

---

**ECE471-571 – Pattern Recognition**

---

**Lecture 17: Support Vector Machine**

Hairong Qi, Gonzalez Family Professor  
 Electrical Engineering and Computer Science  
 University of Tennessee, Knoxville  
<http://www.eecs.utk.edu/faculty/qi>  
 Email: hqi@utk.edu

---

---

---

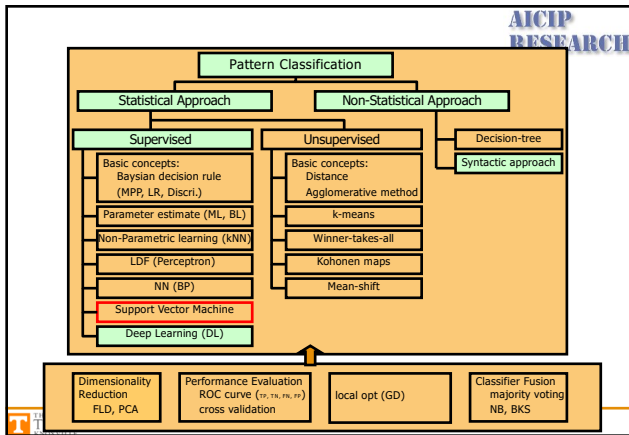
---

---

---

---

---




---

---

---


---

---


---

---

---



- Reference: Christopher J.C. Burges, "A tutorial on support vector machines for pattern recognition," *Data Mining and Knowledge Discovery*, 2, 121-167, 1998


3

---

---

---

---

---

---

---

---

## A bit about Vapnik

- Started SVM study in late 70s
- Fully developed in late 90s
- While at AT&T lab

---

---

---

---

---

---

---

---

## Generalization and Capacity

- For a given learning task, with a given finite amount of training data, the best generalization performance will be achieved if the right balance is struck between the accuracy attained on that particular training set, and the “capacity” of the machine
- Capacity – the ability of the machine to learn any training set without error
  - Too much capacity - overfitting

---

---

---

---

---

---

---

---

## Bounds on the Balance

- ◆ Under what circumstances, and how quickly, the mean of some empirical quantity converges uniformly, as the number of data point increases, to the true mean
- ◆ True mean error (or actual risk)

$$R(\alpha) = \int \frac{1}{2} |y - f(\mathbf{x}, \alpha)| p(\mathbf{x}, y) d\mathbf{x} dy$$

- ◆ One of the bounds

$$R(\alpha) \leq R_{emp}(\alpha) + \sqrt{\frac{h \ln(2/h\eta) + \ln(1/\eta)}{L}} \quad R_{emp}(\alpha) = \frac{1}{L} \sum_{i=1}^L |y_i - f(\mathbf{x}_i, \alpha)|$$

$f(\mathbf{x}, \alpha)$ : a machine that defines a set of mappings,  $\mathbf{x} \rightarrow f(\mathbf{x}, \alpha)$   
 $\alpha$ : parameter or model learned  
 $h$ : VC dimension that measures the capacity, non-negative integer  
 $R_{emp}$ : empirical risk  
 $\eta$ :  $1-\eta$  is confidence about the loss,  $\eta$  is between  $[0, 1]$   
 $L$ : number of observations,  $y$ : label,  $\{+1, -1\}$ ,  $\mathbf{x}$  is n-D vector

Principled method: choose a learning machine that minimizes the RHS with a sufficiently small  $\eta$

