

Statistical Machine Learning (BE4M33SSU)

Lecture 4: Support Vector Machines

Czech Technical University in Prague
V.Franc

Linear classifier with minimal classification error

- ◆ \mathcal{X} is a set of observations and $\mathcal{Y} = \{+1, -1\}$ a set of hidden labels
- ◆ $\phi: \mathcal{X} \rightarrow \mathbb{R}^n$ is fixed feature map embedding \mathcal{X} to \mathbb{R}^n
- ◆ **Task:** find linear classification strategy $h: \mathcal{X} \rightarrow \mathcal{Y}$

$$h(x; \mathbf{w}, b) = \text{sign}(\langle \mathbf{w}, \phi(x) \rangle + b) = \begin{cases} +1 & \text{if } \langle \mathbf{w}, \phi(x) \rangle + b \geq 0 \\ -1 & \text{if } \langle \mathbf{w}, \phi(x) \rangle + b < 0 \end{cases}$$

with minimal expected risk

$$R^{0/1}(h) = \mathbb{E}_{(x,y) \sim p} \left(\ell^{0/1}(y, h(x)) \right) \quad \text{where} \quad \ell^{0/1}(y, y') = [y \neq y']$$

- ◆ We are given a set of training examples

$$\mathcal{T}^m = \{(x^i, y^i) \in (\mathcal{X} \times \mathcal{Y}) \mid i = 1, \dots, m\}$$

drawn from i.i.d. with the distribution $p(x, y)$.

ERM learning for linear classifiers

- ◆ The Empirical Risk Minimization principle leads to solving

$$(\mathbf{w}^*, b^*) \in \underset{(\mathbf{w}, b) \in (\mathbb{R}^n \times \mathbb{R})}{\text{Argmin}} R_{\mathcal{T}^m}^{0/1}(h(\cdot; \mathbf{w}, b)) \quad (1)$$

where the empirical risk is

$$R_{\mathcal{T}^m}^{0/1}(h(\cdot; \mathbf{w}, b)) = \frac{1}{m} \sum_{i=1}^m [y^i \neq h(x^i; \mathbf{w}, b)]$$

In this lecture we address the following issues:

1. The statistical consistency of the ERM for hypothesis space containing linear classifiers.
2. Algorithmic issues: in general, there is no known algorithm solving the task (1) in time polynomial in m .

Training linear classifier from separable examples

Definition 1. The examples $\mathcal{T}^m = \{(x^i, y^i) \in (\mathcal{X} \times \mathcal{Y}) \mid i = 1, \dots, m\}$ are linearly separable w.r.t. feature map $\phi: \mathcal{X} \rightarrow \mathbb{R}^n$ if there exists $(\mathbf{w}, b) \in \mathbb{R}^{n+1}$ such that

$$y^i(\langle \mathbf{w}, \phi(x^i) \rangle + b) > 0, \quad i \in \{1, \dots, m\} \quad (2)$$

Perceptron algorithm:

Input: linearly separable examples \mathcal{T}^m

Output: linear classifier with $R_{\mathcal{T}^m}^{0/1}(h(\cdot; \mathbf{w}, b)) = 0$

step 1: $\mathbf{w} \leftarrow \mathbf{0}, b \leftarrow 0$

step 2: find (x^i, y^i) such that $y^i(\langle \mathbf{w}, \phi(x^i) \rangle + b) \leq 0$.

If not found exit, the current (\mathbf{w}, b) solves the problem.

step 3: $\mathbf{w} \leftarrow \mathbf{w} + y^i \phi(x^i), b \leftarrow b + y^i$ and goto to step 2.

Training linear classifier from NON-separable examples

- ◆ The intractable ERM problem we wish to solve

$$(\mathbf{w}^*, b^*) \in \underset{(\mathbf{w}, b) \in (\mathbb{R}^n \times \mathbb{R})}{\text{Argmin}} \frac{1}{m} \sum_{i=1}^m \underbrace{[y^i \neq h(x^i; \mathbf{w}, b)]}_{\ell^{0/1}(y^i, h(x^i; \mathbf{w}, b))}$$

where $h(x; \mathbf{w}, b) = \text{sign}(\langle \mathbf{w}, \phi(x) \rangle + b)$.

Training linear classifier from NON-separable examples

- ◆ The intractable ERM problem we wish to solve

$$(\mathbf{w}^*, b^*) \in \underset{(\mathbf{w}, b) \in (\mathbb{R}^n \times \mathbb{R})}{\text{Argmin}} \frac{1}{m} \sum_{i=1}^m \underbrace{[y^i \neq h(x^i; \mathbf{w}, b)]}_{\ell^{0/1}(y^i, h(x^i; \mathbf{w}, b))}$$

where $h(x; \mathbf{w}, b) = \text{sign}(\langle \mathbf{w}, \phi(x) \rangle + b)$.

- ◆ The ERM problem is approximated by a tractable **convex problem**

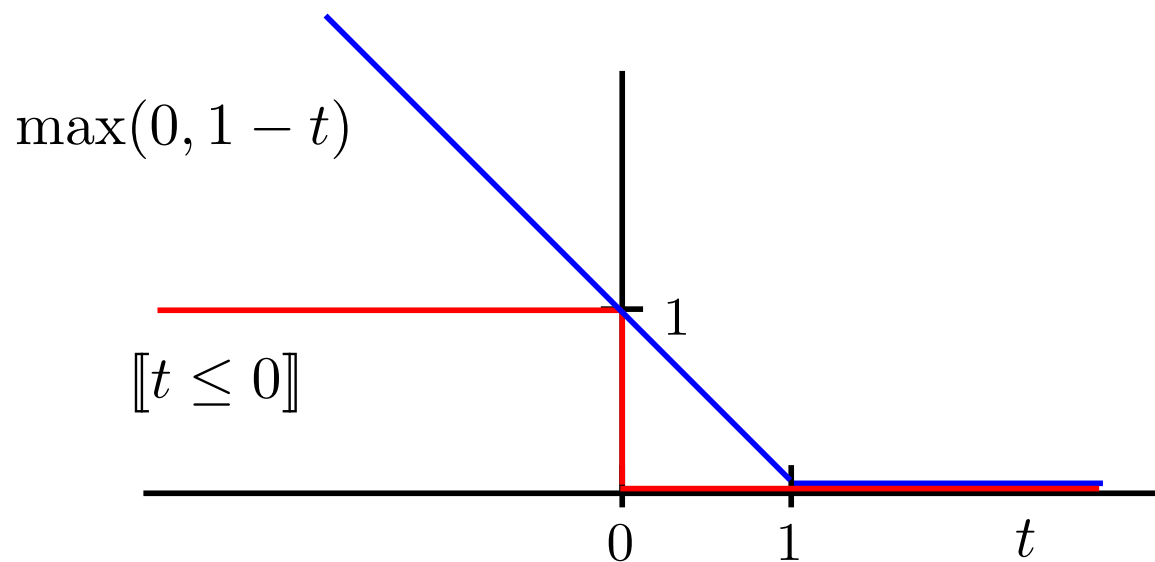
$$(\mathbf{w}^*, b^*) \in \underset{(\mathbf{w}, b) \in (\mathbb{R}^n \times \mathbb{R})}{\text{Argmin}} \frac{1}{m} \sum_{i=1}^m \underbrace{\max\{0, 1 - y^i f(x^i; \mathbf{w}, b)\}}_{\psi(y^i, f(x^i; \mathbf{w}, b))}$$

where $f(x; \mathbf{w}, b) = \langle \mathbf{w}, \phi(x) \rangle + b$ and $\psi(y, f(x))$ is so called Hinge-loss.

The hinge-loss upper bounds the 0/1-loss

- ◆ The hinge-loss is an upper bound of the 0/1-loss evaluated for the predictor $h(x) = \text{sign}(f(x))$:

$$\underbrace{[\text{sign}(f(x)) \neq y]}_{\ell^{0/1}(y, f(x))} = [y f(x) \leq 0] \leq \underbrace{\max\{0, 1 - y f(x)\}}_{\psi(y, f(x))}$$



Support Vector Machines

- ◆ Find linear classifier $h(x; \mathbf{w}, b) = \text{sign}(\langle \phi(x), \mathbf{w} \rangle + b)$ by solving

$$(\mathbf{w}^*, b^*) = \underset{\mathbf{w} \in \mathbb{R}^n, b \in \mathbb{R}}{\text{argmin}} \left(\underbrace{\frac{1}{2} \|\mathbf{w}\|^2}_{\text{penalty term}} + C \underbrace{\sum_{i=1}^m \max\{0, 1 - y^i (\langle \mathbf{w}, \phi(x^i) \rangle + b)\}}_{\text{empirical error}} \right)$$

- ◆ The regularization constant $C \geq 0$ helps to prevent overfitting (i.e. high estimation error) by constraining the parameter space.

- $C_1 < C_2$ implies $\|\mathbf{w}_1^*\| \leq \|\mathbf{w}_2^*\|$

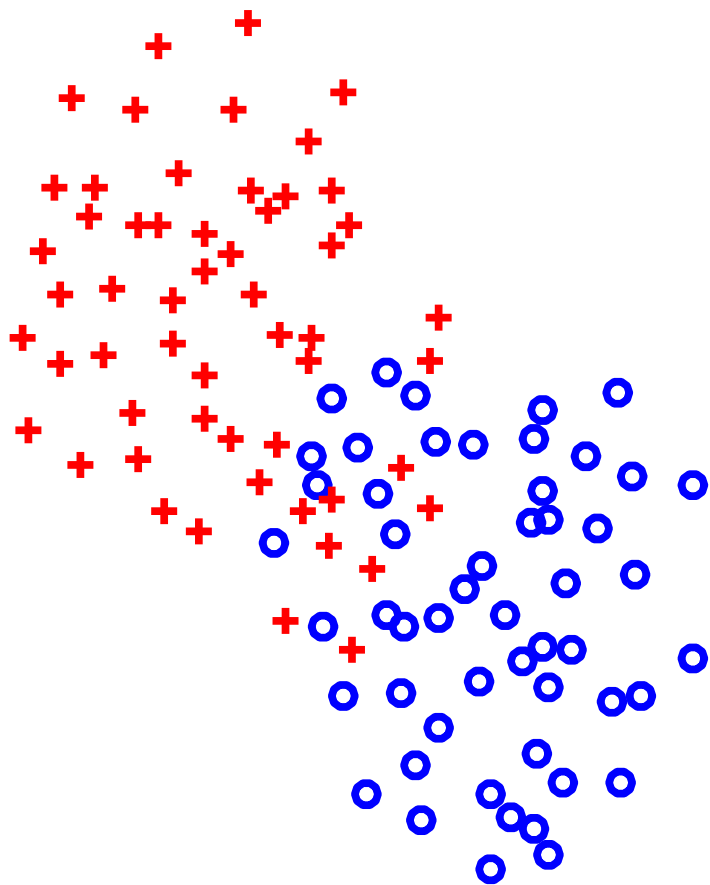
- ◆ Small $\|\mathbf{w}\|$ implies score $f(x; \mathbf{w}, b) = \langle \mathbf{w}, \phi(x) \rangle + b$ varies slowly.

- Cauchy inequality:

$$(\langle \phi(x), \mathbf{w} \rangle - \langle \phi(x'), \mathbf{w} \rangle)^2 \leq \|\phi(x) - \phi(x')\|^2 \|\mathbf{w}\|^2$$

