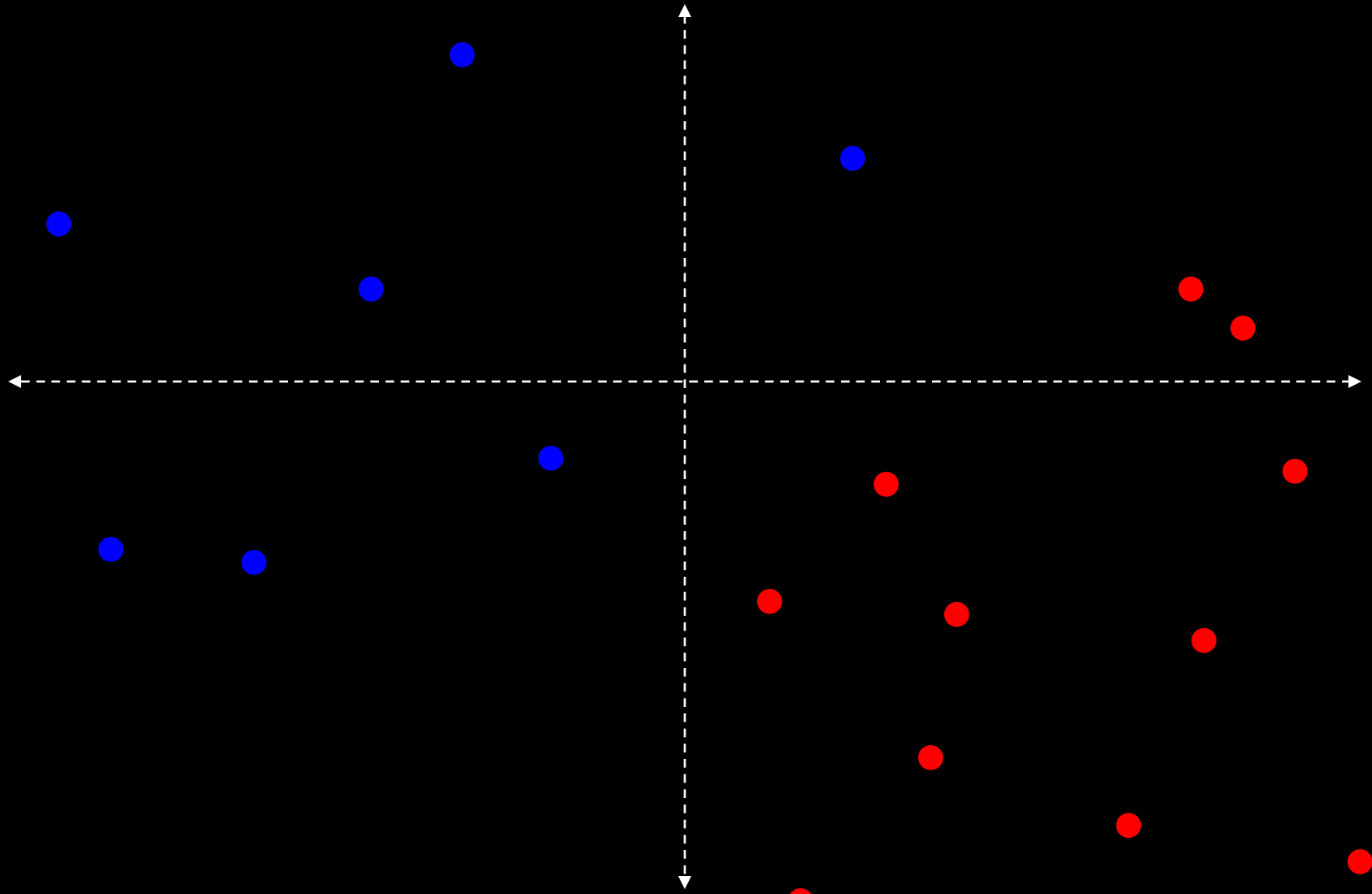
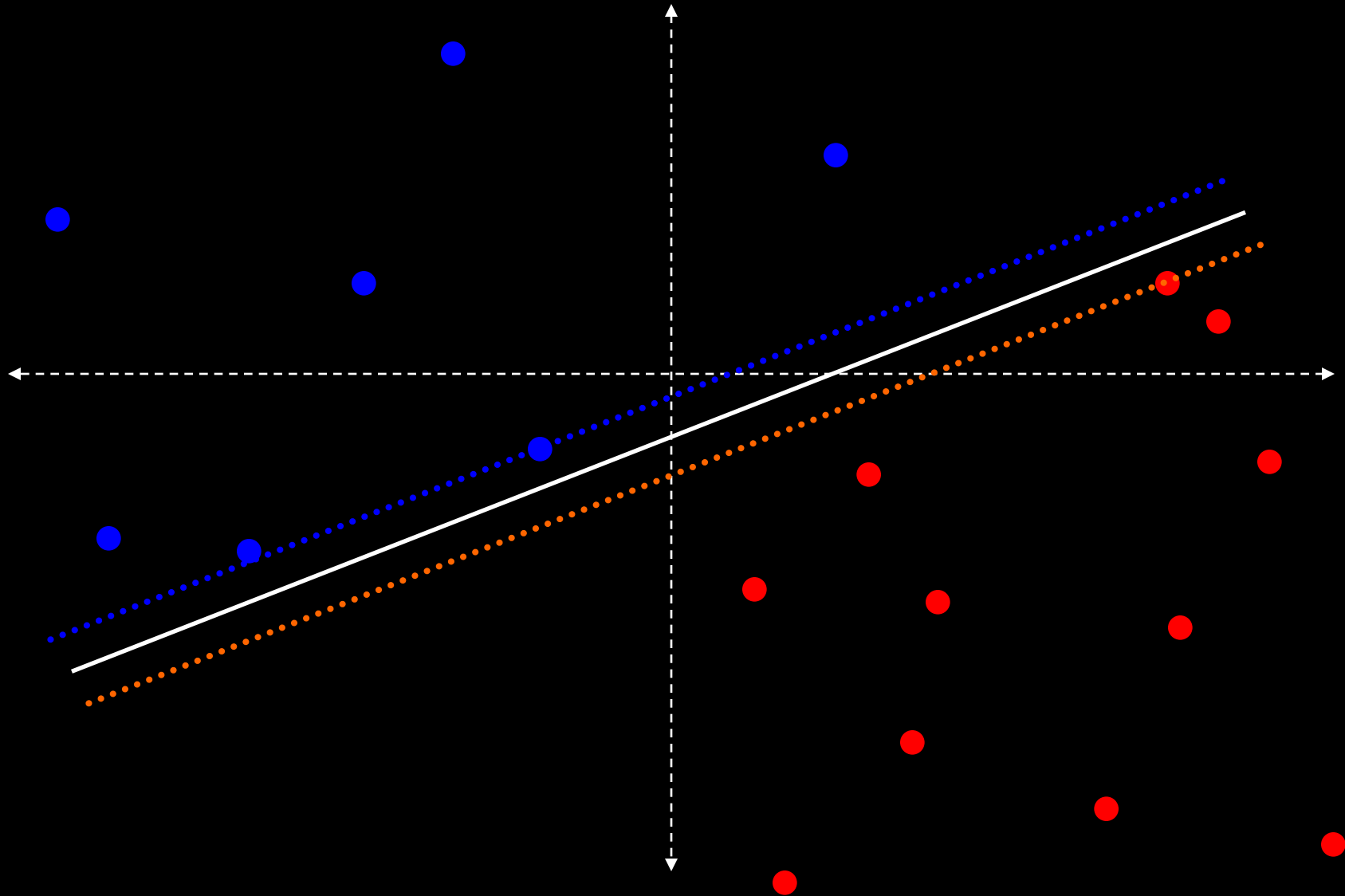


