

Distance Weighted
Discrimination

(improving SVM with SOCP)

Steve Marron
Statistics, UNC

Mike Todd
Operations Research,
Cornell

[www.unc.edu/depts/statistics/
postscript/papers/marron/HDD/DWD](http://www.unc.edu/depts/statistics/postscript/papers/marron/HDD/DWD)

Introduction, problem

Two class discrimination

Have probability distrib. on $X \times Y$

$X \subseteq \mathbb{R}^d$, $Y = \{+1, -1\}$, unknown

Have sample $\{(x_i, y_i) : i=1, 2, \dots, n\}$,
training set.

Want to deduce a method to predict,
for a new observation (x, y) from the
dist., with x known, what y is.

Example: components of x :
medical history, results of tests.

$y = \begin{cases} +1 & \text{have certain disease} \\ -1 & \text{don't} \end{cases}$

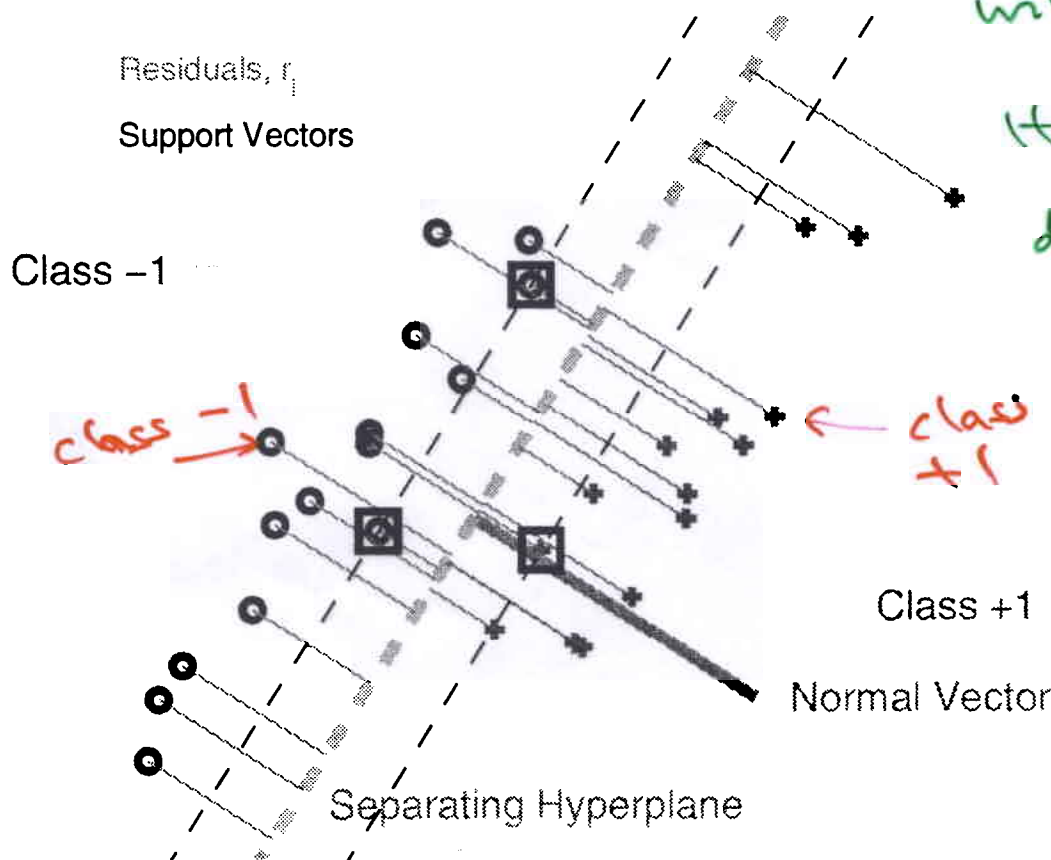
Rule: $y = f(x)$
 ↑
 constructed from
 training set

Linear rule: for some $w \in \mathbb{R}^d$, $\beta \in \mathbb{R}$ 2

$$f(x) = \text{sign}(w^T x + \beta)$$

Seems substantial limitation. But can expand vector x , eg. to $(x_1, \dots, x_d, x_1^2, x_1 x_2, \dots, x_d^2)$, so not bad.

Support Vector Machine (Vapnik, 1982)



What's wrong with this??

