mathbf{x_i}$ 



HyunJung (Helen) Shin, Max Planck Society, European School of Genetic Medicine, 03. 2006

1

SVs are distributed around decision boundary !



## **SVM Decision Function**

 $f(x) = \begin{cases} +1 \ (class 1) & if \ f(x) \ge 0\\ -1 \ (class 2) & if \ f(x) < 0 \end{cases}$ 



Wrap-up

SVM QP Problem: (<u>Non-linear</u> & <u>Non-Separable</u>)

$$\min \cdot \frac{1}{2} \| \vec{w} \|^2 + C \sum_{i=1}^M \xi_i$$
  
s.t.  $y_i (\vec{w} \cdot \Phi(\vec{x}_i) + b) \ge 1 - \xi_i$ ,  
 $i = 1, \dots, M$ 

SVM Decision Function: 
$$f(\vec{x}) = sign\left(\sum_{i \in SV} y_i \alpha_i \Phi(\vec{x}_i) \cdot \Phi(\vec{x}) + b\right)$$

Kernels: 
$$\Phi(\vec{x}) \cdot \Phi(\vec{x}') = k(\vec{x}, \vec{x}') = \begin{cases} \exp\left(-\left\|\vec{x} - \vec{x}'\right\|^2 / 2\sigma^2\right) \\ \tanh(\kappa(\vec{x} \cdot \vec{x}') + \Theta\right) \\ (\vec{x} \cdot \vec{x}' + 1)^P \end{cases}$$

Wrap-up

$$\min_{0 \le \alpha_{i} \le C} W(\alpha_{i}, b) = \frac{1}{2} \sum_{i, j=1}^{M} \alpha_{i} \alpha_{j} y_{i} y_{j} K(\vec{x}_{i}, \vec{x}_{j}) - \sum_{i=1}^{M} \alpha_{i} + b \sum_{i=1}^{M} y_{i} \alpha_{i}$$

$$\frac{\partial W(\alpha_{i}, b)}{\partial \alpha_{i}} = \sum_{j=1}^{M} y_{i} y_{j} K(\vec{x}_{i}, \vec{x}_{j}) \alpha_{j} + y_{i} b - 1 = y_{i} \bar{f}(x_{j}) - 1$$

$$\frac{\partial \mathbf{W}(\boldsymbol{\alpha}_{\mathrm{i}},\mathbf{b})}{\partial b} = \sum_{\mathrm{j=1}}^{\mathrm{M}} y_{\mathrm{j}} \boldsymbol{\alpha}_{\mathrm{j}} = 0$$

where 
$$\overline{f}(\vec{x}) = \sum_{i=1}^{M} y_i \alpha_i K(\vec{x}_i, \vec{x}) + b$$

# Application I

Recognition of Alternatively Spliced Exons in C.elegans

Task : Classification, Data IntegrationModel : Semi-Supervised LearningApplication: C.elegans Genes - Alternative Splicing

HyunJung (Helen) Shin, Max Planck Society, European School of Genetic Medicine, 03. 2006

6

### Splicing



## Splicing



#### Splice sites are

- the exon/intron boundaries
- recognized by five snRNAs
- assembled in snRNPs
- flanked by regulatory elements

#### Spliceosomal Proteins

- interact with snRNPs and mRNA
- regulate recognition of splice sites
- can lead to <u>alternative transcripts</u>

One gene may correspond to several transcripts/proteins !!

## Alternative Splicing



#### Alternative Splicing

#### Alternative Splicing (AS) ..

- can produce several mRNA transcript per gene (sometimes leading to more than 100 slightly different proteins)
- is highly regulated
- greatly increases the proteome diversity in eukaryotes (about 70% of human genes are alternatively spliced!)

#### Alternative Splicing

Methods for identifying alternative splicing ...