Support Vector Machines

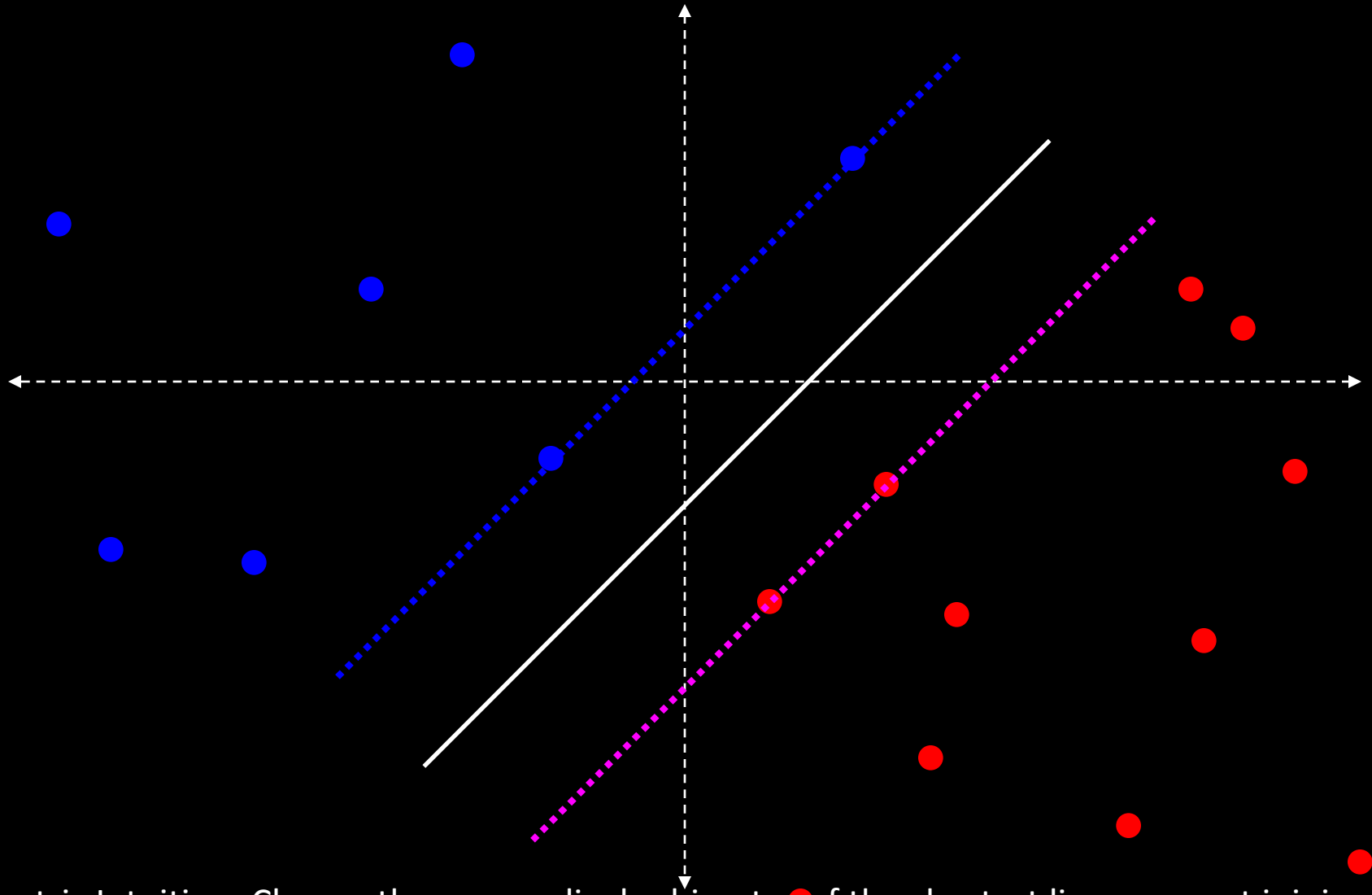
Binary Classification



A Separating Hyperplane



Maximum Margin Hyperplane



Geometric Intuition: Choose the perpendicular bisector of the shortest line segment joining the convex hulls of the two classes

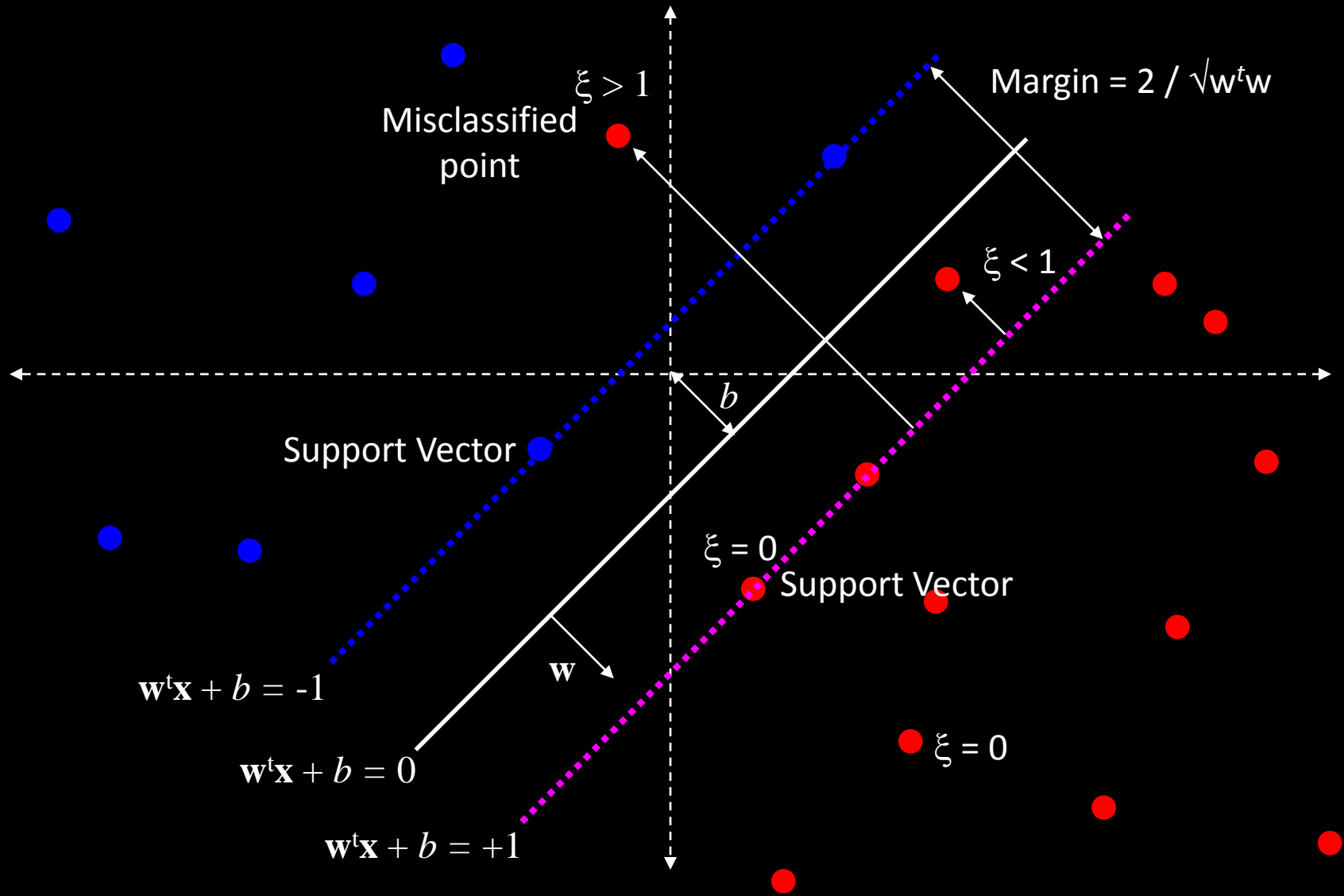
Hard Margin SVM Primal

- Maximize $2/|\mathbf{w}|$
such that $\mathbf{w}^t \mathbf{x}_i + b \geq +1$ if $y_i = +1$
 $\mathbf{w}^t \mathbf{x}_i + b \leq -1$ if $y_i = -1$
- Difficult to optimize directly
- Convex Quadratic Program (QP) reformulation
- Minimize $\frac{1}{2} \mathbf{w}^t \mathbf{w}$
such that $y_i (\mathbf{w}^t \mathbf{x}_i + b) \geq 1$
- Convex QPs can be easy to optimize

