# Support Vector Machine

# Microarray Gene Expression Data

# Support Vector Machine

# Microarray Gene Expression Data

**2001    12**

2

,    ,    ,

,

,                              ,

,    ,              ,         ,

.

2002                    ...

# Support Vector Machine　　　　Microarray
# Gene Expression Data

.

microarray　　　　　SVM(Support Vector Machine)

,

microarray　　　　　　　　　　　(kernel-function)　　　　SVM

.　　　　　　　　　　　　　　　　　SVM

microarray　　　　　　　　　,　　　microarray

(S-PLUS)　　　　　　　.　　　SVM

Chen　　　　R-Package　　　SVM

.

,　　　　　　　　　　　　(　　,　　　,

,　　　,　　　,　　　)　　　　　　　　　,

microarray　　　　　　Cy5

.

.　　,

,　　　　microarray

.

: Support Vector Machine, Microarry, Kernel-function, SVM
classifier

# 1

Human Genome Project 10

30 DNA

. A B
C

(hypothesis-driven)

. ,

(discovery-driven)

, DNA microarray SAGE(serial analysis
of gene expression) .

. DNA microarray

,

(phenotype)

. DNA microarray

, target

.

SVM 1995 Vladmir Vapnik

(Vapnik 1999). SVM

,

microarray .

SVM $k$

SVM

.

SVM

, Microarray SVM

.

SVM (paramet-
er) SVM .

(trade-off) (penalty) $C$

SVM . $C$

$\nu$ , RBF

SVM .

## 2 SVM

### 2. 1

$x_i \in R^n \quad i = 1, \cdots, l$

$(x_i, y_i), \cdots, (x_l, y_l) \qquad y \in \{+1, -1\}$

$x_i \qquad n \qquad\qquad\qquad , \; y_i \qquad\qquad\qquad\qquad .$

$\qquad\qquad\qquad\qquad\qquad f \qquad\qquad\qquad\qquad x \quad \{+1, -1\}$

$\qquad\qquad\qquad$ VC(Vapnik Chernonenkis) $\qquad\qquad\qquad .$

$$R(\alpha) = \int \frac{1}{2} |y - f(x, \alpha)| dP(x, y)$$

$\qquad P(x, y) \qquad x \qquad y \qquad\qquad\qquad\qquad , \; R(\alpha)$

$\qquad\qquad\qquad$ (test error) $\qquad\qquad\qquad$ (risk function)

$\qquad$ (expected risk) $\qquad\qquad .$

$\qquad\qquad\qquad\qquad f$

$, \qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad .$

$\qquad\qquad$ (trained machine) $\qquad\qquad$ (test error)

$R(\alpha) \qquad\qquad\qquad\qquad \alpha \qquad .$

$$P(x, y)$$

$$f(x)$$

(Vapnik 1999).

## 2. 1. 1 (empirical risk)

$$R(\alpha)$$

$$P(x, y) \qquad R(\alpha)$$

.    ,    (empirical risk) $R_{emp}(\alpha)$    (training set)    (error rate)    .

$$R_{emp}(\alpha) = \int \frac{1}{2l} \sum_{i=1}^{l} |y_i - f(x_i, \alpha)|$$

$$R_{emp}(\alpha) \qquad P(x, y)$$

$$\frac{1}{2} |y_i - f(x_i, \alpha)| \qquad ,$$

(Vapnik 1999).

$$R_{emp}(\alpha) \qquad R(\alpha) \qquad ,$$

$$R_{emp}(\alpha) \qquad R(\alpha)$$

.

$$\lim_{l \to \infty} R_{emp}(\alpha) = R(\alpha)$$

$$\lim_{l \to \infty} \min R_{emp}(\alpha) = \min R(\alpha)$$

VC

(Vapnik-Chernonenkis dimension)                    .

## 2. 1. 2 VC    (Vapnik-Chernonenkis dimension)

VC(Vapnik-Chernonenkis)

(capacity)

.

(bounds)                    ,

VC                                                    .

$$R(\alpha) \leq R_{emp}(\alpha) + \sqrt{(\frac{h(\log(2l/h) + 1) - \log(\eta/4)}{l})}$$

h         (non-negative)    VC          ,    l                              ,

$\eta$         (confidence)              $0 \leq \eta \leq 1$                    .

(risk bound)              ,

$\sqrt{(\frac{h(\log(2l/h) + 1) - \log(\eta/4)}{l})}$    VC                          .

,

VC                    $R_{emp}(\alpha)$    0

(test set)              (error)                    .              $R(\alpha)$

$$R_{emp}(\alpha) \qquad \frac{h}{l} \qquad\qquad ,$$

VC                                                                (Vapnik 1999).

## 2. 1. 3    (structural risk)

(empirical error)                                                (empirical risk minimizat-ion : ERM)                                      .

.

Vapnik    Chervonenkis(Vapnik and Chervonenkis, 1974)

VC

(quality)

(complexity)              (trade-off)

(structural risk minimization : SRM)                .

SRM

$R_{emp}(\alpha)$

.

$$S_1 \subset S_2 \subset \cdots \subset S_n \cdots$$

VC                                    .

$$h_1 \le h_2 \le \cdots \le h_n \cdots$$

.



그림 2-1.

　　　그림 2-1                                                                ,
                                                                    .    , VC
                                                                                  .

                                    VC

,                                                        VC         (          $h^*$)
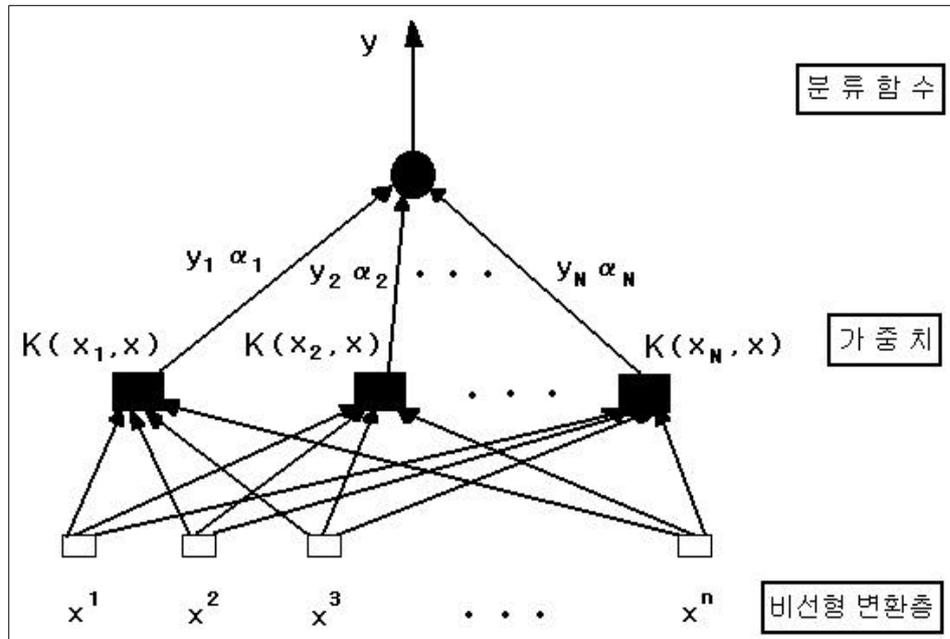(Vapnik  1999).