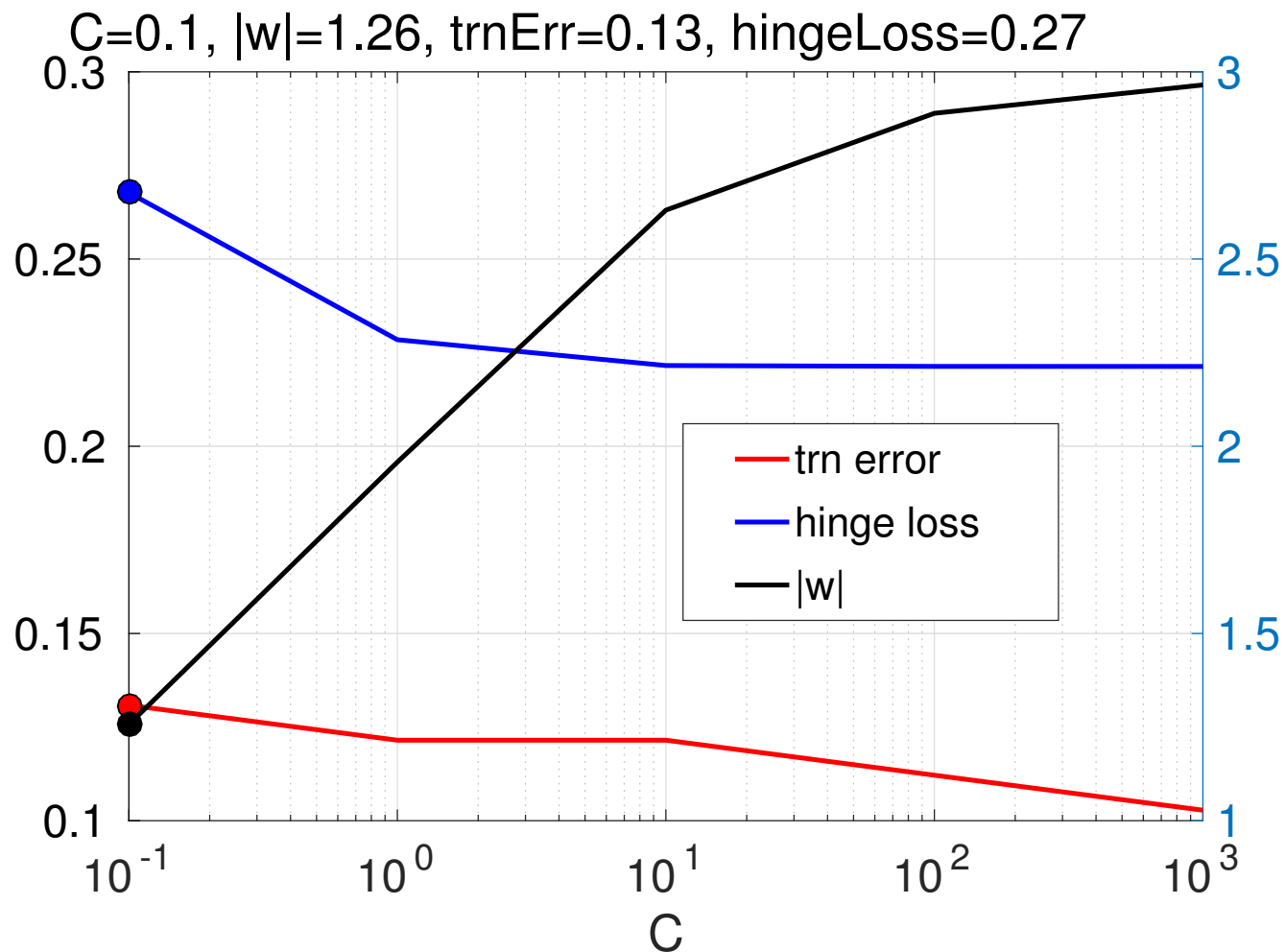
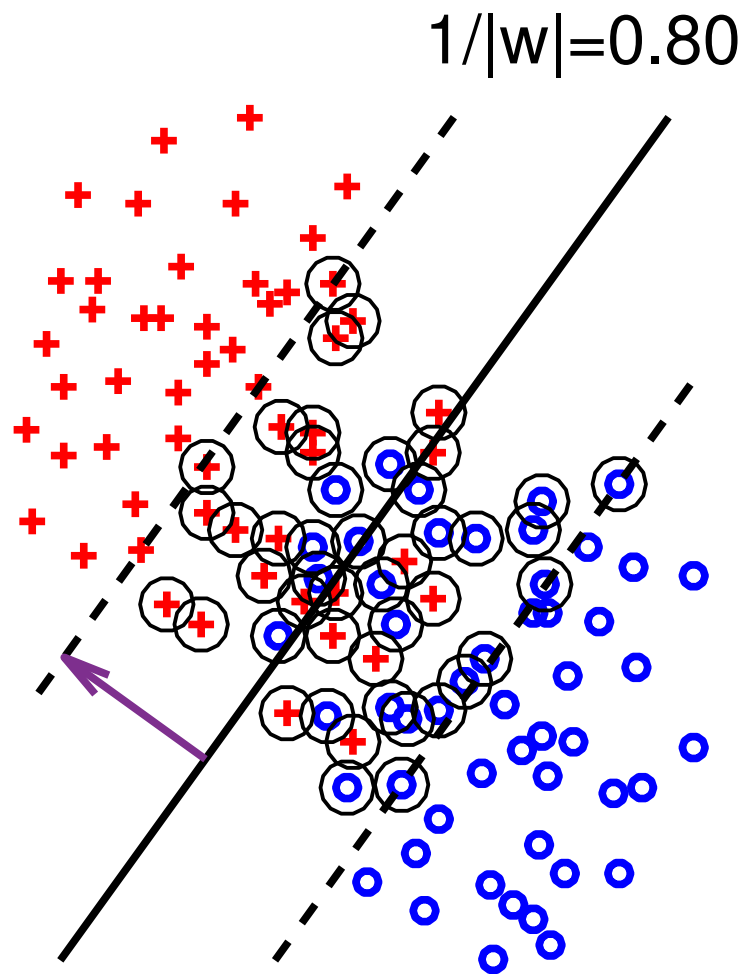
Example: Primal SVM problem

$$(\mathbf{w}^*, b^*) = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^n, b \in \mathbb{R}} \left(\underbrace{\frac{1}{2} \|\mathbf{w}\|^2}_{\text{penalty term}} + C \underbrace{\sum_{i=1}^m \max\{0, 1 - y^i (\langle \mathbf{w}, \phi(x^i) \rangle + b)\}}_{\text{empirical error}} \right)$$



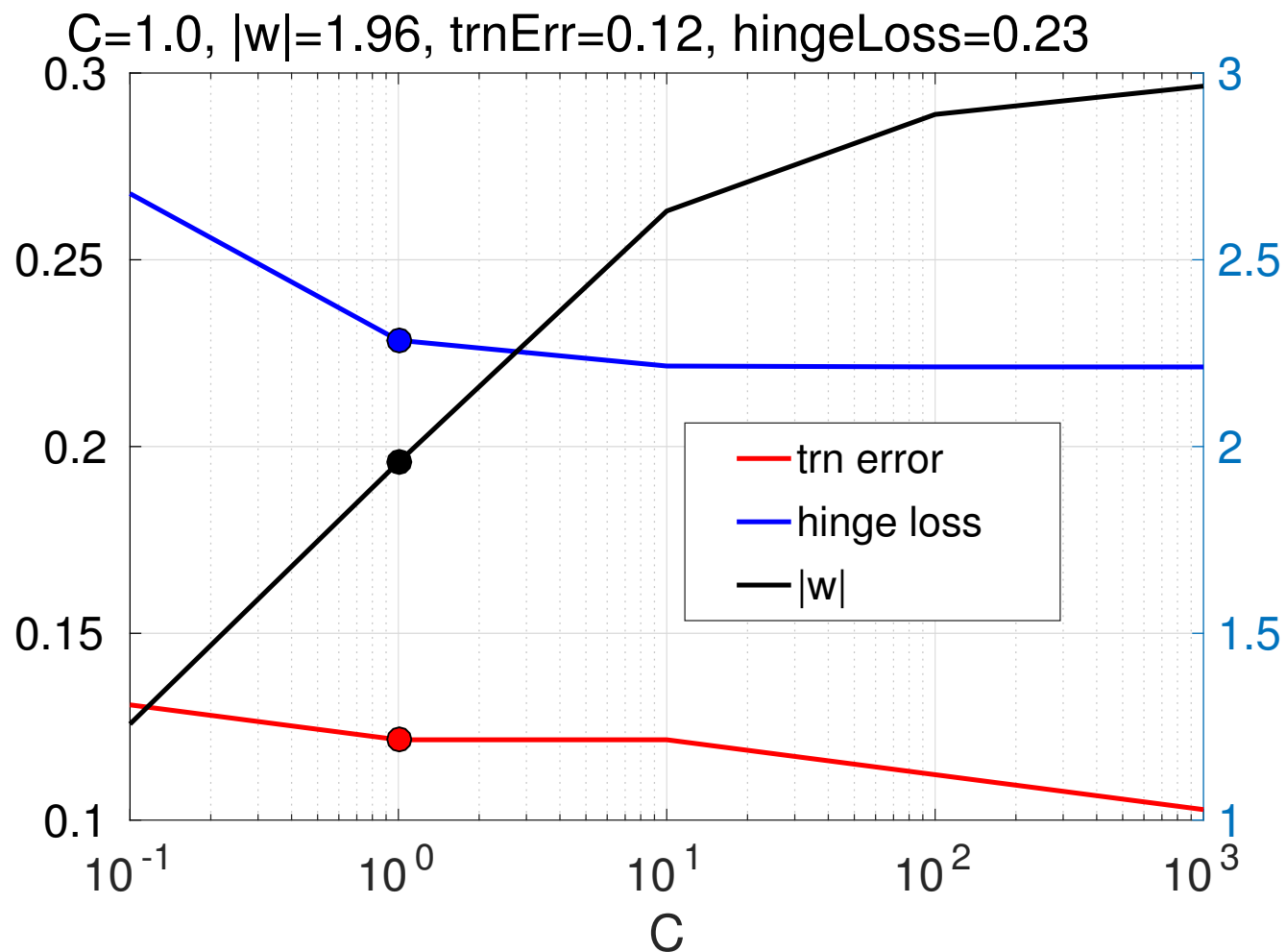
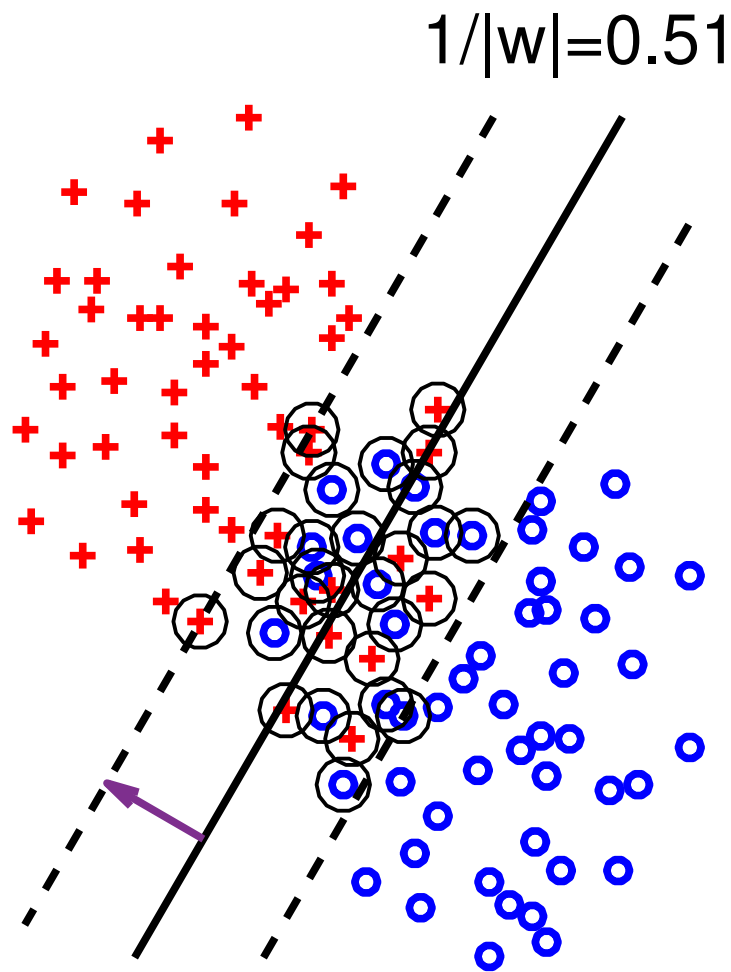
Example: Primal SVM problem

$$(\mathbf{w}^*, b^*) = \underset{\mathbf{w} \in \mathbb{R}^n, b \in \mathbb{R}}{\operatorname{argmin}} \left(\underbrace{\frac{1}{2} \|\mathbf{w}\|^2}_{\text{penalty term}} + C \underbrace{\sum_{i=1}^m \max\{0, 1 - y^i(\langle \mathbf{w}, \phi(x^i) \rangle + b)\}}_{\text{empirical error}} \right)$$