Linearly Inseparable Data

- Minimize $\frac{1}{2}\mathbf{w}^t\mathbf{w} + C \#(\text{Misclassified points})$
such that $y_i(\mathbf{w}^t\mathbf{x}_i + b) \geq 1$ (for “good” points)
- The optimization problem is NP Hard in general
- Disastrous errors are penalized the same as near misses

Inseparable Data – Hinge Loss



The C-SVM Primal Formulation

- Minimize $\frac{1}{2}\mathbf{w}^t\mathbf{w} + C \sum_i \xi_i$
such that $y_i(\mathbf{w}^t\mathbf{x}_i + b) \geq 1 - \xi_i$
 $\xi_i \geq 0$

- The optimization is a convex QP
- The globally optimal solution will be obtained
- Number of variables = $D + N + 1$
- Number of constraints = $2N$
- Solvers can train on 800K points in 47K (sparse) dimensions in less than 2 minutes on a standard PC

The C-SVM Dual Formulation

- Maximize $\mathbf{1}^t \boldsymbol{\alpha} - \frac{1}{2} \boldsymbol{\alpha}^t \mathbf{Y} \mathbf{K} \mathbf{Y} \boldsymbol{\alpha}$
such that $\mathbf{1}^t \mathbf{Y} \boldsymbol{\alpha} = 0$
 $\mathbf{0} \leq \boldsymbol{\alpha} \leq \mathbf{C}$
- \mathbf{K} is a kernel matrix such that $\mathbf{K}_{ij} = K(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^t \mathbf{x}_j$
- $\boldsymbol{\alpha}$ are the dual variables (Lagrange multipliers)
- Knowing $\boldsymbol{\alpha}$ gives us \mathbf{w} and b
- The dual is also a convex QP
 - Number of variables = N
 - Number of constraints = $2N + 1$

Duality

- Primal $P = \text{Min}_{\mathbf{x}} \quad f_0(\mathbf{x})$
s. t. $f_i(\mathbf{x}) \leq 0 \quad 1 \leq i \leq N$
 $h_i(\mathbf{x}) = 0 \quad 1 \leq i \leq M$
- Lagrangian $L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = f_0(\mathbf{x}) + \sum_i \lambda_i f_i(\mathbf{x}) + \sum_i \mu_i h_i(\mathbf{x})$
- Dual $D = \text{Max}_{\boldsymbol{\lambda}, \boldsymbol{\mu}} \quad \text{Min}_{\mathbf{x}} L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu})$
s. t. $\boldsymbol{\lambda} \geq \mathbf{0}$

Duality

- The Lagrange dual is always concave (even if the primal is not convex) and might be an easier problem to optimize
- Weak duality : $P \geq D$
 - Always holds
- Strong duality : $P = D$
 - Does not always hold
 - Usually holds for convex problems
 - Holds for the SVM QP

根据定义, 有

$$\min_x L(x, \lambda, \mu) \leq L(x^*, \lambda, \mu) \quad (5)$$

(其中 x^* 为 Primal 问题的最优解)

因为 $\lambda_i \geq 0, \mu_j \geq 0, f_i(x^*) \leq 0, h_j(x^*) = 0,$

我们有

$$\lambda_i f_i(x^*) \leq 0, \mu_j h_j(x^*) = 0$$

$$\begin{aligned} \text{因此, } L(x^*, \lambda, \mu) &= f_0(x^*) + \sum_i \lambda_i f_i(x^*) \\ &\quad + \sum_j \mu_j h_j(x^*) \\ &\leq f_0(x^*) \quad (6) \end{aligned}$$

由(5)和(6), 有

$$\min_x L(x, \lambda, \mu) \leq f_0(x^*)$$

(对 $\forall \lambda, \forall \mu$ 均成立)

因此

$$\begin{aligned} \text{Dual} &= \max_{\lambda, \mu} \min_x L(x, \lambda, \mu) \\ &\leq f_0(x^*) = \text{Primal} \end{aligned}$$

Karush-Kuhn-Tucker (KKT) Conditions

- If strong duality holds, then for \mathbf{x}^* , λ^* and μ^* to be optimal the following KKT conditions must necessarily hold
- Primal feasibility : $f_i(\mathbf{x}^*) \leq 0$ & $h_i(\mathbf{x}^*) = 0$ for $1 \leq i$
- Dual feasibility : $\lambda^* \geq 0$
- Stationarity : $\nabla_{\mathbf{x}} L(\mathbf{x}^*, \lambda^*, \mu^*) = 0$
- Complimentary slackness : $\lambda_i^* f_i(\mathbf{x}^*) = 0$
- If \mathbf{x}^+ , λ^+ and μ^+ satisfy the KKT conditions for a convex problem then they are optimal