- usually need many EST sequences or- exploit conservation between several organisms

Novel AS prediction method only using the pre-mRNA

#### Alternatively Spliced Exons



#### Idea: Use Machine Learning to

- understand differences between alternative and constitutive splicing
- exploit and identify regulative elements
- predict unknown alternative splicing events

#### Alternatively Spliced Exons



Previous work Analysis of conserved alternatively spliced exons (Sorek et al., Yeo et al. and others) - consider conserved alternative spliced exons (ACE) - exploit that ACE and flanking introns are more conserved between mouse and human

> *Problem* only works for conserved exons

Derive the features from the "pre-mRNA" in order to find "novel" exons !!
### Task Formulation

#### Two-class Classification Problem



A (or B) is true splice site or not?

#### Use Support Vector Machines!

### Remind the Procedure of Kernel Methods !



#### Procedure - Data Set



*True sites* (y=1): fixed window around a true splice site Decoys sites (y=-1): generated by shifting the window

### Procedure - Data Set



**Procedure - Kernel Function** 

Kernels measure similarities between sequences Weighted Degree Kernel (Sonnenburg et al., 2002)



Given two sequences S1 and S2 of equal length, the kernel consists of a weighted sum to which each match in the sequences makes a contribution. The longer matches contribute more significantly.

### **Procedure - Kernel Function**



# **Procedure - Modeling**



# **Procedure - Modeling**



#### Results

#### Exons Known



- 21,000 exons and 28,000 introns (single EST confirmed)

22

- 25 random exons & introns from 1-2% top ranks

- RT-PCR with primers in flanking exons

#### Results

Based on wetlab experiments and accuracy estimates:

- ~ 1% of known exons are alternatively spliced (AS)
- ~ 0.25% of AS exons are yet completely unknown
- 280 AS spliced exons (total)
- 13 confirmed by RT-PCR
- Additional 80 AS exons can be found with less than 200 additional RT-PCRs

Semi-Supervised Learning Methods: Graph-Based SSL

# Semi-Supervised Learning

Semi-Supervised Learning...

- Utilizes every possible information in hand for prediction
- Seems to show a better performance than Supervised-Learning



- Adjacency (similarity) matrix of the network: W
- *Known Labels* :  $y_1, ..., y_l \in \{-1, 1\}$
- *Unknown Labels* :  $y_{l+1}, ..., y_n \in \{0\}$
- Predicted outputs :  $f_1, ..., f_n$ 
  - $f_i$  should be close to those of adjacent nodes,  $f_j$ 's where  $i \sim j$ .  $f_i$  should be close to the given label  $y_i$  at training nodes



Learning Problem

min 
$$\mu \sum_{i \sim j} w_{ij} (f_i - f_j)^2 + \sum_{i} (f_i - y_i)^2$$

#### Equivalent Vector Form

$$\min_{f} \mu f^{T}L f + (f - y)^{T} (f - y)$$

$$f$$

$$L \text{ is called the graph Laplacian matrix where}$$

$$L = D - W, \quad D = diag(d_{i}), \quad d_{i} = \sum_{i} w_{ij}$$

#### HyunJung (Helen) Shin, Max Planck Society, European School of Genetic Medicine, 03. 2006

J



Objective Function $\min_{f} \mu f^T L f + (f - y)^T (f - y)$ f $f = \{I + \mu L\}^{-1} y$ 

# Application II-1

### Functional Class Prediction on a Protein Network

Task : Classification

Model : Semi-Supervised Learning

Application: Protein – Protein Function Prediction





+1/-1 : Labeled proteins with/without a specific function

? : Unlabeled proteins

Similarities between Proteins : Edges



Edges in Physical Interaction Network: Two proteins physically interact (e.g., docking)
Edges in Metabolic Network of Enzymes: Two enzymes catalyzing successive reactions

The task is to predict labels of unlabeled proteins using similarities.

# Example: Metabolic Gene Network





Substrate of PFK1

Enzyme or <u>Protein</u>, PFK1

Product of PFK1









HyunJung (Helen) Shin, Max Planck Society, European School of Genetic Medicine, 03. 2006 41



HyunJung (Helen) Shin, Max Planck Society, European School of Genetic Medicine, 03. 2006 42







# Application II-2

### Functional Class Prediction with <u>Multiple Networks</u>

Task : Classification, Data IntegrationModel : Semi-Supervised LearningApplication : Yeast Protein – Protein Function Classification

If Multiple Graphs are Given ?