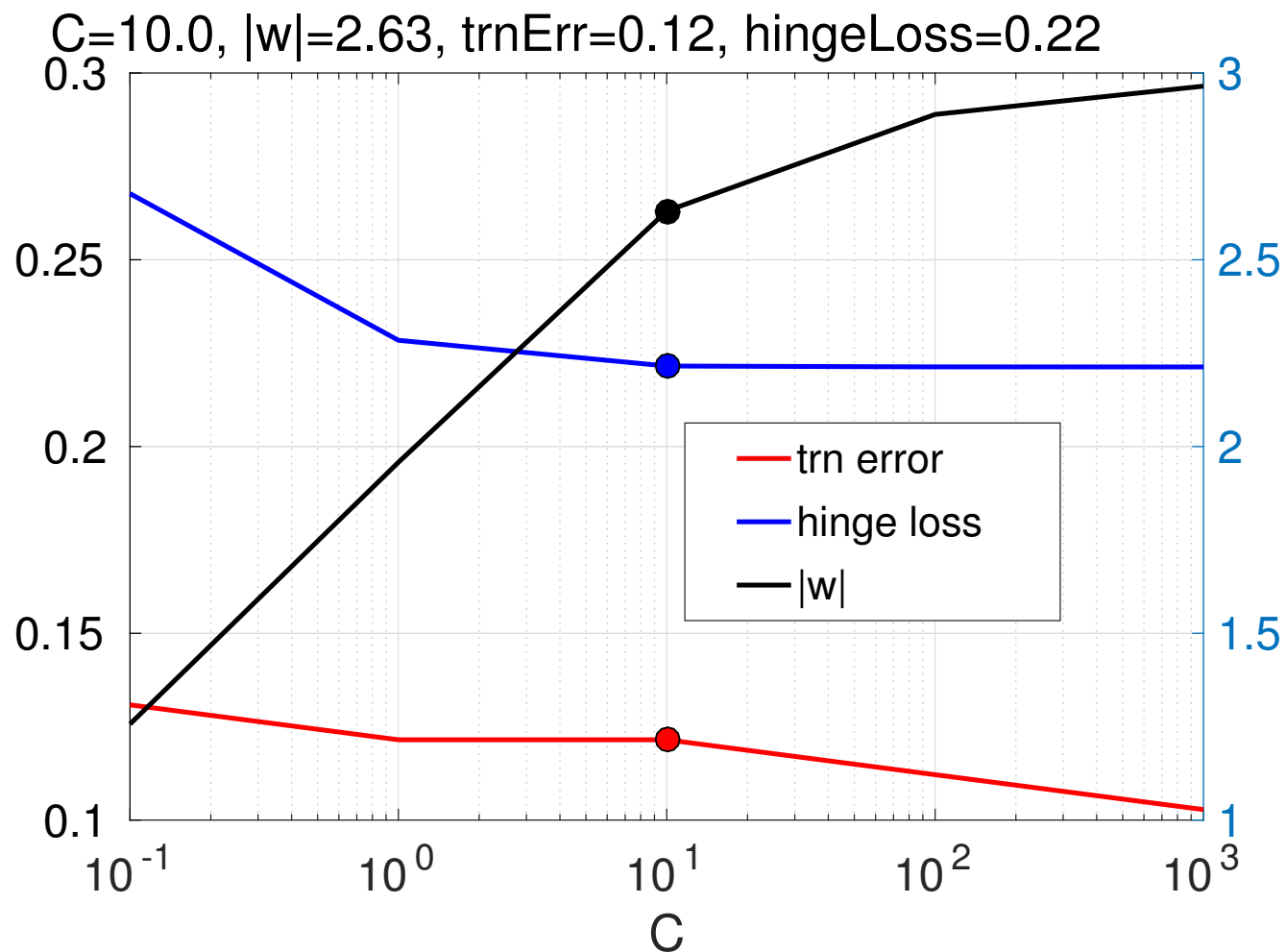
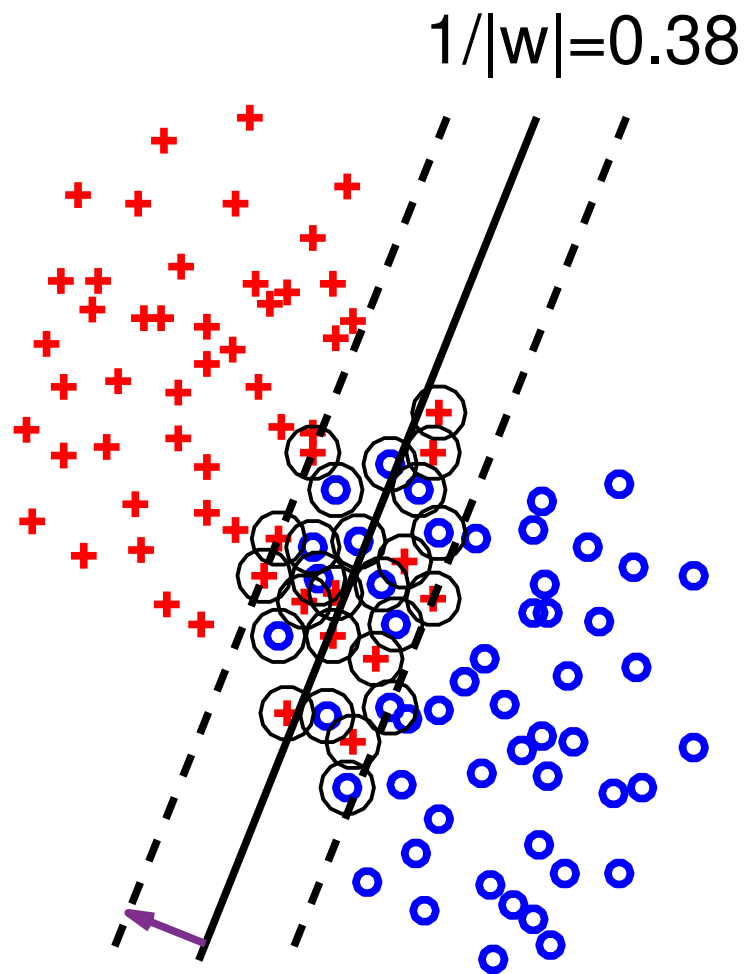
Example: Primal SVM problem

$$(\mathbf{w}^*, b^*) = \underset{\mathbf{w} \in \mathbb{R}^n, b \in \mathbb{R}}{\operatorname{argmin}} \left(\underbrace{\frac{1}{2} \|\mathbf{w}\|^2}_{\text{penalty term}} + C \underbrace{\sum_{i=1}^m \max\{0, 1 - y^i (\langle \mathbf{w}, \phi(x^i) \rangle + b)\}}_{\text{empirical error}} \right)$$



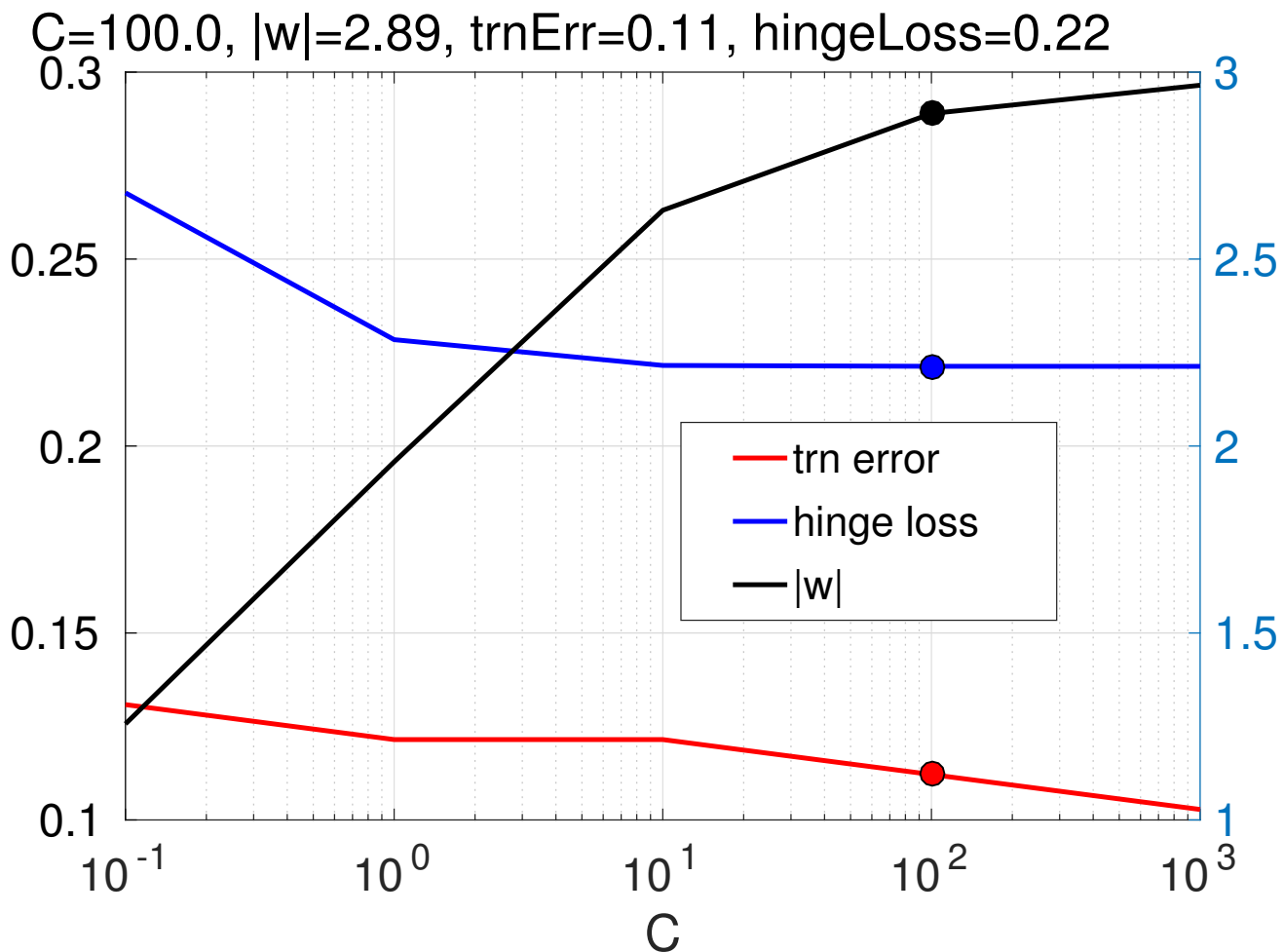
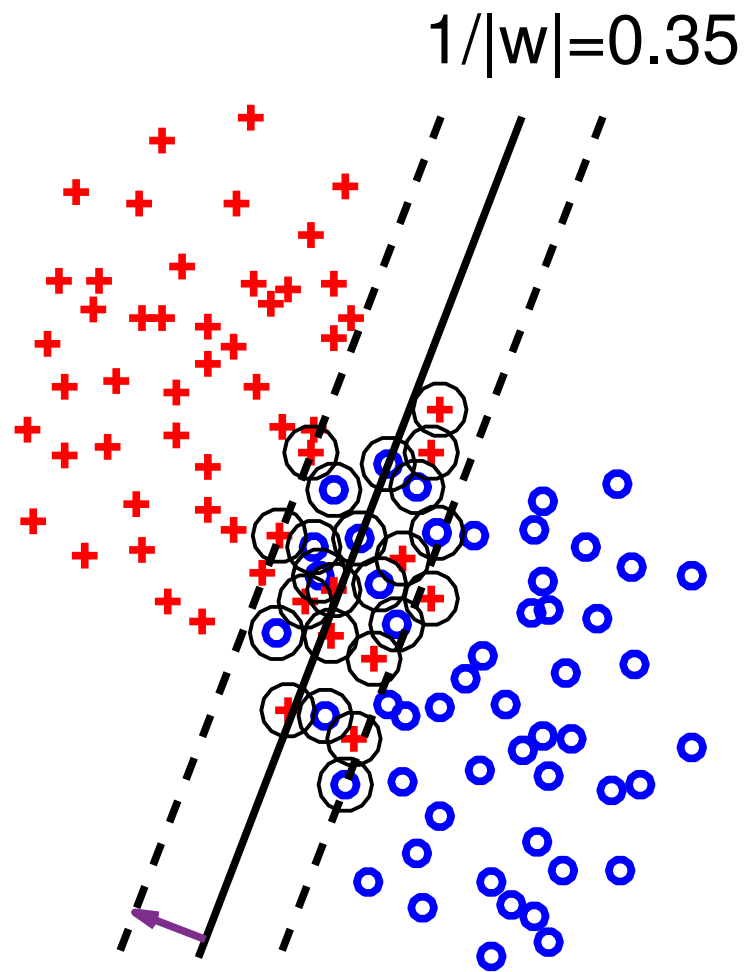
Example: Primal SVM problem

$$(\mathbf{w}^*, b^*) = \underset{\mathbf{w} \in \mathbb{R}^n, b \in \mathbb{R}}{\operatorname{argmin}} \left(\underbrace{\frac{1}{2} \|\mathbf{w}\|^2}_{\text{penalty term}} + C \underbrace{\sum_{i=1}^m \max\{0, 1 - y^i(\langle \mathbf{w}, \phi(x^i) \rangle + b)\}}_{\text{empirical error}} \right)$$



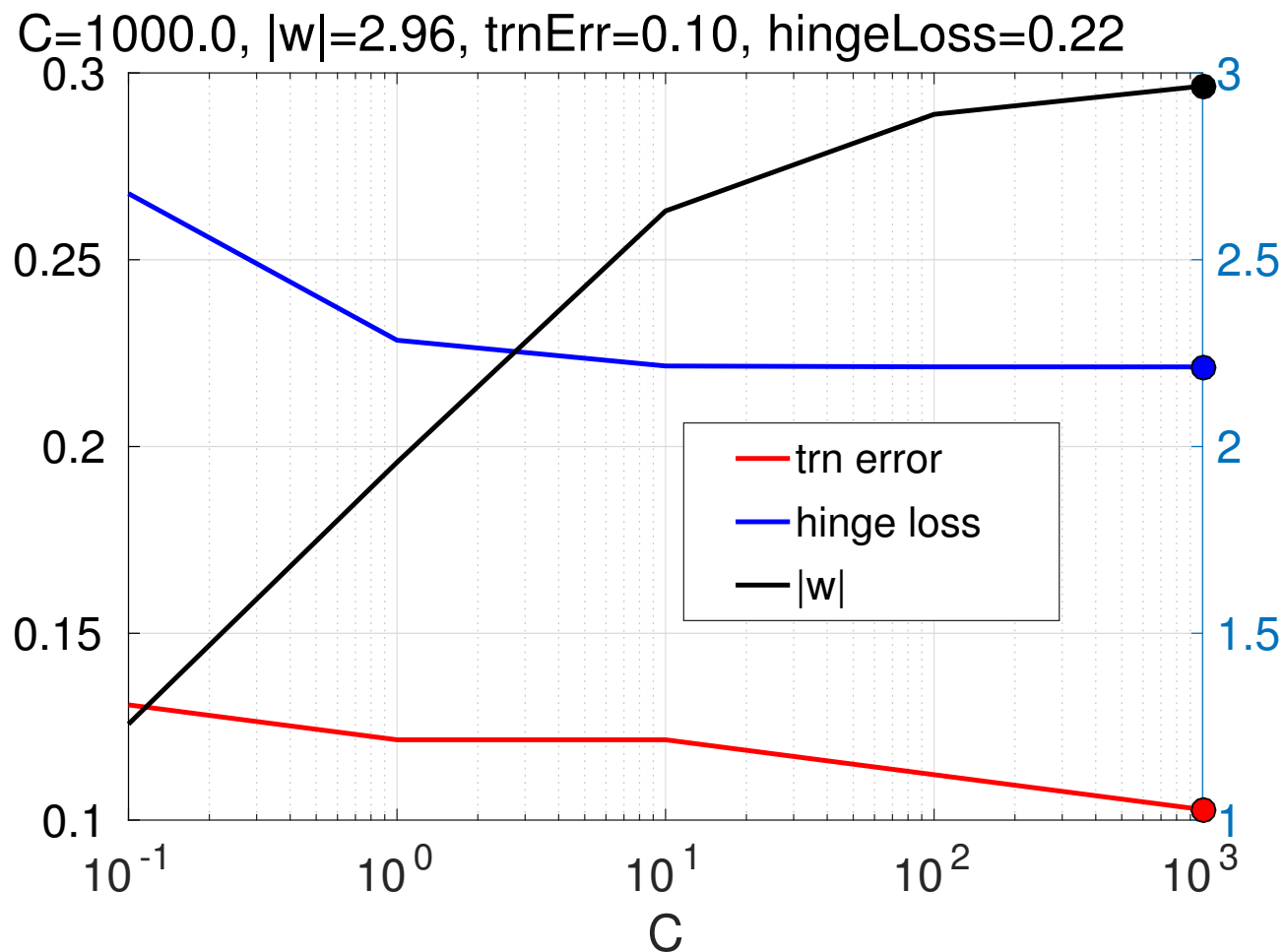
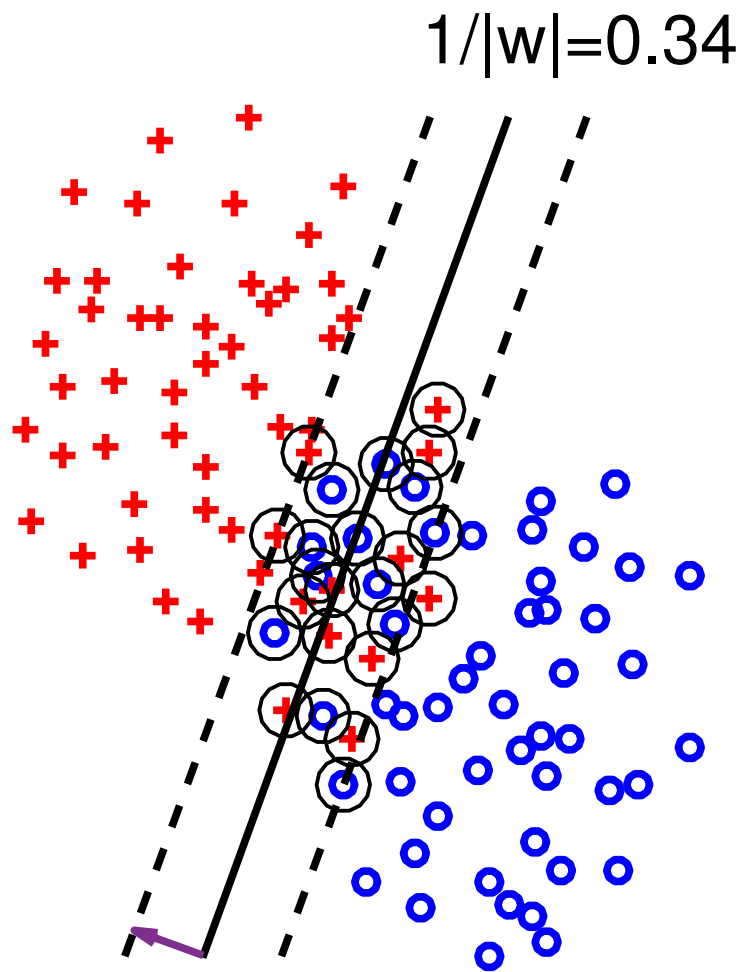
Example: Primal SVM problem