SVM问题形式

$$\min_{w, b, \varepsilon_i} \frac{1}{2} w^T w + c \sum_{i=1}^N \varepsilon_i$$

$$\text{s.t. } y_i (w^T x_i + b) \geq 1 - \varepsilon_i \quad (1)$$

$$\varepsilon_i \geq 0$$

$$(i = 1, 2, \dots, N)$$

把(1)变成原始问题的标准形式:

$$\min_{w, b, \varepsilon_i} \frac{1}{2} w^T w + c \sum_{i=1}^N \varepsilon_i$$

$$\text{s.t. } 1 - \varepsilon_i - y_i (w^T x_i + b) \leq 0 \quad (i = 1, 2, \dots, N)$$

$$-\varepsilon_i \leq 0 \quad (2)$$

根据(2)得到的拉格朗日函数是:

$$L(w, b, \varepsilon_i, \alpha_i, \beta_i)$$

$$= \frac{1}{2} w^T w + c \sum_{i=1}^N \varepsilon_i \quad (3)$$

$$+ \sum_{i=1}^N \alpha_i [1 - \varepsilon_i - y_i (w^T x_i + b)]$$

$$+ \sum_{i=1}^N \beta_i (-\varepsilon_i)$$

$$= \frac{1}{2} w^T w + \sum_{i=1}^N (c - \alpha_i - \beta_i) \varepsilon_i + \sum_{i=1}^N \alpha_i$$

$$+ \sum_{i=1}^N (-\alpha_i y_i) (w^T x_i + b)$$

假设 w^* , b^* , ε_i^* 是问题(2)的最优解, 则根据 KKT 条件, 有

$$1 - \varepsilon_i^* - \gamma_i (w^{*T} x_i + b^*) \leq 0$$

$$-\varepsilon_i^* \leq 0$$

$$\alpha_i^* \geq 0$$

$$\beta_i^* \geq 0$$

(α_i^* , β_i^* 是对偶问题的最优解)

$$\frac{\partial L}{\partial w} = w - \sum_{i=1}^N \alpha_i \gamma_i x_i = 0$$

$$\frac{\partial L}{\partial b} = -\sum_{i=1}^N \gamma_i \alpha_i = 0$$

$$\frac{\partial L}{\partial \varepsilon_i} = c - \alpha_i - \beta_i = 0 \quad (i=1, 2, \dots, N)$$

$$\alpha_i [1 - \varepsilon_i - \gamma_i (w^T x_i + b)] = 0$$

$$\beta_i (-\varepsilon_i) = 0$$

(3) 可化简为

$$\begin{aligned} & \frac{1}{2} \left(\sum_{i=1}^N \alpha_i \gamma_i x_i^T \right) \left(\sum_{j=1}^N \alpha_j \gamma_j x_j \right) + 0 + \sum_{i=1}^N \alpha_i \\ & - \left(\sum_{i=1}^N \alpha_i \gamma_i x_i^T \right) \left(\sum_{j=1}^N \alpha_j \gamma_j x_j \right) + b \times 0 \quad (4) \\ & = \sum_{i=1}^N \alpha_i - \frac{1}{2} \left(\sum_{i=1}^N \alpha_i \gamma_i x_i^T \right) \left(\sum_{j=1}^N \alpha_j \gamma_j x_j \right) \end{aligned}$$

(4) 中只有一组变量 α_i , 通过 KKT 条件中的不等式有

$$\begin{cases} \alpha_i \geq 0 \\ C - \alpha_i = \beta_i \geq 0 \end{cases}$$

因此, 有

$$0 \leq \alpha_i \leq C \quad (i=1, 2, \dots, N)$$

最后得到对偶问题的形式

$$\max_{\alpha_i} \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \gamma_i \gamma_j \alpha_i \alpha_j$$

$$\text{s.t. } 0 \leq \alpha_i \leq C \quad (i=1, 2, \dots, N)$$

$$\sum_{i=1}^N \gamma_i \alpha_i = 0$$

SVM – Duality

- Primal P = $\text{Min}_{\mathbf{w}, \xi, b}$ $\frac{1}{2}\mathbf{w}^t\mathbf{w} + \mathbf{C}^t\xi$
s. t. $\mathbf{Y}(\mathbf{X}^t\mathbf{w} + b\mathbf{1}) \geq \mathbf{1} - \xi$
 $\xi \geq \mathbf{0}$
- Lagrangian $L(\alpha, \beta, \mathbf{w}, \xi, b)$ = $\frac{1}{2}\mathbf{w}^t\mathbf{w} + \mathbf{C}^t\xi - \beta^t\xi$
 $-\alpha^t[\mathbf{Y}(\mathbf{X}^t\mathbf{w} + b\mathbf{1}) - \mathbf{1} + \xi]$
- Dual D = Max_{α} $\mathbf{1}^t\alpha - \frac{1}{2}\alpha^t\mathbf{Y}\mathbf{K}\mathbf{Y}\alpha$
s. t. $\mathbf{1}^t\mathbf{Y}\alpha = 0$
 $\mathbf{0} \leq \alpha \leq \mathbf{C}$