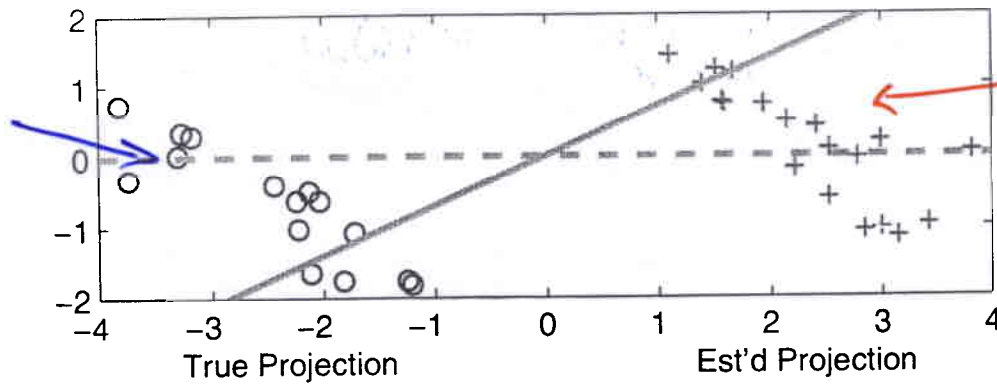
HDLSS setting

$d \gg 1$, $n \ll d$

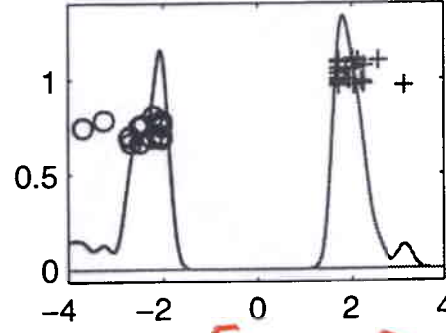
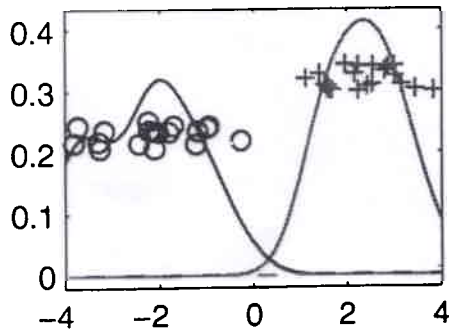
sensitivity!

Linear SVM, C = 1000, dimension = 39

20
class
-1
pts



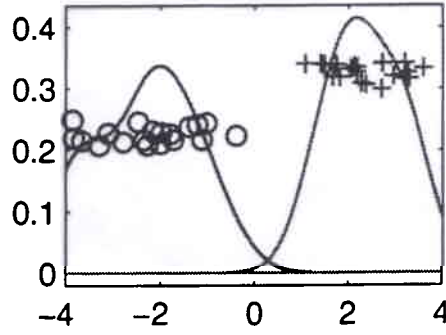
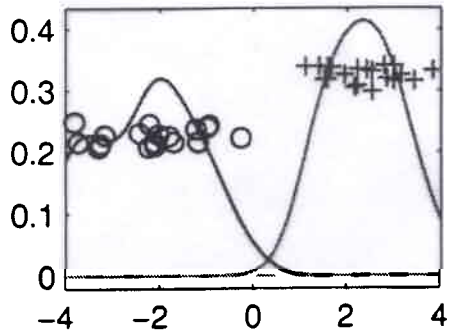
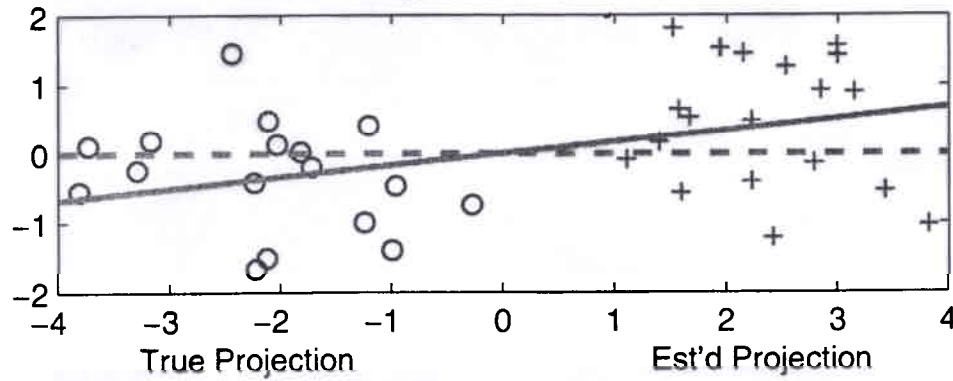
20
class
+1
pts



If $y = \pm 1, x \sim N(\pm 2.2e_1, I)$

"data piling"

Distance Weighted Disc., dimension = 39



Optimization Formulations

4

Have $(x_i, y_i) \in \mathbb{R}^d \times \{\pm 1\}$, $i=1, \dots, n$.