$$(\mathbf{w}^*, b^*) = \underset{\mathbf{w} \in \mathbb{R}^n, b \in \mathbb{R}}{\operatorname{argmin}} \left(\underbrace{\frac{1}{2} \|\mathbf{w}\|^2}_{\text{penalty term}} + C \underbrace{\sum_{i=1}^m \max\{0, 1 - y^i(\langle \mathbf{w}, \phi(x^i) \rangle + b)\}}_{\text{empirical error}} \right)$$



Example: Primal SVM problem

$$(\mathbf{w}^*, b^*) = \underset{\mathbf{w} \in \mathbb{R}^n, b \in \mathbb{R}}{\operatorname{argmin}} \left(\underbrace{\frac{1}{2} \|\mathbf{w}\|^2}_{\text{penalty term}} + C \underbrace{\sum_{i=1}^m \max\{0, 1 - y^i(\langle \mathbf{w}, \phi(x^i) \rangle + b)\}}_{\text{empirical error}} \right)$$



SVM as Quadratic Program

- ◆ Find linear classifier $h(x; \mathbf{w}, b) = \text{sign}(\langle \phi(x), \mathbf{w} \rangle + b)$ by solving

$$(\mathbf{w}^*, b^*) = \underset{\mathbf{w} \in \mathbb{R}^n, b \in \mathbb{R}}{\text{argmin}} \left(\underbrace{\frac{1}{2} \|\mathbf{w}\|^2}_{\text{penalty term}} + C \underbrace{\sum_{i=1}^m \max\{0, 1 - y^i (\langle \mathbf{w}, \phi(x^i) \rangle + b)\}}_{\text{empirical error}} \right)$$

where $C > 0$ is the regularization constant.

SVM as Quadratic Program

- Find linear classifier $h(x; \mathbf{w}, b) = \text{sign}(\langle \phi(x), \mathbf{w} \rangle + b)$ by solving

$$(\mathbf{w}^*, b^*) = \underset{\mathbf{w} \in \mathbb{R}^n, b \in \mathbb{R}}{\text{argmin}} \left(\underbrace{\frac{1}{2} \|\mathbf{w}\|^2}_{\text{penalty term}} + C \underbrace{\sum_{i=1}^m \max\{0, 1 - y^i(\langle \mathbf{w}, \phi(x^i) \rangle + b)\}}_{\text{empirical error}} \right)$$

where $C > 0$ is the regularization constant.

- It can be re-formulated as a convex *quadratic program*

$$(\mathbf{w}^*, b^*, \boldsymbol{\xi}^*) = \underset{\substack{(\mathbf{w}, b) \in \mathbb{R}^{n+1} \\ \boldsymbol{\xi} \in \mathbb{R}^m}}{\text{argmin}} \left(\frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m \xi_i \right)$$

subject to

$$\begin{aligned} y^i(\langle \mathbf{w}, \phi(x^i) \rangle + b) &\geq 1 - \xi_i, & i \in \{1, \dots, m\} \\ \xi_i &\geq 0, & i \in \{1, \dots, m\} \end{aligned}$$

From Primal SVM to Dual SVM problem

- ◆ Lagrangian of the primal SVM problem:

$$\begin{aligned}
 L(\mathbf{w}, b, \boldsymbol{\xi}, \boldsymbol{\alpha}, \boldsymbol{\mu}) = & \underbrace{\frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m \xi_i}_{\text{original objective}} \\
 & - \underbrace{\sum_{i=1}^m \alpha_i (y^i (\langle \mathbf{w}, \phi(x^i) \rangle + b) - 1 + \xi_i) - \sum_{i=1}^m \mu_i \xi_i}_{\text{constraint violation penalty}}
 \end{aligned}$$

- ◆ Strong duality:

$$\underbrace{\min_{\substack{\mathbf{w} \in \mathbb{R}^n \\ b \in \mathbb{R} \\ \boldsymbol{\xi} \in \mathbb{R}^m}} \max_{\substack{\boldsymbol{\alpha} \in \mathbb{R}_+^m \\ \boldsymbol{\mu} \in \mathbb{R}_+^m}} L(\mathbf{w}, b, \boldsymbol{\xi}, \boldsymbol{\alpha}, \boldsymbol{\mu})}_{\text{primal problem}} = \underbrace{\max_{\substack{\boldsymbol{\alpha} \in \mathbb{R}_+^m \\ \boldsymbol{\mu} \in \mathbb{R}_+^m}} \min_{\substack{\mathbf{w} \in \mathbb{R}^n \\ b \in \mathbb{R} \\ \boldsymbol{\xi} \in \mathbb{R}^m}} L(\mathbf{w}, b, \boldsymbol{\xi}, \boldsymbol{\alpha}, \boldsymbol{\mu})}_{\text{dual problem}}$$

Dual SVM problem

- ◆ The dual SVM formulation is a convex quadratic program

$$\begin{aligned}
 \boldsymbol{\alpha}^* &= \operatorname{argmax}_{\boldsymbol{\alpha} \in \mathbb{R}^m} \left(\sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y^i y^j \langle \boldsymbol{\phi}(x^i), \boldsymbol{\phi}(x^j) \rangle \right) \\
 \text{s.t.} \quad & \sum_{i=1}^m \alpha_i y^i = 0, \quad 0 \leq \alpha_i \leq C, \quad i \in \{1, \dots, m\}
 \end{aligned}$$

Dual SVM problem

- ◆ The dual SVM formulation is a convex quadratic program

$$\begin{aligned}
 \boldsymbol{\alpha}^* = \operatorname{argmax}_{\boldsymbol{\alpha} \in \mathbb{R}^m} & \left(\sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y^i y^j \langle \boldsymbol{\phi}(x^i), \boldsymbol{\phi}(x^j) \rangle \right) \\
 \text{s.t.} & \quad \sum_{i=1}^m \alpha_i y^i = 0, \quad 0 \leq \alpha_i \leq C, \quad i \in \{1, \dots, m\}
 \end{aligned}$$

- ◆ The primal variables (\boldsymbol{w}, b) are obtained from the dual variables $\boldsymbol{\alpha}$ by

$$\boldsymbol{w} = \sum_{i=1}^m y^i \boldsymbol{\phi}(x^i) \alpha_i = \sum_{i \in \mathcal{I}_{\text{SV}}} y^i \boldsymbol{\phi}(x^i) \alpha_i$$

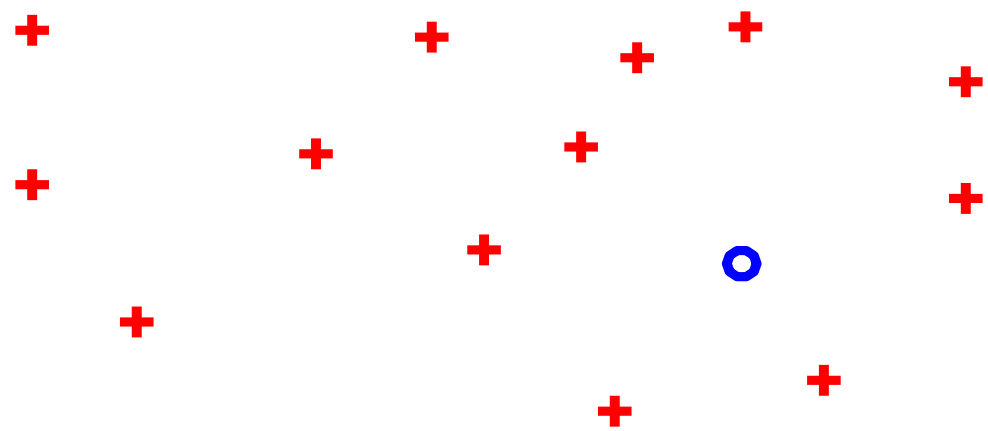
$$b = y^i - \langle \boldsymbol{w}, \boldsymbol{\phi}(x^i) \rangle, \quad \forall i \in \mathcal{I}_{\text{SV}}^b = \{j \mid 0 < \alpha_j < C\}$$

- ◆ $\boldsymbol{\alpha}$ is sparse; \boldsymbol{w} is lin. combination of Support Vectors $\mathcal{I}_{\text{SV}} = \{j \mid \alpha_j > 0\}$

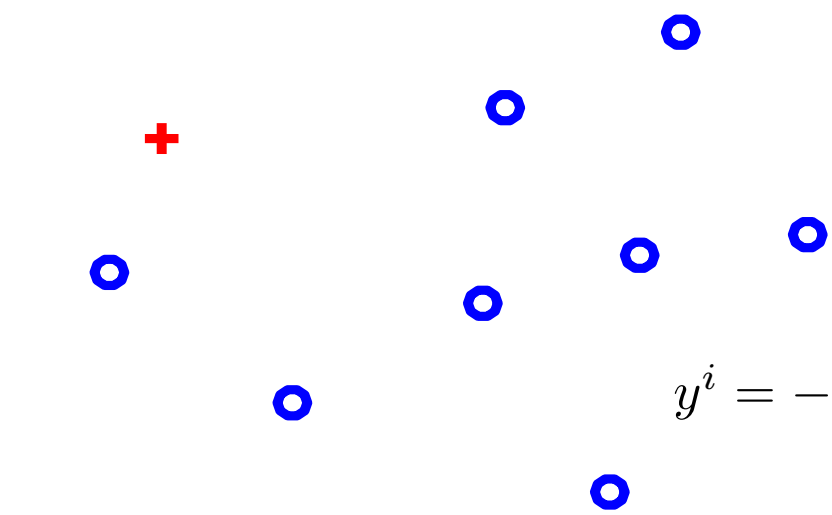
Example: SVM classifier

$$f(x) = \langle \mathbf{w}, \phi(x) \rangle + b = \langle \underbrace{\sum_{i=1}^m y^i \alpha_i \phi(x^i)}_{\mathbf{w}}, \phi(x) \rangle + b$$