SVM – KKT Conditions

- Lagrangian $L(\alpha, \beta, \mathbf{w}, \xi, b) = \frac{1}{2}\mathbf{w}^t\mathbf{w} + \mathbf{C}^t\xi - \beta^t\xi - \alpha^t[\mathbf{Y}(\mathbf{X}^t\mathbf{w} + b\mathbf{1}) - \mathbf{1} + \xi]$
- Stationarity conditions
 - $\nabla_{\mathbf{w}} L = 0 \Rightarrow \mathbf{w}^* = \mathbf{X}\mathbf{Y}\alpha^*$ (Representer Theorem)
 - $\nabla_{\xi} L = 0 \Rightarrow \mathbf{C} = \alpha^* + \beta^*$
 - $\nabla_b L = 0 \Rightarrow \alpha^{*t}\mathbf{Y}\mathbf{1} = 0$
- Complimentary Slackness conditions
 - $\alpha_i^* [y_i(\mathbf{x}_i^t\mathbf{w}^* + b^*) - 1 + \xi_i^*] = 0$
 - $\beta_i^* \xi_i^* = 0$

Support Vector Classification

- Training data $(\mathbf{x}_i, y_i), i = 1, \dots, l, \mathbf{x}_i \in R^n, y_i = \pm 1$

$$\min_{\mathbf{w}} \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^l \max(0, 1 - y_i \mathbf{w}^T \phi(\mathbf{x}_i))$$

- C : regularization parameter
- High dimensional (maybe infinite) feature space

$$\phi(\mathbf{x}) = [\phi_1(\mathbf{x}), \phi_2(\mathbf{x}), \dots]^T.$$

- We omit the bias term b
- \mathbf{w} : may have infinite variables



Support Vector Classification (Cont'd)

- The **dual** problem (**finite** # variables)

$$\begin{aligned} \min_{\alpha} f(\alpha) &= \frac{1}{2} \alpha^T Q \alpha - \mathbf{e}^T \alpha \\ \text{subject to} \quad & 0 \leq \alpha_i \leq C, i = 1, \dots, l, \end{aligned}$$

where $Q_{ij} = y_i y_j \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$ and $\mathbf{e} = [1, \dots, 1]^T$

- At optimum

$$\mathbf{w} = \sum_{i=1}^l \alpha_i y_i \phi(\mathbf{x}_i)$$

- Kernel: $K(\mathbf{x}_i, \mathbf{x}_j) \equiv \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$



Large Dense Quadratic Programming

- $Q_{ij} \neq 0$, Q : an l by l **fully dense** matrix

$$\begin{aligned} \min_{\alpha} f(\alpha) &= \frac{1}{2} \alpha^T Q \alpha - \mathbf{e}^T \alpha \\ \text{subject to} \quad & 0 \leq \alpha_i \leq C, i = 1, \dots, l \end{aligned}$$

- 50,000 training points: 50,000 variables:
($50,000^2 \times 8/2$) bytes = 10GB RAM to store Q
- Traditional methods:
Newton, Quasi Newton **cannot** be directly applied
- Now most use **decomposition methods**
[Osuna et al., 1997, Joachims, 1998, Platt, 1998]



Decomposition Methods

- We consider a **one-variable** version
Similar to **coordinate descent** methods
- Select the i th component for update:

$$\begin{aligned} \min_d \quad & \frac{1}{2}(\boldsymbol{\alpha} + d\mathbf{e}_i)^T Q(\boldsymbol{\alpha} + d\mathbf{e}_i) - \mathbf{e}^T(\boldsymbol{\alpha} + d\mathbf{e}_i) \\ \text{subject to} \quad & 0 \leq \alpha_i + d \leq C \end{aligned}$$

where

$$\mathbf{e}_i \equiv \underbrace{[0 \dots 0]_{i-1}}_{i-1} [1 \ 0 \dots 0]^T$$

- $\boldsymbol{\alpha}$: current solution; the i th component is changed



Avoid Memory Problems

- The new objective function

$$\frac{1}{2} Q_{ii} d^2 + (Q\alpha - \mathbf{e})_i d + \text{constant}$$

- To get $(Q\alpha - \mathbf{e})_i$, only Q 's *ith row* is needed

$$(Q\alpha - \mathbf{e})_i = \sum_{j=1}^l Q_{ij} \alpha_j - 1$$

- Calculated when needed. Trade time for space
- Used by popular software (e.g., *SVM^{light}*, LIBSVM)
They update 10 and 2 variables at a time



Decomposition Methods: Algorithm

- Optimal d :

$$-\frac{(Q\alpha - \mathbf{e})_i}{Q_{ii}} = -\frac{\sum_{j=1}^I Q_{ij}\alpha_j - 1}{Q_{ii}}$$

- Consider lower/upper bounds: $[0, C]$
- Algorithm:

While α is not optimal

1. Select the i th element for update

2. $\alpha_i \leftarrow \min \left(\max \left(\alpha_i - \frac{\sum_{j=1}^I Q_{ij}\alpha_j - 1}{Q_{ii}}, 0 \right), C \right)$



Select an Element for Update

Many ways

- Sequential (easiest)
- Permuting $1, \dots, l$ every l steps
- Random
- Existing software check gradient information

$$\nabla_1 f(\boldsymbol{\alpha}), \dots, \nabla_l f(\boldsymbol{\alpha})$$

But is $\nabla f(\boldsymbol{\alpha})$ available?