## 2. 2 SVM

SVM $x$ (support

vector) $x_i$ $\Phi$ $\Phi(x)$ , $\Phi(x_i)$ ,

(feature space) ( $\Phi(\cdot) : R^n \rightarrow R^p$ ( $p \gg n$ ) ). ,

$K(x, x_i)$ (input space)

(Scholkopf et al. 1999).

SVM

SVM ,

SVM .

2-2. SVM

## 2. 2. 1    SVM (Linear Support Vector Machines)

SVM    Vladimir  Vapnik                                    [ 2-3 ]

(margin)

$$w \cdot x + b = 0$$                    .

$w$                    , $x$            , $b$                    .

$$D = \{ (x_i, y_i) \}$$                    ,                    $x_i$

(class)            $y_i$    +1          ,                    - 1                    .
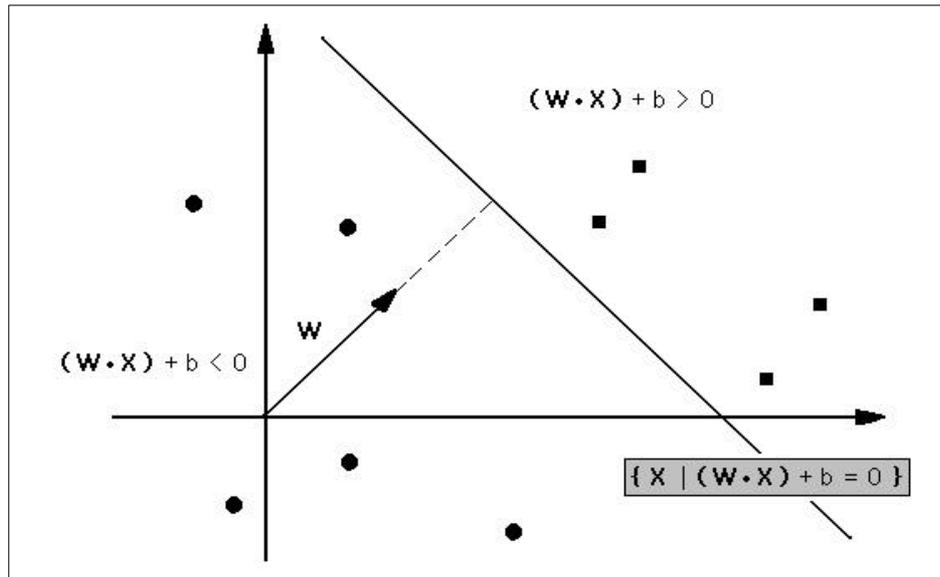
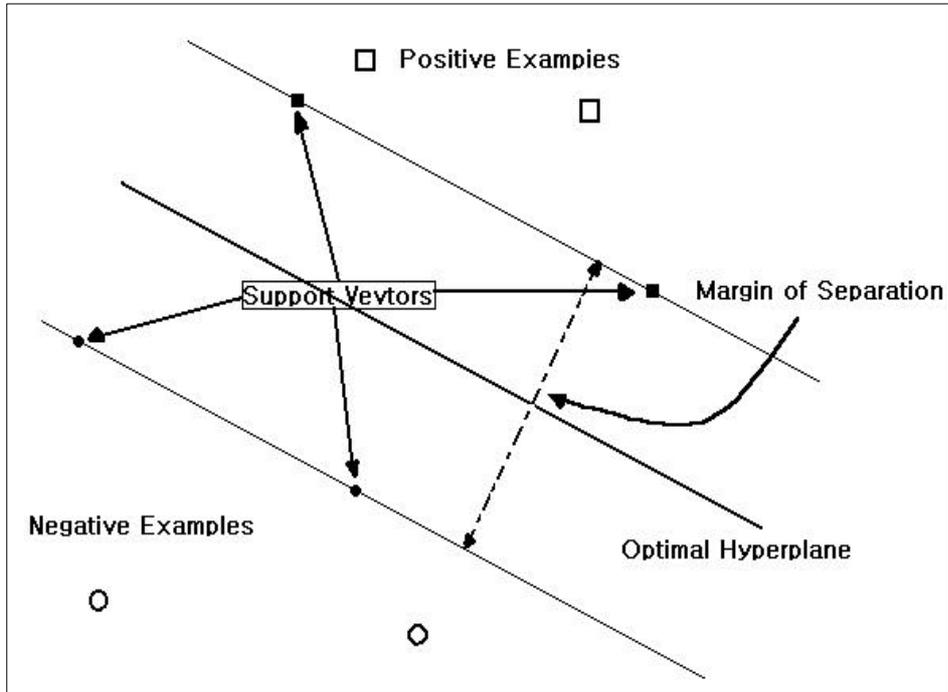SVM                                                    $w$    $b$                    .

$$x_i \cdot w + b \geq + 1 \quad ( \ y_i = \ + 1 \qquad )$$

$$x_i \cdot w + b \leq - 1 \quad ( \ y_i = \ - 1 \qquad )$$



2-3.

,

.

2-4. SVM

## 2. 2. 1. 1 (maximum margin classifier)

$$y_i(x_i \cdot w + b) \geq 1 \ , \quad i = 1, \cdots, i$$

$$\Psi(w) = ||w||^2 \qquad ||w||$$

.

(Constrained Optimization)

1 (primal) 2 (dual) .