Let $X = [x_1 \dots x_n]$, $y = (y_i) \in \mathbb{R}^n$,

$Y = \text{Diag}(y) \in \mathbb{R}^{n \times n}$ ($X \in \mathbb{R}^{d \times n}$).

Given $w \in \mathbb{R}^d$, $\beta \in \mathbb{R}$, residuals

$$\bar{r}_i := y_i (x_i^T w + \beta) \quad \text{or}$$

$$\bar{r} := Y(X^T w + \beta e) = YX^T w + \beta y$$

Want all components of \bar{r} "large and pos."

May be impossible to even make all pos. Introduce $\xi \in \mathbb{R}_+^n$ errors, to give perturbed residuals

$$r := \bar{r} + \xi = YX^T w + \beta y + \xi$$

SVM: "max_{w, \beta} min_i r_i" $\|r^{-1}\|_{\infty}$

DWD: "min_{w, \beta} \sum_i 1/r_i" $\|r^{-1}\|_1$

Support vector machine

5

First, assume no errors.

$$\max_{\delta, w, \beta} \delta, \quad \bar{r} = YX^T w + \beta y, \quad \bar{r} \geq \delta e, \quad \frac{1}{2} w^T w \leq \frac{1}{2}$$

$$\min_{w, \beta} \frac{1}{2} w^T w, \quad YX^T w + \beta y \geq e.$$

Now allow perturbations, penalized via L_1 norm

$$\begin{aligned} \min_{w, \beta, \xi} \quad & \frac{1}{2} w^T w + C e^T \xi \\ (P_{SVM}) \quad & YX^T w + \beta y + \xi \geq e, \\ & \xi \geq 0. \end{aligned}$$

QP with dual

$$\begin{aligned} (D_{SVM}) \quad & \max_{\alpha} \quad -\frac{1}{2} \alpha^T YX^T X Y \alpha + e^T \alpha \\ & y^T \alpha = 0 \\ & 0 \leq \alpha \leq C e. \end{aligned}$$

Optimality Conditions:

$$(X_+ \alpha_+ - X_- \alpha_-) X^T \alpha = w, \quad y^T \alpha = 0, \quad e_+^T \alpha_+ = e_-^T \alpha_-$$

$$s = Y X^T w + \beta y + \gamma - e \geq 0, \quad \alpha \geq 0, \quad s^T \alpha = 0$$

$$C e - \alpha \geq 0, \quad \gamma \geq 0, \quad (C e - \alpha)^T \gamma = 0.$$

$$X = [X_+ \quad X_-], \quad y = \begin{pmatrix} \vdots \\ \vdots \\ \vdots \\ \vdots \\ \vdots \\ \vdots \\ \vdots \\ \vdots \\ \vdots \\ \vdots \end{pmatrix}, \quad \alpha = \begin{pmatrix} \alpha_+ \\ \alpha_- \end{pmatrix}$$

Geometric interpretation:

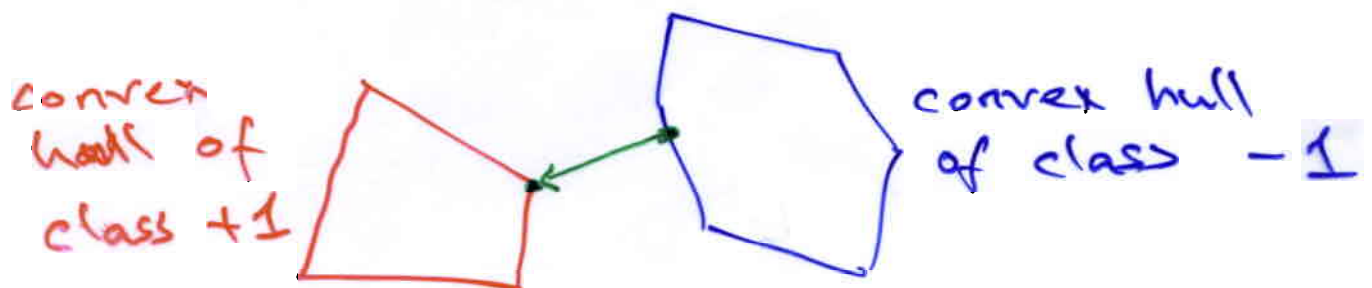
Assume optimal $\alpha < C e$.