---

---

---

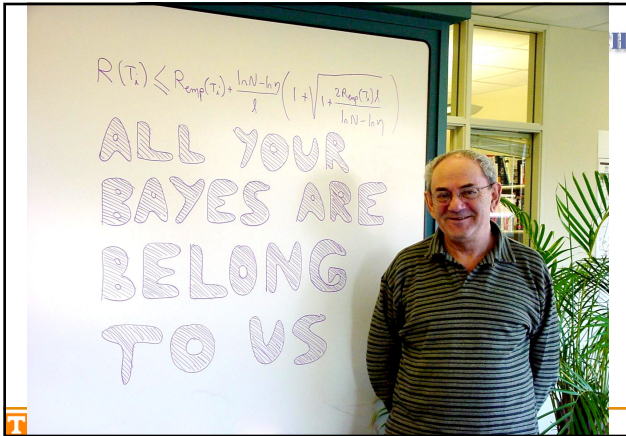
---

---

---

---

---




---

---

---

---

---

---

---

---

**AICIP RESEARCH**

### VC Dimension

- For a given set of  $l$  points, there can be  $2^l$  ways to label them. For each labeling, if a member of the set  $\{f(\alpha)\}$  can be found that correctly classifies them, we say that set of points is **shattered** by that set of functions.
- VC dimension of that set of functions  $\{f(\alpha)\}$  is defined as the maximum number of training points that can be shattered by  $\{f(\alpha)\}$
- We should minimize  $h$  in order to minimize the bound

THE UNIVERSITY OF TENNESSEE

8

---

---

---

---

---

---

---

---

**AICIP RESEARCH**

### Example ( $f(\alpha)$ is perceptron)

*Figure 1. Three points in  $\mathbb{R}^2$ , shattered by oriented lines.*

THE UNIVERSITY OF TENNESSEE

9

---

---

---

---

---

---

---

---



## Non-separable Case – Soft Margin

$$\begin{cases} \mathbf{x}_i \cdot \mathbf{w} + b \geq 1 - \xi_i & \text{for } y_i = +1 \\ \mathbf{x}_i \cdot \mathbf{w} + b \leq -1 + \xi_i & \text{for } y_i = -1 \end{cases} \quad \text{for } \xi_i \geq 0$$

Minimizing  $\|\mathbf{w}\|^2$

$$\text{s.t. } y_i (\mathbf{x}_i \cdot \mathbf{w} + b) - 1 + \xi_i \geq 0$$

$$\text{Minimize } L_p = \frac{1}{2} \|\mathbf{w}\|^2 - C \left( \sum_i \xi_i \right)^k$$

$$\text{Maximize } L_D = -\frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \mathbf{x}_i \cdot \mathbf{x}_j + \sum_i \alpha_i$$

$$\text{s.t. } 0 \leq \alpha_i \leq C, \quad \sum_i \alpha_i y_i = 0$$

---

---

---

---

---

---

---

---

---

---

## Non-separable Cases – Kernel Trick

- If there were a “kernel function”,  $K$ , s.t.

$$K(\mathbf{x}_i, \mathbf{x}_j) = \Phi(\mathbf{x}_i) \cdot \Phi(\mathbf{x}_j) = e^{-\|\mathbf{x}_i - \mathbf{x}_j\|^2 / 2\sigma^2}$$

Gaussian Radial Basis Function (RBF)