1                    1           *w*    *b*                                        ,

(Lagrange  multipliers)  $\alpha_i$                                                  .

$$L(w, b, \alpha) = \frac{1}{2}(w \cdot w) - \sum_{i=1}^{l} \alpha_i \{ [(x_i \cdot w) - b] y_i - 1 \}$$

$( \alpha_i \geq 0 , \forall_i )$

(saddle  point)

.

$$\frac{\partial L(w_0, b_0, \alpha^0)}{\partial b} = \sum_{i=1}^{l} \alpha_i^0 y_i = 0$$

$$\frac{\partial L(w_0, b_0, \alpha^0)}{\partial w} = w_0 - \sum_{i=1}^{l} y_i \alpha_i^0 x_i = 0$$

.

$$W(\alpha) = \sum_{i=1}^{l} \alpha_i - \frac{1}{2} \sum_{i,j}^{l} \alpha_i \alpha_j y_i y_j (x_i \cdot x_j)$$

$\alpha$

2                (Quadratic  Optimization)                                   .

.

$$w_0 = \sum_{support\ vectors} y_i \alpha_i^0 x_i$$

.

$$f(x) = sign\ (\sum_{support\ vectors} y_i \alpha_i^0 (x_i \cdot x) -\ b_o\ )$$

$x_i$                 (support vector)

$\alpha$             "0"

$b_0$                            .

$$b_0 = \frac{1}{2}\ [\ (\ w_0 \cdot x^*(1)\ ) +\ (\ w_0 \cdot x^*(-1)\ )\ ]$$

$x^*(1)$                          , $x^*(-1)$

(Vapnik 1999).

2-5.


## 2. 2. 1. 2                                **(soft margin classifier)**


SVM


(slack)       $\xi_i$                                      (generalized

hyperplane)                                   .

   ( $w$, $b$)            $r$                     ($x_{i,}y_i$)                $\xi_i$

         .


$$\xi_i\,(\,(\,x_{i,}y_i\,)\,,\,(\,w,\,b)\,,\,r\,) = \;\xi_i = \;\max\,(\,0\,,\,r\,-\;\,y_i\,(\,w\cdot x_i\,+\;\,b)\,)$$


- 14 -

,  $\xi_i > r$ $(x_i, y_i)$ ,

$\xi_i$ $r$

(Cristianini et al. 2000).

$\xi_i \geq 0$

$(i = 1, \cdots, l)$ $y_i \cdot ((w, x_i) + b) \geq 1 - \xi_i$

$C$ (Cortes, and Vapnik 1995).

$$\tau(w, \xi) = \frac{1}{2} \|w\|^2 + C \sum_{i=1}^{l} \xi_i$$

$\alpha_i$ $k(x_i, x_j)$

. $0 \leq \alpha_i \leq C (i = 1, \cdots, l)$ $\sum_{i=1}^{l} \alpha_i y_i = 0$

.

$$\max \ W(\alpha) = \sum_{i=1}^{l} \alpha_i - \frac{1}{2} \sum_{i,j=1}^{l} y_i y_j \alpha_i \alpha_j (x_i \cdot x_j)$$

$b$ $\alpha_i < C$ $\xi_i$ "0"

(Scholkopf et al. 1999).

$$f(x) = sign \left( \sum_{i=1}^{l} y_i \alpha_i \cdot k(x, x_i) + b \right)$$

$\xi_i$               $C$

(Coretes, and Vapnik 1995).        $C$

(trade-off)                          ,

.

$r$

$\xi_i$    $\xi_j$           $r$

,                                    "0"

.



2-6.

## 2. 2. 2     S V M

(nonlinear decision surface)                    .

$K ( \cdot , \cdot )$

.

.

$$K (x , y) = ( \Phi(x) \cdot \Phi(y) )$$

,        $x$      $\Phi( \cdot )$

.



2-7.

(polynomial)　　　　　, RBF
(Radial Basis Function)　　　　　, 　　　　　(multi-layer perceptron)

　　　　　,

.

### 2. 2. 2. 1

(dot product)

,

.

$$K(x, x_i) = ((x \cdot x_i) + 1)^d$$

$d$　　　　　　　　　　　　　.

### 2. 2. 2. 2 RBF

$$f(x) = sign\left( \sum_{i=1}^{N} a_i K_\gamma(|x - x_i|) - b \right)$$

$K_\gamma(|x - x_i|)$　　　　　　　　　　$|x - x_i|$

,　　　　　　　　　　　　　　.

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

RBF

, RBF . $\sigma$

(smoother) .

# 3    Microarray

## 3. 1 DNA Microarray

Array                    3                            filter    array,  Oligonucleotide
array        cDNA    microarray            .                filter    array
Oligonucleotide array    cDNA microarray
cDNA    microarray                        .  cDNA
microarray                        Northern  blot
(Reverse  Northern).    , filter      mRNA
(probe)            Northern
mRNA
.                Northern

.

,                                      filter  paper
. Filter  paper
.

,
.

,              mRNA
Northern  blot                                    tagging

, cDNA microarray                          mRNA      labeling

        .                                                    ,
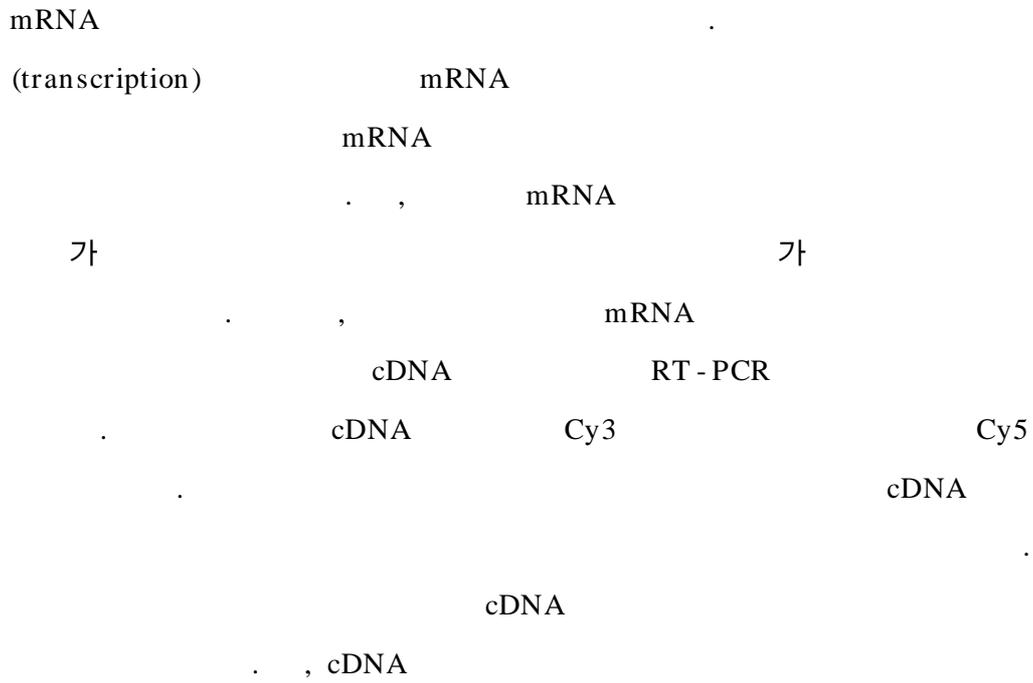
                                    ,                              cDNA

microarray                                    labeling           .