$y^i = +1$

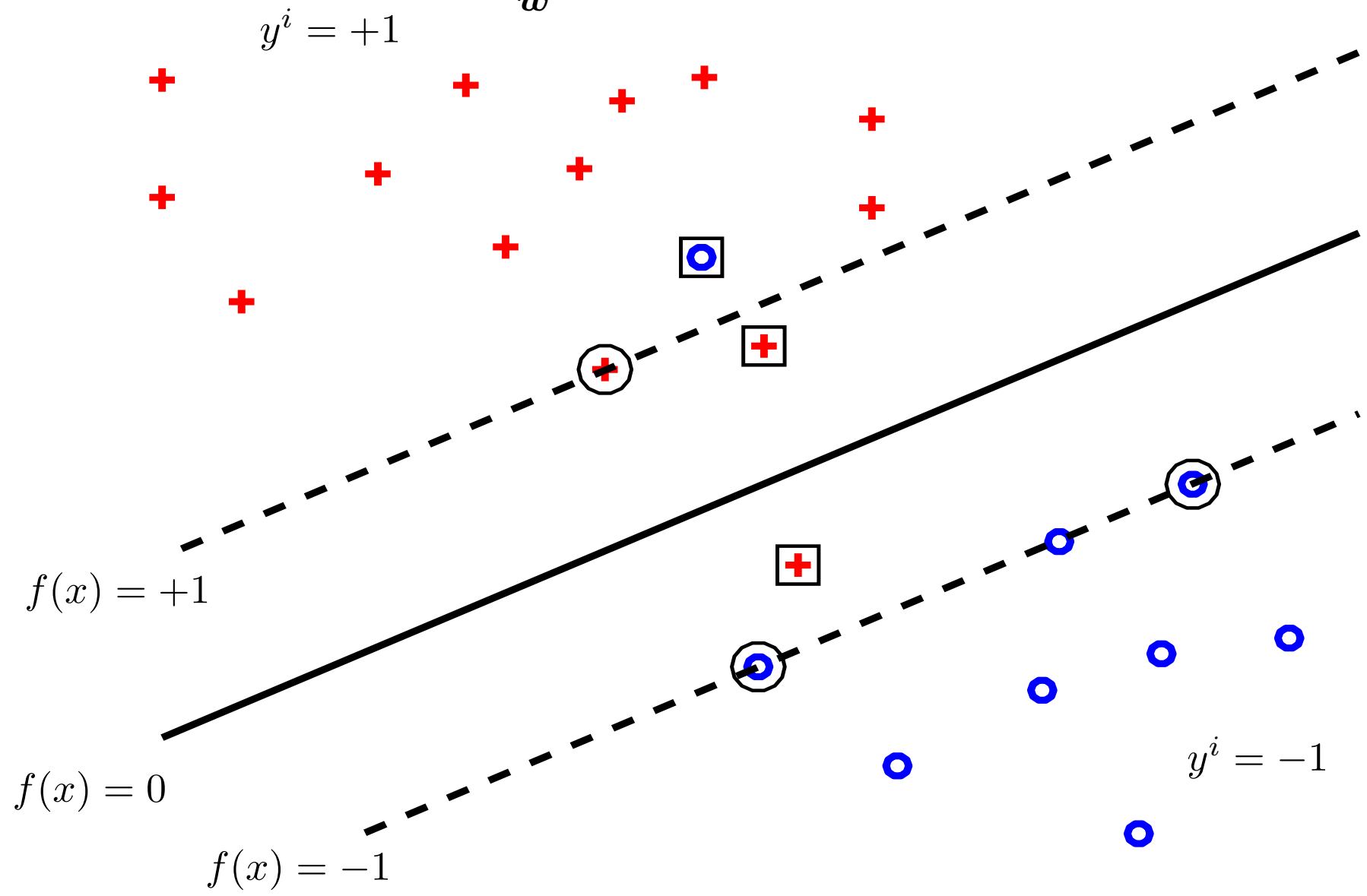


$y^i = -1$



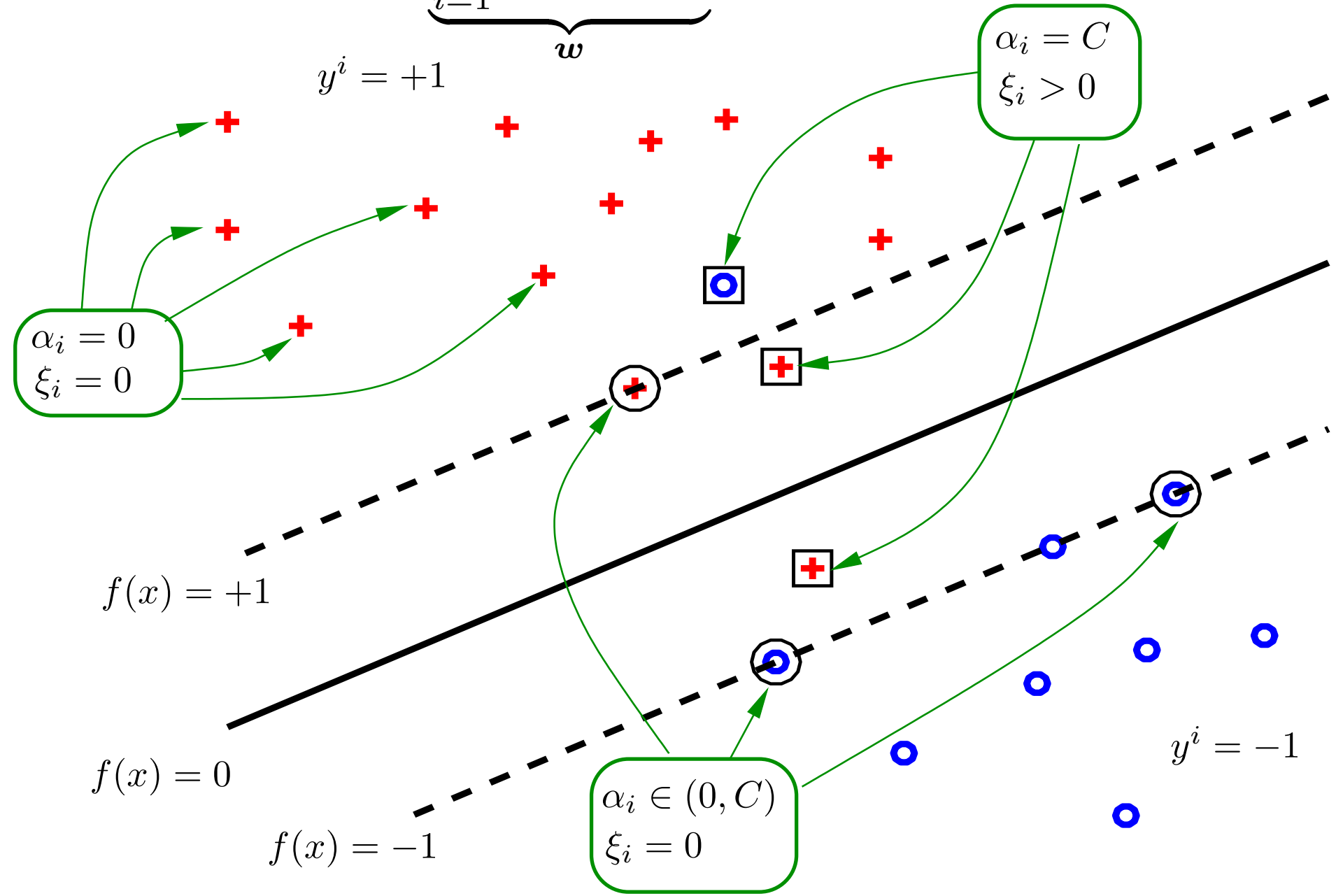
Example: SVM classifier

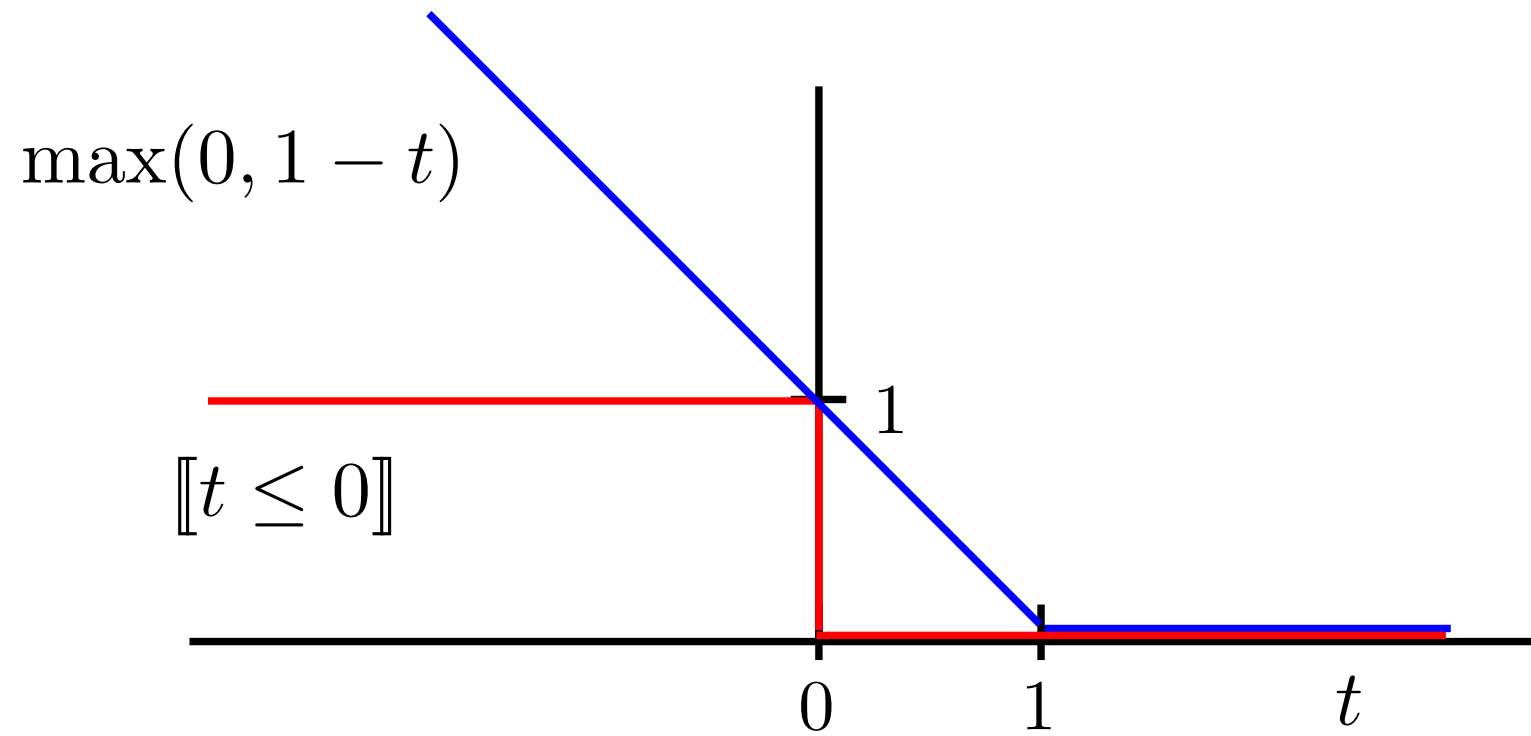
$$f(x) = \langle \mathbf{w}, \phi(x) \rangle + b = \langle \underbrace{\sum_{i=1}^m y^i \alpha_i \phi(x^i)}_{\mathbf{w}}, \phi(x) \rangle + b$$