Select an Element for Update (Cont'd)

- We can easily maintain gradient

$$\nabla f(\boldsymbol{\alpha}) = Q\boldsymbol{\alpha} - \mathbf{e}$$

$$\nabla_s f(\boldsymbol{\alpha}) = (Q\boldsymbol{\alpha})_s - 1 = \sum_{j=1}^l Q_{sj}\alpha_j - 1$$

- Initial $\boldsymbol{\alpha} = \mathbf{0}$

$$\nabla f(\mathbf{0}) = -\mathbf{e}$$

- α_j updated to $\bar{\alpha}_j$

$$\nabla_s f(\boldsymbol{\alpha}) \leftarrow \nabla_s f(\boldsymbol{\alpha}) + Q_{sj}(\bar{\alpha}_j - \alpha_j), \quad \forall s$$

- $O(l)$ if $Q_{sj} \forall s$ (i th column) are available



Select an Element for Update (Cont'd)

- No matter maintaining $\nabla f(\alpha)$ or not
 Q 's i th row (column) always needed

$$\bar{\alpha}_i \leftarrow \min \left(\max \left(\alpha_i - \frac{\sum_{j=1}^l Q_{ij} \alpha_j - 1}{Q_{ii}}, 0 \right), C \right)$$

Q is **symmetric**

- Using $\nabla f(\alpha)$ to select i : faster convergence
i.e., fewer iterations



Decomposition Methods: Using Gradient

The new procedure

- $\alpha = \mathbf{0}, \nabla f(\alpha) = -\mathbf{e}$
- While α is not optimal
 1. Select the i th element using $\nabla f(\alpha)$
 2. $\bar{\alpha}_i \leftarrow \min \left(\max \left(\alpha_i - \frac{\sum_{j=1}^l Q_{ij} \alpha_j - 1}{Q_{ii}}, 0 \right), C \right)$
 3. $\nabla_s f(\alpha) \leftarrow \nabla_s f(\alpha) + Q_{si}(\bar{\alpha}_i - \alpha_i), \forall s$

Cost per iteration

- $O(ln)$, l : # instances, n : # features
- Assume each $Q_{ij} = y_i y_j K(\mathbf{x}_i, \mathbf{x}_j)$ takes $O(n)$



Linear SVM

- Primal without the bias term b

$$\min_{\mathbf{w}} \quad \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^l \max(0, 1 - y_i \mathbf{w}^T \mathbf{x}_i)$$

- Dual

$$\begin{aligned} \min_{\alpha} \quad f(\alpha) &= \frac{1}{2} \alpha^T Q \alpha - \mathbf{e}^T \alpha \\ \text{subject to} \quad & 0 \leq \alpha_i \leq C, \forall i \end{aligned}$$

- $Q_{ij} = y_i y_j \mathbf{x}_i^T \mathbf{x}_j$



Revisit Decomposition Methods

- While α is not optimal
 1. Select the i th element for update
 2. $\alpha_i \leftarrow \min \left(\max \left(\alpha_i - \frac{\sum_{j=1}^l Q_{ij} \alpha_j - 1}{Q_{ii}}, 0 \right), C \right)$
- $O(ln)$ per iteration; n : # features, l : # data
- For linear SVM, define

$$\mathbf{w} \equiv \sum_{j=1}^l y_j \alpha_j \mathbf{x}_j \in R^n$$

- $O(n)$ per iteration

$$\sum_{j=1}^l Q_{ij} \alpha_j - 1 = \sum_{j=1}^l y_i y_j \mathbf{x}_i^T \mathbf{x}_j \alpha_j - 1 = y_i \mathbf{w}^T \mathbf{x}_i - 1$$



- All we need is to maintain \mathbf{w} . If

$$\bar{\alpha}_i \leftarrow \alpha_i$$

then $O(n)$ for

$$\mathbf{w} \leftarrow \mathbf{w} + (\bar{\alpha}_i - \alpha_i)y_i\mathbf{x}_i$$

- Initial \mathbf{w}

$$\alpha = \mathbf{0} \quad \Rightarrow \quad \mathbf{w} = \mathbf{0}$$

- Give up maintaining $\nabla f(\alpha)$

- Select i for update

Sequential, random, or

Permuting $1, \dots, l$ every l steps



Algorithms for Linear and Nonlinear SVM

Linear:

- While α is not optimal
 1. Select the i th element for update
 2. $\bar{\alpha}_i \leftarrow \min \left(\max \left(\alpha_i - \frac{y_i \mathbf{w}^T \mathbf{x}_i - 1}{Q_{ii}}, 0 \right), C \right)$
 3. $\mathbf{w} \leftarrow \mathbf{w} + (\bar{\alpha}_i - \alpha_i) y_i \mathbf{x}_i$

Nonlinear:

- While α is not optimal
 1. Select the i th element using $\nabla f(\alpha)$
 2. $\bar{\alpha}_i \leftarrow \min \left(\max \left(\alpha_i - \frac{\sum_{j=1}^l Q_{ij} \alpha_j - 1}{Q_{ii}}, 0 \right), C \right)$
 3. $\nabla_s f(\alpha) \leftarrow \nabla_s f(\alpha) + Q_{si} (\bar{\alpha}_i - \alpha_i), \forall s$



