

Introduction to Machine Learning

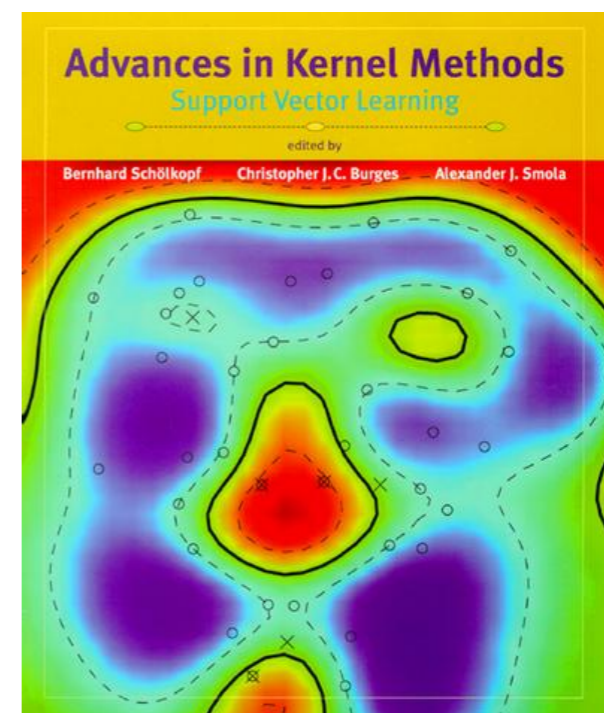
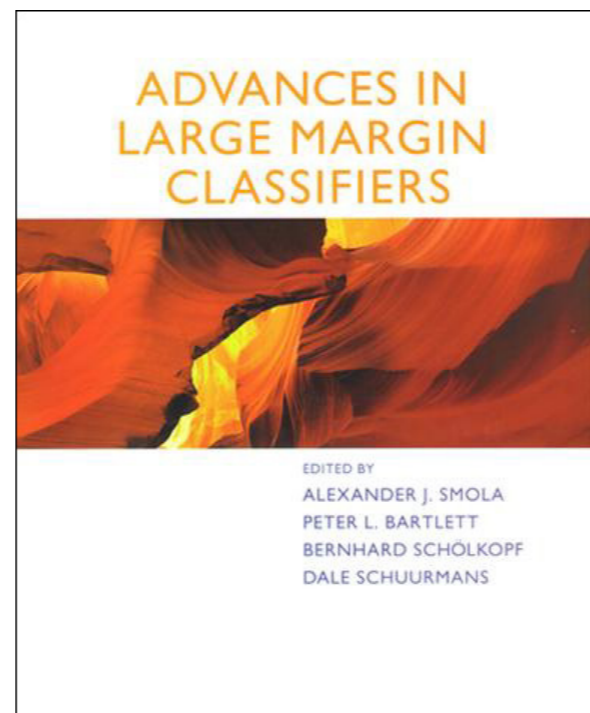
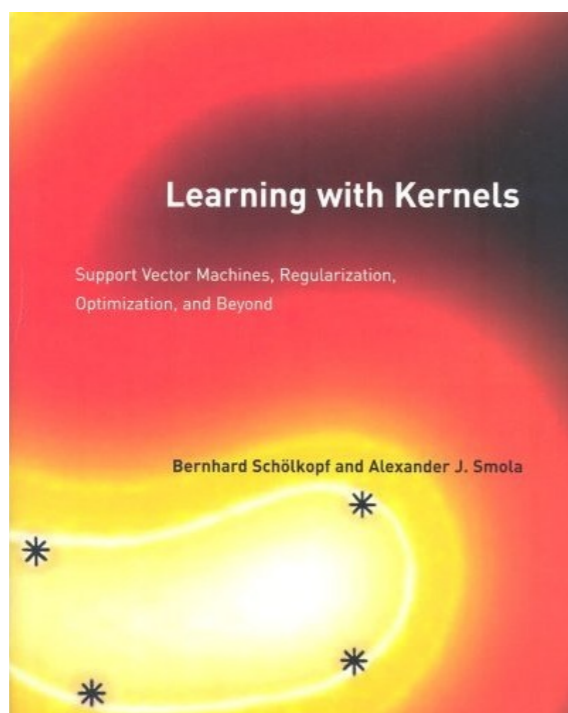
6. Support Vector Classification

Geoff Gordon and Alex Smola
Carnegie Mellon University

<http://alex.smola.org/teaching/cmu2013-10-701x>
10-701

Outline

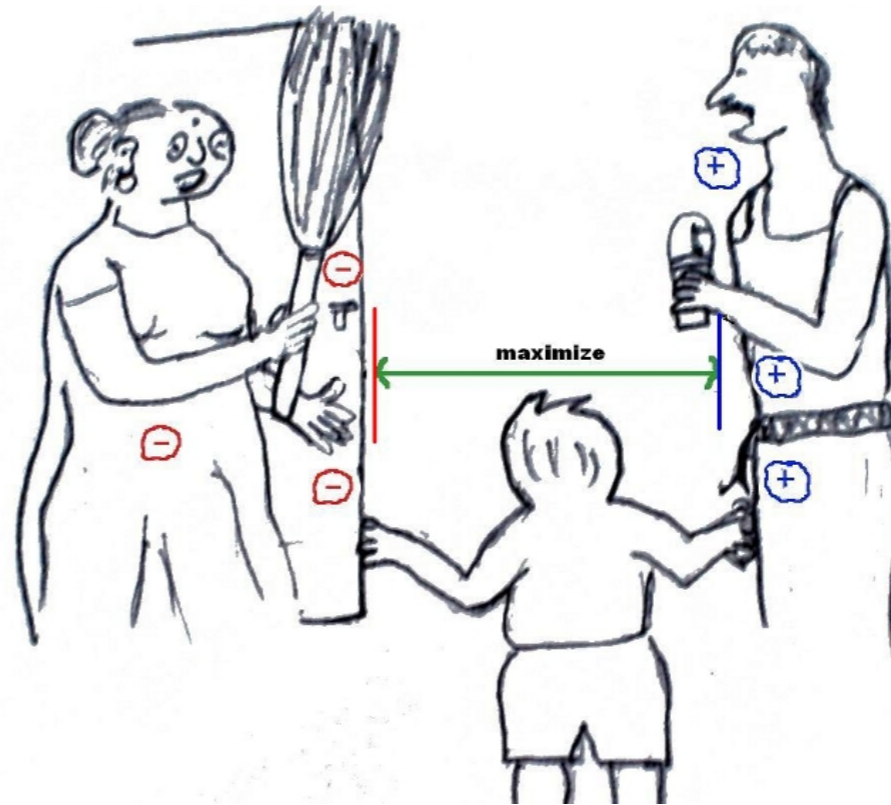
- Support Vector Classification
Large Margin Separation, optimization problem
- Properties
Support Vectors, kernel expansion
- Soft margin classifier
Dual problem, robustness





MAGIC Etch A Sketch[®] SCREEN

Support
Vector
Machines



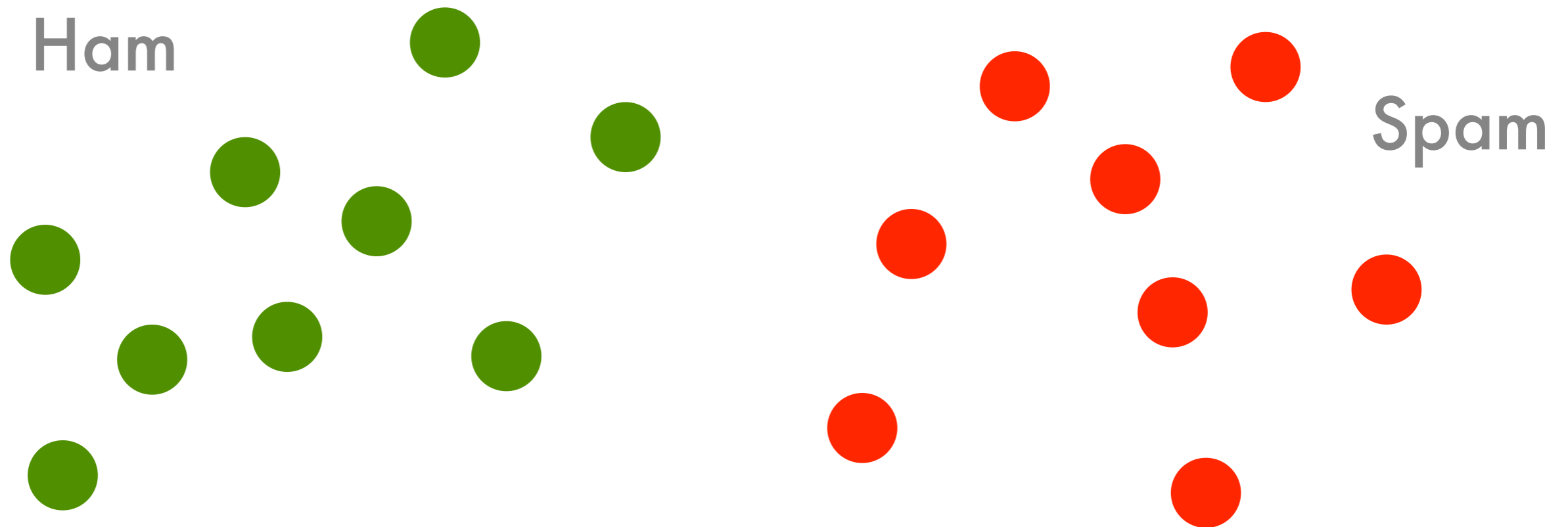
Horizontal
Grid

OHIO ART "The World of Toys"