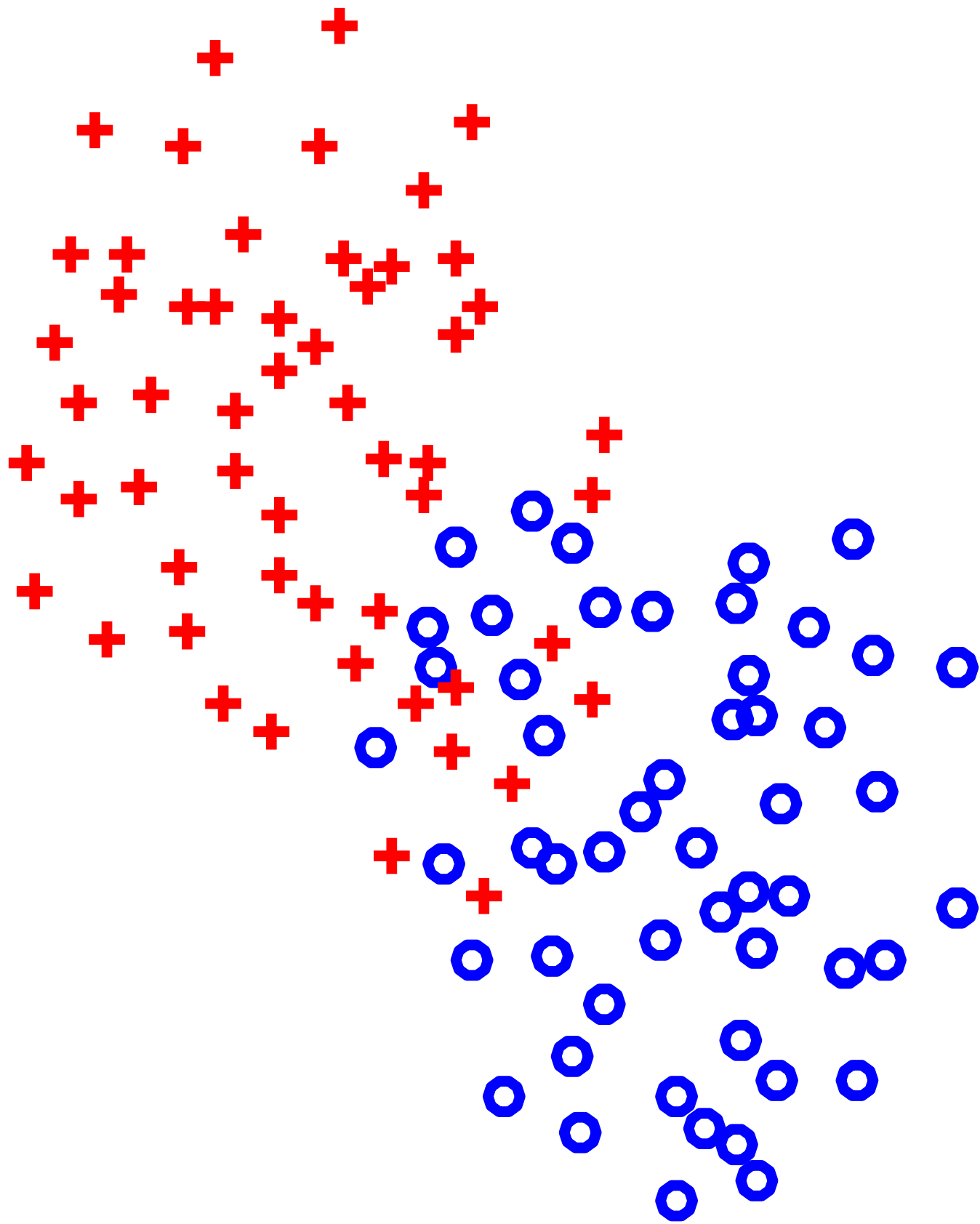


Example: SVM classifier

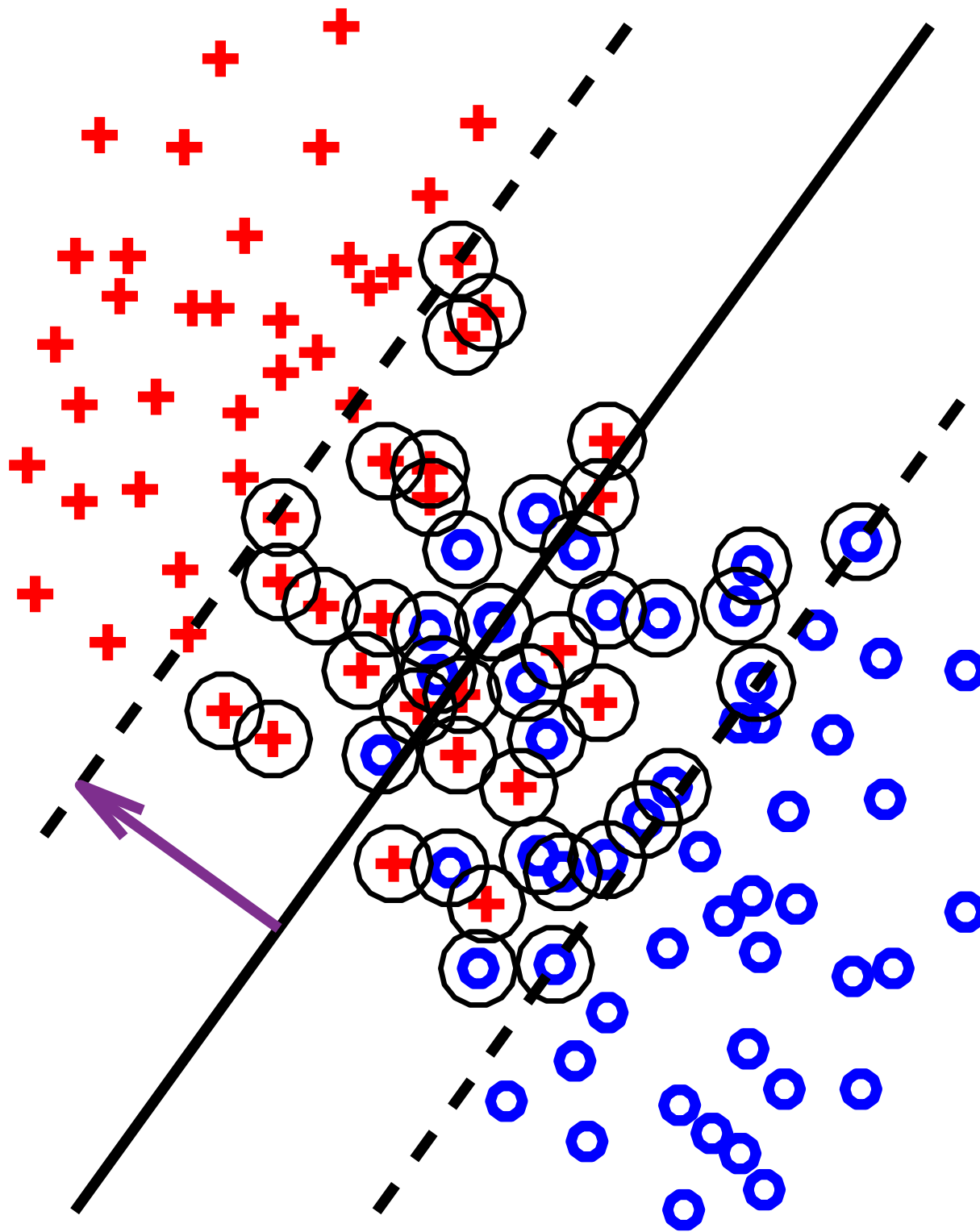
$$f(x) = \langle \mathbf{w}, \phi(x) \rangle + b = \langle \underbrace{\sum_{i=1}^m y^i \alpha_i \phi(x^i)}_{\mathbf{w}}, \phi(x) \rangle + b$$