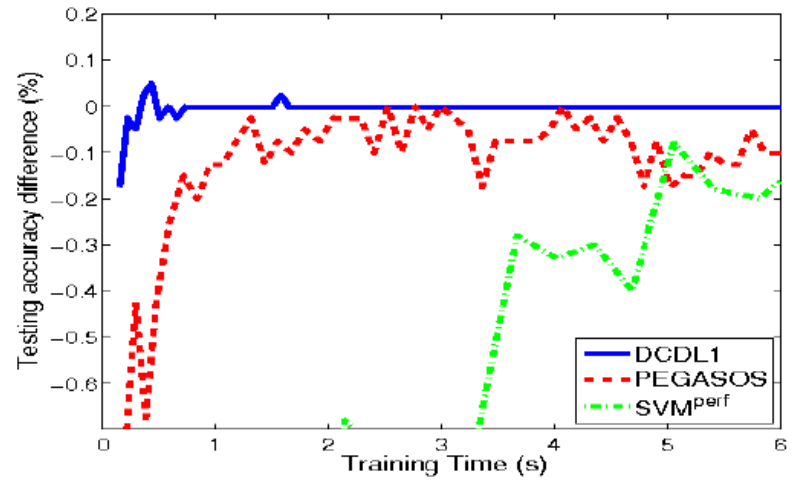
Analysis

- Decomposition method for nonlinear (also linear):
 $O(ln)$ per iteration (used in LIBSVM)
- New way for linear:
 $O(n)$ per iteration (used in LIBLINEAR)
- Faster if $\#$ iterations not l times more
- Experiments

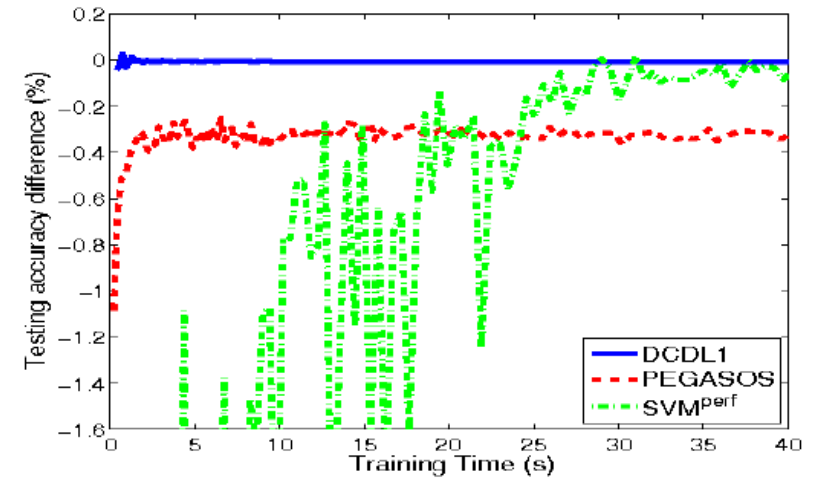
Problem	l : # data	n : # features
news20	19,996	1,355,191
yahoo-japan	176,203	832,026
rcv1	677,399	47,236
yahoo-korea	460,554	3,052,939



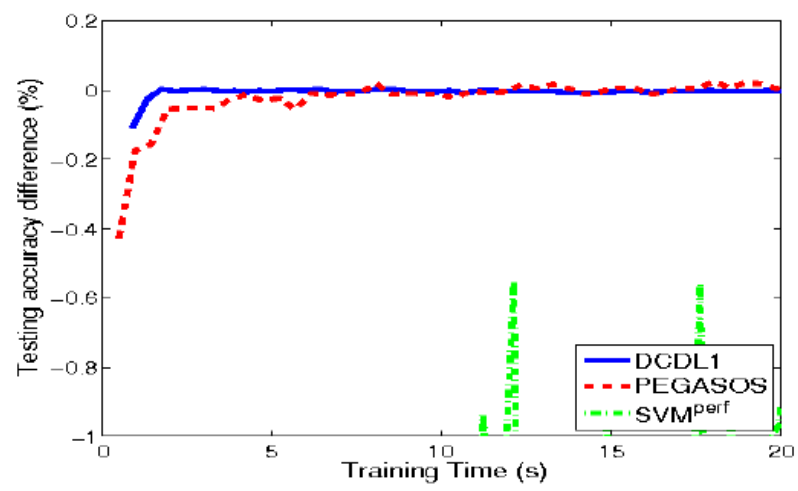
Testing Accuracy versus Training Time



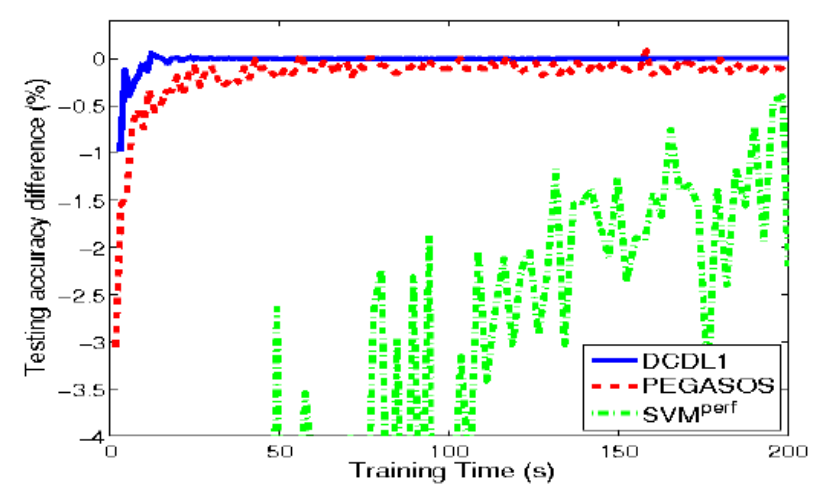
news20



yahoo-japan



rcv1



yahoo-korea