        Cy3    Cy5    Alexa dye                         ,

                                    end  labeling

Cy3-dUTP,  Cy5-dUTP                    (reverse  transcription)

cDNA                            .

        ,                                            DNA

                    (hybridization)

                    .          cDNA  microarray

(treatment  group)        (control group) 2        mRNA        Cy3

Cy5    labeling

                                                    .


                                                    .

    DNA  chip


                    .   ,



            .          DNA chip          cDNA chip          Affimatrix

    oligochip              .          Affimatrix      oligochip



                    .     , cDNA chip

.

DNA chip                                    DNA microarray          ,

microarray                      .

.

(normalization)                        microarray

.                    cDNA microarray

Chen   et   al.(1997)    Cy3(          )              (intensity)

Cy5(          )

.                                          (image

analysis  program)                                                    .

,  Yang  et  al.(2001)    2001    1    SPIE  BiOE

LOWESS

.         Dudoit(2000)          Yang

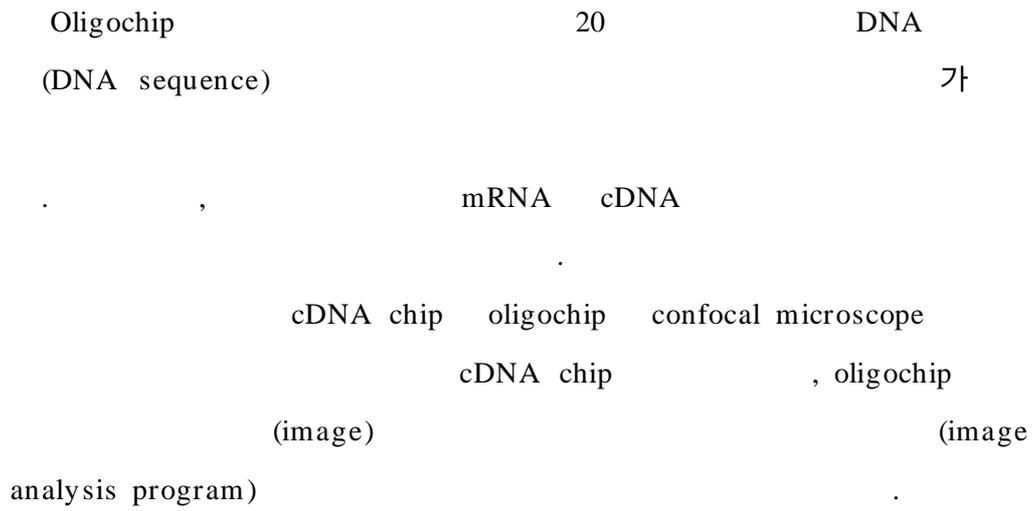LOWESS                                    .

microarray                                        (Newton  2001)

(Eisen  1998).          ,

.


## 3. 2 DNA  chip


### 3. 2. 1 DNA  chip


cDNA  chip

mRNA                                                    .

(transcription)                        mRNA


                        mRNA


                .    ,                    mRNA



                .      ,                        mRNA

                    cDNA                RT‐PCR

        .            cDNA            Cy3                        Cy5

            .                                            cDNA

                                                            .


                        cDNA

        .    , cDNA


                                                        .


## 3. 2. 2 Oligochip


 Oligochip                            20                    DNA

    (DNA  sequence)



        .            ,                    mRNA    cDNA

                                        .

                    cDNA chip    oligochip    confocal microscope

                        cDNA  chip                , oligochip

            (image)                                        (image

analysis program)                                        .

## 3. 2. 3 　　　(Normalization)

microarray

.　　　,

,

,

.

## 3. 2. 3. 1

.

housekeeping 　　　(gene)

spiked .

## 3. 2. 3. 2

housekeeping Cy3

. ,

.

.

## 3. 2. 4 　　　　(Normalization)

　　　　　　　　　　　　　　　　　　　　　　Cy3　　　　　　$G_j$

Cy5　　　　$R_j$　　　　　.　　, $j$　　　　　　　　　　　　$id$　　.

$$M_j \;=\; \log \frac{R_j}{G_j} \;=\; \log R_j \;-\; \log G_j$$

$$A_j \;=\; \log \sqrt{R\,G} \;=\; \frac{(\log R_j \;+\; \log G_j)}{2}$$

　　　　　　　　　　　.　　　　　　　　　　　　　$R_j$　　$G_j$

　　　　　　　　　　　.

## 3. 3 Microarray

　　cDNA　microarray

　　　　　　　　　　　　　　　　　　　　　　　　　　　.

　　Dudoit et al.(2000)　　　　　　,

　　　　　　　　3　　　　　　　　.

　　　　,

　　　　　　　　　.　　, 　　　　　(cluster analysis)

　　　　　　(unsupervised learning)　　　　.

,

.                                    (discriminant  analysis)

,                    (supervised  learning)                .

,                            class                      ,   '

(marker  gene)'                              .                    (variable  selection)

.

# 4                                                      SVM

## 4.1          Microarray

SVM

                              microarray                    .

  microarray                          ,             Cy3    Cy5

                  .     ,                            20                    ,

              100               .                100                ,

      20%                                  ,    , Cy5

      ,          80%    Cy3    Cy5                          0,    ,

                  .