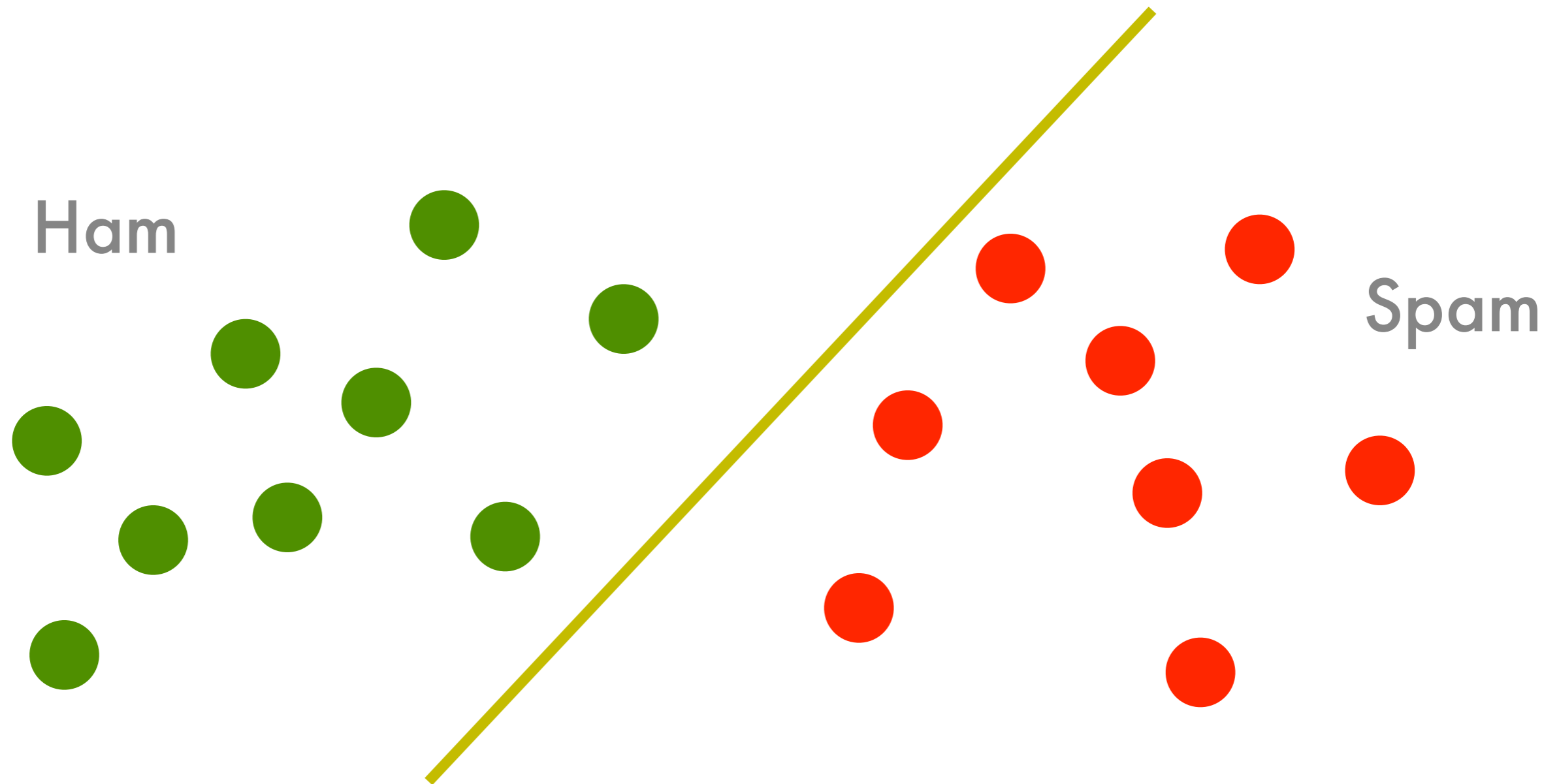
Vertical
Grid

MAGIC SCREEN IS GLASS SET IN DURABLE PLASTIC FRAME
USE WITH CARE

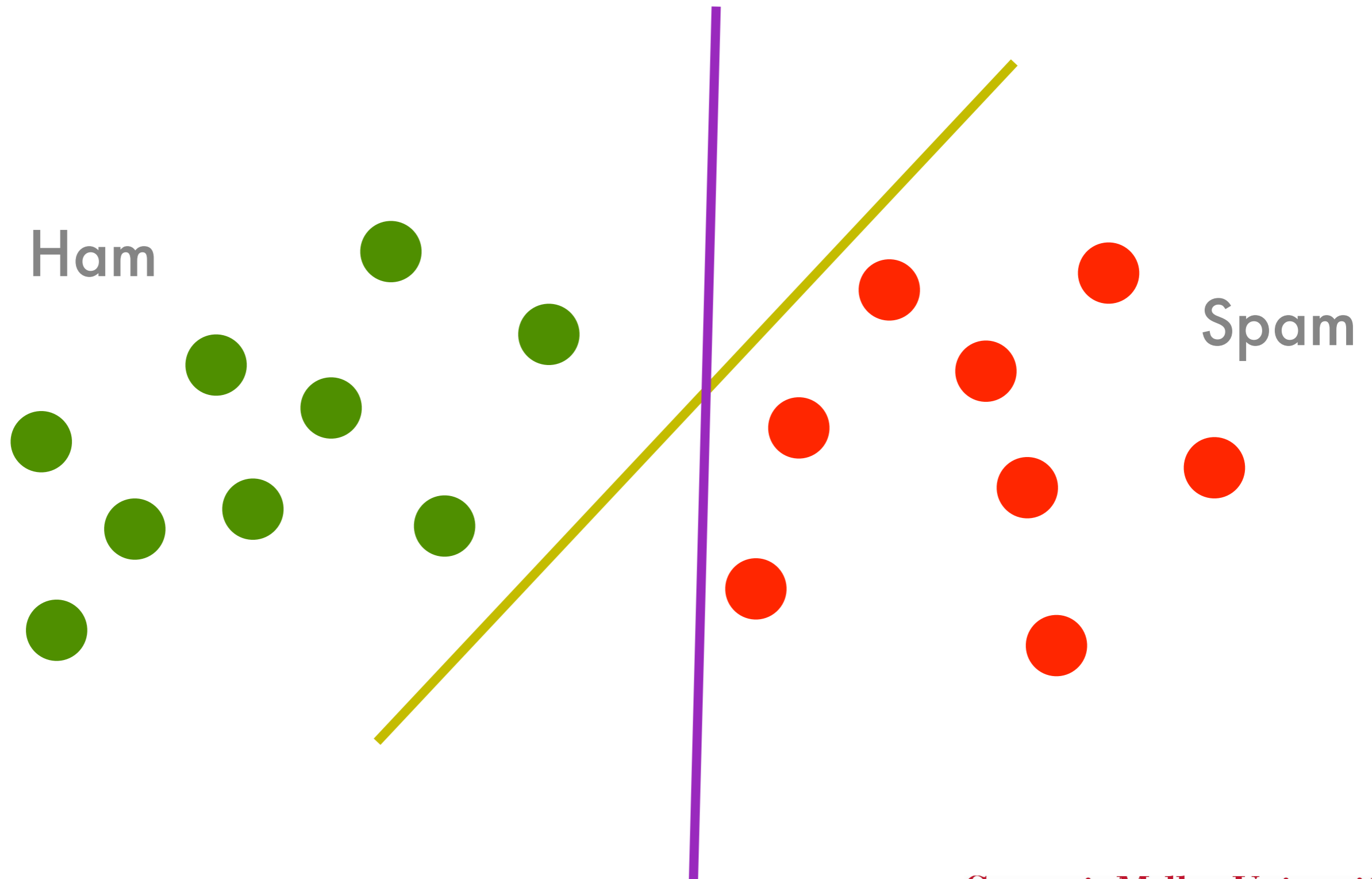
Linear Separator



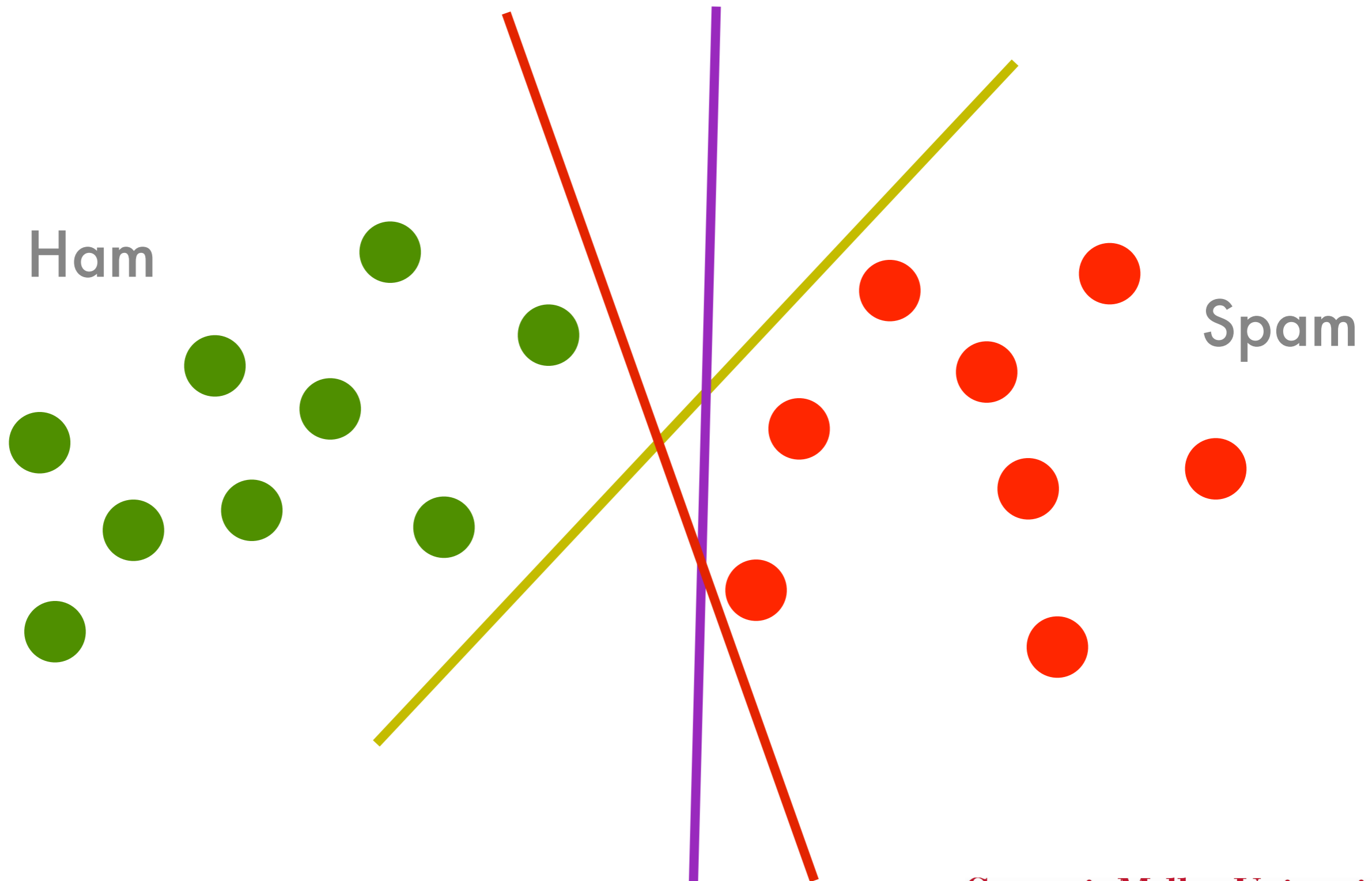
Linear Separator



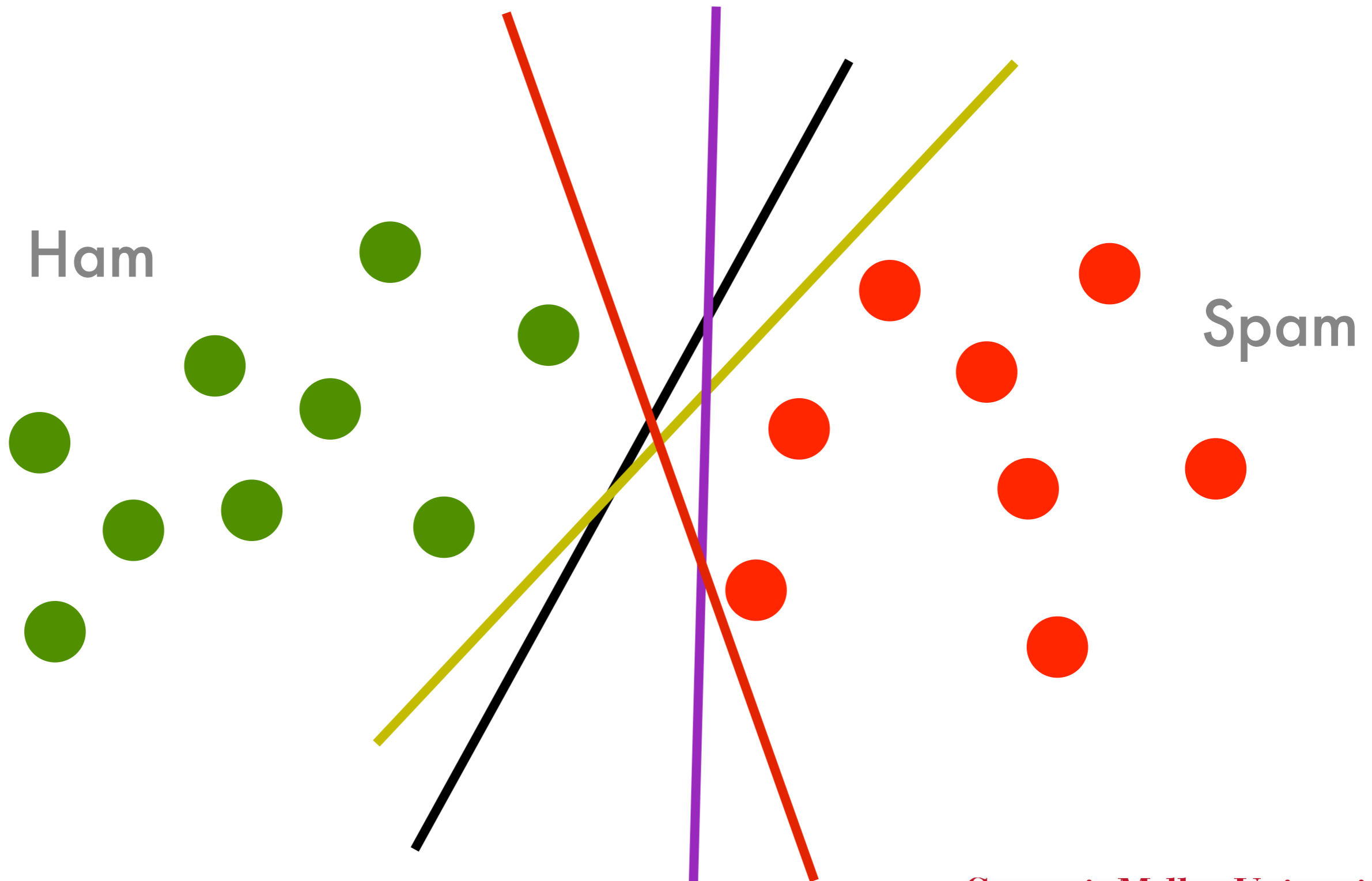
Linear Separator



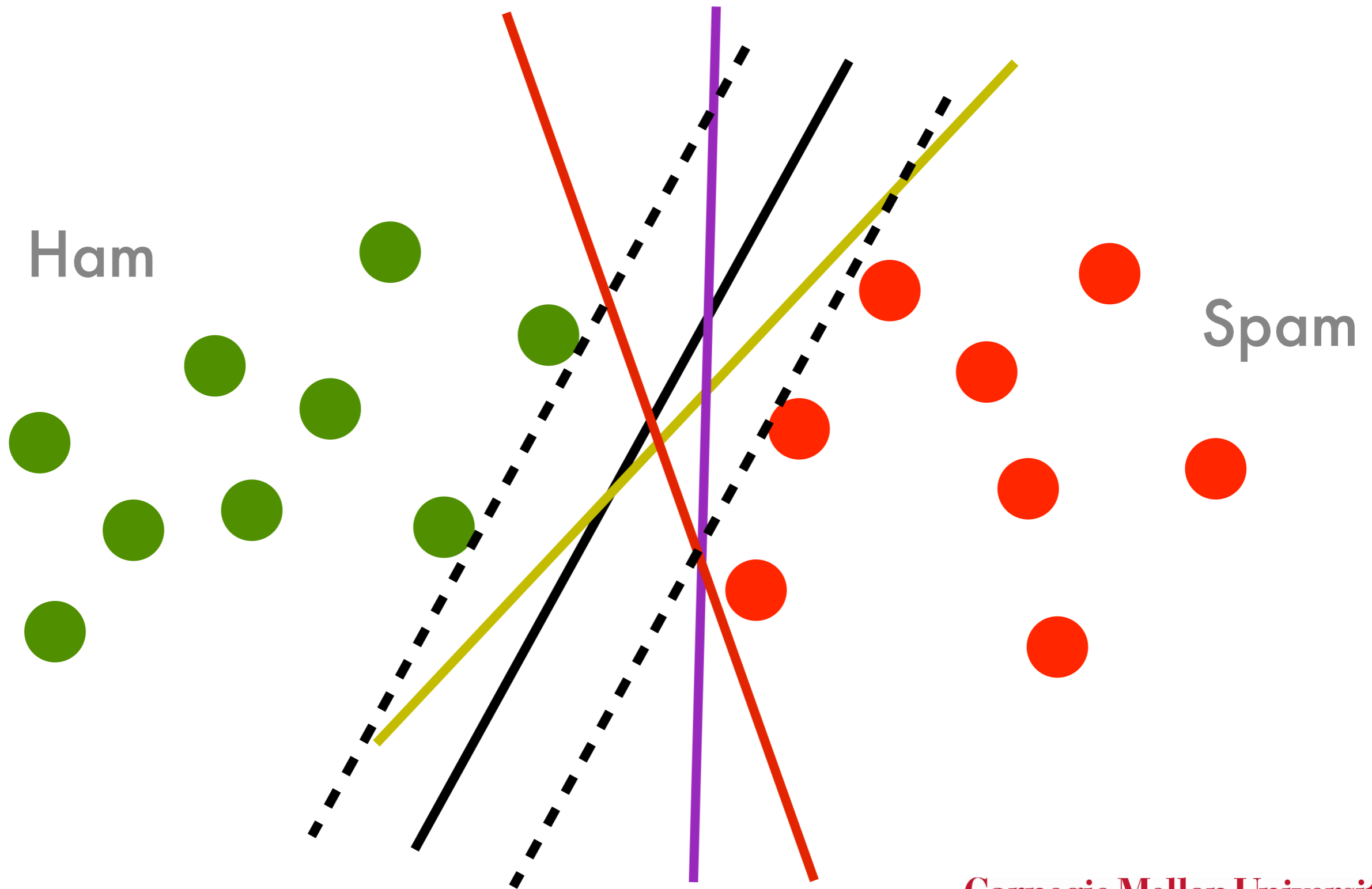
Linear Separator



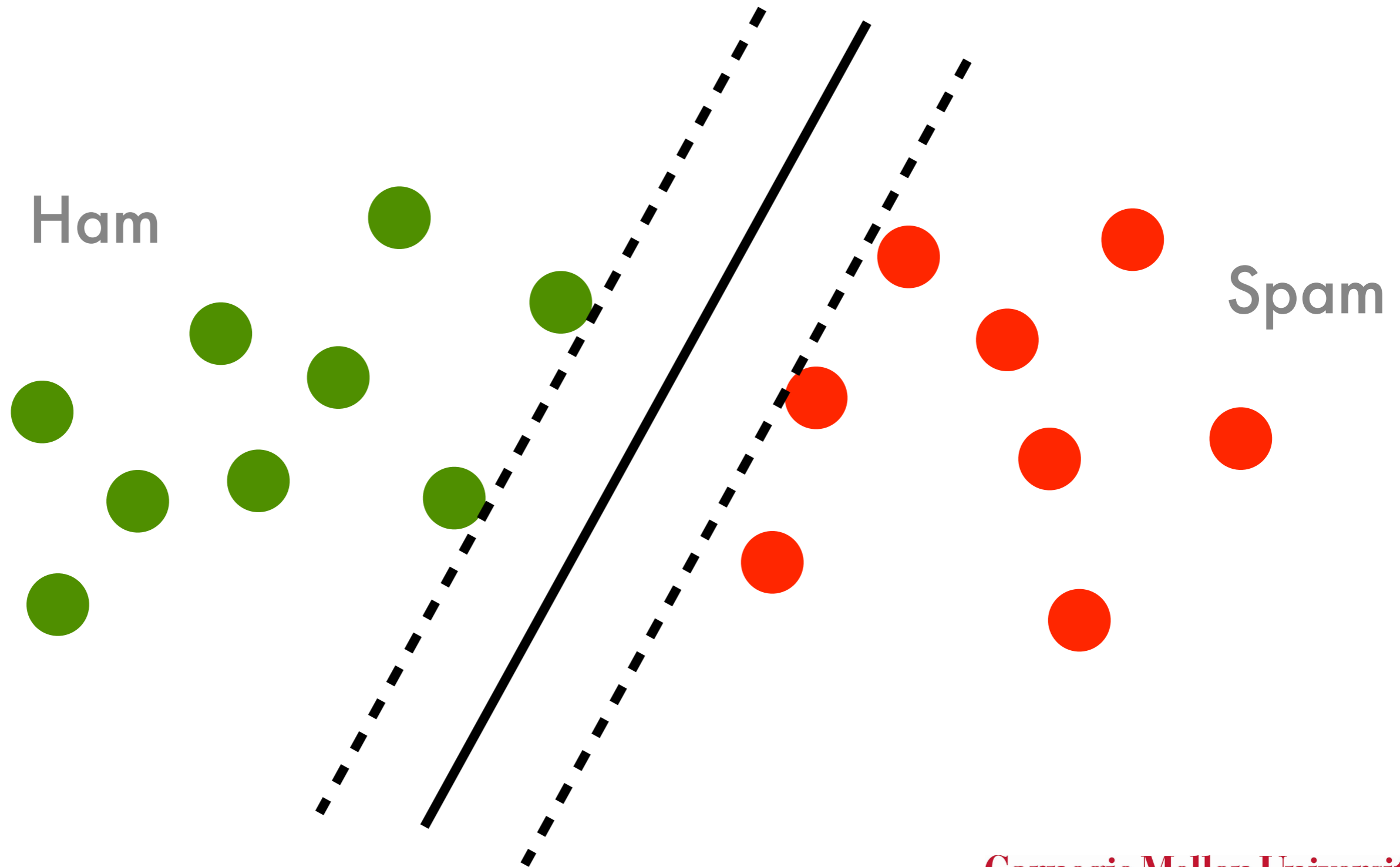
Linear Separator



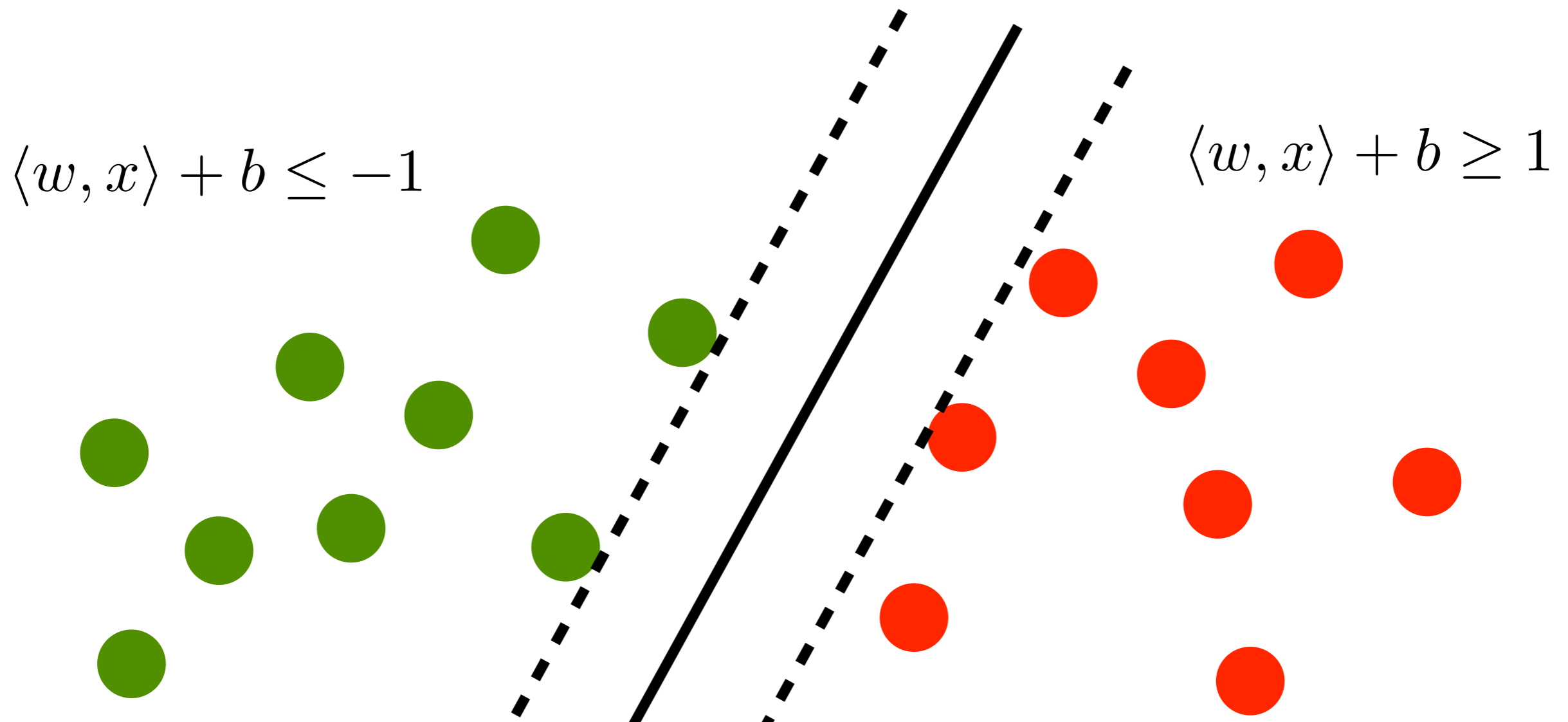
Linear Separator



Linear Separator



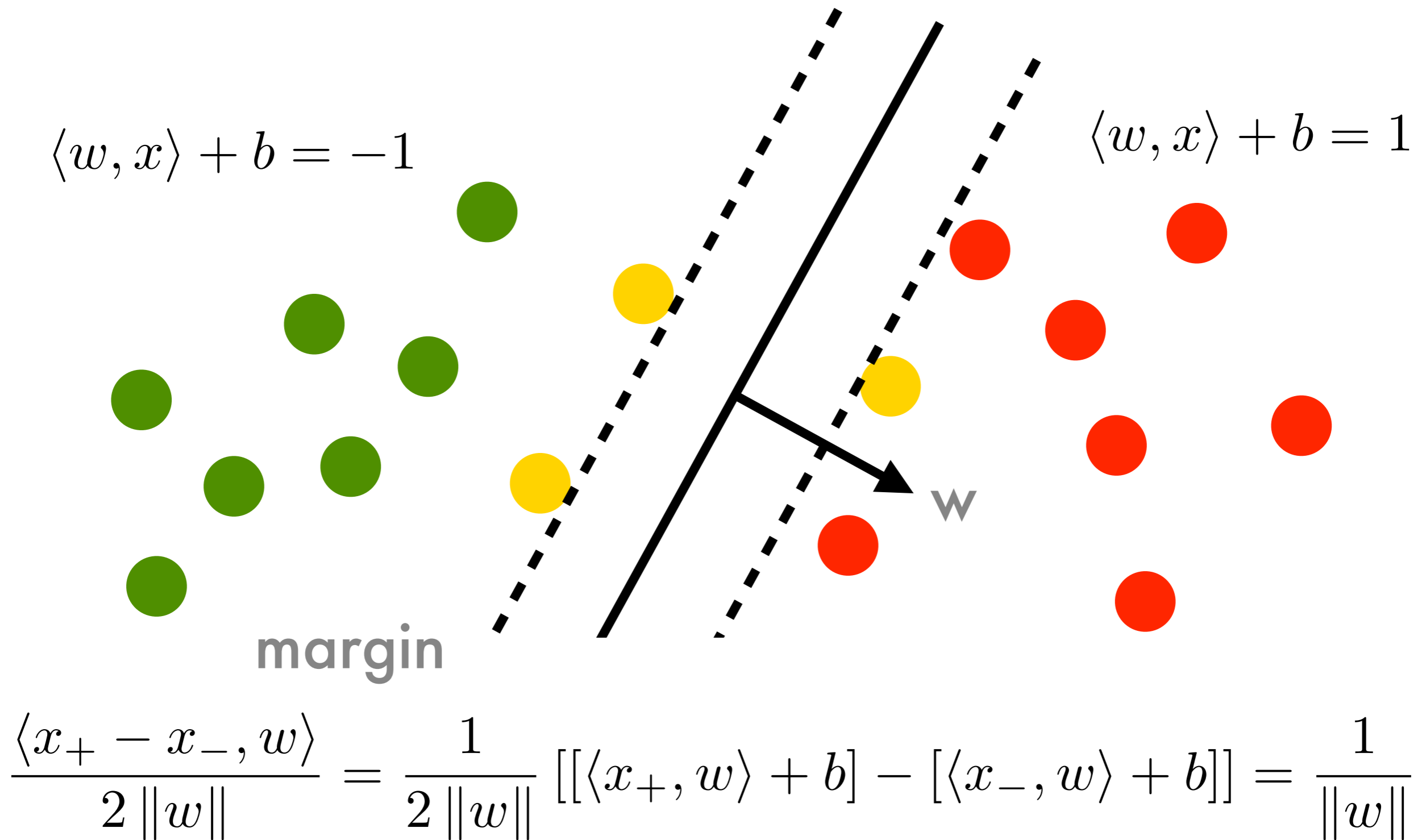
Large Margin Classifier



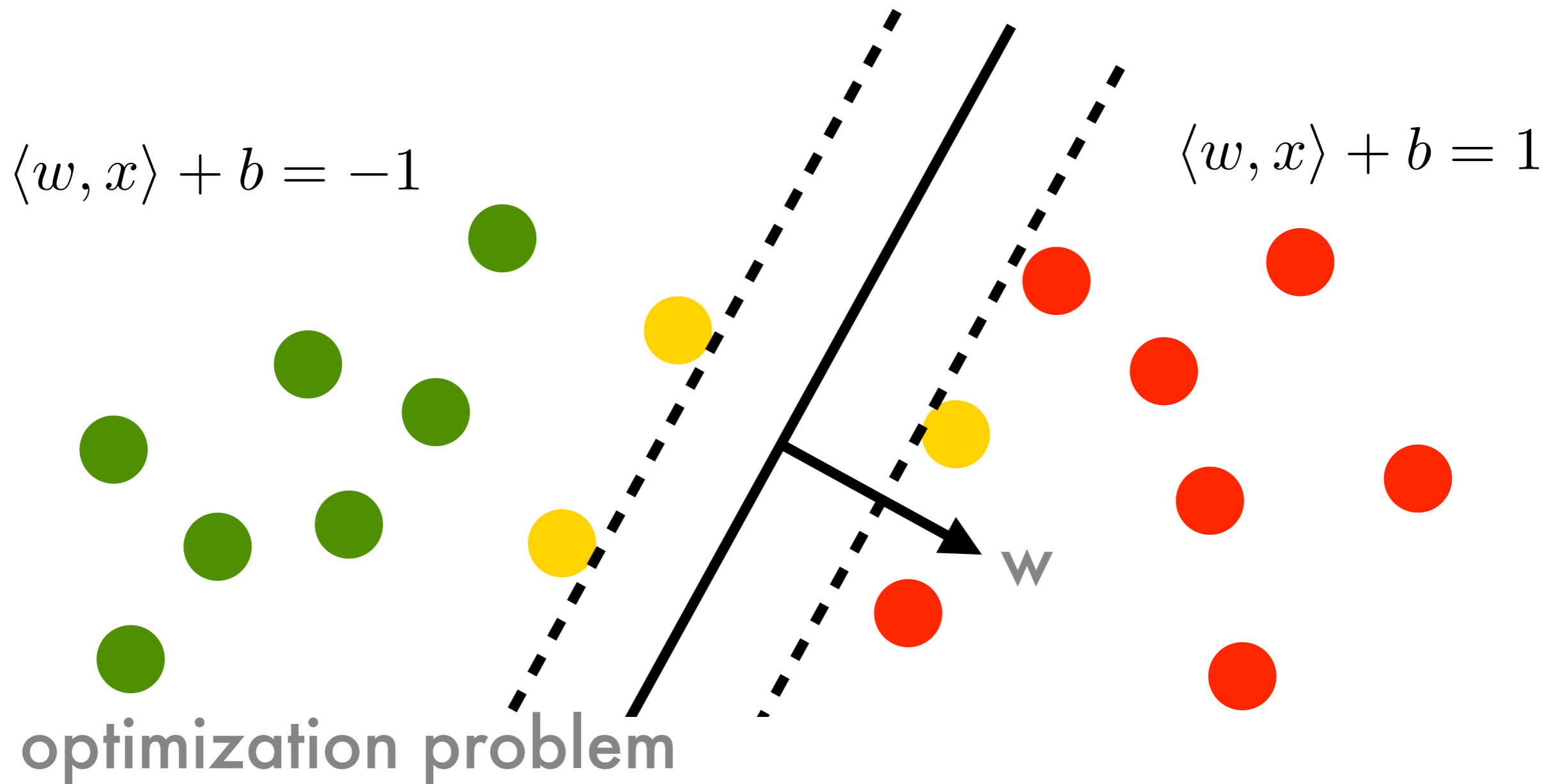
linear function

$$f(x) = \langle w, x \rangle + b$$

Large Margin Classifier

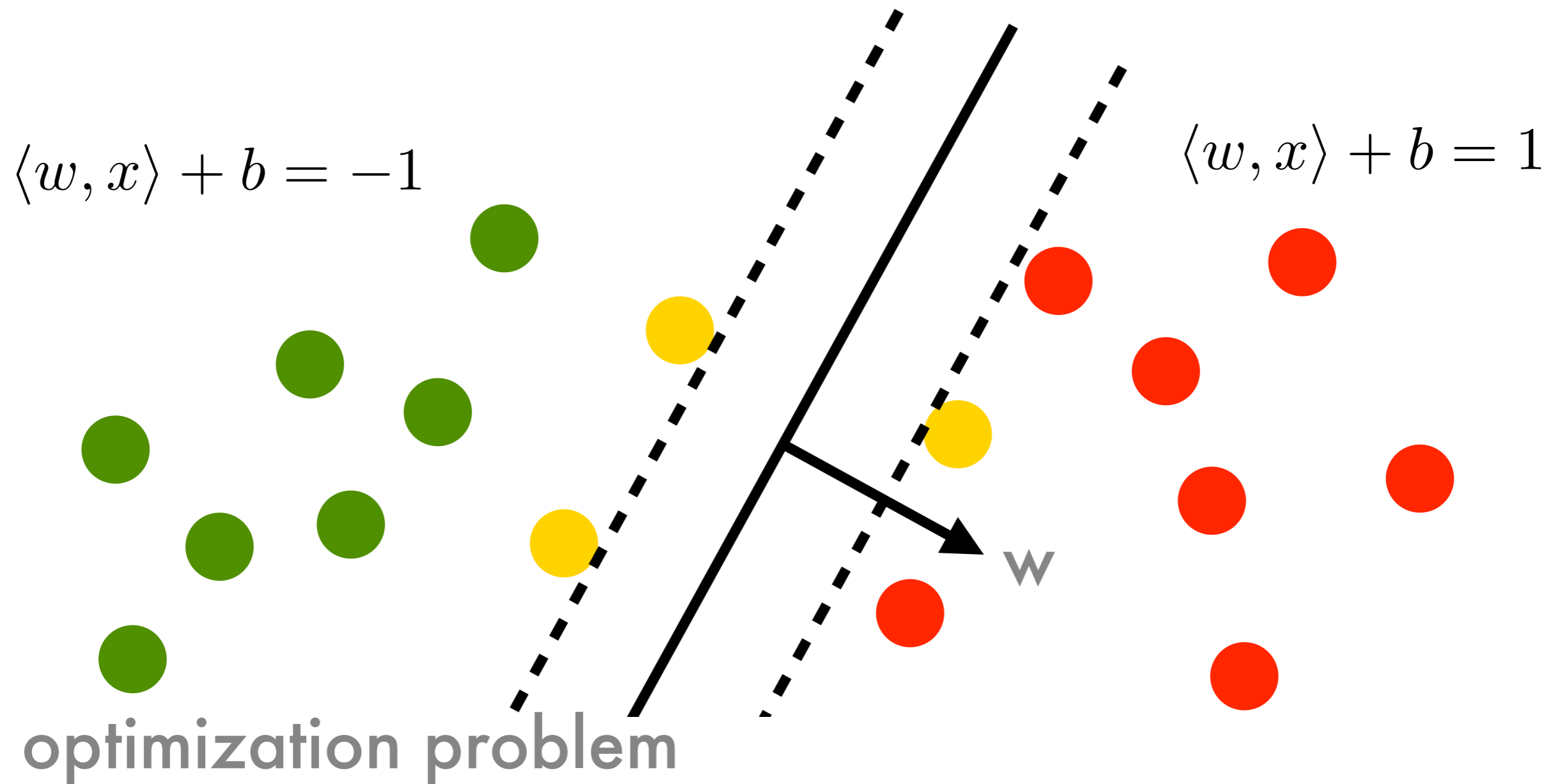


Large Margin Classifier



$$\text{maximize}_{w,b} \frac{1}{\|w\|} \text{ subject to } y_i [\langle x_i, w \rangle + b] \geq 1$$

Large Margin Classifier



$$\text{minimize}_{w,b} \frac{1}{2} \|w\|^2 \quad \text{subject to } y_i [\langle x_i, w \rangle + b] \geq 1$$

Dual Problem

- Primal optimization problem

$$\text{minimize}_{w,b} \frac{1}{2} \|w\|^2 \text{ subject to } y_i [\langle x_i, w \rangle + b] \geq 1$$

- Lagrange function

$$L(w, b, \alpha) = \frac{1}{2} \|w\|^2 - \sum_i \alpha_i [y_i [\langle x_i, w \rangle + b] - 1]$$

constraint

Optimality in w, b is at saddle point with α

- Derivatives in w, b need to vanish

Dual Problem

- **Lagrange function**

$$L(w, b, \alpha) = \frac{1}{2} \|w\|^2 - \sum_i \alpha_i [y_i [\langle x_i, w \rangle + b] - 1]$$

- **Derivatives in w , b need to vanish**