---

---

---

---

---

---

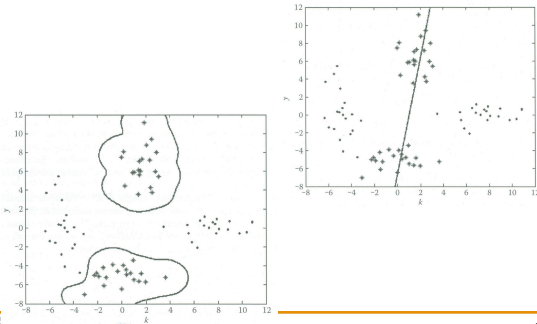
---

---

---

---

## Comparison - XOR




---

---

---

---

---

---

---

---

---

---

**AICIP  
RESEARCH**

## Limitation

- Need to choose parameters

---

---

---


---

---

---

---

---

 THE UNIVERSITY OF TENNESSEE 16

---

---

---

---

---

---

---

---

**AICIP  
RESEARCH**

## Packages

- libSVM
  - Use one-against-one (1a)
- SVM<sup>light</sup>

---

---

---


---

---

---

---

---

 THE UNIVERSITY OF TENNESSEE 17

---

---

---

---

---

---

---

---

**AICIP  
RESEARCH**

## Package Installation

- Download:  
<http://www.csie.ntu.edu.tw/~cjlin/libsvm/>
- Installation (Three choices)
  - On Unix systems, type 'make' to build the 'svm-train' and 'svm-predict' programs.
  - On other systems, consult 'Makefile' to build them
  - Use the pre-built binaries (Windows binaries are in the directory 'windows').
  - More details pls refer to the README file

---

---

---


---

---

---

---

---

 THE UNIVERSITY OF TENNESSEE

---

---

---

---

---

---


---

---

**AICIP  
RESEARCH**

### Steps

- ◆ Step 1: Transform the data to the format of an SVM package
- ◆ Step 2: Conduct simple scaling on the data
- ◆ Step 3: Consider the RBF kernel  $K(x, y) = e^{-\gamma \|x - y\|^2}$
- ◆ Step 4: Select the best parameter  $C$  and  $\gamma$  to train the whole training set
- ◆ Step 5: Test

 THE UNIVERSITY OF TENNESSEE

---

---

---

---

---

---


---

---

**AICIP  
RESEARCH**

### Example

- Dataset: pima.tr and pima.te
- Step 1: Transform the data to the format of an SVM package
  - $P_{tr} \in R^{m \times f}$  (training data: every row is a feature vector)
  - $P_{te} \in R^{n \times f}$  (testing data: every row is a feature vector)
  - $l_{tr}$  (label vector for training data pima.tr)
  - $l_{te}$  (label vector for testing data pima.te)

 THE UNIVERSITY OF TENNESSEE

---

---

---

---

---

---

---


---

**AICIP  
RESEARCH**

### Example

- Step 2: Data scaling
  - Avoid attributes in greater numeric ranges dominating those in smaller numeric ranges
  - Avoid numerical difficulties during the calculation
- How?
  - Calculate the min and max for every feature from the training dataset
  - For every feature (train or test)  $f$ , the scaled feature  $f_s$  can be calculated by  $f_s = (f - \min) / (\max - \min)$

Details pls refer to:  
<http://www.csie.ntu.edu.tw/~cjlin/libsvm/faq.html#f407>  
<http://www.csie.ntu.edu.tw/~cjlin/papers/guide/guide.pdf>

 THE UNIVERSITY OF TENNESSEE

---

---

---

---

---

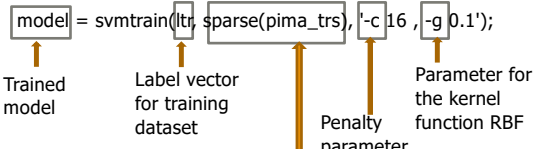
---

---

---

### Example

- Step 3: Train SVM on given parameters



- 1, Pima\_trs is the matrix for the scaled features of training dataset
- 2, Sparse(pima\_trs) is an operation to generate a sparse matrix in matlab, required by the libsvm package




---

---

---

---

---

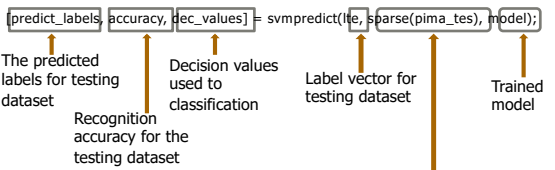
---

---

---

### Example

- Step 4: Test on the trained model



- 1, Pima\_tes is the matrix for the scaled features of testing dataset
- 2, Sparse(pima\_tes) is an operation to generate a sparse matrix in matlab, required by the libsvm package




---

---

---

---

---

---

---

---

### Example

- Result:
  - Scaled:
    - Accuracy = 80.1205% (266/332)
  - Non-scaled:
    - Accuracy = 66.8675% (222/332)




---

---

---

---

---

---

---

---