                                              (gene)                ,    , Cy3

Cy5                          ($y$)                                        .


$$y = \log\left(\frac{R}{G}\right) + \varepsilon \quad , \quad \varepsilon \sim MVN(\widetilde{\mu}, \widehat{\Sigma})$$


      ,                      80%

    , $\left(\dfrac{R}{G}\right) = 1$                                        .


$$y = \varepsilon \quad , \quad \varepsilon \sim MVN(\widetilde{\mu}, \widehat{\Sigma})$$

20%

．

$$y = \log (\frac{R}{G}) + \varepsilon \quad , \quad \varepsilon \sim MVN(\tilde{\mu}, \hat{\Sigma})$$

## 4. 2 　　　　　　　　　　　　　　SVM

SVM

(parameter) 　　　　　　　　　　　．

(trade-off)

(penalty) 　　　$C$　　　．

SVM 　　　　　　　(classification) 　　　　　　$c$- 　　　　$\nu$-

．

$c$- 　　　(c-classification) 　　　(cost)

(Sensitivity), 　　　(Specificity), 　　　　　(Positive Predict-

ed Value), 　　　(Negative Predicted Value), 　　　(Correct Prop-

ortion), 　　　(Miss Correct Proportion) 　　　　　　　．

(tuninig)

$k$-fold 　　　　　　(cross validation) 　　LOO(Leave-One

-Out) 　　　(Duan et al. 2001), 　　　　　LOO

．

$k$-fold 　　　　　　　　　　(training data)

(mutually exclusive) 　　　　　　　$k$

- 28 -

. LOO $l-1$

1 (test) . $l$

(expected generalization error) . (Duan et al. 2001)

$\nu\text{-classification}$ $\nu$

$0 < \nu < 1$ .

## 4.3

$C$

, $C$ (infinite)

(training data) (smola,

1997).

SVM $C$ $\nu$ .

$C$ , LOO(Leave-One-Out)

[ 4-1 ], [ 4-2 ]

$C = 1.4$ .

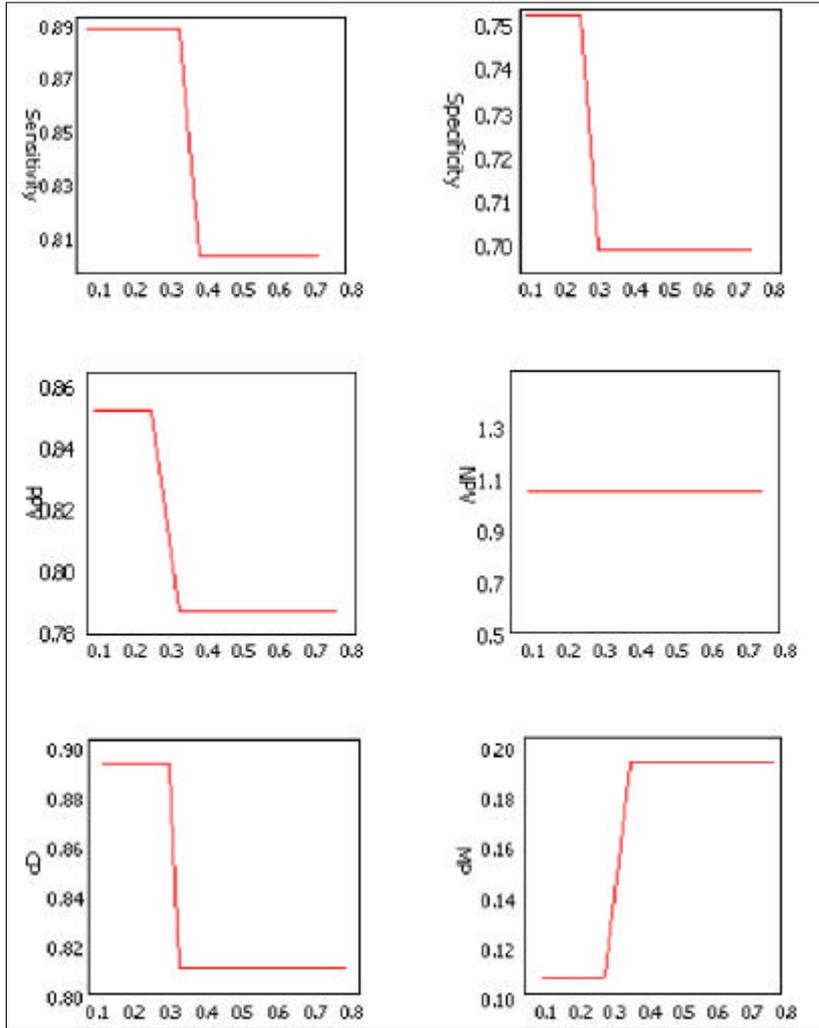, $\nu$ $0 < \nu < 1$ [ 4-3 ], [

4-4 ] 0.3 .

4- 1.                                                      (radial)

4-2.                                                    (polynomial)

4-3. ν                                    (radial)

4-4. ν (polynomial)

表 1 . 　　　　　　　　SV 　　　　　(c-classification)

| | cost | $\sigma=0.5$ | $\sigma=1.0$ | $\sigma=1.5$ |
|---|---|---|---|---|
| radial | 1.0 | 30 | 30 | 30 |
| | 1.3 | 30 | 30 | 30 |
| | 1.5 | 30 | 30 | 30 |
| | 1.7 | 30 | 30 | 30 |
| | 1.9 | 30 | 30 | 30 |
| | 2.0 | 30 | 30 | 30 |
| polynomial (degree=3) | 1.0 | 24 | 25 | 26 |
| | 1.3 | 24 | 25 | 26 |
| | 1.5 | 24 | 25 | 26 |
| | 1.7 | 24 | 25 | 26 |
| | 1.9 | 24 | 25 | 26 |
| | 2.0 | 24 | 25 | 26 |

$c$-　　　　　　　　　　　　　Support　Vector

　　　　　[　1 ]　　　　　　　.

　　　　　　　　　radial

Support　Vector　　　　　　　　　　　　　,

　　　　　　　　　　　　　　　　Support　Vector

　　　　.

## 4．4

　$c$-　　　　　　　1.4 ，$\nu$-　　　$\nu$　　0.3　　　　　，

radial　　3　　　　　　　Cy5

　　　　　　　　　　　100　　　　　，

，　　　　　，　　　　，　　　，

[　2,3 ]　　　　　　．

2 .                                                    (c-classification)

| parameter mean±sd | log 1.5 | log 2.0 | log 2.5 | log 3.0 |
|---|---|---|---|---|
| **radial** | | | | |
| Sensitivity | 0.729±0.218 | 0.841±0.188 | 0.896±0.163 | 0.968±0.083 |
| Specificity | 0.761±0.220 | 0.812±0.194 | 0.882±0.193 | 0.942±0.128 |
| PPV | 0.746±0.239 | 0.821±0.187 | 0.892±0.168 | 0.958±0.082 |
| NPV | 0.726±0.231 | 0.845±0.170 | 0.917±0.128 | 0.970±0.082 |
| CP | 0.700±0.152 | 0.804±0.130 | 0.883±0.121 | 0.957±0.062 |
| MP | 0.300±0.152 | 0.196±0.130 | 0.117±0.121 | 0.043±0.062 |
| **polynomial (degree=3)** | | | | |
| Sensitivity | 0.765±0.206 | 0.883±0.166 | 0.949±0.118 | 0.992±0.043 |
| Specificity | 0.727±0.207 | 0.771±0.199 | 0.837±0.201 | 0.921±0.145 |
| PPV | 0.697±0.215 | 0.753±0.212 | 0.845±0.178 | 0.936±0.110 |
| NPV | 0.788±0.201 | 0.896±0.148 | 0.962±0.089 | 0.994±0.031 |
| CP | 0.718±0.140 | 0.800±0.138 | 0.885±0.118 | 0.958±0.069 |
| MP | 0.282±0.140 | 0.200±0.138 | 0.115±0.118 | 0.043±0.069 |