$$\partial_w L(w, b, a) = w - \sum_i \alpha_i y_i x_i = 0$$

$$\partial_b L(w, b, a) = \sum_i \alpha_i y_i = 0$$

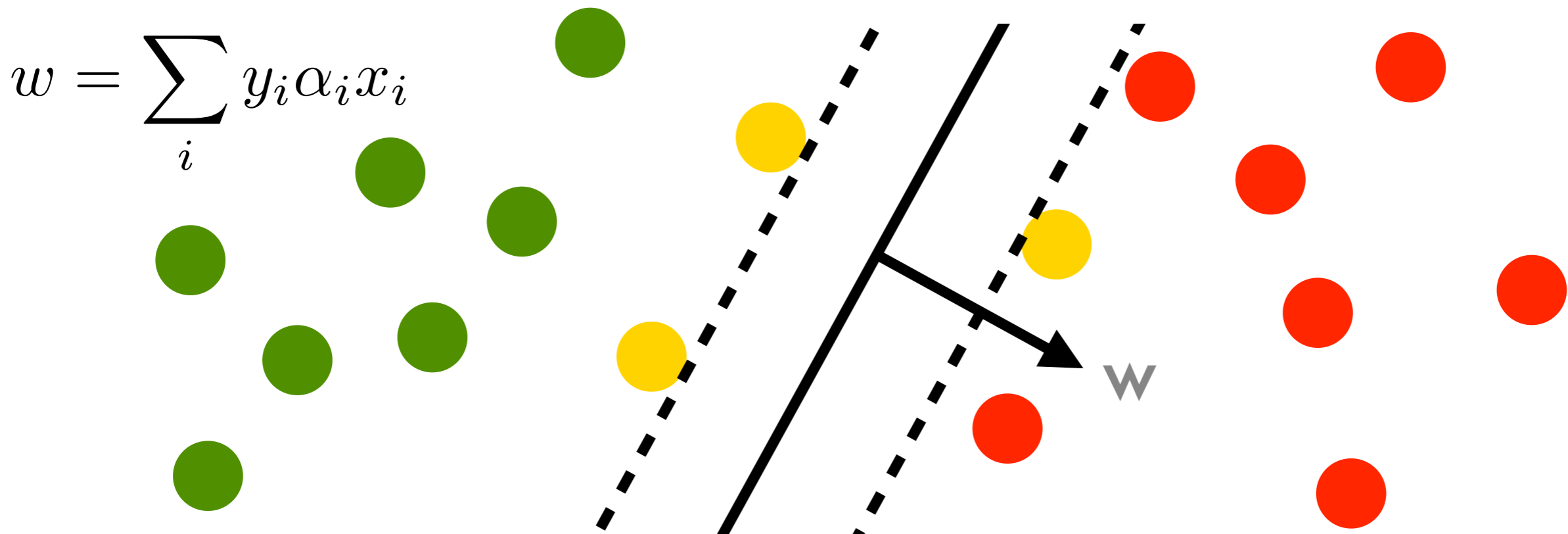
- **Plugging terms back into L yields**

$$\underset{\alpha}{\text{maximize}} - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \langle x_i, x_j \rangle + \sum_i \alpha_i$$

$$\text{subject to } \sum_i \alpha_i y_i = 0 \text{ and } \alpha_i \geq 0$$

Support Vector Machines

$$\text{minimize}_{w,b} \frac{1}{2} \|w\|^2 \quad \text{subject to } y_i [\langle x_i, w \rangle + b] \geq 1$$



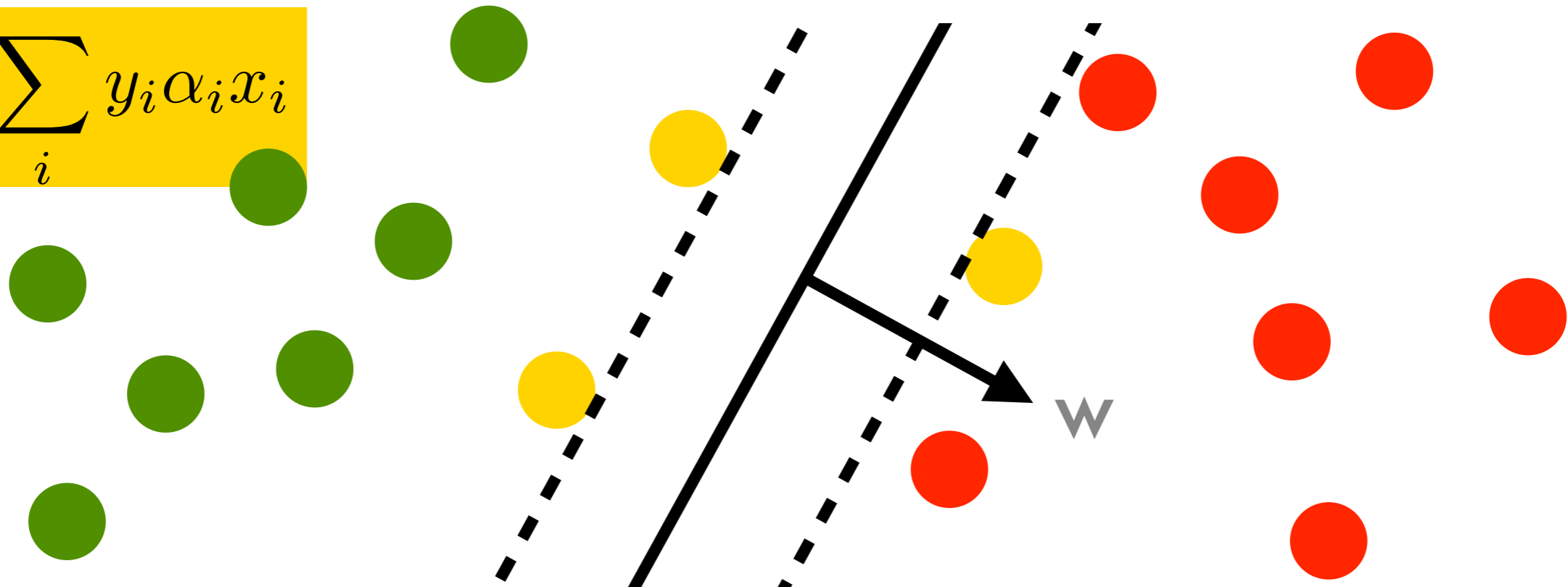
$$\text{maximize}_{\alpha} -\frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \langle x_i, x_j \rangle + \sum_i \alpha_i$$

$$\text{subject to } \sum \alpha_i y_i = 0 \text{ and } \alpha_i \geq 0$$

Support Vectors

$$\text{minimize}_{w,b} \frac{1}{2} \|w\|^2 \quad \text{subject to } y_i [\langle x_i, w \rangle + b] \geq 1$$

$$w = \sum_i y_i \alpha_i x_i$$



Karush Kuhn Tucker

Optimality condition

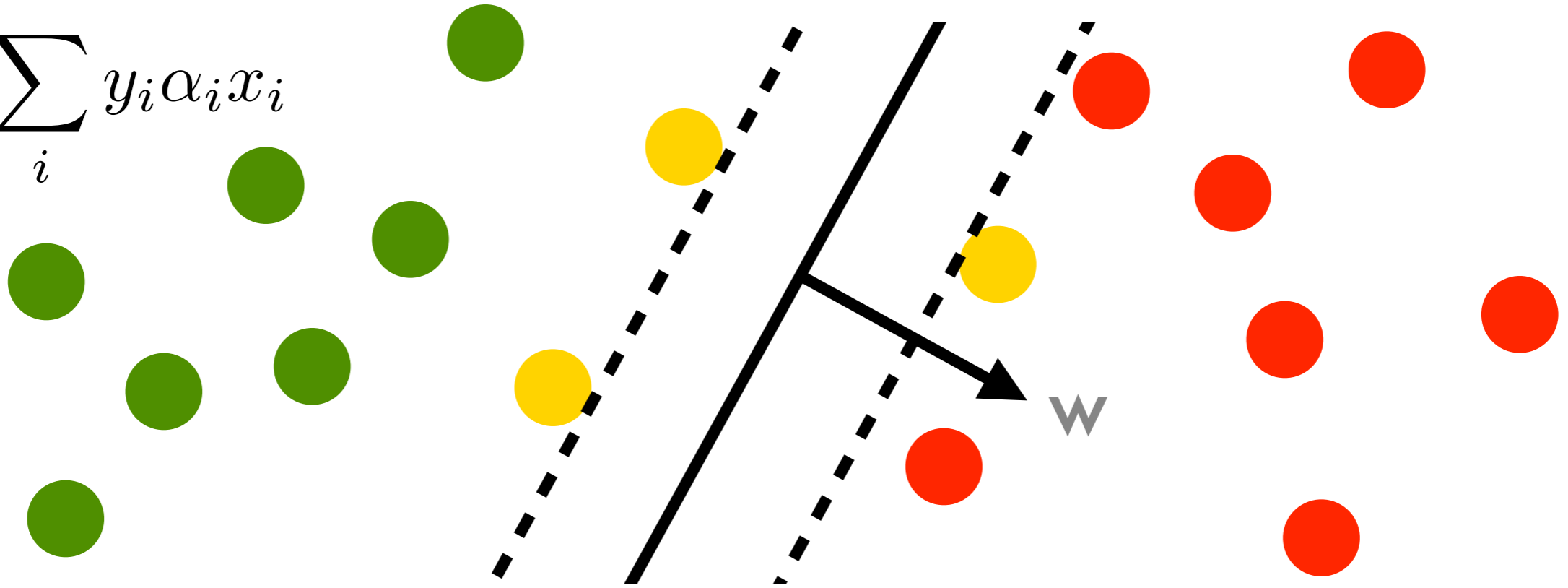
$$\alpha_i [y_i [\langle w, x_i \rangle + b] - 1] = 0$$

$$\alpha_i = 0$$

$$\alpha_i > 0 \implies y_i [\langle w, x_i \rangle + b] = 1$$

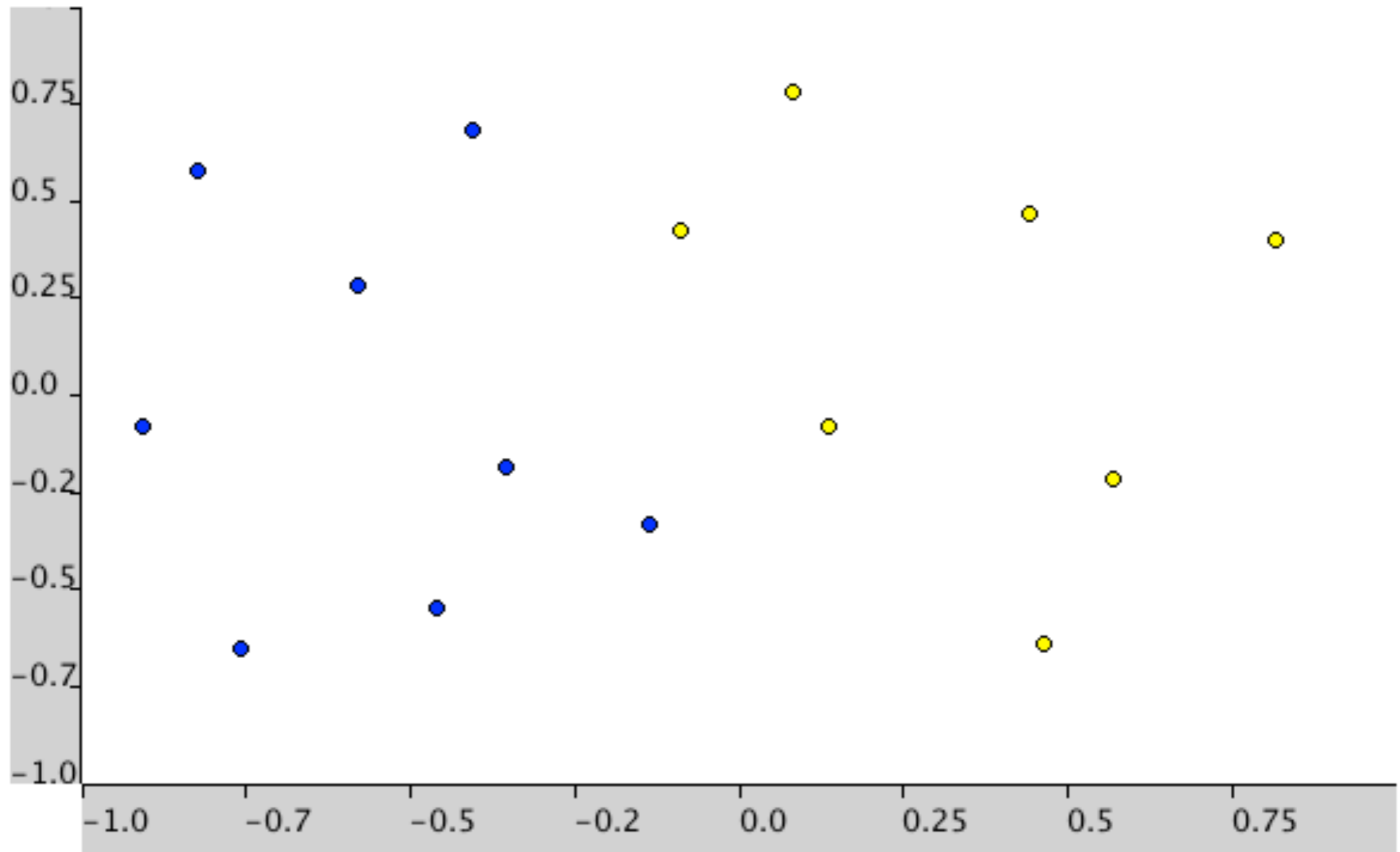
Properties

$$w = \sum_i y_i \alpha_i x_i$$



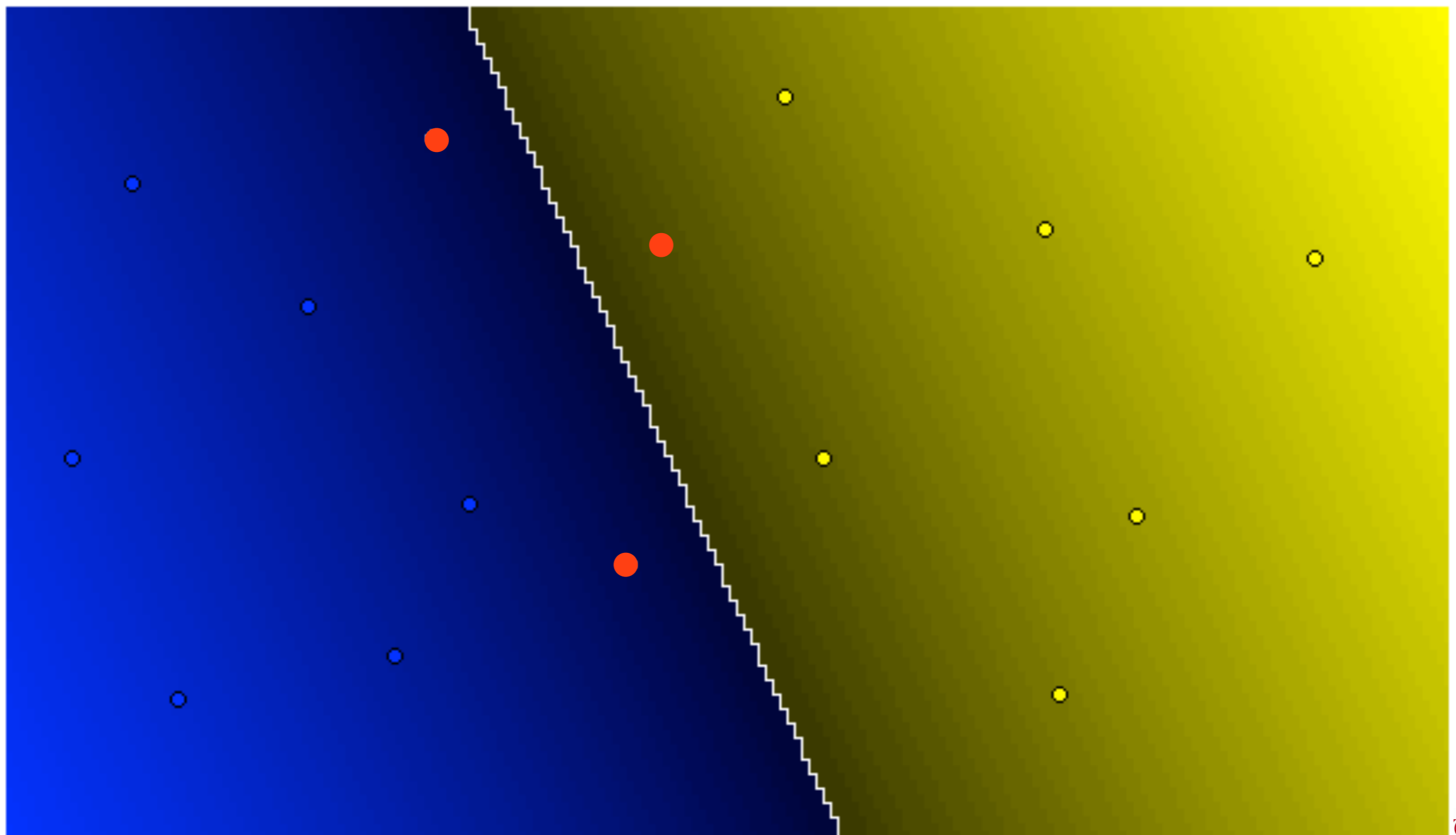
- Weight vector w as weighted linear combination of instances
- Only points on margin matter (ignore the rest and get same solution)
- Only inner products matter
 - Quadratic program
 - We can replace the inner product by a kernel
- Keeps instances away from the margin