$$\text{Let } \alpha = \gamma \hat{\alpha}, \quad e_+^T \hat{\alpha}_+ = e_-^T \hat{\alpha}_- = 1,$$

$$\text{so } w = \gamma (X_+ \hat{\alpha}_+ - X_- \hat{\alpha}_-)$$

By optimizing over γ for fixed $\hat{\alpha}$,

$$(D_{\text{svm}}) \equiv \max 2 / \|X_+ \hat{\alpha}_+ - X_- \hat{\alpha}_-\|$$



Distance weighted discrimination

7

Now minimize sum of reciprocals of perturbed residuals, + penalty term:

$$\min_{r, w, \beta, \xi} \sum_i \frac{1}{r_i} + C e^T \xi$$

$$r = YX^T w + \beta y + \xi \geq 0$$

$$\frac{1}{2} w^T w \leq \frac{1}{2}, \quad \xi \geq 0.$$

(Could consider $\sum f(r_i)$, f smooth, convex, $\rightarrow +\infty$ as $\arg \downarrow 0$.)

Reformulate as SOCP. Use second-order

Lorentz cone

$$S_{d+1} := \{ (\psi; u) \in \mathbb{R}^{d+1} : \psi \geq \|u\| \}.$$

Write $r_i = \rho_i - \sigma_i$, where $\rho_i = \frac{r_i + \frac{1}{r_i}}{2}$,

$\sigma_i = \frac{\frac{1}{r_i} - r_i}{2}$. Then $\rho_i + \sigma_i = \frac{1}{r_i}$, and

$\rho_i^2 - \sigma_i^2 = 1$, or $(\rho_i; \sigma_i; 1) \in S_3$. Also,

$(1, w) \in S_{d+1}$.

$$\min_{\omega, w, \beta, \xi, \rho, \sigma, \tau} C e^T \xi + e^T \rho + e^T \sigma$$

(P DWD)

$$\begin{aligned} YK^T w + \beta z + \xi - \rho + \sigma &= 0 \\ \omega &= 1 \\ \tau &= \rho \end{aligned}$$

$$(\omega, w) \in S_{d+1}, \xi \geq 0, (\rho_i, \sigma_i, \tau_i) \in S_3, \text{ all } i.$$

SOCTs have duals: get

$$\max_{\alpha, \gamma, \tau, \pi, \rho, x, \lambda, \mu, \nu}$$

$$\gamma + e^T \gamma$$

(D DWD)

$$\begin{aligned} \gamma + \pi &= 0 \\ X\gamma\alpha + \rho &= 0 \\ y^T \alpha &= 0 \\ \alpha + x &= \rho e \\ -\alpha &= e \\ \alpha + \lambda &= e \\ \alpha + \mu &= e \\ \gamma + \nu &= 0 \end{aligned}$$

$$(\pi, \rho) \in S_{d+1}, x \geq 0, (\lambda_i, \mu_i, \nu_i) \in S_3 \text{ all } i.$$

Can greatly simplify: equivalent to ⁹

$$\max_{\alpha} - \|X^T \alpha\| + 2e^T \sqrt{\alpha},$$
$$y^T \alpha = 0, \quad 0 \leq \alpha \leq Ce.$$

Can check (P_{DWD}) and (D_{DWD}) have strictly feasible solutions, so have optimal solns and no duality gap.