## If Multiple Graphs are Given?



# If Multiple Graphs are Given?




Each graph can solely predict the label of the unlabeled nodes depending on its own similarity.



Since different graphs contain <u>partly independent</u> and <u>partly complementary</u> pieces of information about the problem at hand,

one thus can enhance the total information about the problem by <u>combining those graphs</u>.

Example: Multiple Graph Sources on Proteins

*Physical interactions of the proteins* [Schwikowski,et al., 2000, Uetz et al., 2000, von Mering et al., 2002]

*Gene regulatory relationships* [Lee et al., 2002, Ihmels et al., 2002, Segal et al., 2003]

Edges in a metabolic pathway [Kanehisa et al., 2004]

Similarities between protein sequences [Yona et al., 1999]

etc.

2

### Lee et. al., 2004. A Probabilistic Functional Network of Yeast Genes, Science, vol. 306

Lee et. al., 2004. A Probabilistic Functional Network of Yeast Genes, Science, vol. 306

QP.







#### **Previous Approach**

#### SDP/SVM : Semi-Definite Programming based Support Vector Machine [Lanckriet et al., Bioinfomatics, 2004]



Each graph is <u>converted to a kernel matrix</u>



Kernel matrices are combined with weights which are automatically learned by Semi-Definite Programming

SVM

Labels are predicted based on the combined kernel matrix





Each graph is converted to a kernel matrix

SDP

Kernel matrices are combined with weights which are automatically learned by Semi-Definite Programming

SVM

Labels are predicted based on the combined kernel matrix

Good accuracy which is much better than Markov Random Field

> But Very Slow

SDP/SVM : Semi-Definite Programming based SVM

In SDP/SVM, multiple kernel matrices corresponding to each of data sources are combined with weights obtained by solving an SDP.

However, when trying to apply SDP/SVM to large problems, the computational cost can become prohibitive, since both Converting the data to a kernel matrix for the SVM and Solving the SDP are <u>time and memory demanding</u>

Diffusion Kernel

[Kondor and Lafferty, 2002].

$$K_{\beta} = e^{\beta \boldsymbol{L}} = \lim_{s \to \infty} (I + \frac{\beta \boldsymbol{L}}{s})^s = I + \beta \boldsymbol{L} + \frac{\beta^2}{2} \boldsymbol{L}^2 + \frac{\beta^3}{6} \boldsymbol{L}^3 + \dots$$

L : graph Laplacian.

 $\beta$  : diffusion rate

Semi-Definite Programming

[Vandenberg and Boyd, 1996] [Boyd and Vandenberg, 2003]

$$\min_{\boldsymbol{u}} \quad \boldsymbol{c}^{T}\boldsymbol{u}$$

$$\boldsymbol{u}$$
s.t.  $F^{j}(\boldsymbol{u}) = F_{o} + \sum_{k=1}^{K} u_{k}F_{k} \ge 0, \quad j = 1, ..., J.$ 

where  $c \in \mathbb{R}^{K}$ ,  $F_{k} \in \mathbb{R}^{n \times n}$ ,  $F(u) \in \mathbb{R}^{n \times n}$ : symmetric, positive semi - definite

*Convex optimization problem since its objective and constrains are convex* 

$$SDP/SVM$$
[Lanckriet et al., 2004]  
Single Kernel Case  
SVM dual problem
$$SVM cast as an SDP$$

$$\max 2a^{T}e - a^{T}(G(K) + \tau I)a$$

$$a$$

$$s.t. \quad 0 \le a \le C, \quad a^{T}y = 0,$$

$$trace(K) = c.$$

$$(G(K_{tr}) + \tau I_{ntr} \quad e + v - \delta + \lambda y)$$

$$(e + v - \delta + \lambda y)^{T} \quad t - 2C\delta^{T}e$$

$$v \ge 0,$$

$$\delta \ge 0.$$



Calculating a Diffusion Kernel from a Graph

 $O(n^3)$ , A dense matrix of  $n \ge n$ 

Solving SDP

 $O((m+n)^2 n^{2.5})$ 

*m: the number of kernel matrices n: number of nodes (data)* 

Computationally Expensive both in Time and Memory

Why not use a more direct approach for combining graphs based on significant progress of graph-based semi-supervised learning methods ?