Example

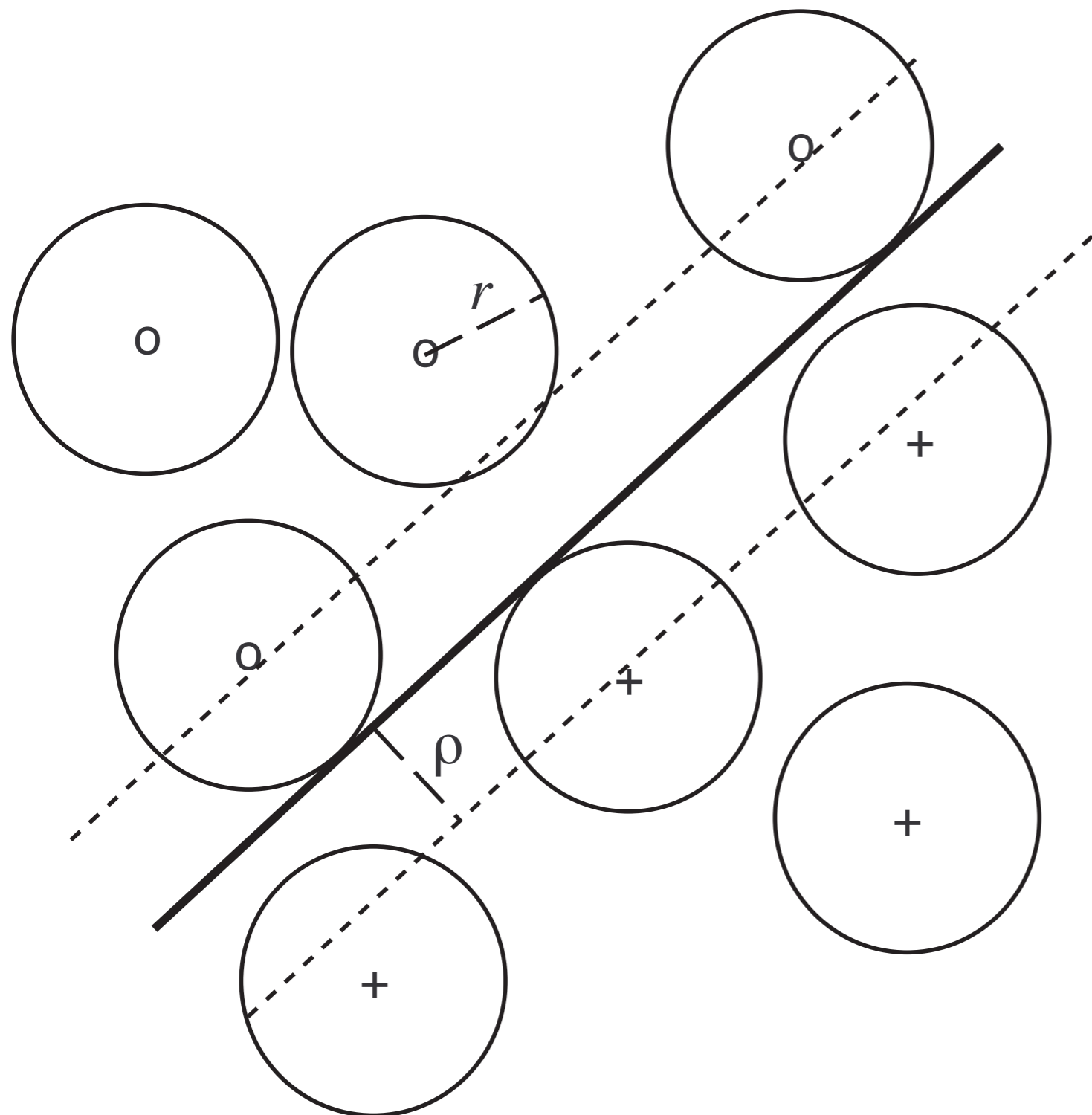


Example

Number of Support Vectors: **3** (-ve: 2, +ve: 1) Total number of points: 15



Why large margins?



- **Maximum robustness relative to uncertainty**
- **Symmetry breaking**
- **Independent of correctly classified instances**
- **Easy to find for easy problems**



MAGIC Etch A Sketch[®] SCREEN



CLASSIFIERS

Horizontal
Grid

OHIO ART The World of Toys[®]

Mexico
1963

MAGIC SCREEN IS GLASS SET IN DURABLE PLASTIC FRAME
USE WITH CARE

Leaderboard

10-701: Machine Learning (f13)

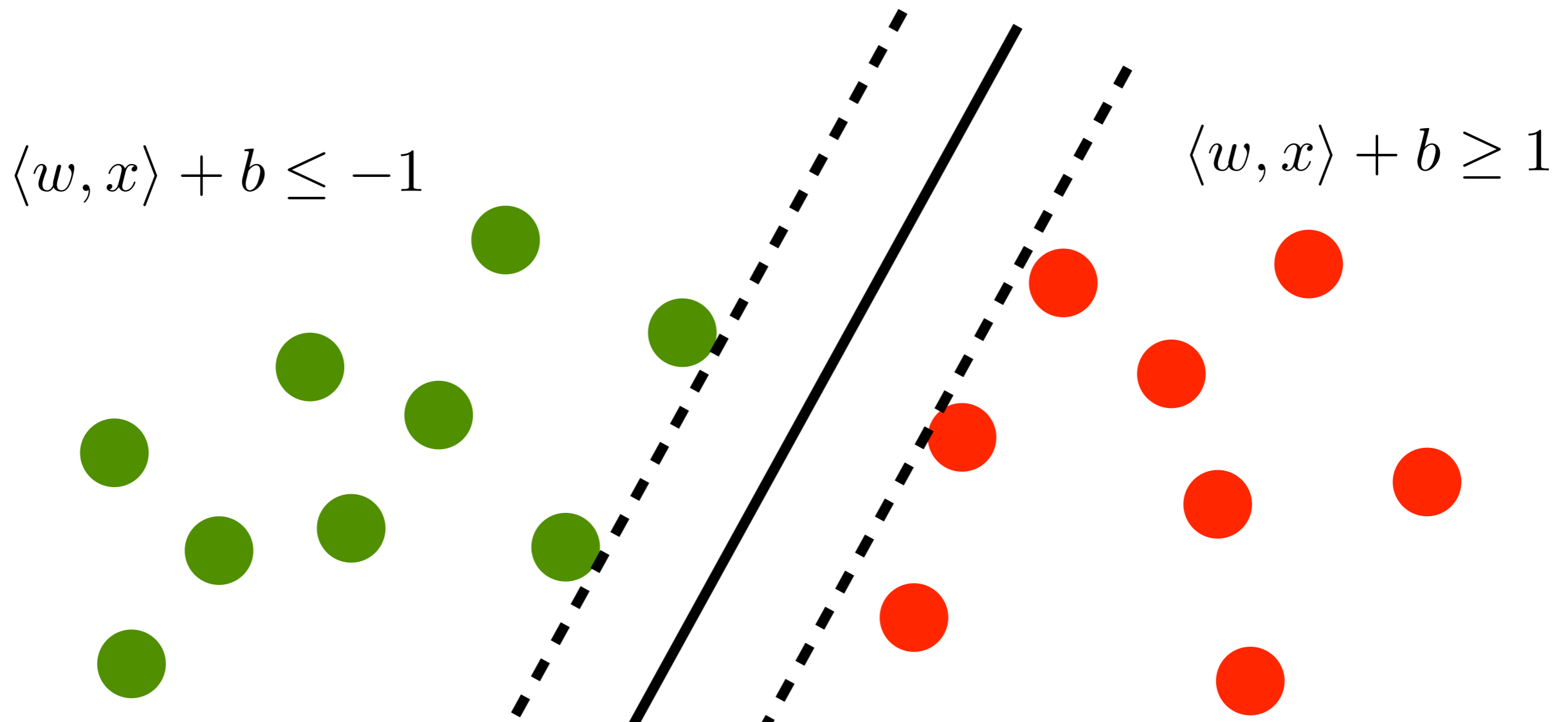
[Home](#) | [Gradebook](#) | [Account](#) | [Jobs](#) | [Admin](#)

Class Scoreboard for homework2

SPAM Classification Contest!

0	NICKNAME	VERSION	TIME	CLASSIFICATION
1	cywu	18	2013-10-14 09:55:54	93.2
2	AndHobbes	13	2013-10-14 19:28:24	90.1
3	pxie	17	2013-10-15 10:22:30	90.8
4	lélouch	12	2013-10-13 22:27:24	89
5	YHK	26	2013-10-13 21:26:59	88
6	teach smola ML	24	2013-10-15 14:09:38	88
7	DD	9	2013-10-14 09:15:16	86.9
8	cmalings	27	2013-10-15 09:06:15	76.8
9	dontworry	18	2013-10-15 16:50:19	74.3
10	tcarlone	17	2013-10-15 21:53:08	71.3