* cost=1.4, number of iteration=100

3 .　　　　　　　　　　　　　　　　　　（ $\nu$-classification）

| parameter mean±sd | | log 1.5 | log 2.0 | log 2.5 | log 3.0 |
|---|---|---|---|---|---|
| | Sensitivity | 0.727±0.214 | 0.842±0.188 | 0.898±0.159 | 0.968±0.083 |
| | Specificity | 0.762±0.219 | 0.817±0.189 | 0.883±0.191 | 0.942±0.128 |
| | PPV | 0.749±0.245 | 0.827±0.182 | 0.893±0.164 | 0.958±0.082 |
| radial | NPV | 0.721±0.226 | 0.844±0.171 | 0.919±0.127 | 0.970±0.082 |
| | CP | 0.700±0.144 | 0.808±0.129 | 0.885±0.118 | 0.957±0.062 |
| | MP | 0.300±0.144 | 0.192±0.129 | 0.115±0.118 | 0.043±0.062 |
| | Sensitivity | 0.765±0.206 | 0.883±0.166 | 0.949±0.118 | 0.992±0.043 |
| | Specificity | 0.727±0.207 | 0.771±0.199 | 0.837±0.201 | 0.921±0.145 |
| | PPV | 0.697±0.215 | 0.753±0.212 | 0.845±0.178 | 0.936±0.110 |
| polynomial (degree=3) | NPV | 0.788±0.201 | 0.896±0.148 | 0.962±0.089 | 0.994±0.031 |
| | CP | 0.718±0.140 | 0.800±0.138 | 0.885±0.118 | 0.956±0.069 |
| | MP | 0.282±0.140 | 0.200±0.138 | 0.115±0.118 | 0.043±0.069 |

\* $\nu$=0.3, number of iteration=100

( ,

, , , , )

, microarray Cy5

.

. log 0.5 , *c*- radial

0.729±0.218 , $\nu$- 0.727±0.214

.

*c*- 1.4 , $\nu$- $\nu$ 0.3 , radial

3

100 , , , ,

, [ 4,5 ] .

4 .                                                                          (c-classification)

| parameter<br>mean±sd | $\sigma$=0.5 | $\sigma$=1.0 | $\sigma$=1.5 |
|---|---|---|---|
| | Sensitivity | 0.920±0.157 | 0.808±0.202 | 0.763±0.237 |
| | Specificity | 0.941±0.113 | 0.784±0.209 | 0.767±0.232 |
| | PPV | 0.941±0.109 | 0.766±0.241 | 0.749±0.264 |
| radial | NPV | 0.927±0.136 | 0.814±0.195 | 0.764±0.235 |
| | CP | 0.925±0.100 | 0.756±0.151 | 0.708±0.173 |
| | MP | 0.075±0.100 | 0.244±0.151 | 0.292±0.173 |
| | Sensitivity | 0.976±0.061 | 0.885±0.175 | 0.851±0.178 |
| | Specificity | 0.917±0.135 | 0.761±0.237 | 0.687±0.202 |
| | PPV | 0.928±0.114 | 0.760±0.238 | 0.644±0.235 |
| polynomial<br>(degree=3) | NPV | 0.969±0.081 | 0.879±0.181 | 0.860±0.189 |
| | CP | 0.944±0.079 | 0.790±0.154 | 0.722±0.147 |
| | MP | 0.056±0.079 | 0.210±0.154 | 0.278±0.147 |

* log 2.0, cost=1.4, number of iteration=100

表 5 .　　　　　　　　　　　　　　　　　　（$\nu$-classification）

| parameter<br>mean±sd | | $\sigma$=0.5 | $\sigma$=1.0 | $\sigma$=1.5 |
|---|---|---|---|---|
| radial | Sensitivity | 0.922±0.145 | 0.806±0.204 | 0.770±0.225 |
| | Specificity | 0.937±0.113 | 0.782±0.210 | 0.770±0.227 |
| | PPV | 0.932±0.116 | 0.765±0.241 | 0.753±0.259 |
| | NPV | 0.925±0.137 | 0.812±0.194 | 0.765±0.236 |
| | CP | 0.925±0.089 | 0.756±0.153 | 0.712±0.170 |
| | MP | 0.075±0.089 | 0.244±0.153 | 0.288±0.170 |
| polynomial<br>(degree=3) | Sensitivity | 0.976±0.061 | 0.885±0.175 | 0.851±0.178 |
| | Specificity | 0.917±0.135 | 0.761±0.237 | 0.687±0.202 |
| | PPV | 0.928±0.114 | 0.760±0.238 | 0.644±0.235 |
| | NPV | 0.969±0.081 | 0.879±0.181 | 0.860±0.189 |
| | CP | 0.944±0.079 | 0.790±0.154 | 0.722±0.147 |
| | MP | 0.056±0.079 | 0.210±0.154 | 0.278±0.147 |

\* log 2.0, $\nu$=0.3, number of iteration=100

(sensitivity, specificity, PPV, NPV, CP, EP)                    ,

    microarray                 (noise)

                        .

        ,

        .              (σ)    0.5        ,    *c*-              radial

0.920±0.157      ,    ν-                  0.922±0.145

                            .        ,                    *c*-              radial

        0.920±0.157        ,                        0.976±0.661

                    .

**5**

microarray        SVM

. microarray             (generating)

SVM             .

SVM    Vapnik

,

.

SVM

microarray            ,     microarray

(S-PLUS)         .     SVM

Chen        R-Package     SVM

.

, $c$-        radial

Support Vector                  ,

Support

Vector          .

(   ,    ,      ,

,     ,     )               ,     microarray

Cy5

.    ,

.    ,                  (   ,

,      ,      ,     ,     )

,    microarray

.