- Zhou et al., 2004
- Belkin and Niyogi, 2003
- Zhu et al., 2003
- Chapelle et al., 2003

# Semi-Supervised Learning Extension to Multiple Graphs

- Combining weights are automatically assigned to Graphs
- Comparable Accuracy to SDP/SVM
- Very Fast

*Linear Combination of Laplacians* 

$$L(\beta) = \sum_{k=1}^{K} \beta_k L_k$$



How to Find Combining Weights ?

Single Graph

$$\min_{\boldsymbol{f}} \quad (\boldsymbol{f} - \boldsymbol{y})^T (\boldsymbol{f} - \boldsymbol{y}) + c \boldsymbol{f}^T L \boldsymbol{f}$$

Without loss of generality, the problem is rewritten by penalizing the upper-bound

$$\min_{\boldsymbol{f},\boldsymbol{\gamma}} \quad (\boldsymbol{f}-\boldsymbol{y})^T (\boldsymbol{f}-\boldsymbol{y}) + c\boldsymbol{\gamma}, \qquad \boldsymbol{f}^T L \, \boldsymbol{f} \leq \boldsymbol{\gamma}.$$

Multiple Graphs

$$\min_{\boldsymbol{f}} \quad (\boldsymbol{f} - \boldsymbol{y})^T (\boldsymbol{f} - \boldsymbol{y}) + c \{\beta_1 \boldsymbol{f}^T \boldsymbol{L}_1 \boldsymbol{f} + \beta_2 \boldsymbol{f}^T \boldsymbol{L}_2 \boldsymbol{f} + \dots + \beta_k \boldsymbol{f}^T \boldsymbol{L}_k \boldsymbol{f} \}$$

Without loss of generality, the problem is rewritten by penalizing the upper-bound

$$\min_{\substack{f,\gamma}} (\boldsymbol{f} - \boldsymbol{y})^T (\boldsymbol{f} - \boldsymbol{y}) + c\gamma, \quad \boldsymbol{f}^T L_k \boldsymbol{f} \leq \gamma, \quad k = 1, \dots K.$$

### Extension to Multiple Graphs: Optimization

Prime  

$$\begin{array}{l} \min \quad (f - y)^{T} (f - y) + c\gamma, \\ f, \gamma \\ \text{s.t.} \quad f^{T} L_{k} \quad f \leq \gamma, \\ \gamma \geq 0, \quad k = 1, \dots K. \end{array}$$

$$\begin{array}{l} \min \quad d(\boldsymbol{\beta}) \equiv y^{T} (\boldsymbol{I} + \sum_{k=1}^{K} \beta_{k} \quad L_{k})^{-1} y, \\ \boldsymbol{\beta} \\ \text{s.t.} \quad \sum_{k=1}^{K} \beta_{k} \leq c \\ k = 1 \end{array}$$

 $\beta_k$ : Weight for Network k, Lagrange Multiplier.

### Extension to Multiple Graphs: Solution

Solution

$$f = \left\{ I + \sum_{k=1}^{K} \beta_k L_k \right\}^{-1} y$$

Matrix Inversion

 $\beta_k$ : Weight for Network k, Lagrange Multiplier.

Sparse Linear System



Linear Systems

### Extension to Multiple Graphs: Meaning of Weights

By KKT complementarity condition, we have the following relationship at the optimal solution,

$$\beta_{k}(f^{T}L_{k}f-\delta) = 0$$

$$\beta_{k} = 0 \quad \text{iff } f^{T}L_{k}f < \delta \qquad \beta_{k} > 0 \quad \text{iff } f^{T}L_{k}f = \delta$$

The score vector **f** would not be changed much with those graphs, thus those are considered as **redundant** 

*Those graphs are considered important.* 

Computational Efficiency

1. <u>Repetition</u> of an Identical Form of Inverse Matrix

2. <u>Implicit</u> Calculation of Matrix Inversion

Computational Efficiency

1. <u>Repetition</u> of an Identical Form of Inverse Matrix: in the objective function and the derivative, (and the network output).