Optimality conditions:

Want feasibility + orthog:

$$(\pi; \rho)^T (w; w) = x^T \xi = (\lambda; \mu; \nu)^T (\rho; \sigma; \tau) = 0$$

Note: for SOC: $(\psi; u)^T (\phi; v) = 0 \iff$

$(\psi; u) = 0$ or $(\phi; v) = 0$ or $\|u\| = \psi, \|v\| = \phi, u \perp v$.

$$YX^T w \rightarrow \beta y + \xi - \rho + \sigma = 0, \quad y^T \alpha = 0$$

$$\alpha > 0, \quad \alpha \leq Ce, \quad \xi \geq 0, \quad (Ce - \alpha)^T \xi = 0,$$

Either $X^T \alpha = 0$ and $\|w\| \leq 1$ or

$$w = \frac{X^T \alpha}{\|X^T \alpha\|} = \frac{X_+ \alpha_+ - X_- \alpha_-}{\|X_+ \alpha_+ - X_- \alpha_-\|}$$

$$\rho_i = \frac{\alpha_i + 1}{2\sqrt{\alpha_i}}, \quad \sigma_i = \frac{\alpha_i - 1}{2\sqrt{\alpha_i}} \quad \text{all } i.$$

Geometric interpretation :

Again, assume all optimal $\alpha < C\epsilon$.

Write $\alpha = \gamma \hat{\alpha}$, $e_+^T \hat{\alpha}_+ = e_-^T \hat{\alpha}_- = 1$.

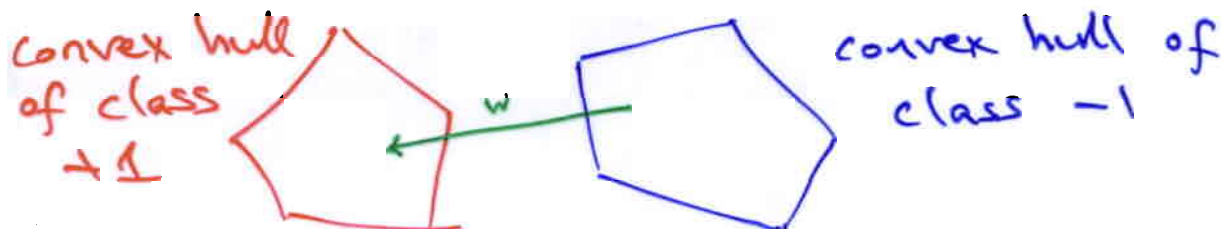
Then (D_{opt}) equiv to

$$\max_{\gamma, \hat{\alpha}} -\gamma \|XY\hat{\alpha}\| + 2\sqrt{\gamma} e^T \sqrt{\hat{\alpha}}, \quad \hat{\alpha} \dots$$

for fixed $\hat{\alpha}$, choose $\gamma = (e^T \sqrt{\hat{\alpha}} / \|XY\hat{\alpha}\|)^2$,

$$\Rightarrow \max \frac{(e_+^T \sqrt{\hat{\alpha}_+} + e_-^T \sqrt{\hat{\alpha}_-})^2}{\|X_+ \hat{\alpha}_+ - X_- \hat{\alpha}_-\|} \quad \begin{matrix} e_+^T \hat{\alpha}_+ = 1, \\ e_-^T \hat{\alpha}_- = 1, \\ \hat{\alpha}_+, \hat{\alpha}_- \geq 0. \end{matrix}$$

So we want to minimize the norm between convex combinations of the class +1 and the class -1 points, but scaled down by the square of the sum of the square roots of all convex weights (this forces all $\alpha_i > 0$).



Dimension reduction

(P_{DWD}) has $2n+1$ eq. constraints,
 $d+2+4n$ variables (S_{d+1}).

Bad in HDLSS case, $d \gg n$.

Core: $X = QR$, $Q \in \mathbb{R}^{d \times n}$, orthog cols
 $R \in \mathbb{R}^{n \times n}$, upper triang.

Replace $YX^T w$... $(w; w) \in S^{d+1}$
with $YR^T \bar{w}$... $(w; \bar{w}) \in S^{n+1}$.

"Correspondence" of feasible solutions via
 $\bar{w} = Q^T w$. Solve smaller problem to
get optimal \bar{w} , and $w = Q\bar{w}$ sol'n to
original problem.

A similar idea works to treat
the nonlinear case via a kernel function,
similar to with the SVM.