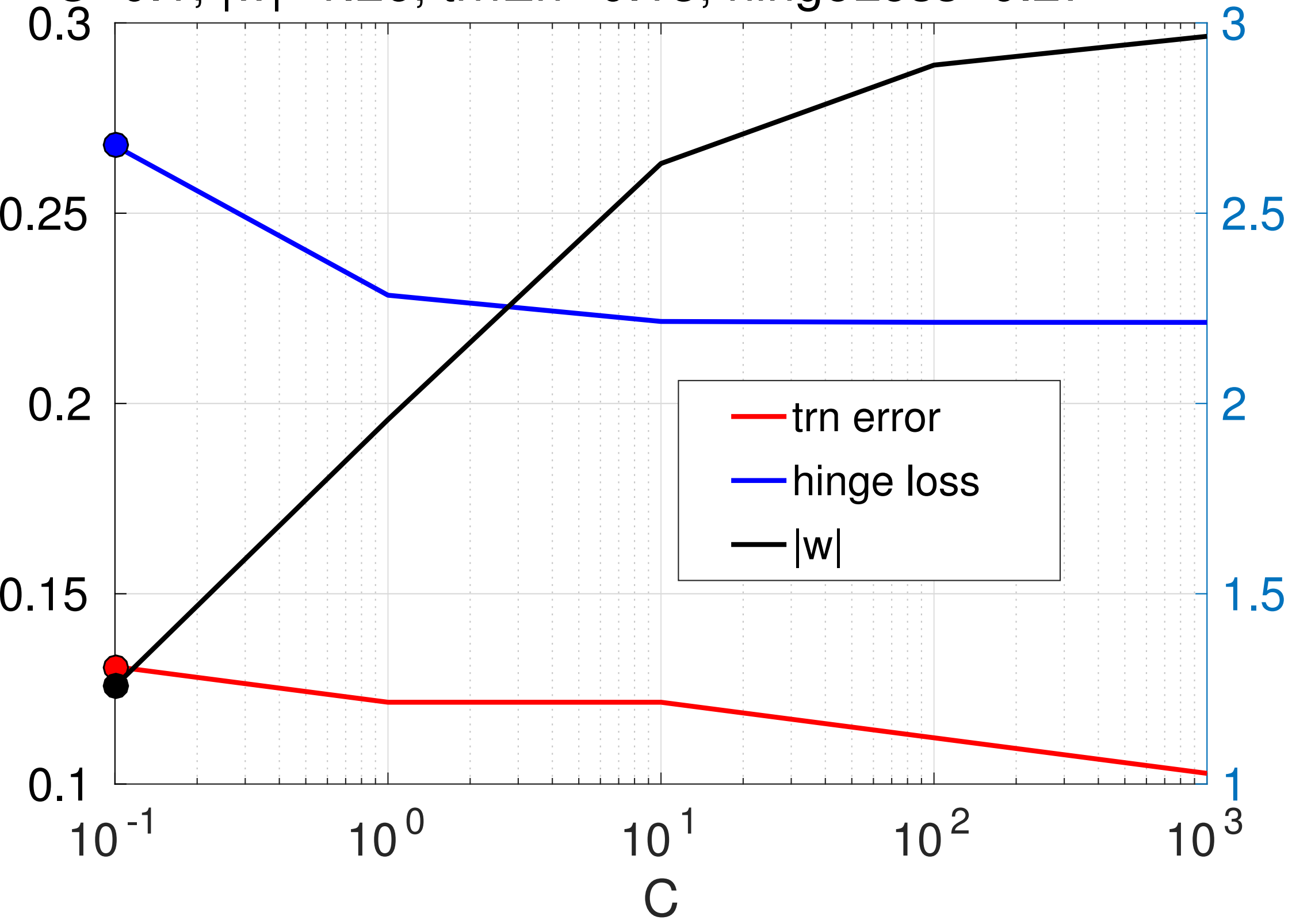




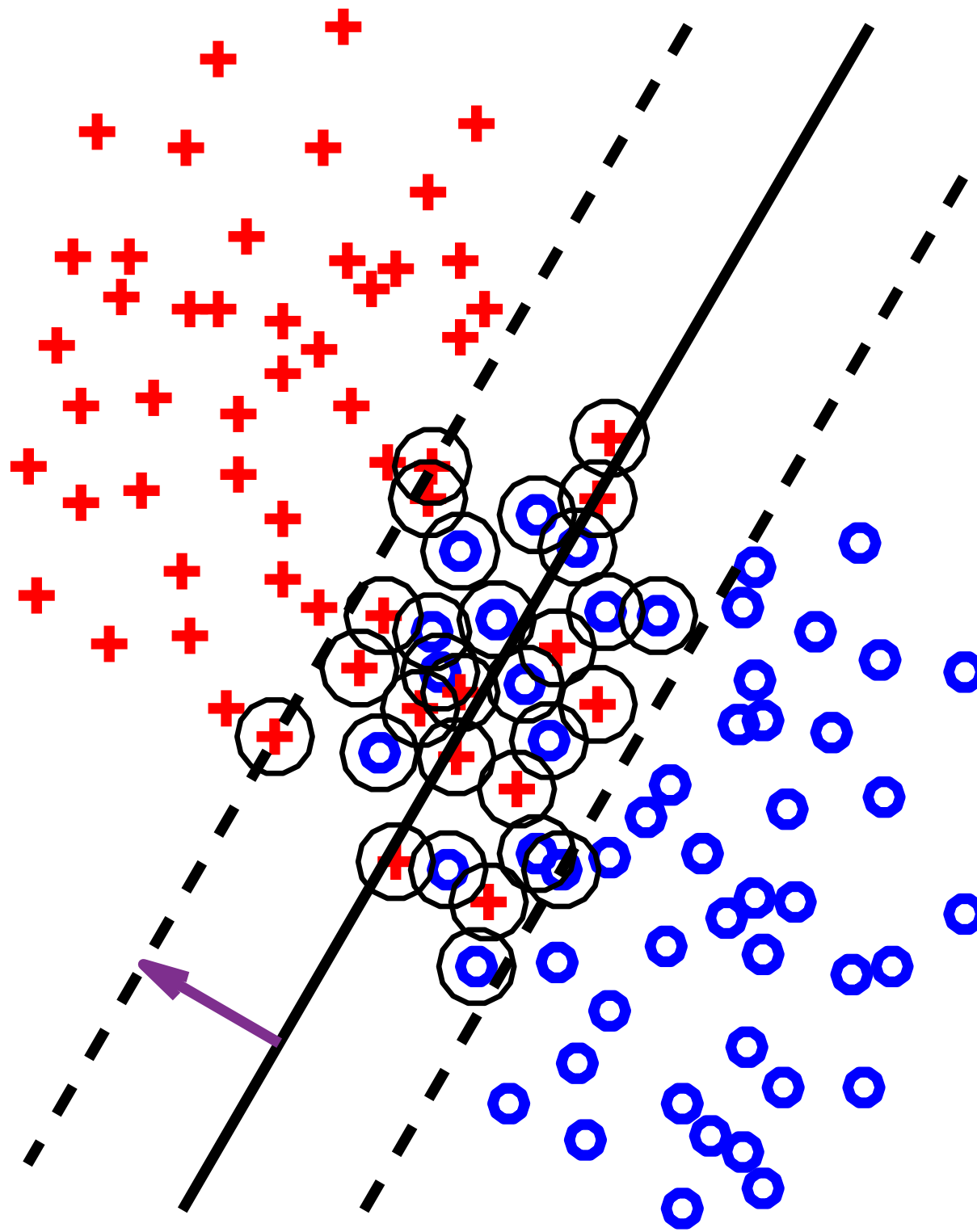
$$1/|w|=0.80$$



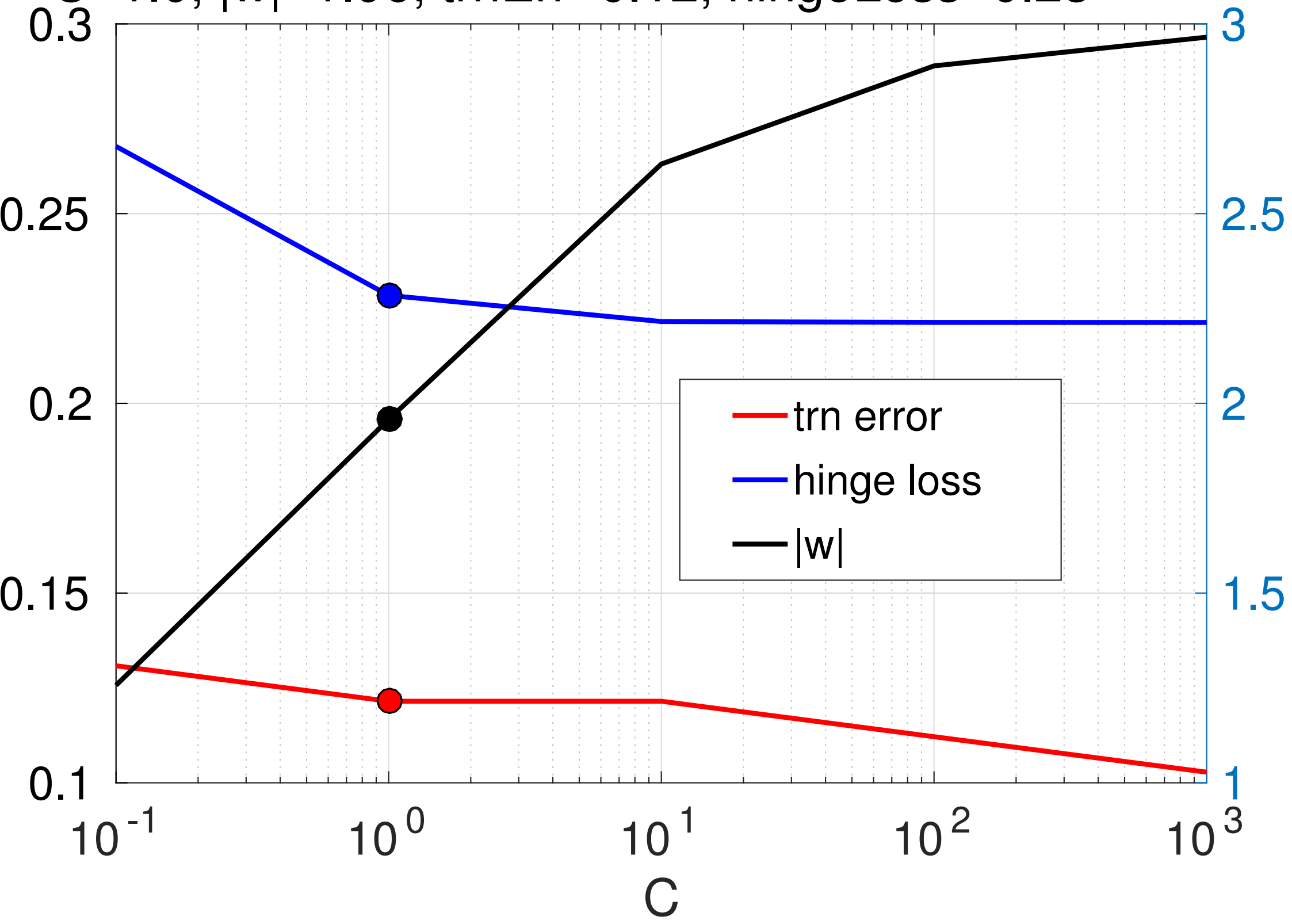
$C=0.1$, $|w|=1.26$, $\text{trnErr}=0.13$, $\text{hingeLoss}=0.27$



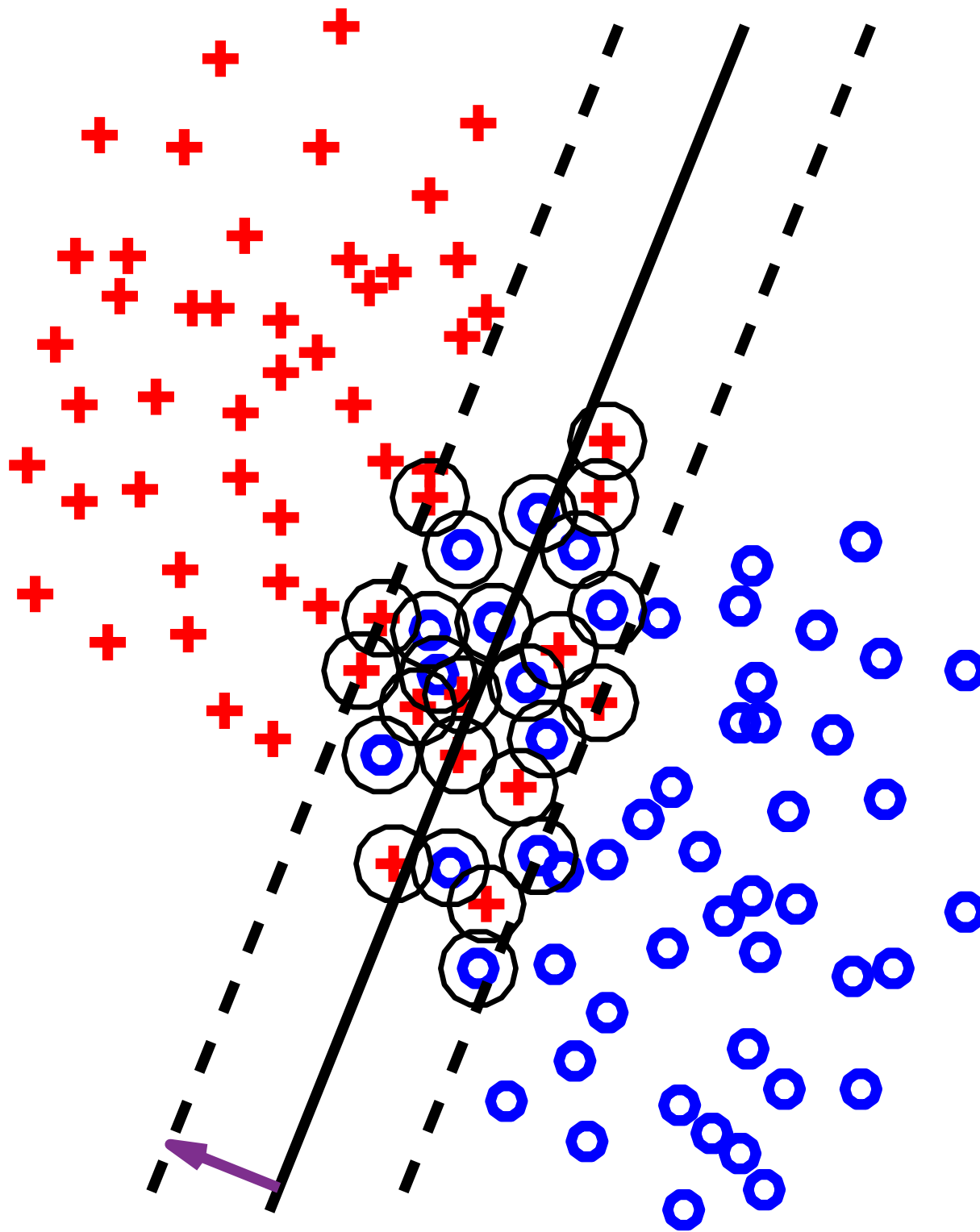
$$1/|w|=0.51$$



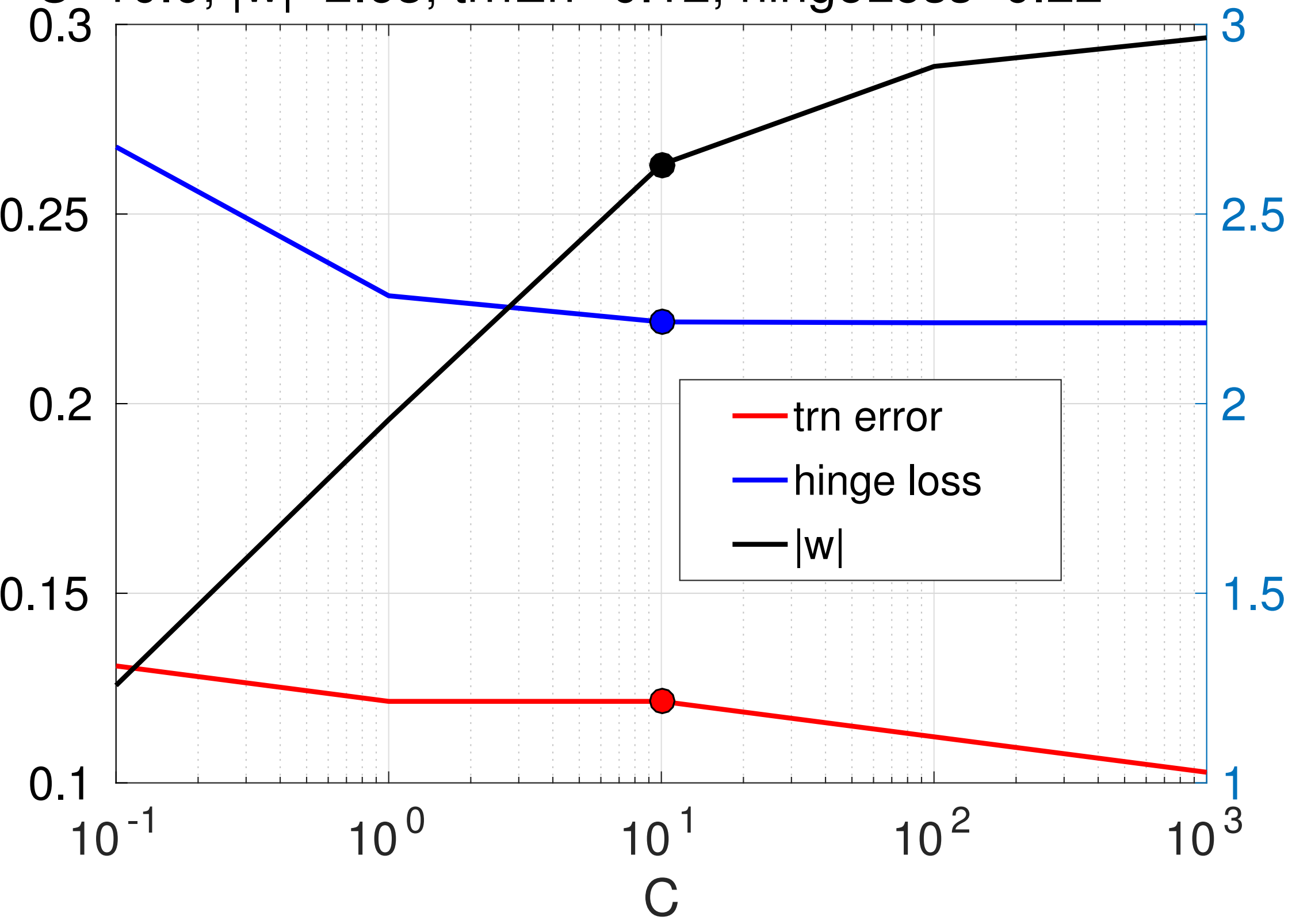
$C=1.0$, $|w|=1.96$, $\text{trnErr}=0.12$, $\text{hingeLoss}=0.23$