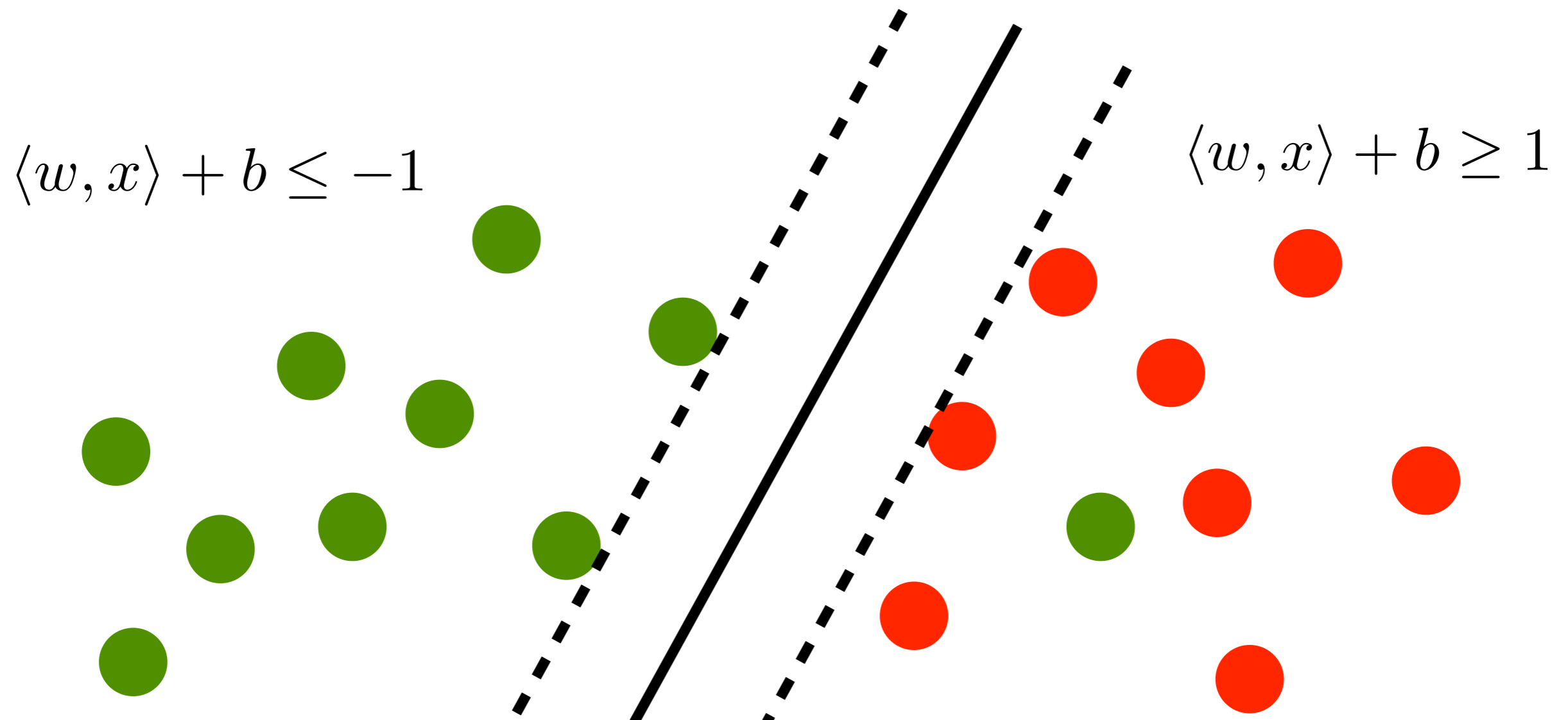
Large Margin Classifier



linear function

$$f(x) = \langle w, x \rangle + b$$

Large Margin Classifier

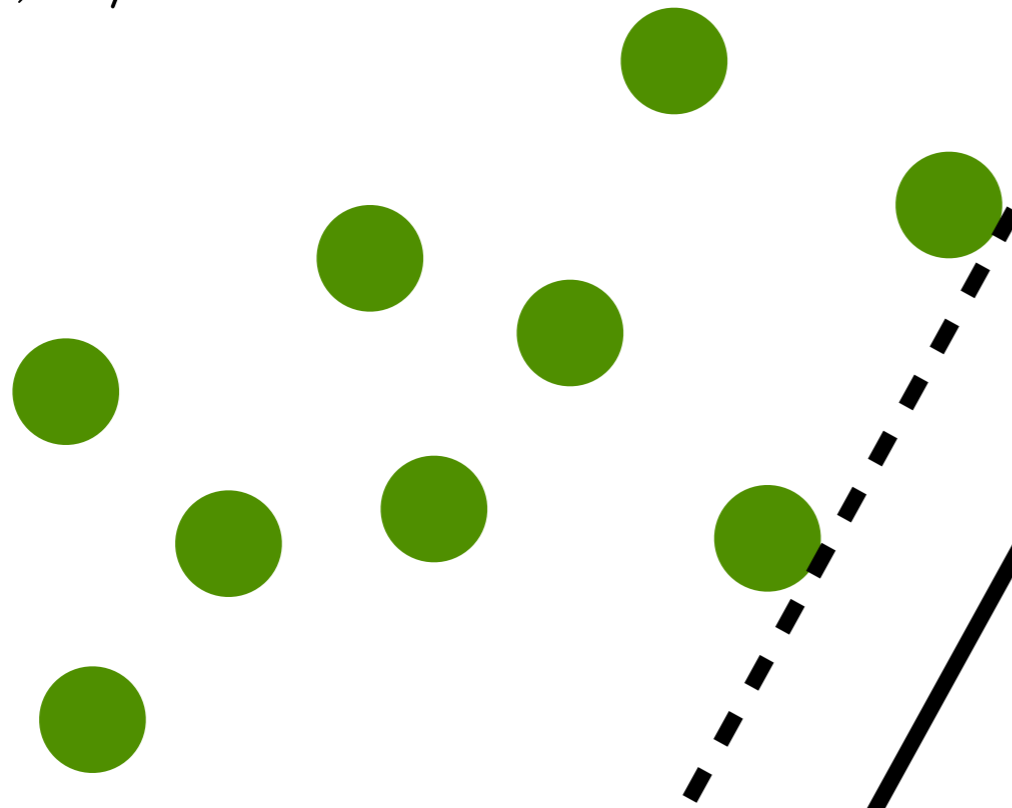


linear function

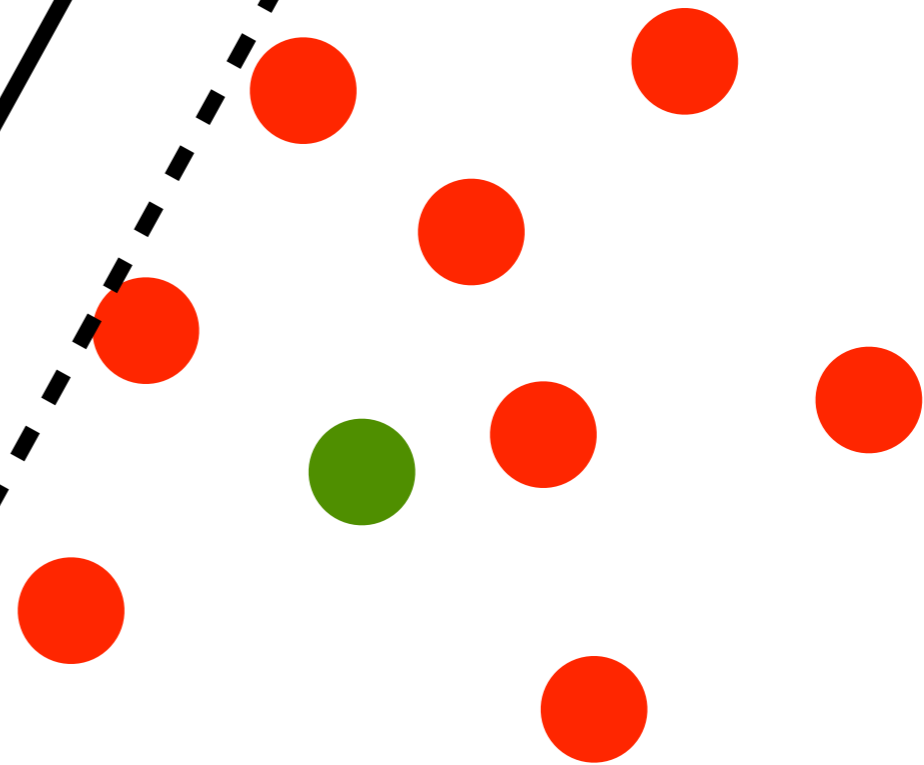
$$f(x) = \langle w, x \rangle + b$$

Large Margin Classifier

$$\langle w, x \rangle + b \leq -1$$



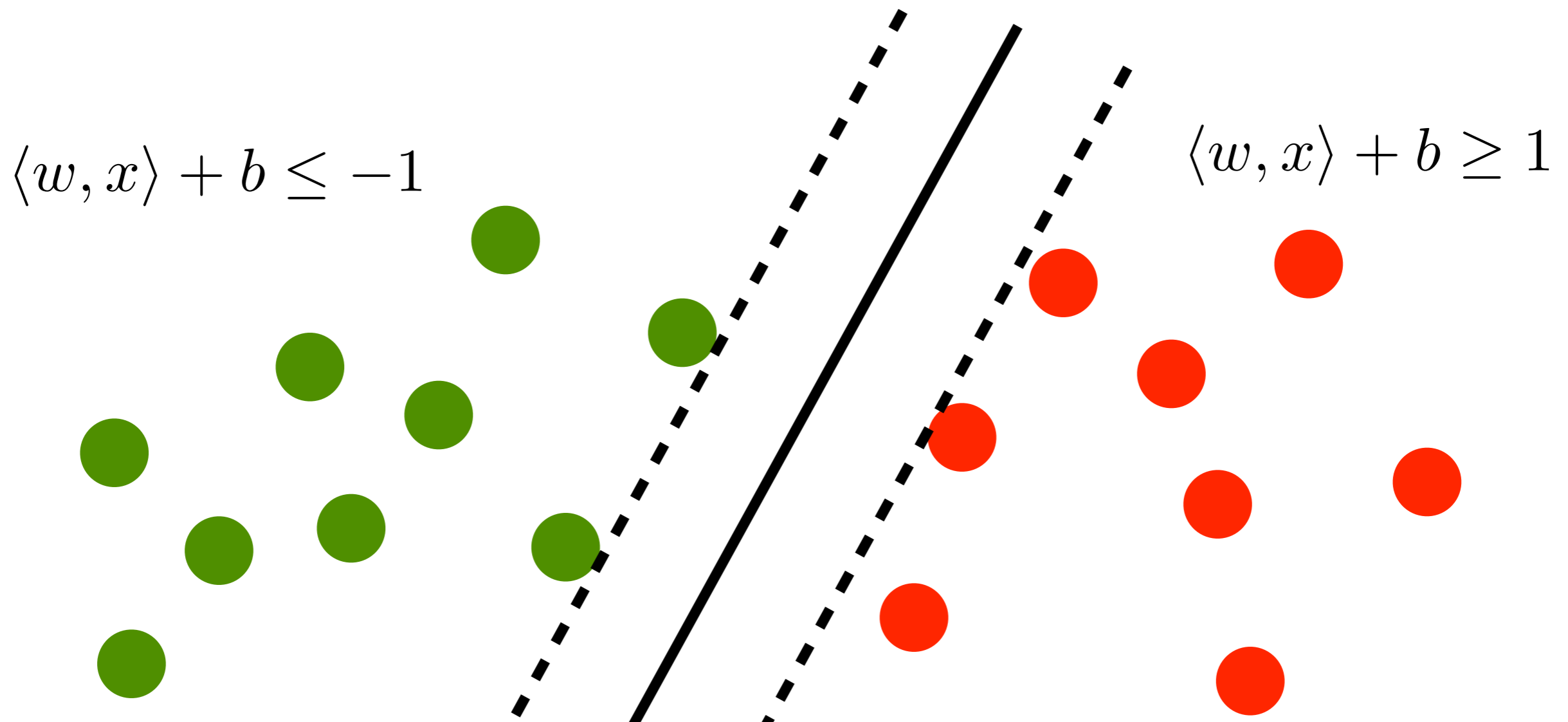
$$\langle w, x \rangle + b \geq 1$$



linear function
 $f(x) = \langle w, x \rangle + b$

**linear separator
is impossible**

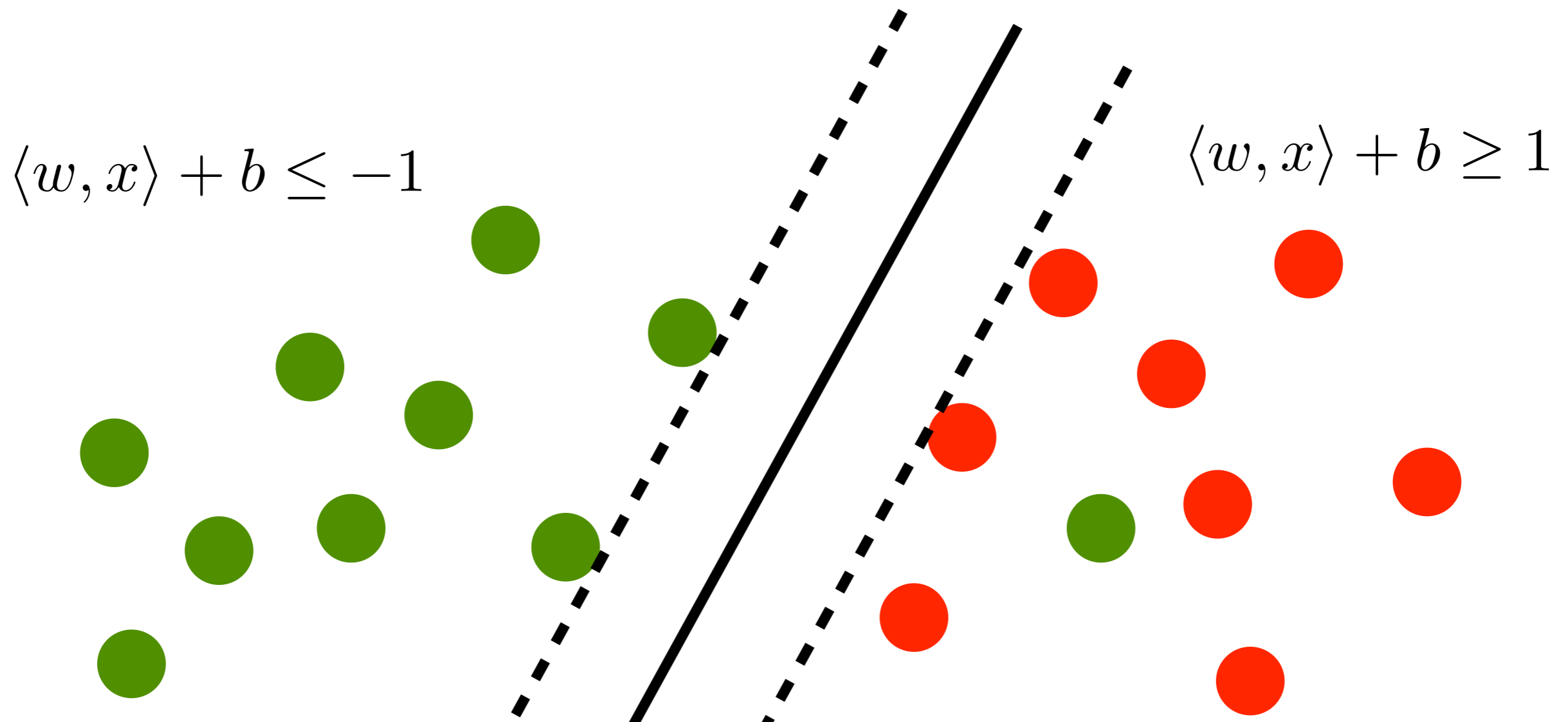
Large Margin Classifier



Theorem (Minsky & Papert)

Finding the minimum error separating hyperplane is NP hard

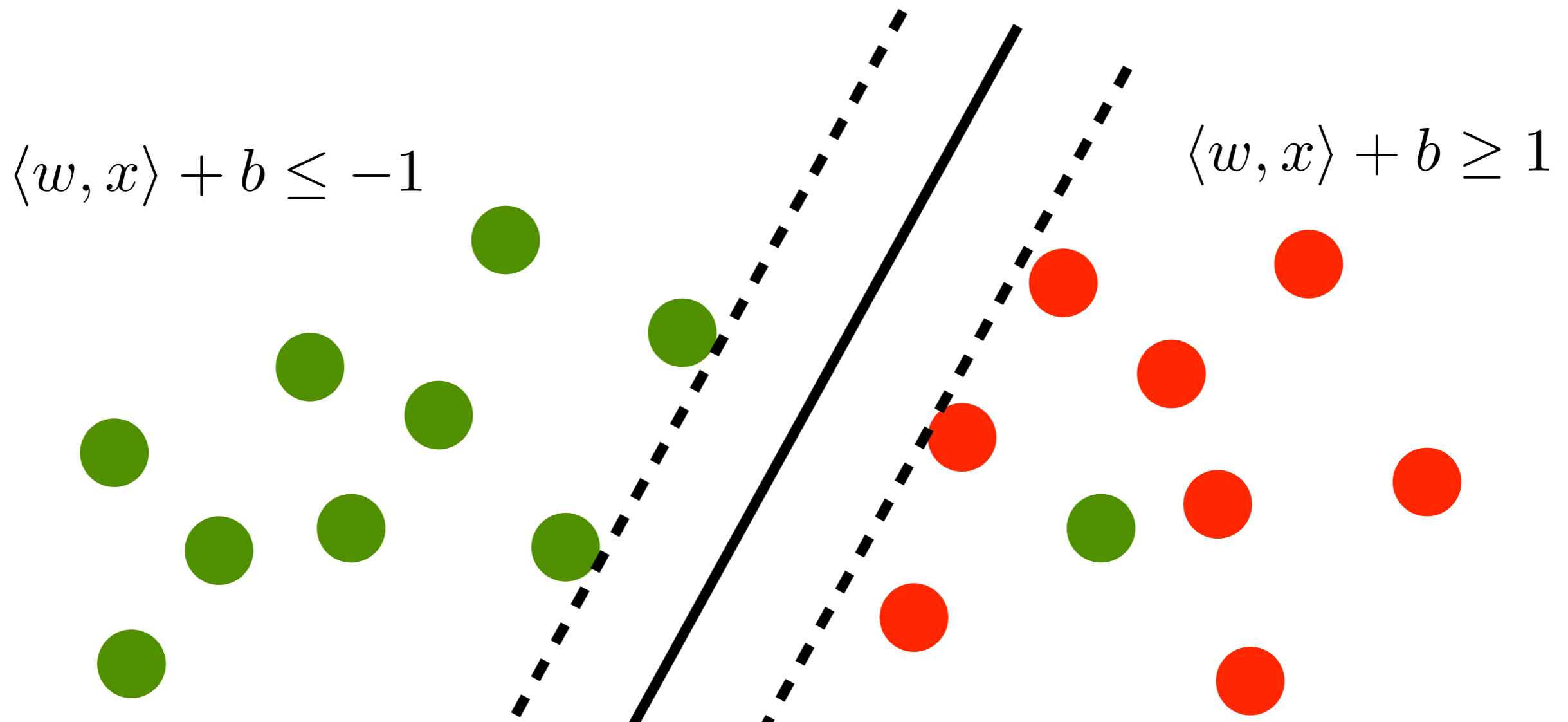
Large Margin Classifier



Theorem (Minsky & Papert)

Finding the minimum error separating hyperplane is NP hard

Large Margin Classifier



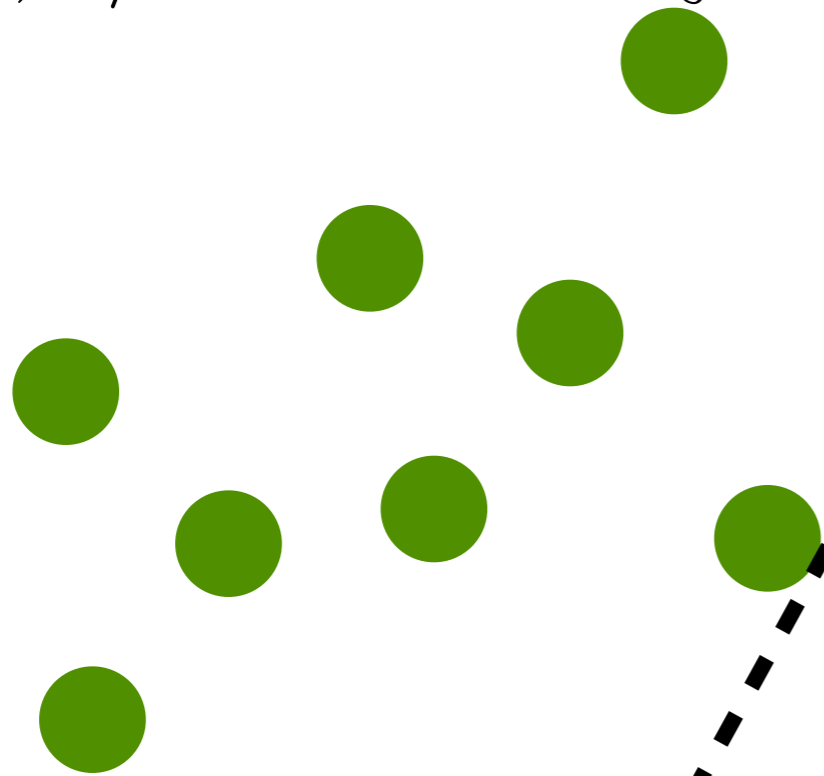
**minimum error separator
is impossible**

Theorem (Minsky & Papert)

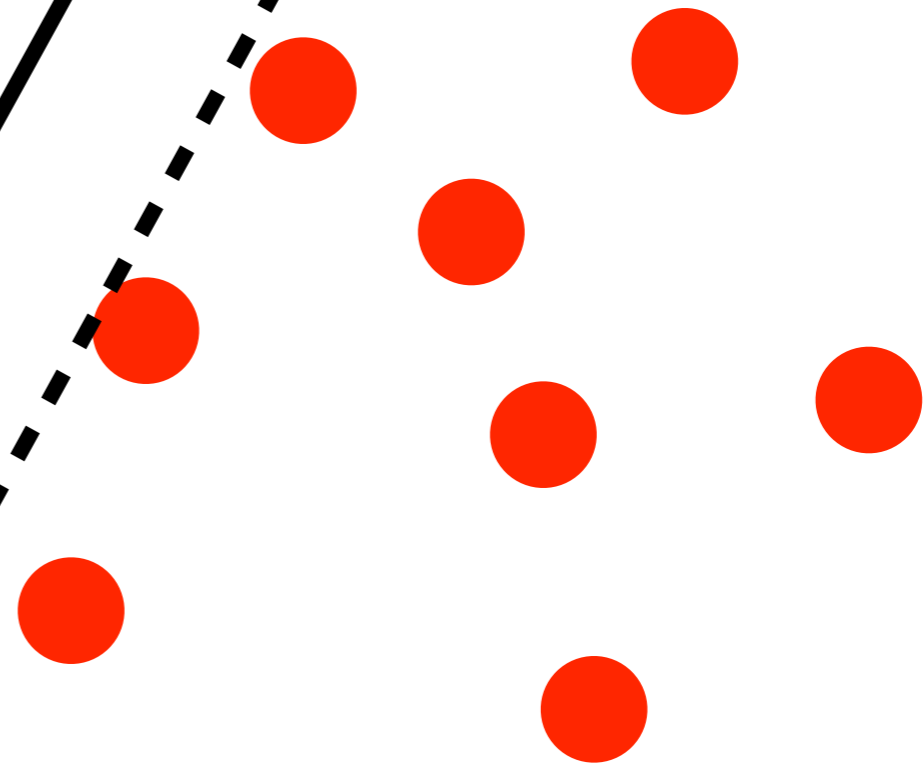
Finding the minimum error separating hyperplane is NP hard

Adding slack variables

$$\langle w, x \rangle + b \leq -1 + \xi$$



$$\langle w, x \rangle + b \geq 1 - \xi$$



Convex optimization problem

Adding slack variables

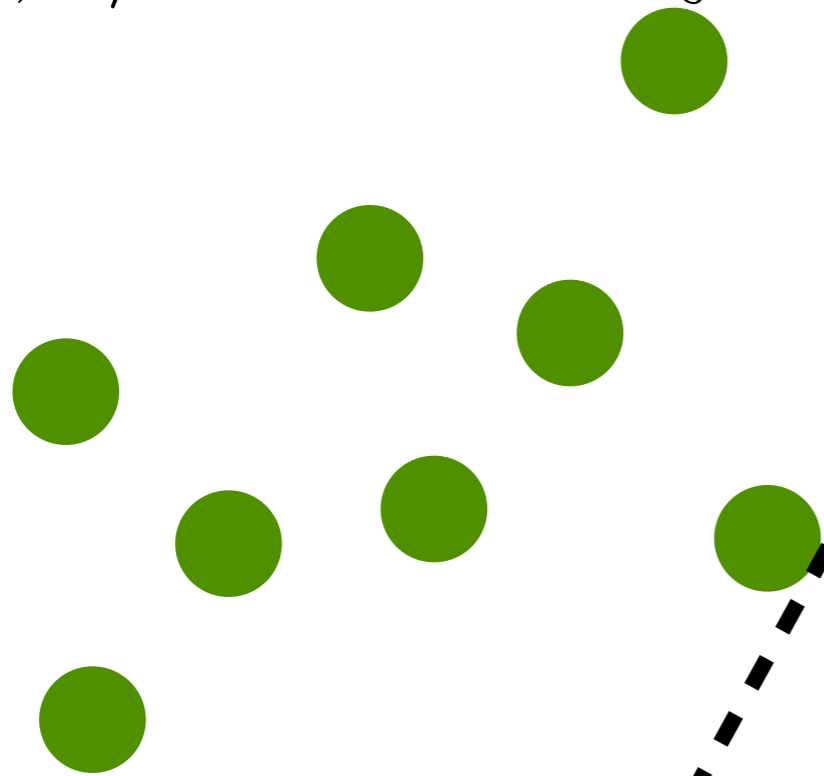
$$\langle w, x \rangle + b \leq -1 + \xi$$

$$\langle w, x \rangle + b \geq 1 - \xi$$

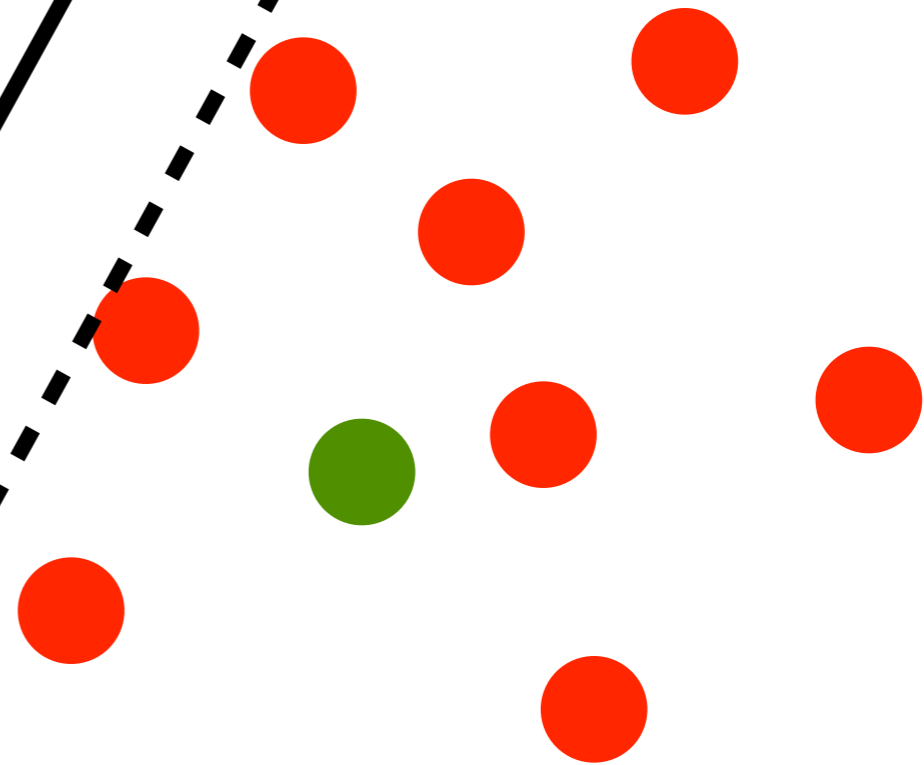
Convex optimization problem

Adding slack variables

$$\langle w, x \rangle + b \leq -1 + \xi$$



$$\langle w, x \rangle + b \geq 1 - \xi$$



Convex optimization problem

minimize amount
of slack

Intermezzo

Convex Programs for Dummies

- **Primal optimization problem**

$$\underset{x}{\text{minimize}} f(x) \text{ subject to } c_i(x) \leq 0$$

- **Lagrange function**

$$L(x, \alpha) = f(x) + \sum_i \alpha_i c_i(x)$$

- **First order optimality conditions in x**

$$\partial_x L(x, \alpha) = \partial_x f(x) + \sum_i \alpha_i \partial_x c_i(x) = 0$$

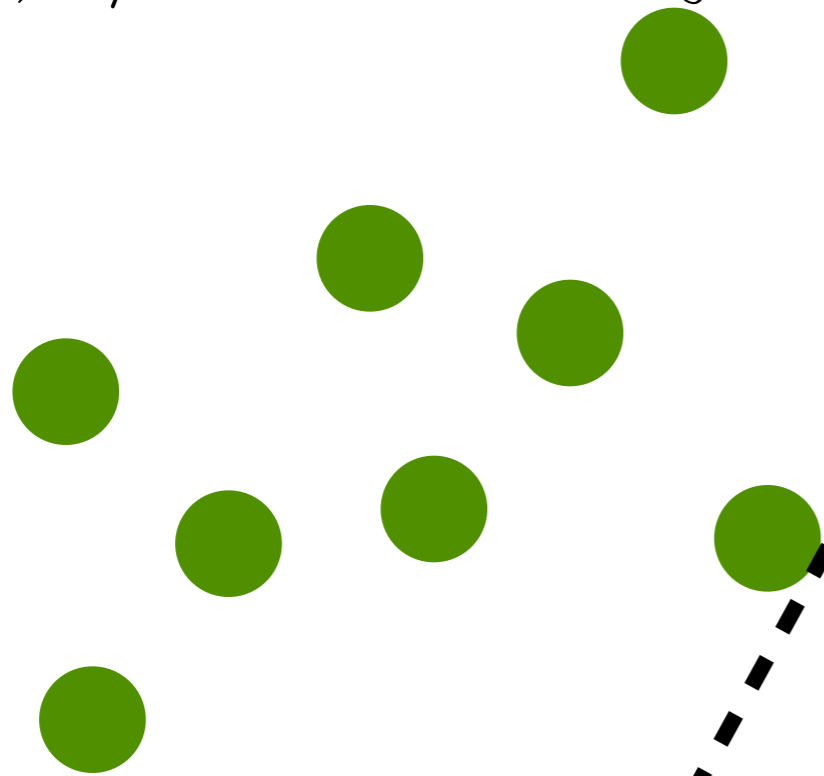
- **Solve for x and plug it back into L**

$$\underset{\alpha}{\text{maximize}} L(x(\alpha), \alpha)$$

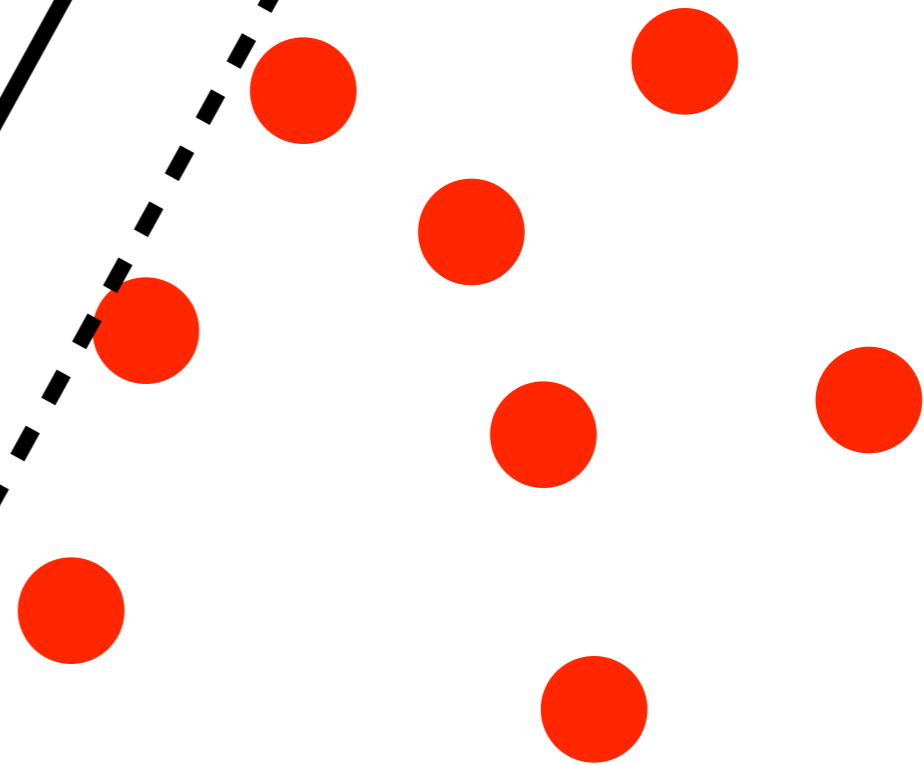
(keep explicit constraints)