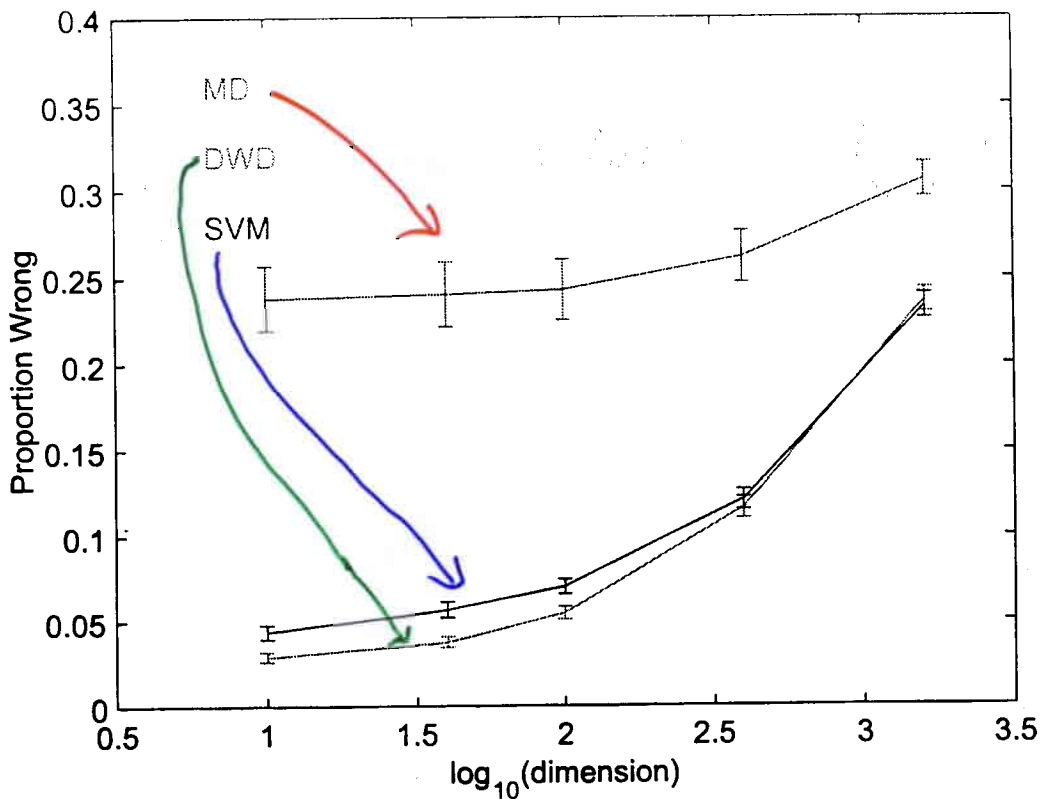


FIGURE 7: *Simulation comparison, for the "wobble" distribution. This is a case where DWD gives superior performance to MD and SVM.*