SVM

.

SVM

.

SV    VC

.

microarray                                                                    ,

.

.

,          , cDNA  microarray.                       , 2001, 21(3):467-476

,          , S-PLUS                           .                , 2000

,          ,  SVM                         .                , 2000,  49(4):
44-48

,          ,          ,          , S-PLUS                          .          , 1997

,          ,          ,  SVM
.                              (    ),  1999,  26(2):60-62

Burges,  Christopher.  J.  C.,  A  Tutorial  on  Support  Vector  Machines
for  Pattern  Recognition. Boston  :  Kluwer  Academic  Publisher,  1998

Chen,  Ratio-Based  Decision,  and  the  Quantitative  Analysis  of  cDNA
Microarray  Images. *Journal of Biomedical Optics*, 1997, 2(4):364-374

Chris  J. C.  Burges,  Federico  Girosi,  Partha  Niyogi,  Tomaso  Poggio,
Vladimir  Vapnik.,  Comparing  Support  Vector  Machines  with
Gaussian  Kernels  to  Radial  Basis  Function  Classifiers. *IEEE
Transactions  on  signal processing,*  1997,  45(11):2758-2765

Cristianini, Nello, and John Shawe-Taylor., An introduction to support Vector Machines and other kernel-based learning method. Cambridge: Cambridge University Press, 2000

Cleveland, Robust locally weighted regression and smoothing scatterplots. *Journal of the American Statistical Association*, 1979, 74(368):829-836

Cortes, C., Vladimir Vapnik, Support Vector Networks. Machine Learning, 1995, 20:273-297

Kaibo Duan, S Sathiya Keerthi, Aun Neow Poo., Evaluation of Simple Performance measures for Tuning SVM Hyperparameters. *Control Division Technical Report*, 2001

Michael B. Eisen, Paul T. Spellman, Patrick O. Brown, and David Bostein, Cluster analysis and display of genome-wide expression patterns. *In Proceedings of the National Academy of Sciences*, 1998, 95:14863-14868

M. A. Newton, C. M. Kendziorski, C. S. Richmond, F. R. Blattner, K. W. Tsui, On Differential Variability of Expression Ratios: Improving Statistical Inference about Gene Expression Changes from Microarray Data. *Journal of the American Computation Biology*, 2001, 8(1):37-52

M. Kathleen Kerr, Mitcheel Martin, Gary A. Churchill, Analysis of Variance for Gene Expression Data. *Journal of the American Computation Biology*, 2001, 7:819-837

Michael P. S. Brown, William Noble Grundy, David Lin, Nello Cristianini, Charles Sugnet, Manuel Ares, Jr., David Haussler, Knowledge-based analysis of microarray gene expression data by using support vector machines. *PNAS*, 2000, 97:262-267

Michael P. S. Brown, William Noble Grundy, David Lin, Nello Cristianini, Charles Sugnet, Manuel Ares, Jr., David Haussler, Support Vector Machine Classification of Microarray Gene Expression Data. *Technical Report* UCSC-CRL-99-09, 1999

Sandrine Dudoit, Jane Fridlyand, and Terry Speed, Comparison of discrimination methods for the classification of tumors using gene expression data. *Technical Report*, 2000
,
Sandrine Dudoit, Yee Hwa Yang, Matthew J. Callow, and Terence P. Speed, Statistical methods for identifying differentially expressed genes in replicated cDNA microarray experiments. *Technical Report* 578, Department of Biochemistry, Stanford University School of Medicine, 2000

Scholkopf, Bernhard, Chris Burges, and Alex J. Smola, Advences In

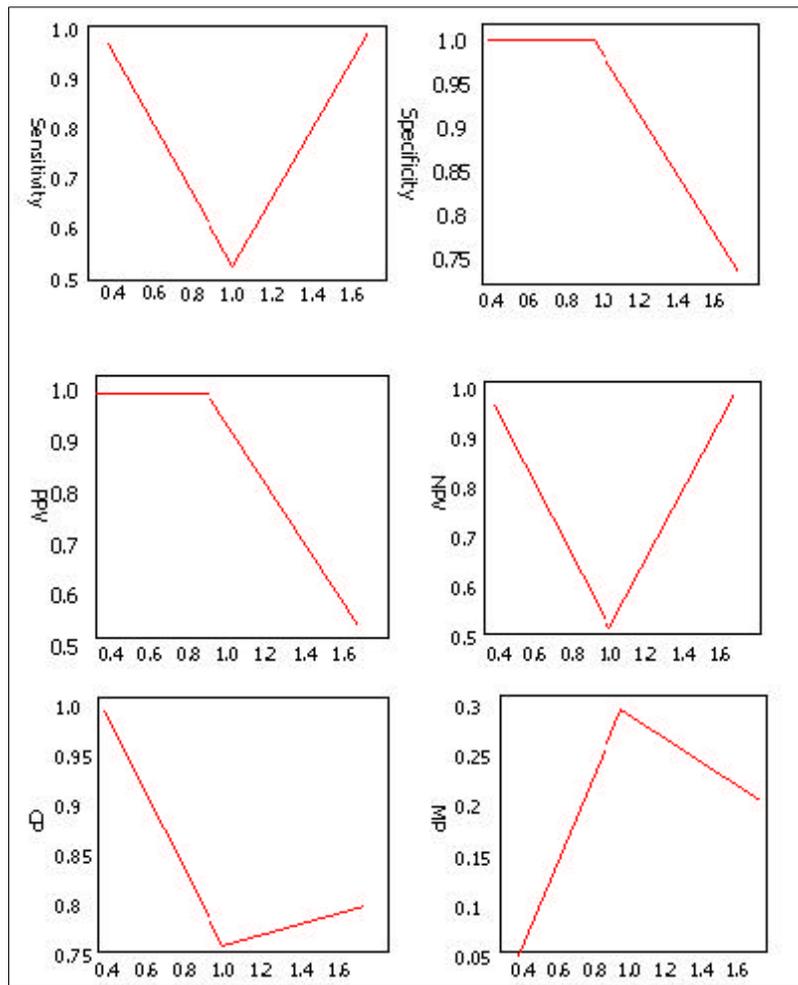kernel methods:Support Vevtor Machines. Cambridge:MIT Press, 1999

A. Smola, and B. Scholkopf, On a kernel-based Method for Pattern Recognition, Regression, Approximation and Operator Inversion. *GMD Technical Report*, 1997