Objective Function

$$\min_{\boldsymbol{\beta}} d(\boldsymbol{\beta}) \equiv \boldsymbol{y}^T (\boldsymbol{I} + \sum_{k=1}^K \beta_k L_k)^{-1} \boldsymbol{y}$$

Solution Update

$$\frac{\partial d}{\partial \beta_k} = -\mathbf{y}^T \left(I + \sum_{j=1}^K \beta_j L_j\right)^{-1} L_k \left(I + \sum_{j=1}^K \beta_j L_j\right)^{-1} \mathbf{y}$$

Network Output

$$\boldsymbol{f} = \left\{ \boldsymbol{I} + \sum_{k=1}^{K} \beta_k \boldsymbol{L}_k \right\}^{-1} \boldsymbol{y}$$

Computational Efficiency

2. Implicit Calculation of Matrix Inversion: The solution can be obtained by solving the "sparse linear systems." Therefore, computational cost is nearly linear in the number of non-zero entries of  $\sum_{k=1}^{K} \beta_k L_k$  – (Spielman and Teng, 2004).



### Function Prediction Experiments

MIPS Comprehensive Yeast Genome Database (CYGD-mips.gsf.de/proj/yeast).



## Protein Functional Categories

MIPS Comprehensive <u>Yeast Genome</u> Database (CYGD-mips.gsf.de/proj/yeast).

13 CYGD functional Classes

- 1. metabolism
  - 2. energy
  - 3. cell cycle and DNA processing
  - 4. transcription
  - 5. protein synthesis
  - 6. protein fate
  - 7. cellular transportation and transportation mechanism
  - 8. cell rescue, defense and virulence
  - 9. interaction with cell environment
- 10. cell fate
- 11. control of cell organization
- 12. transport facilitation
- $\sim$  13. others

# Inputs (5 networks)

Network created from <u>Pfam domain structure</u>. A protein is represented by a 4950-dimensional binary vector, in which each bit represents the presence or absence of one Pfam domain. An edge is created if the inner product between two vectors exceeds 0.06. The edge weight corresponds to the inner product.

- W<sub>2</sub> Co-participation in a protein complex (determined by tandem affinity purification, TAP). An edge is created if there is a bait-prey relationship between two proteins.
- $W_3$  <u>Protein-protein interactions</u> (MIPS physical interactions)
- $W_4$  <u>Genetic interactions</u> (MIPS genetic interactions)

 $\sim W_5$ 

Network created from the <u>cell cycle gene expression measurements</u> [Spellman et al., 1998]. An edge is created if the Pearson coefficient of two profiles exceeds 0.8. The edge weight is set to 1. This is identical with the network used in [Deng et al., 2003]

# Inputs (5 networks)



## Density of Working Matrices



HyunJung (Helen) Shin, Max Planck Society, European School of Genetic Medicine, 03. 2006 32

# Inputs (5 networks)



### Methods in Comparison

$L_k$	<i>Label propagation with an Individual Graphs</i> $(k=15)$
L <sub>opt</sub>	Laplacian of Combined Graph with Optimized Weights
L <sub>fix</sub>	Label propagation with Equal Weights
MRF	Markov Random Field, proposed by Deng et al [2003]
SDP/SVM	Semi-definite Programming based Support Vector Machines, proposed by Lanckriet et al [2004]

### Measurements

*ROC (receiver operating characteristic) score TP1FP*, *TP10FP Computational Time*
## Measurements



## Measurements



## Measurements



**Results** : ROC scores of  $L_{opt}$ ,  $L_{fix}$  vs. the Best Performing Individual  $L_k$ 

*White: the best performing individual network Blue:*  $L_{fix}$ *Black:*  $L_{opt}$ 



Across the 13 classes,  $L_{fix}$  or  $L_{opt}$  outperforms the best performing individual.