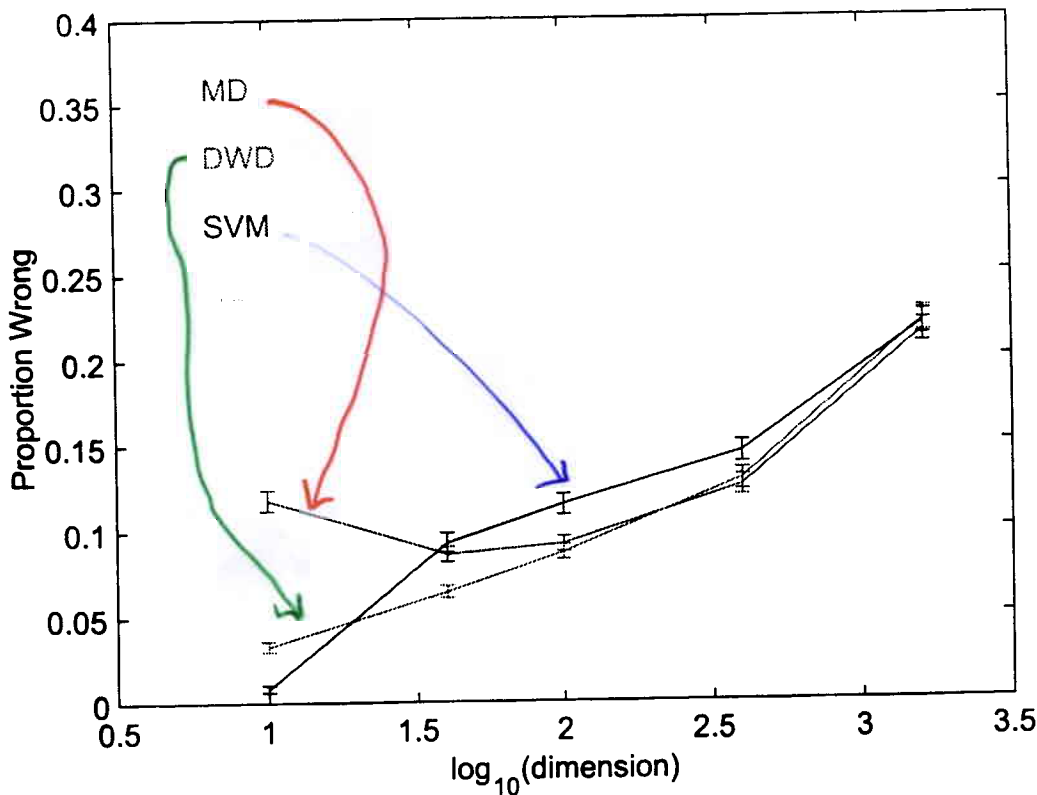


FIGURE 8: *Simulation comparison, for the "nested sphere" distribution. This case shows a fair overall summary, because each method is best for some d , and DWD tends to be near whichever method is best.*

- Group 1 Luminal cancer vs. other cancer types and normals: A first rough classification suggested by clustering of the data in Perou et al. ([12]). Tested using $n_+ = 47$ and $n_- = 38$ training cases, and 51 test cases.
- Group 2 Luminal A vs. Luminal B&C: an important distinction that was linked to survival rate in Perou et al. ([12]). Tested using $n_+ = 35$ and $n_- = 15$ training cases, and 21 test cases.
- Group 3 Normal vs. Erb & Basal cancer types. Tested using $n_+ = 13$ and $n_- = 25$ training cases, and 30 test cases.
- Group 4 Erb vs. Basal cancer types. Tested using $n_+ = 11$ and $n_- = 14$ training cases, and 21 test cases.

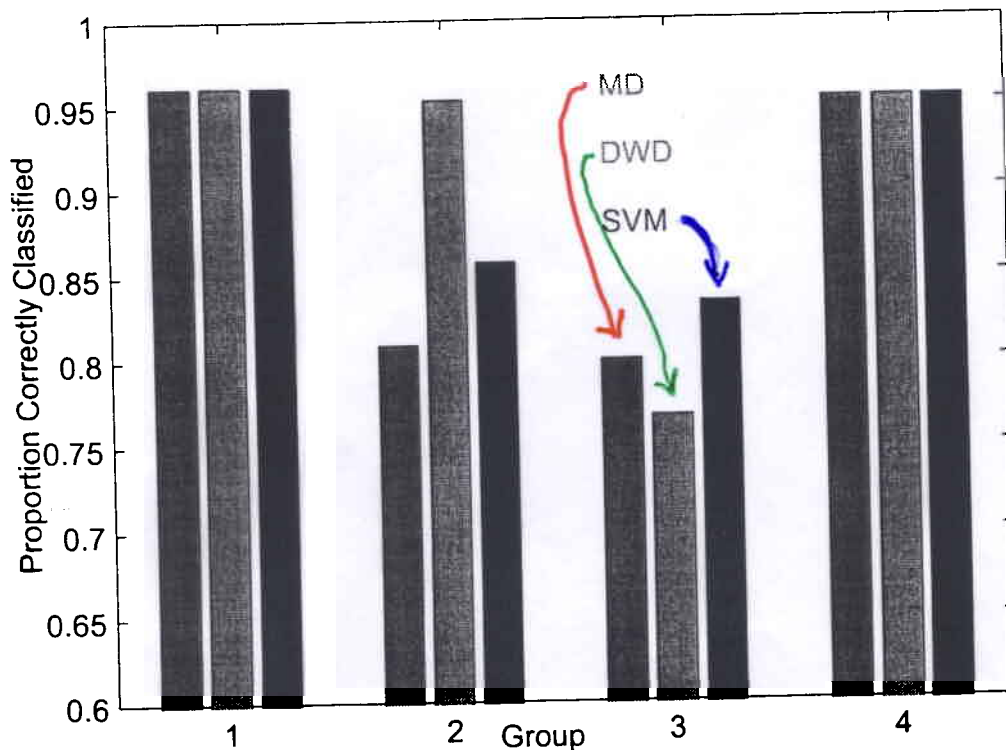


FIGURE 9: Graphical summary of correct classification rates for gene expression data.