T. R. Golub, D. K. Slonim, P.Tamayo, C. Huard, M. Gaasenbeek, J. P. Meirov, H. Coller, M. L. Loh, J. R. Downing, M. A. Caligiuri, C. D. Bloomfield, E. S. Lander, Molecular Classification of Cancer: Class Discovery and Class Prediction by gene Expression Monitoring. *Scinece*, 1999, 286:531-537

Vladimir N. Vapnik, The Nature of Statistical Learning Theory. Springer, 1999

W. N. Venables, B. D. Ripley, Modern Applied Statistics with S-PLUS (3rd edition). Springer, 1999

Yee Hwa Yang, Normalization for cDNA Microarray Data. In *SPIE BioE,* 2001
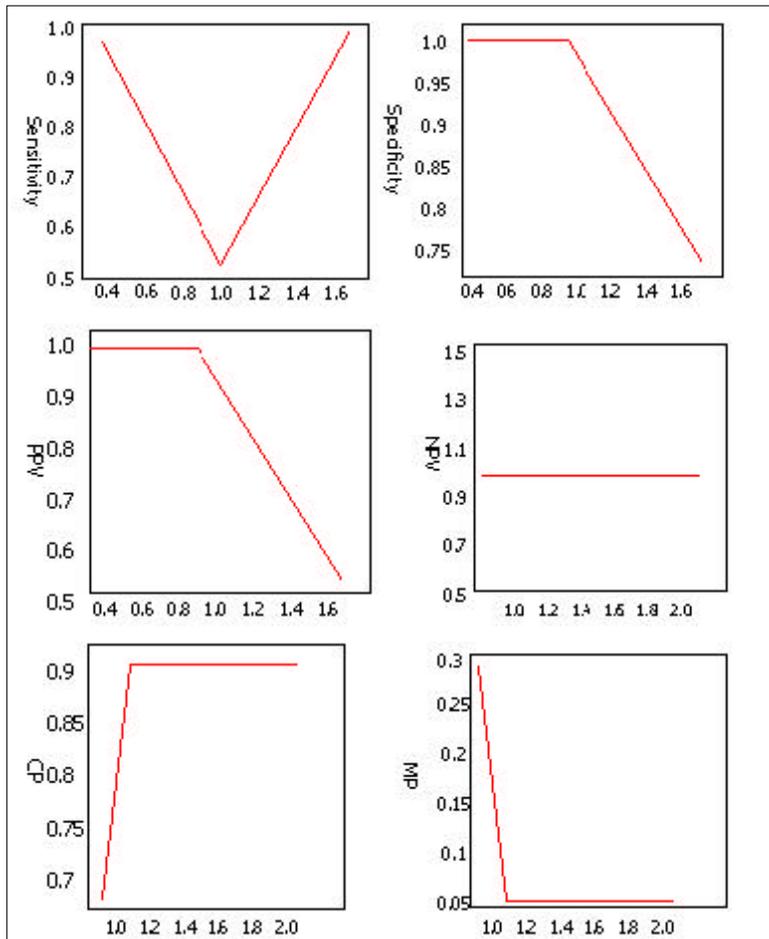
1.                                          radial (cost=1.4)
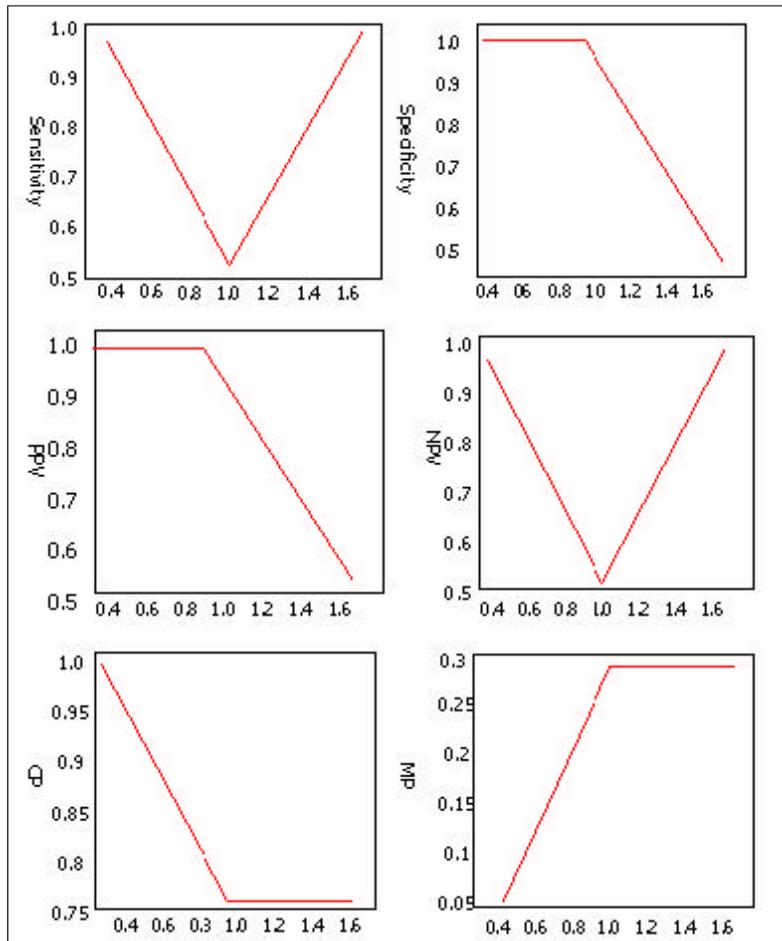
2.                                                                        polynomial (cost=1.4)

3.                                              radial （$\nu$ =0.3)

4. polynomial （ $\nu$ =0.3)

# ABSTRACT

## Microarray Gene Expression Data Classification Using Support Vector Machine

Ku, Kyung Min

Dept. of Biostatistics and Computing

The Graduate School

Yonsei University

In the thesis, we introduce the Support Vector Machine(SVM) classification from microarray data and use simulation of microarry data for kernel-function in order to evaluate SVM. The point of this thesis is to evaluate by SVM classifier using microarray data which is generated by adopting S-PLUS and R-Package.

In conclusion, the simulation result has the following result.

First, the increase of log ratio(Cy5/Cy3), the value of each evaluation item (sensitivity, specificity, Positive Predicted Value, Negative Predicted Value, Correct Proportion, Miss Correct Proportion) was improved. The intensity of Cy5 appeared high in microarray experiment.

Second, classification was more accurate but, there was no significant difference between the kernel-function and classification method. With the increase of standard deviation, the value of each evaluation

item was decreased. And the classification became poorer as the noise in microarray experiment increased.