Adding slack variables

$$\langle w, x \rangle + b \leq -1 + \xi$$



$$\langle w, x \rangle + b \geq 1 - \xi$$



Convex optimization problem

Adding slack variables

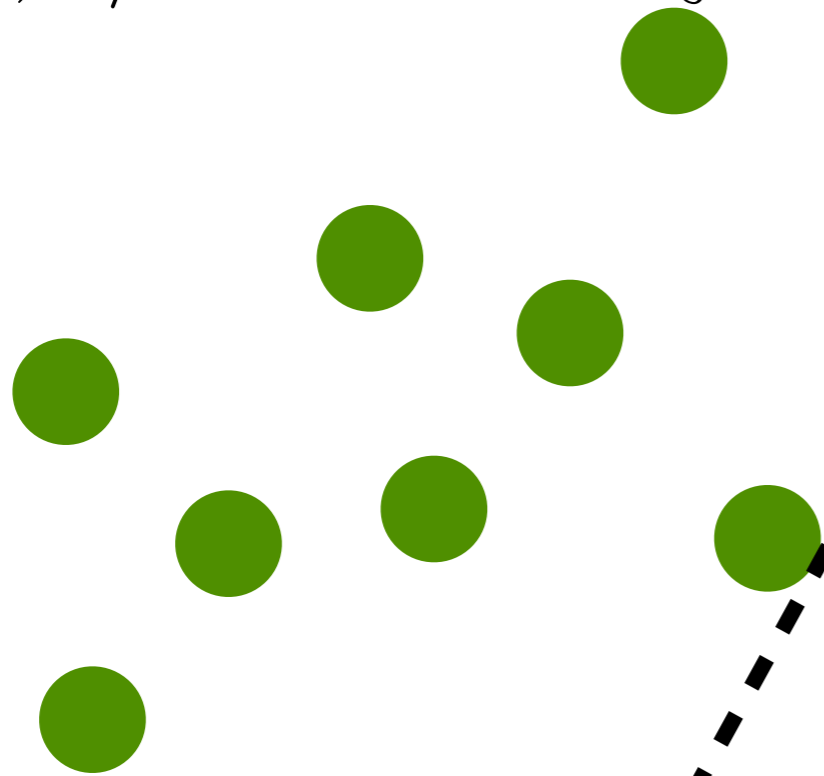
$$\langle w, x \rangle + b \leq -1 + \xi$$

$$\langle w, x \rangle + b \geq 1 - \xi$$

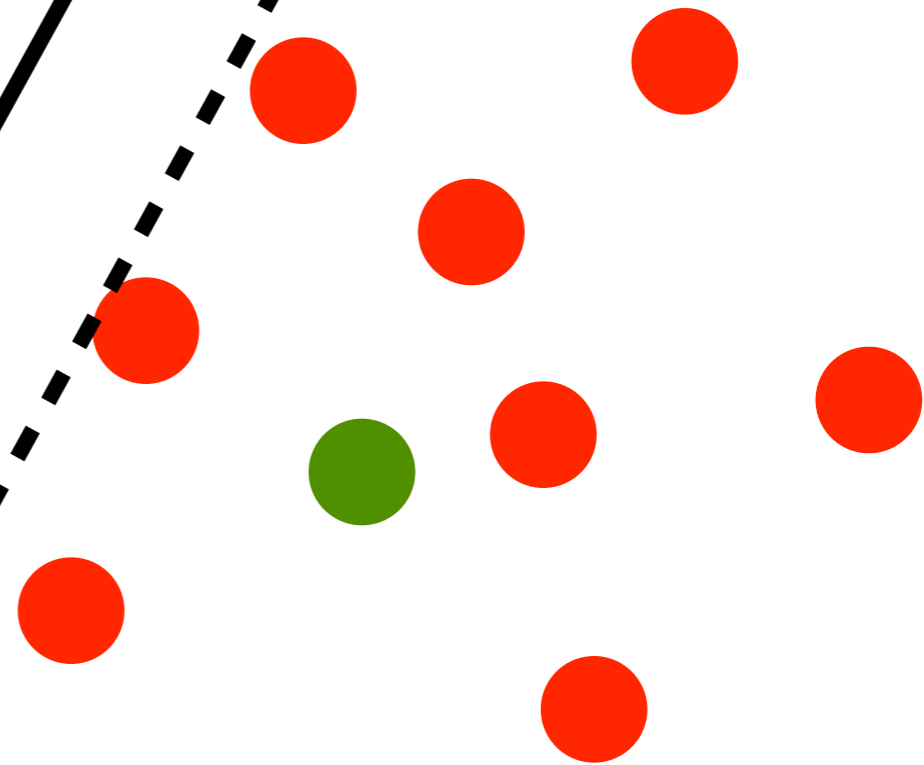
Convex optimization problem

Adding slack variables

$$\langle w, x \rangle + b \leq -1 + \xi$$



$$\langle w, x \rangle + b \geq 1 - \xi$$



Convex optimization problem

minimize amount
of slack

Adding slack variables

- **Hard margin problem**

$$\underset{w,b}{\text{minimize}} \frac{1}{2} \|w\|^2 \quad \text{subject to } y_i [\langle w, x_i \rangle + b] \geq 1$$

- **With slack variables**

$$\underset{w,b}{\text{minimize}} \frac{1}{2} \|w\|^2 + C \sum_i \xi_i$$

$$\text{subject to } y_i [\langle w, x_i \rangle + b] \geq 1 - \xi_i \text{ and } \xi_i \geq 0$$

Problem is always feasible. Proof:

$w = 0$ and $b = 0$ and $\xi_i = 1$ (also yields upper bound)

Dual Problem

- Primal optimization problem

$$\underset{w,b}{\text{minimize}} \quad \frac{1}{2} \|w\|^2 + C \sum_i \xi_i$$

subject to $y_i [\langle w, x_i \rangle + b] \geq 1 - \xi_i$ and $\xi_i \geq 0$

- Lagrange function

$$L(w, b, \alpha) = \frac{1}{2} \|w\|^2 + C \sum_i \xi_i - \sum_i \alpha_i [y_i [\langle x_i, w \rangle + b] + \xi_i - 1] - \sum_i \eta_i \xi_i$$

Optimality in w, b, ξ is at saddle point with α, η

- Derivatives in w, b, ξ need to vanish

Dual Problem

- **Lagrange function**

$$L(w, b, \alpha) = \frac{1}{2} \|w\|^2 + C \sum_i \xi_i - \sum_i \alpha_i [y_i [\langle x_i, w \rangle + b] + \xi_i - 1] - \sum_i \eta_i \xi_i$$

- **Derivatives in w , b need to vanish**

$$\partial_w L(w, b, \xi, \alpha, \eta) = w - \sum_i \alpha_i y_i x_i = 0$$

$$\partial_b L(w, b, \xi, \alpha, \eta) = \sum_i \alpha_i y_i = 0$$

$$\partial_{\xi_i} L(w, b, \xi, \alpha, \eta) = C - \alpha_i - \eta_i = 0$$

- **Plugging terms back into L yields**

$$\text{maximize}_{\alpha} - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \langle x_i, x_j \rangle + \sum_i \alpha_i$$

$$\text{subject to } \sum_i \alpha_i y_i = 0 \text{ and } \alpha_i \in [0, C]$$

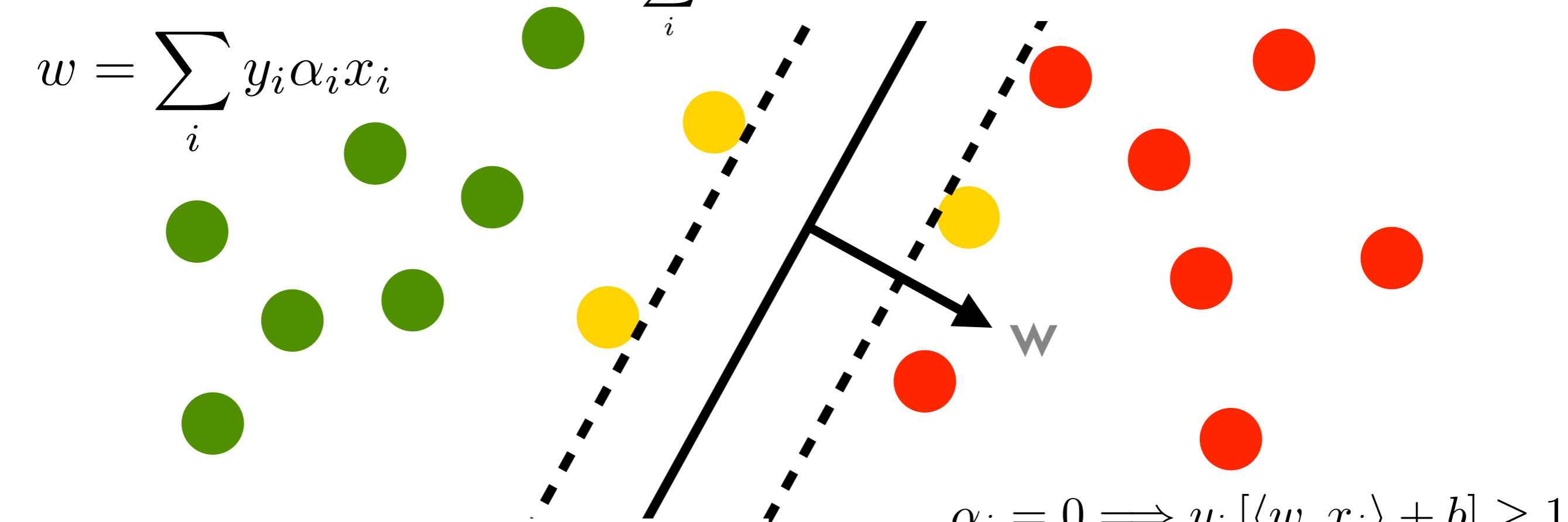
bound
influence

Karush Kuhn Tucker Conditions

$$\text{maximize}_{\alpha} -\frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \langle x_i, x_j \rangle + \sum_i \alpha_i$$

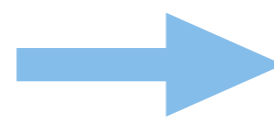
$$\text{subject to } \sum_i \alpha_i y_i = 0 \text{ and } \alpha_i \in [0, C]$$

$$w = \sum_i y_i \alpha_i x_i$$



$$\alpha_i [y_i [\langle w, x_i \rangle + b] + \xi_i - 1] = 0$$

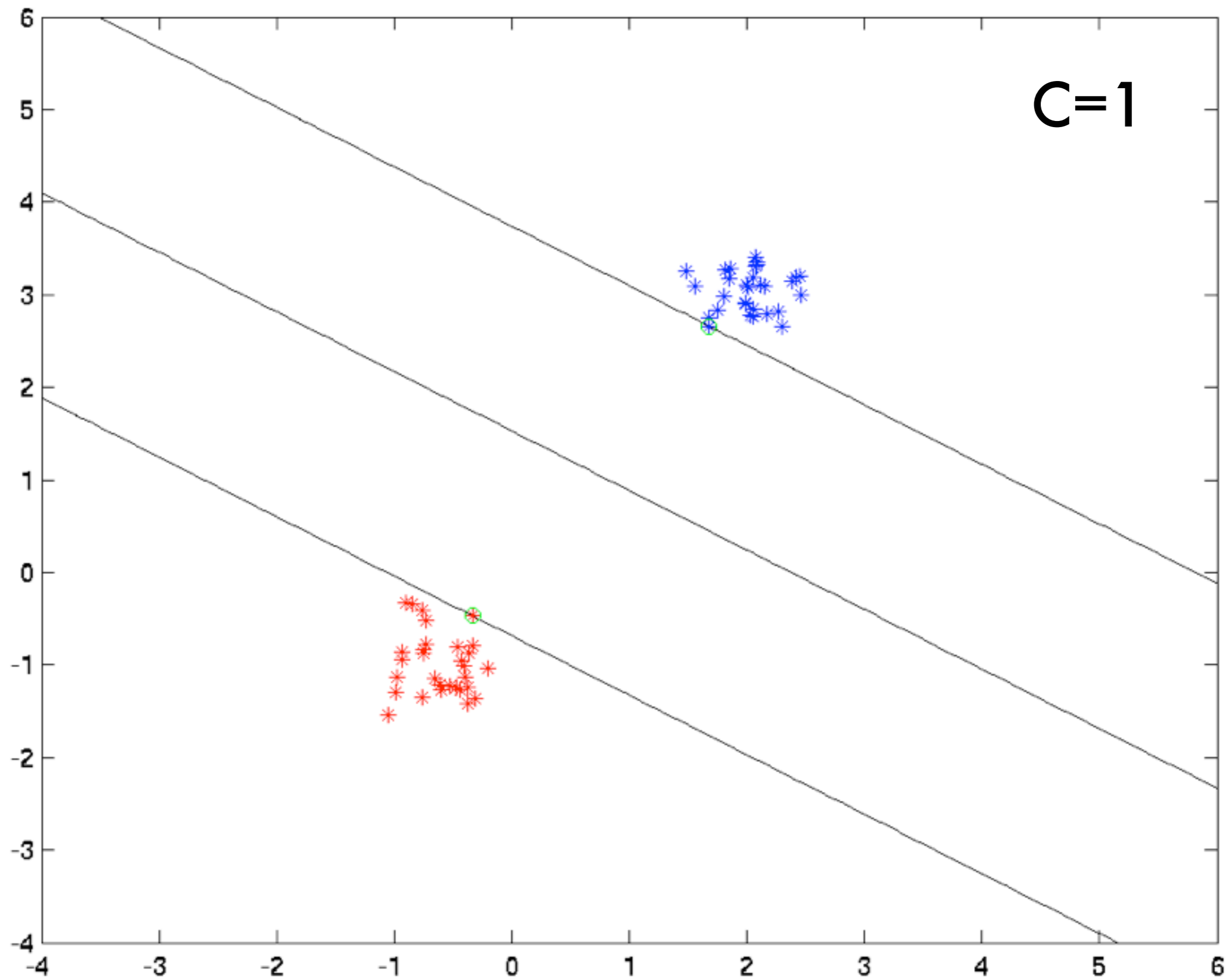
$$\eta_i \xi_i = 0$$

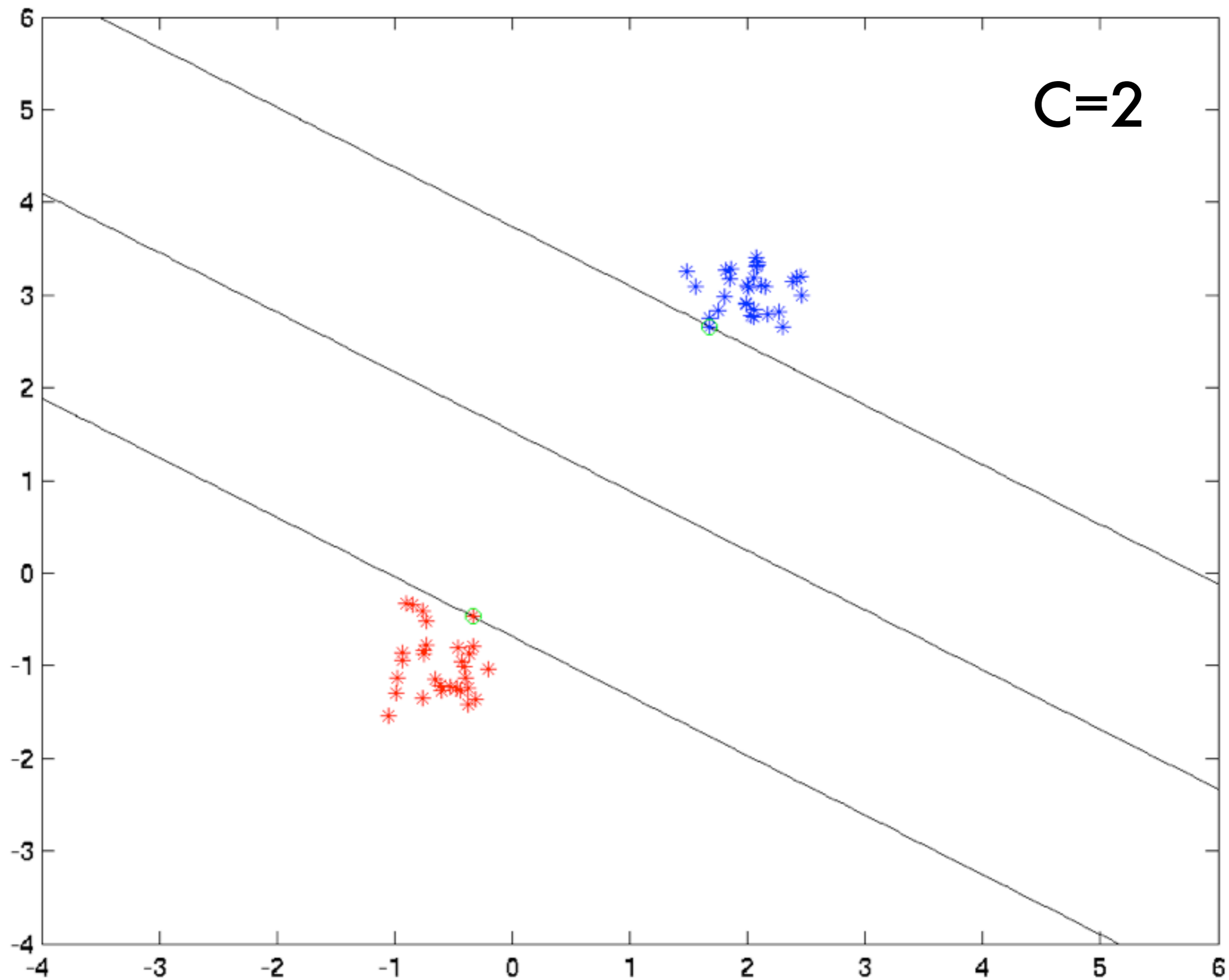


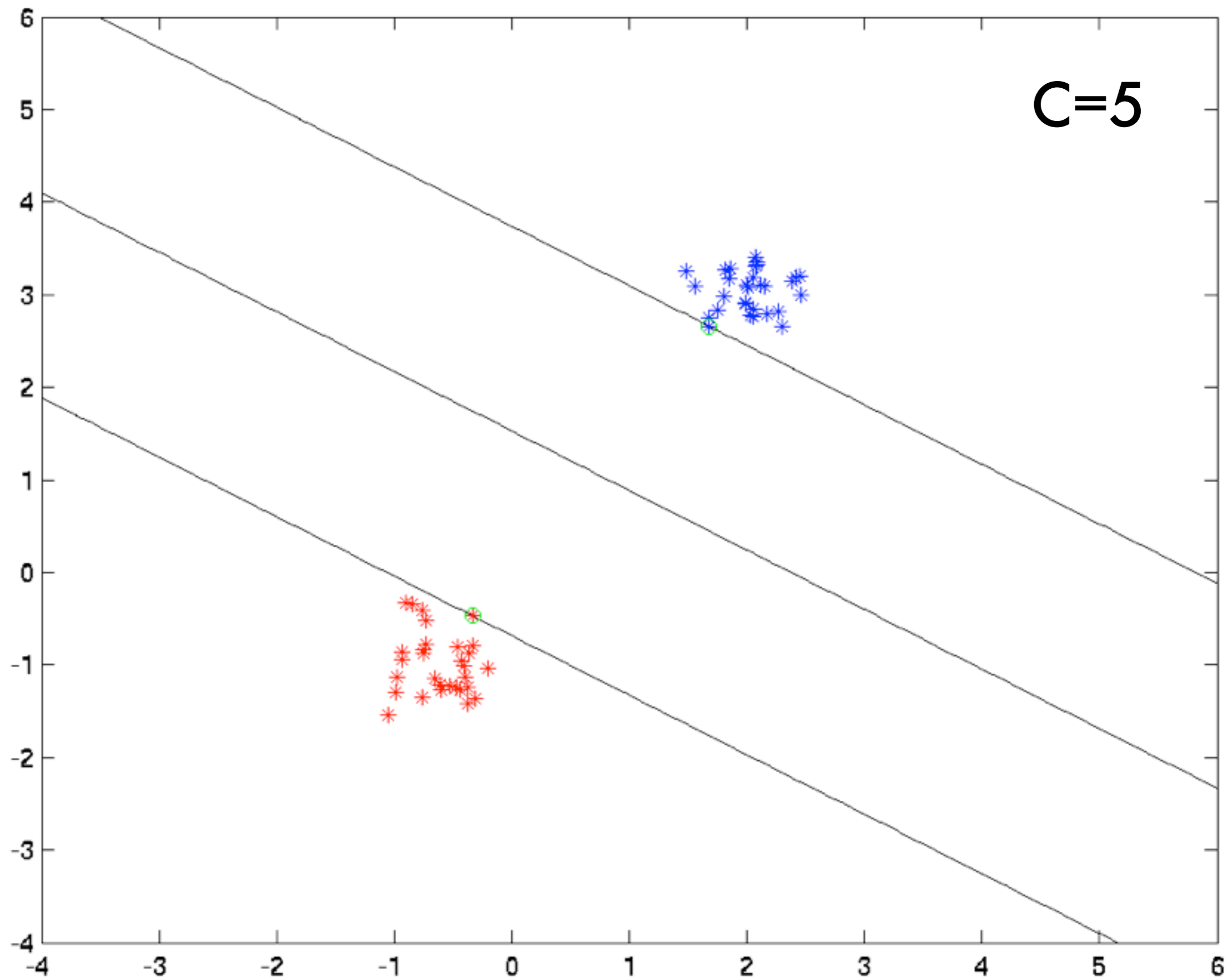
$$\alpha_i = 0 \implies y_i [\langle w, x_i \rangle + b] \geq 1$$