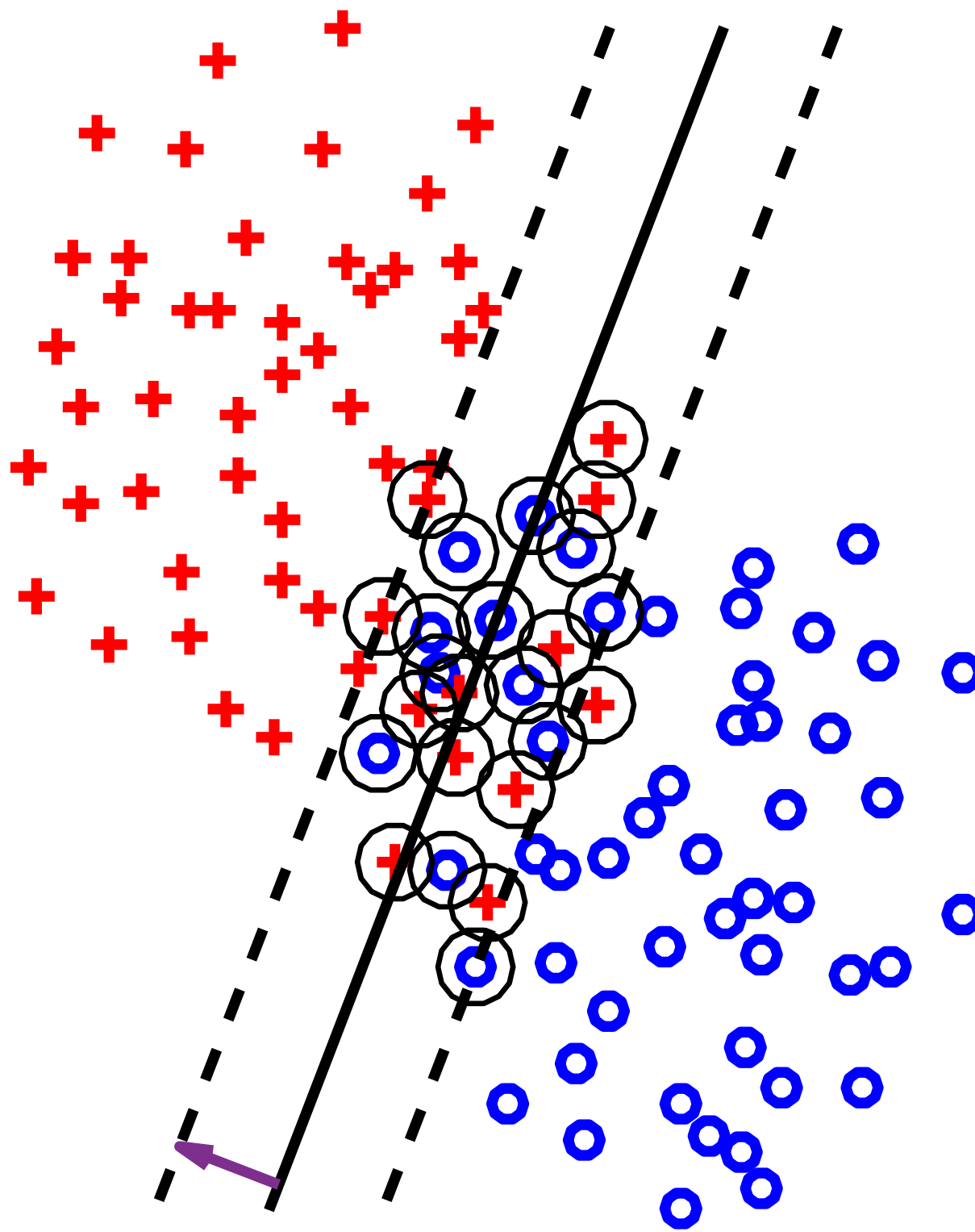
$$1/|w|=0.38$$



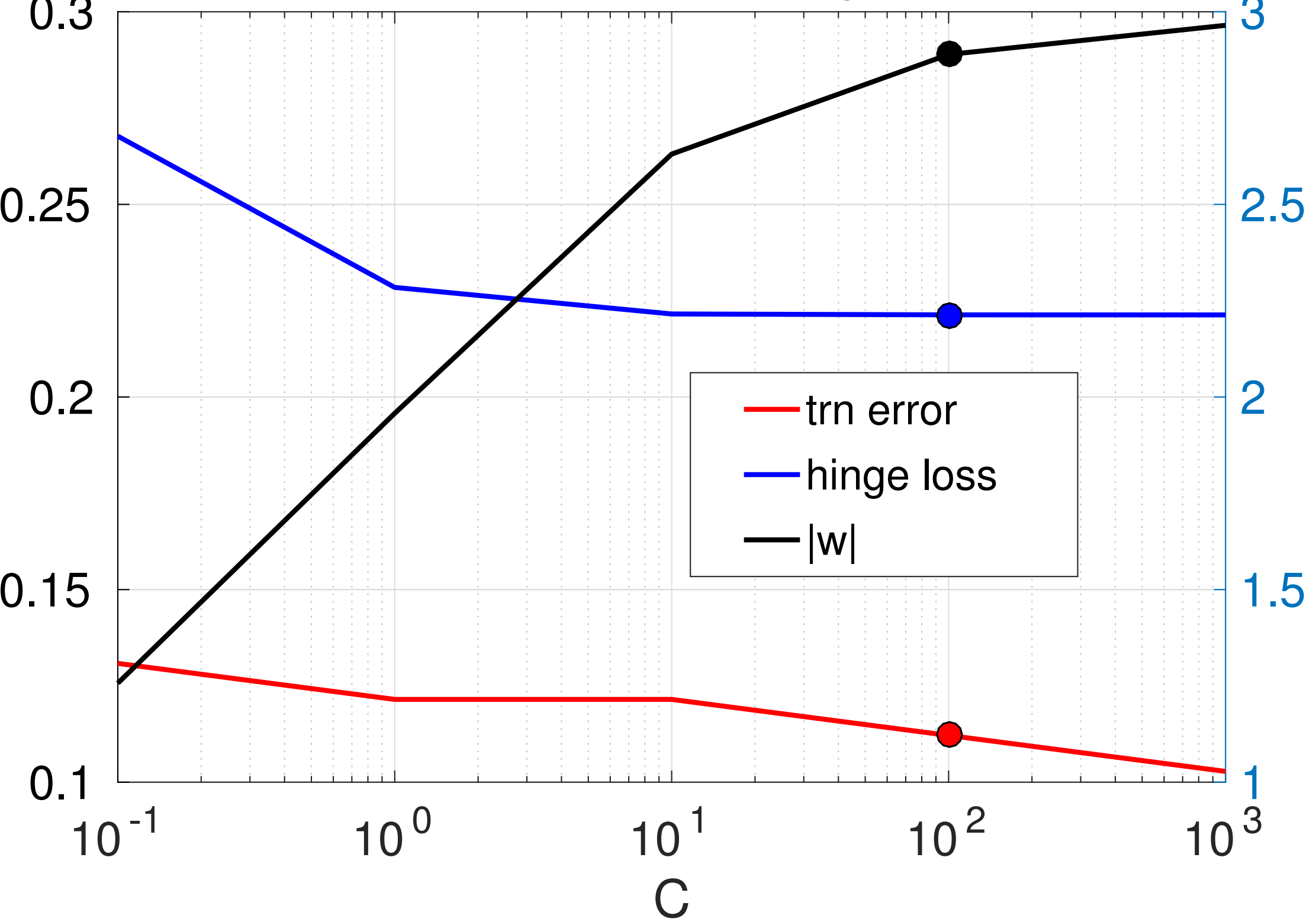
$C=10.0$, $|w|=2.63$, $\text{trnErr}=0.12$, $\text{hingeLoss}=0.22$



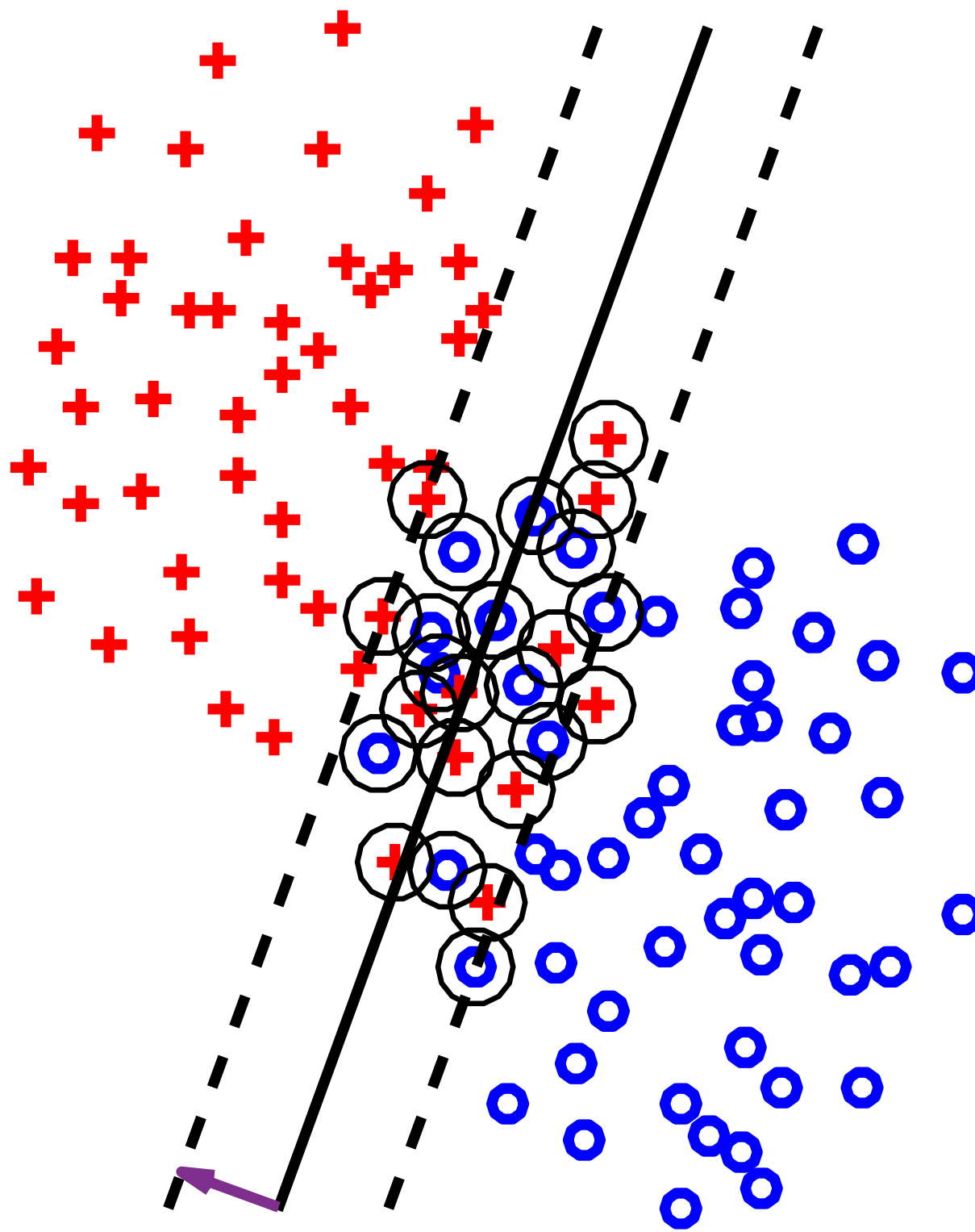
$$1/|w|=0.35$$



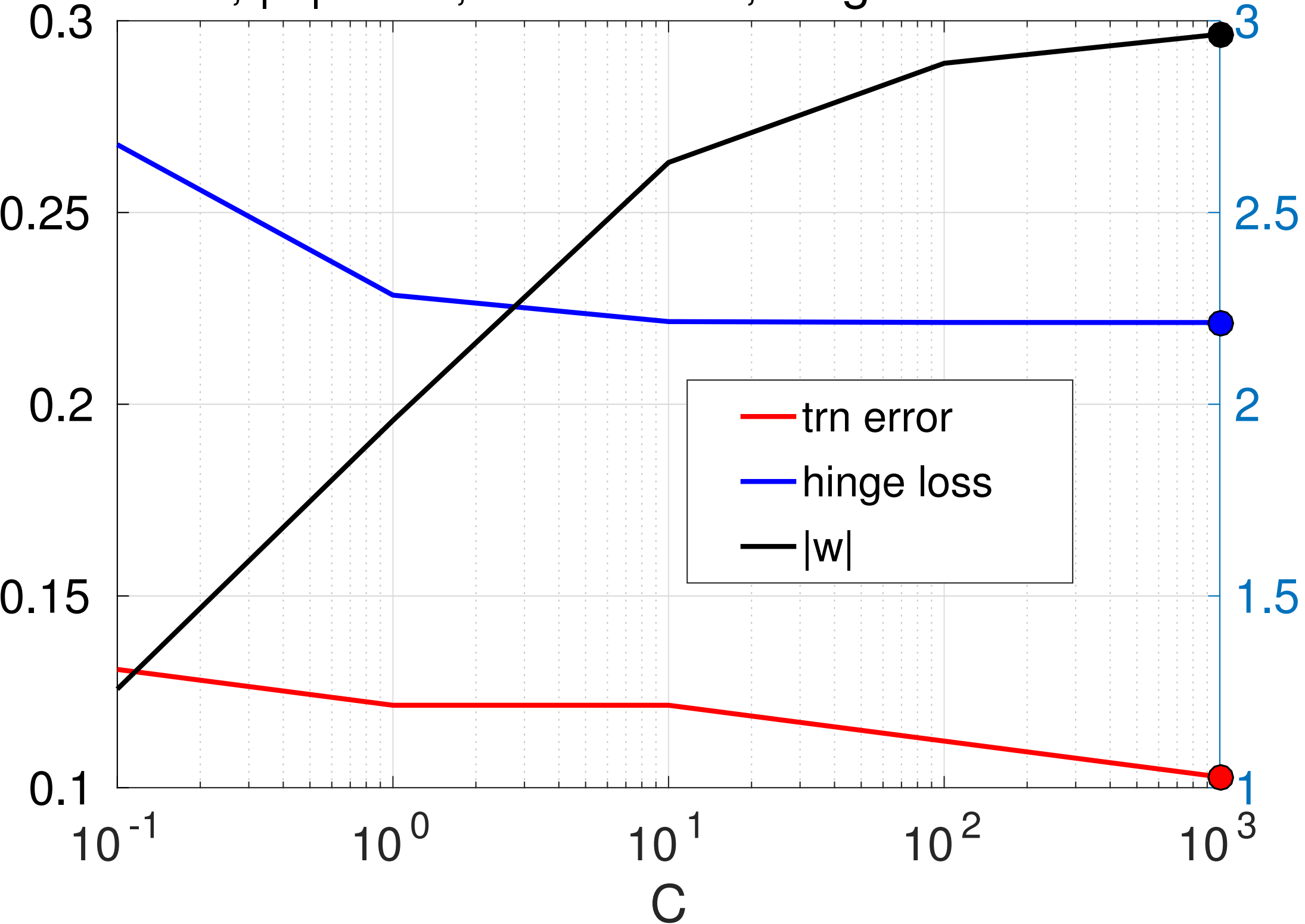
$C=100.0$, $|w|=2.89$, $\text{trnErr}=0.11$, $\text{hingeLoss}=0.22$



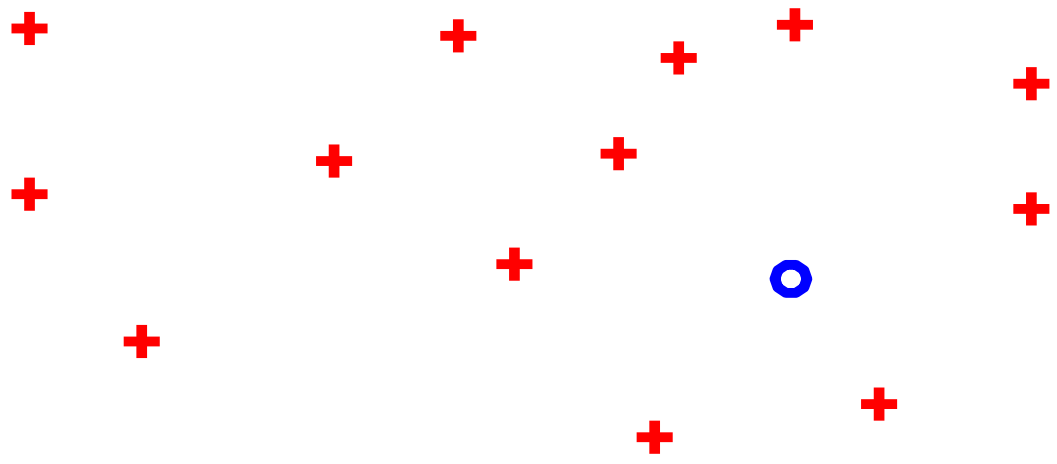
$$1/|w|=0.34$$



$C=1000.0$, $|w|=2.96$, $\text{trnErr}=0.10$, $\text{hingeLoss}=0.22$



$y^i = +1$



$y^i = -1$