**Results** : TP1FP and TP10FP of  $L_{opt}$ ,  $L_{fix}$ , vs. Individual  $L_{k's}$ 



## *Results – McNemar's Test:*

A pairwise test for ROC score difference: the combined graph vs. individual graphs



## Results – McNemar's Test:

#### A smaller p-value indicates a more statistically significant difference



*Results – McNemar's Test:* 

In 61% of the total number of trials, there is a statistically significant difference (at a significance level of alpha=0.05).



# **Results** : Obtained Weights



## **Results** : Comparison between Methods



For most classes, the proposed method achieves high scores, which are similar to the SDP/SVM methods. In classes 11 and 13, the proposed method **performs poor** (but still better than the MRF method), However, taking into account the <u>Simplicity and Efficiency</u> the method shows the promising results

**Results** : Computational Time



SDP/SVM :

*Approx. Several <u>CPU days</u>* (G. Lanckriet, personal communication)

\* Measured in a standard 2.2Ghz PC with 1GByte memory

*Results* : *Computational Time* 

Average Computation Time

Combining Graphs:

*Nearly <u>linearly proportional</u> to the number of non-zero entries of sparse matrices* 

SDP/SVM :

 $O(n^3) + O((m+n)^2 n^{2.5})$ 

**Results : Summary** 

Combining Graphs with <u>Optimized Weights</u> has "MORE"

Selectivity

When Compared with Combining Graphs with *Fixed Weights* 

Combining Graphs has "MORE"

Simplicity, Computational Efficiency, thus Scalablity

When Compared with <u>SDP/SVM</u>

**Results : Summary** 

Combining Graphs with <u>Optimized Weights</u> is "LESS"

Simple

When Compared with Combining Graphs with *Fixed Weights* 

Combining Graphs is "LESS"

Accurate

When Compared with <u>SDP/SVM</u>

# **Results : Summary**

#### Semi-Supervised Learning with Multiple Networks

- ... is Fast and Scalable
- ... provides Selectivity
  - (redundant / irrelevant networks can be excluded)

# Conclusion of Lecture

We explored representative models, from traditional statistical models to recent machine learning models, presenting several up-to-date research projects in bioinfomatics to exemplify how biological questions can benefit from a machine learning approach.

Machine Learning can alleviate the burden of solving many biological problems,

saving the time and cost required for experimentsproviding predictions that guide new experiments.

Multivariate Statistical Methods

R.A. Johnson & D.W. Wichern, Applied multivariate statistical analysis, Prentice-Hall. Inc, 1998.

B.F.J. Manly,

Multivariate statistical methods: A primer,

Chapman & Hall, 1997.

Kernel Methods and SVM

V. Vapnik, Statistical learning theory, Wiley, NY, 1998

### Kernel Methods and SVM

#### C.J.C. Burges,

A tutorial on support vector machines for pattern recognition, Data Mining and Knowledge Discovery, 1998.

N. Cristianini and J. Shawe-Taylor, An introduction to support vector machines, Cambridge University Press, Cambridge, UK, 2000.

> B. Scholkopf and A. J. Smola, Learning with Kernels, MIT press, MA, 2002.

Kernel Methods and Bioinfomatics

B. Scholkopf, K. Tsuda and J-P. Vert, Kernel Methods in Computational Biology, MIT press, London, 2004.

Semi-supervised Learning

Olivier Chapelle, Bernhard Schoelkopf and Alexander Zien, Semi-Supervised Learning, MIT press, 2005

Application I: Alternative Splicing

G. Rätsch, S. Sonnenburg and B. Schölkopf, A RASE: Recognition of Alternatively Spliced Exons in C. elegans, Bioinfomatics, 2004.

http://www.fml.tuebingen.mpg.de/raetsch/projects/RASE

Application II: Protein Function Classification

H. Shin and K. Tsuda, Prediction of Protein Function from Networks, In book: Semi-Supervised Learning, MIT press, London, 2006. <u>http://www.kyb.tuebingen.mpg.de/~shin</u> <u>http://www.fml.tuebingen.mpg.de/~shin</u>

K. Tsuda, H. Shin, and B. Schölkopf, Fast Protein Classification with Multiple Networks, Bioinformatics, 2005.

4