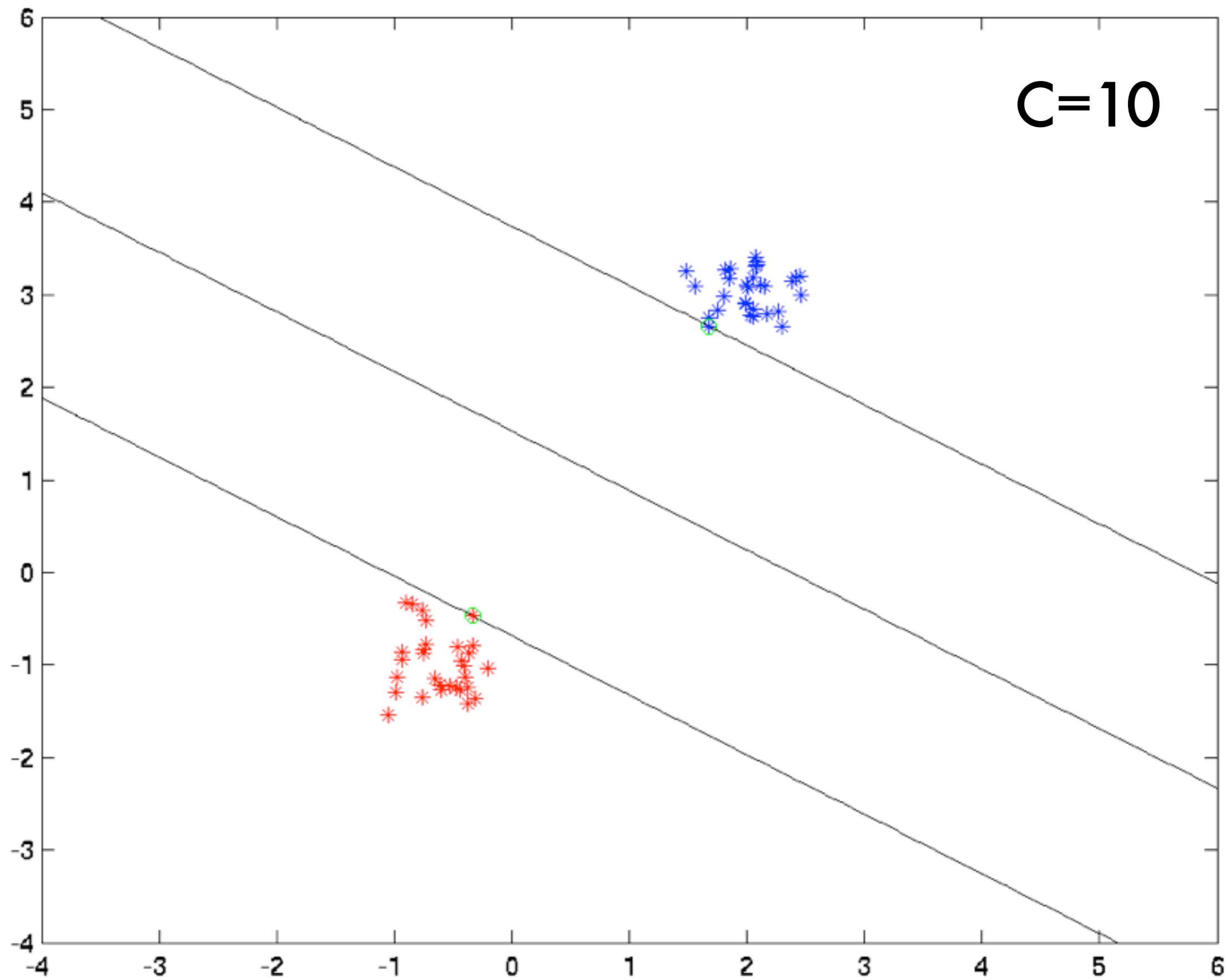
$$0 < \alpha_i < C \implies y_i [\langle w, x_i \rangle + b] = 1$$

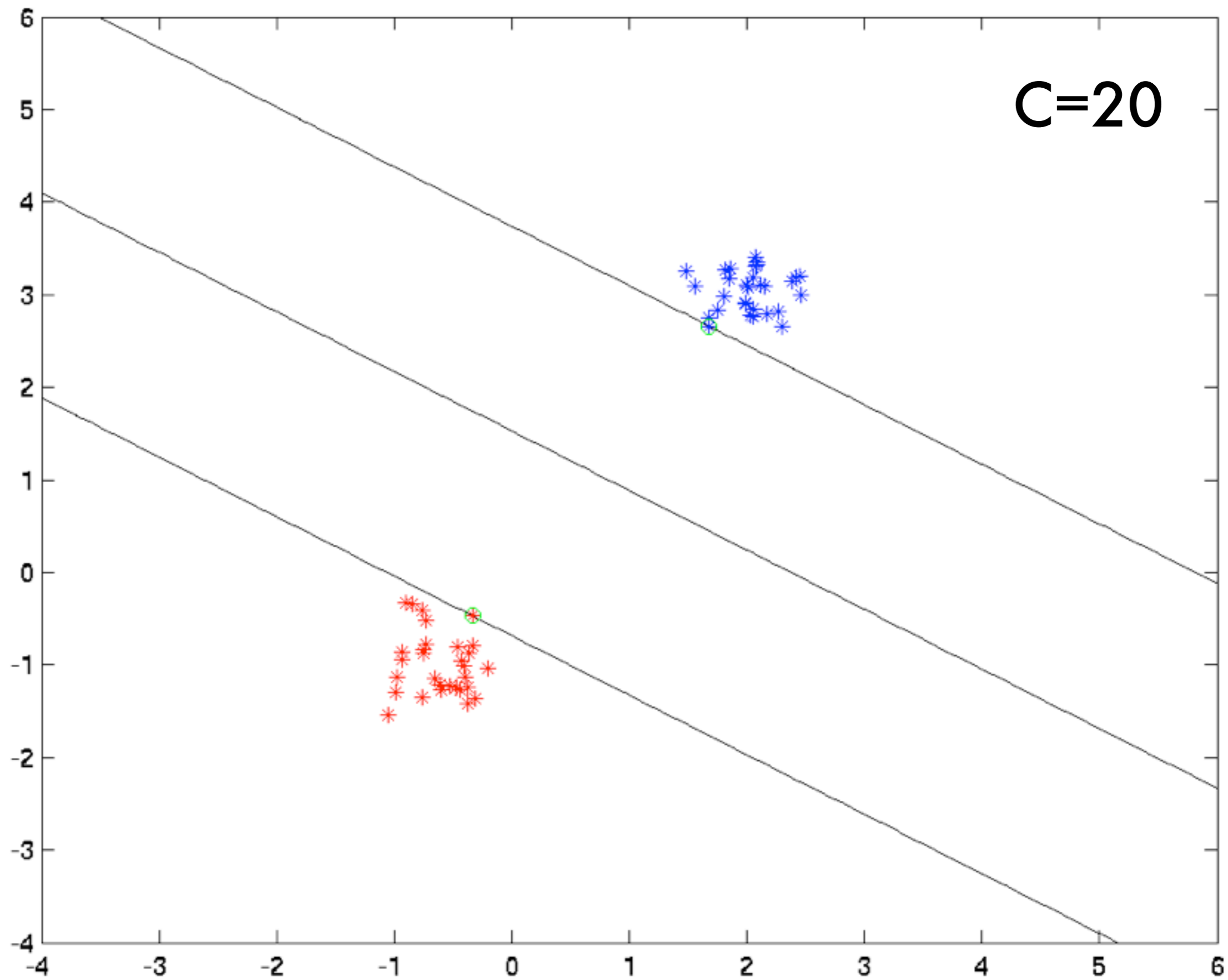
$$\alpha_i = C \implies y_i [\langle w, x_i \rangle + b] \leq 1$$

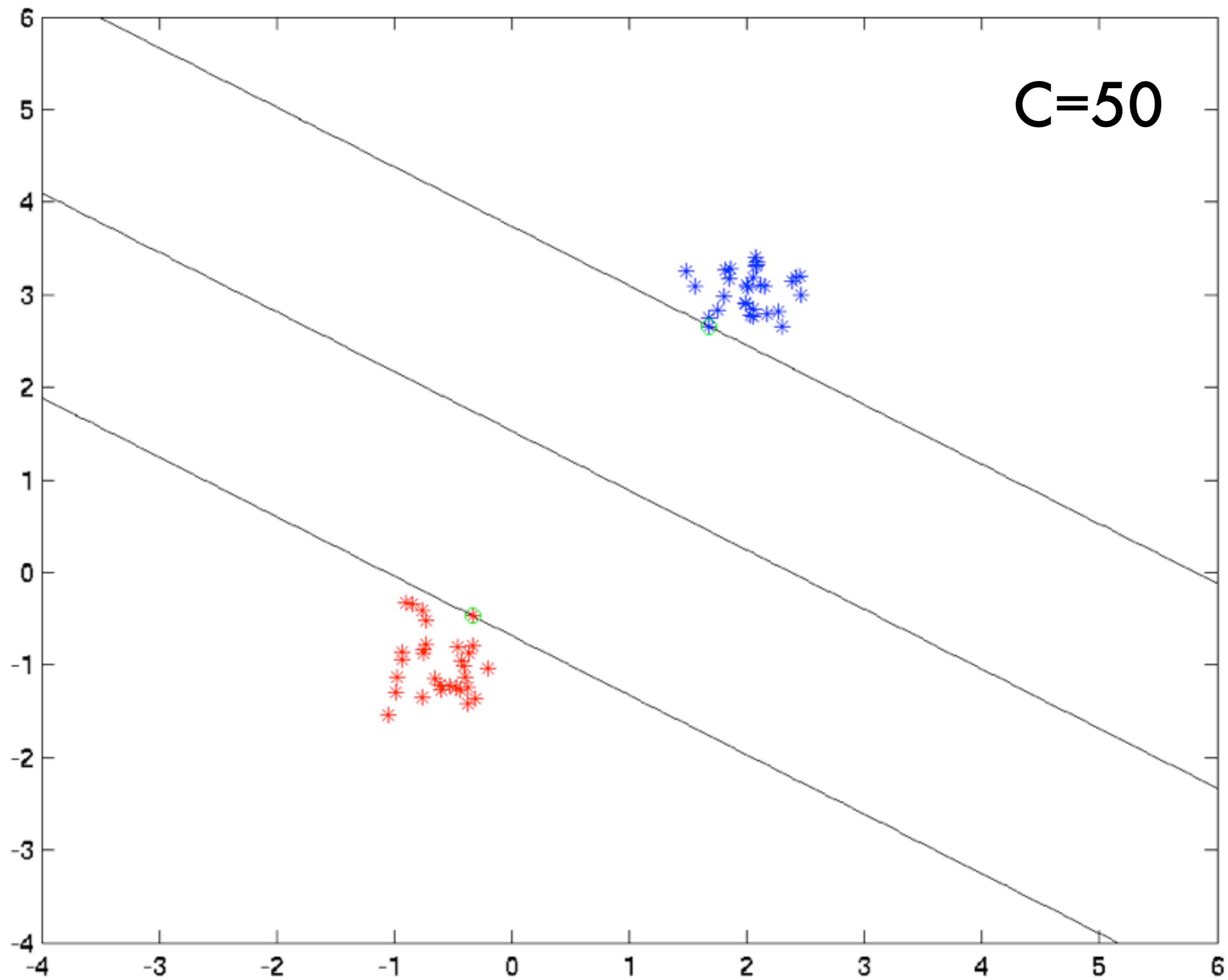


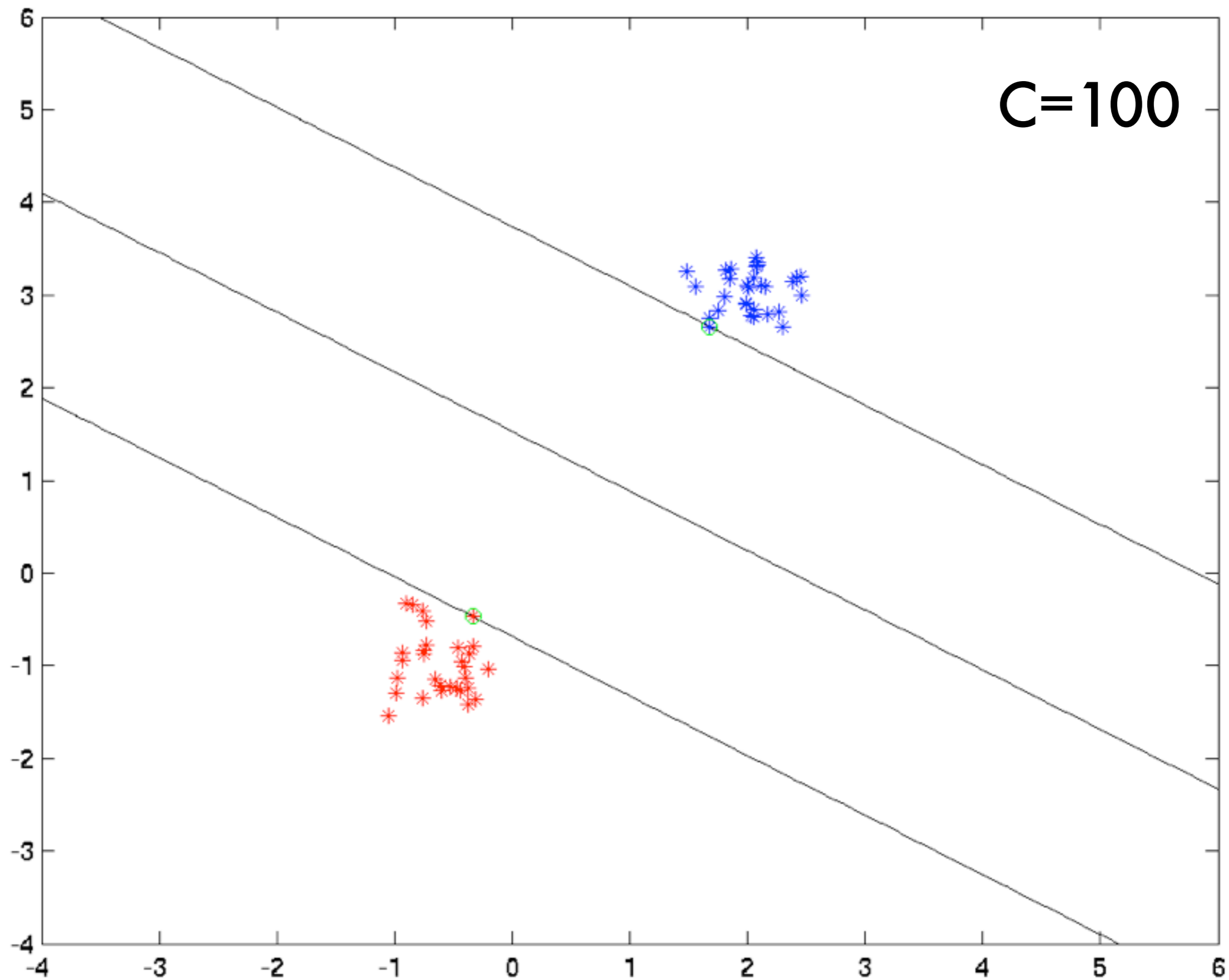


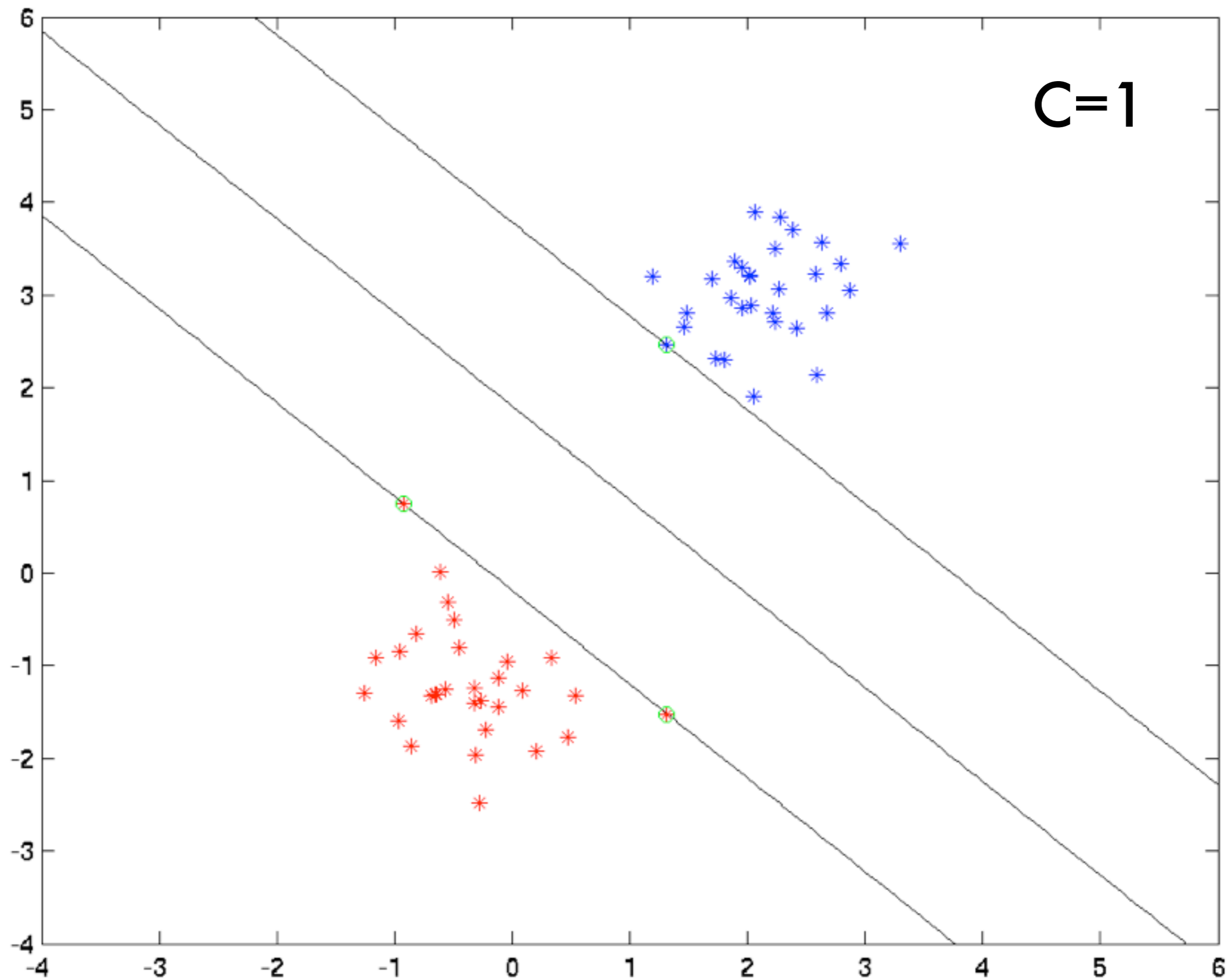


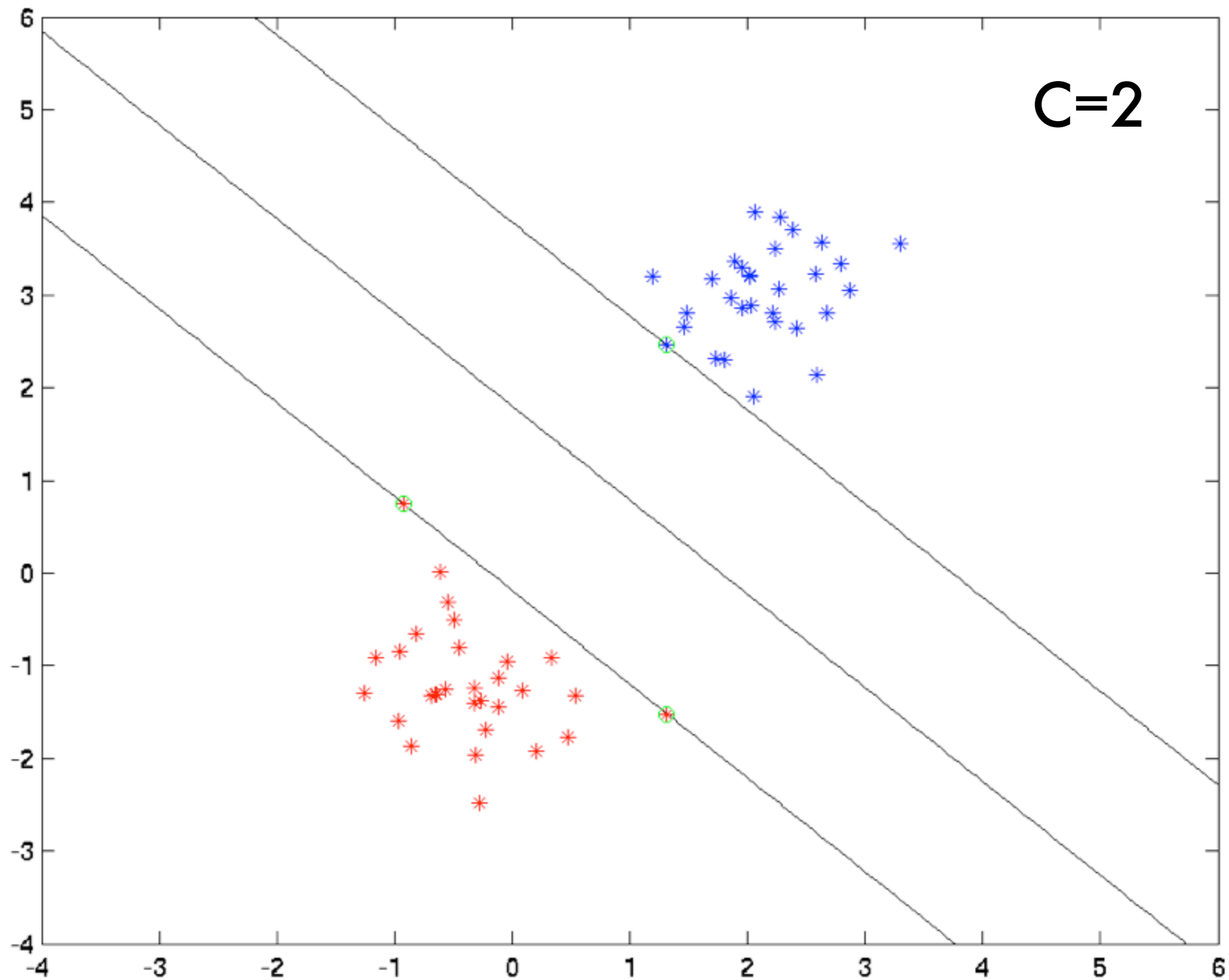


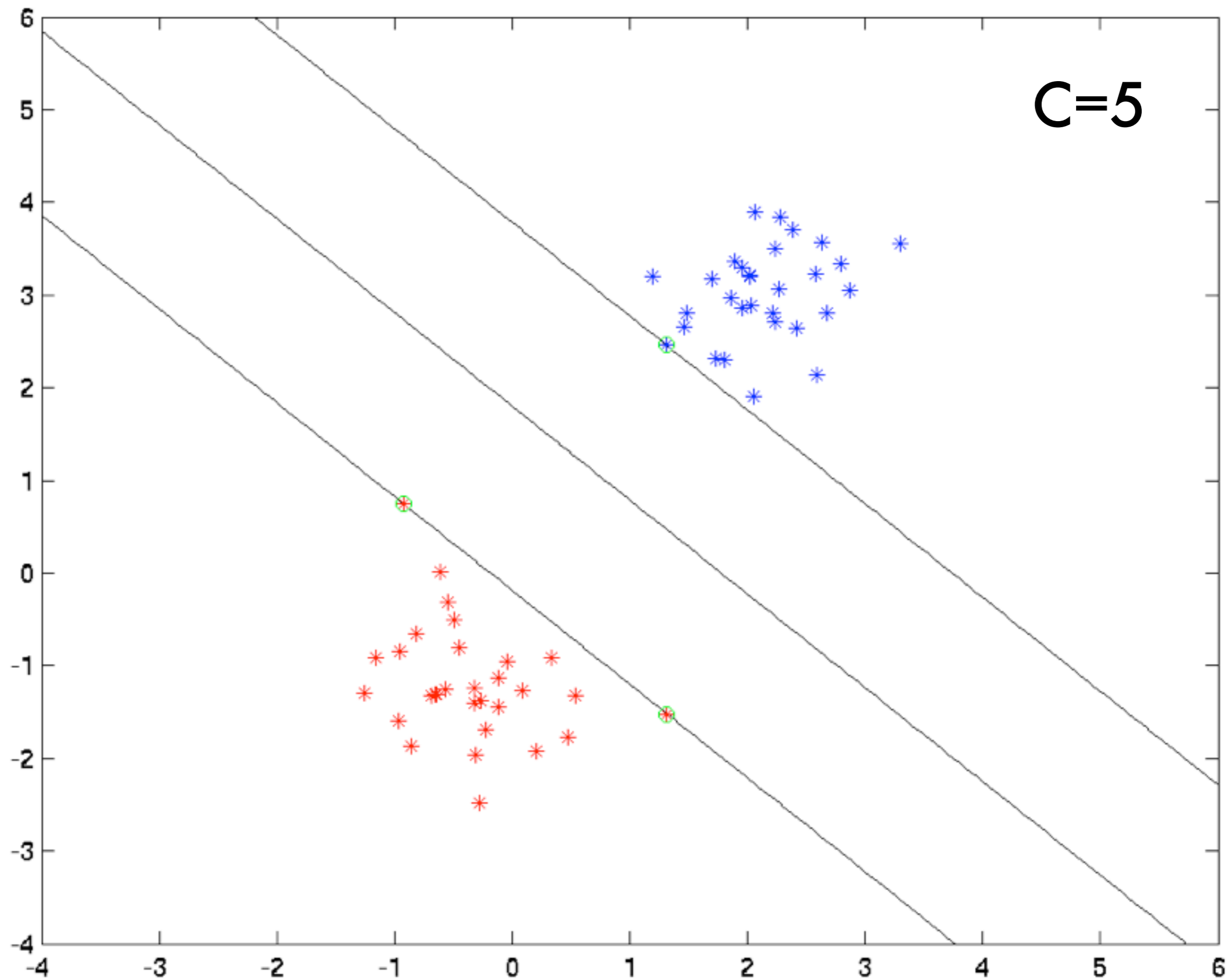


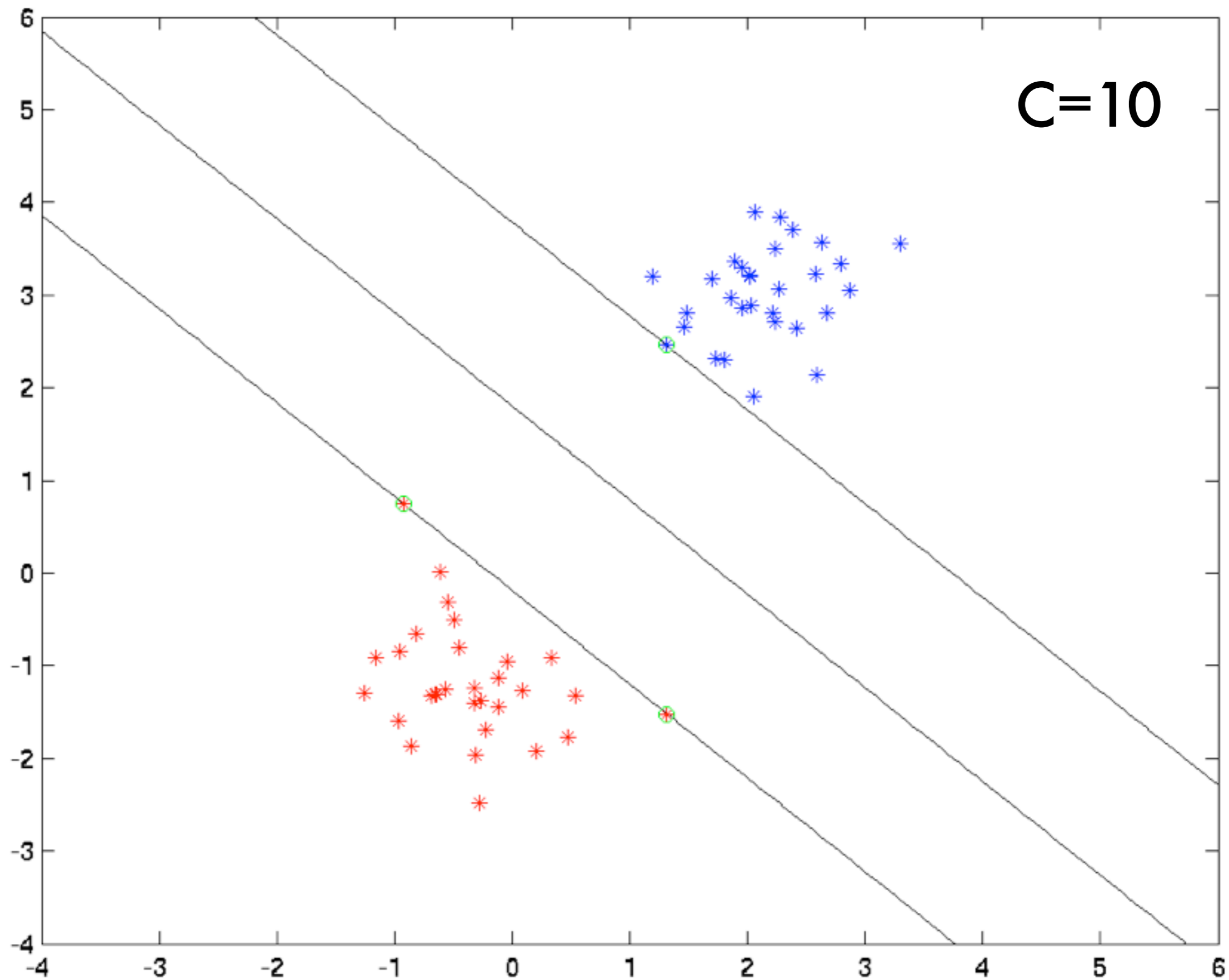


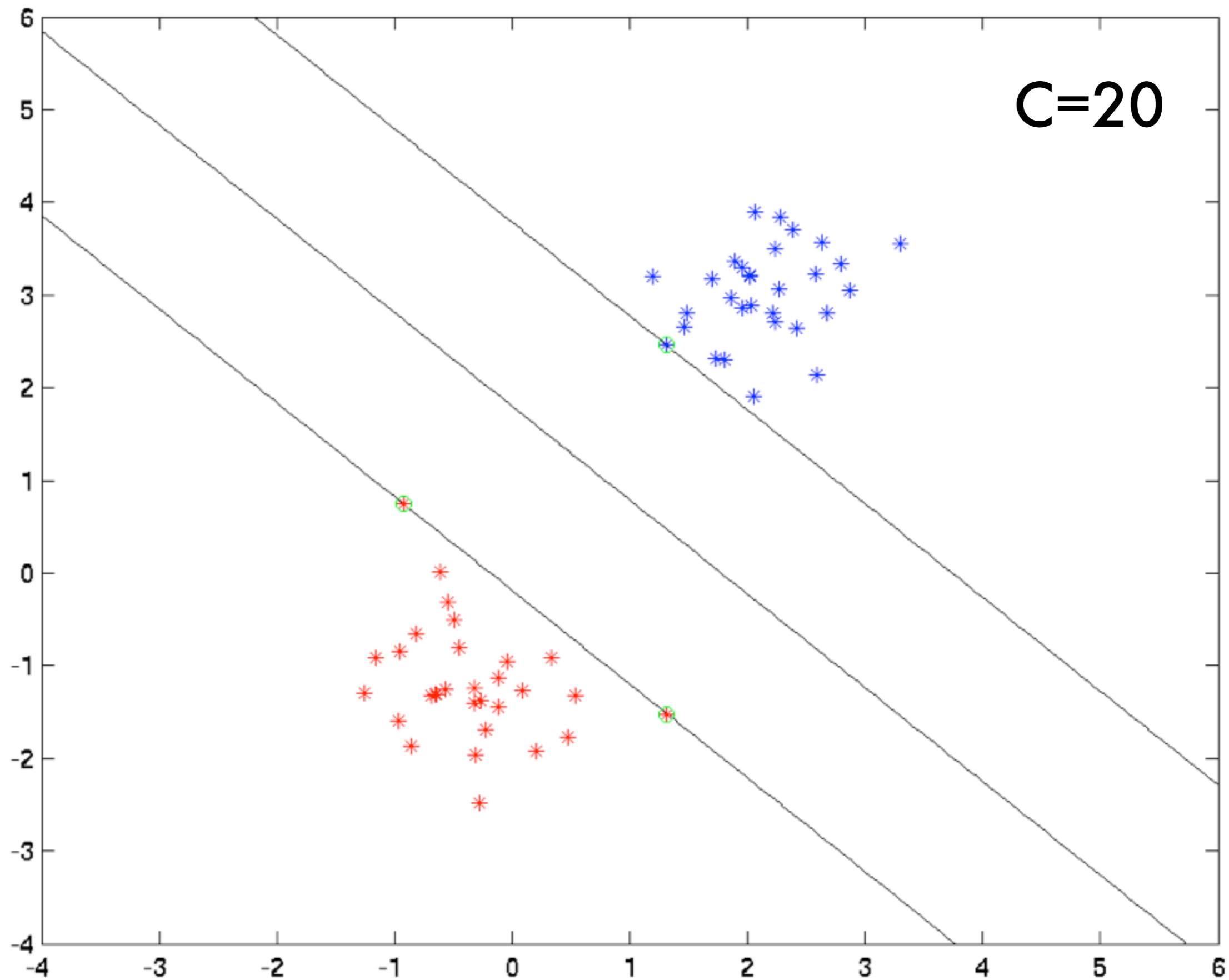


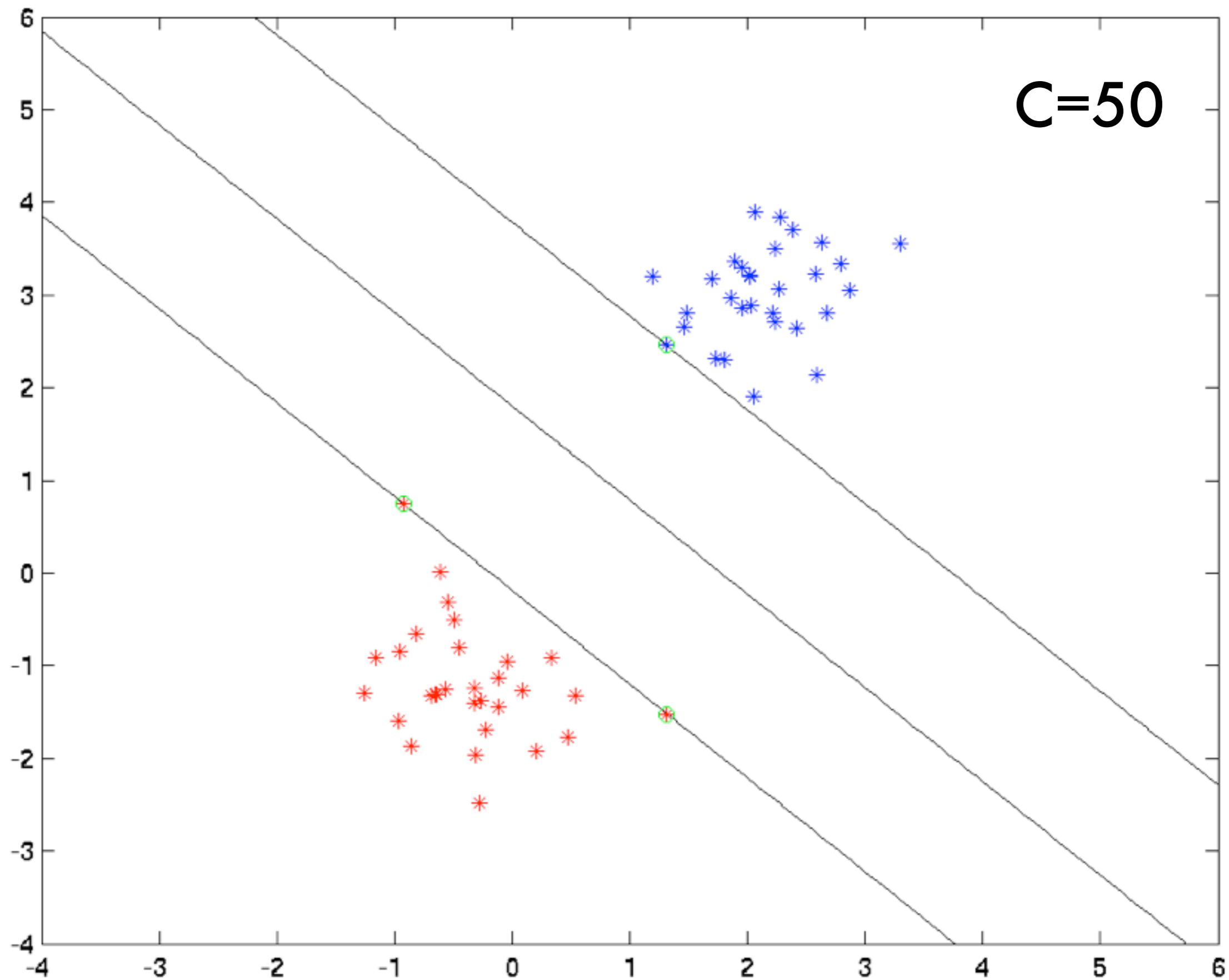


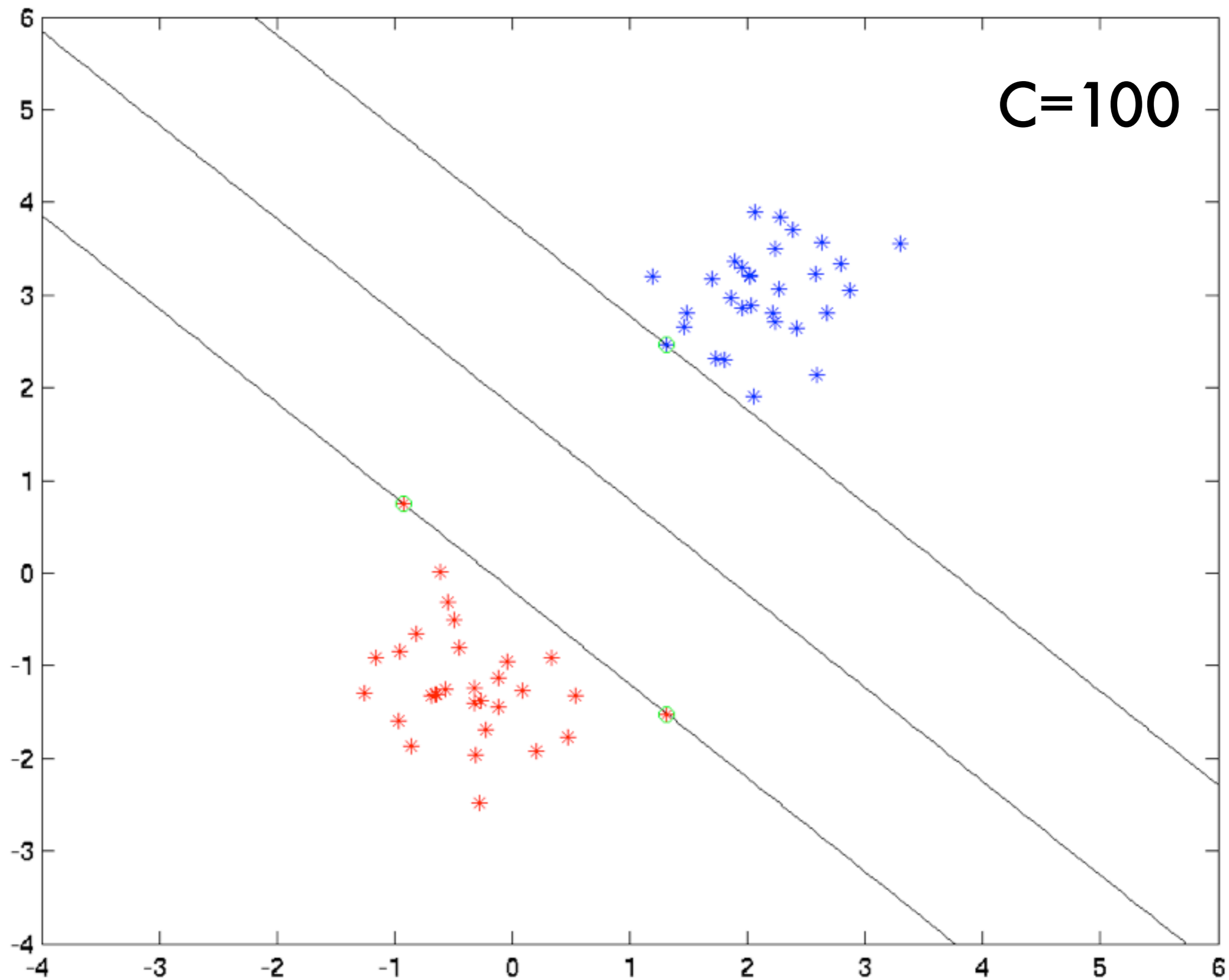


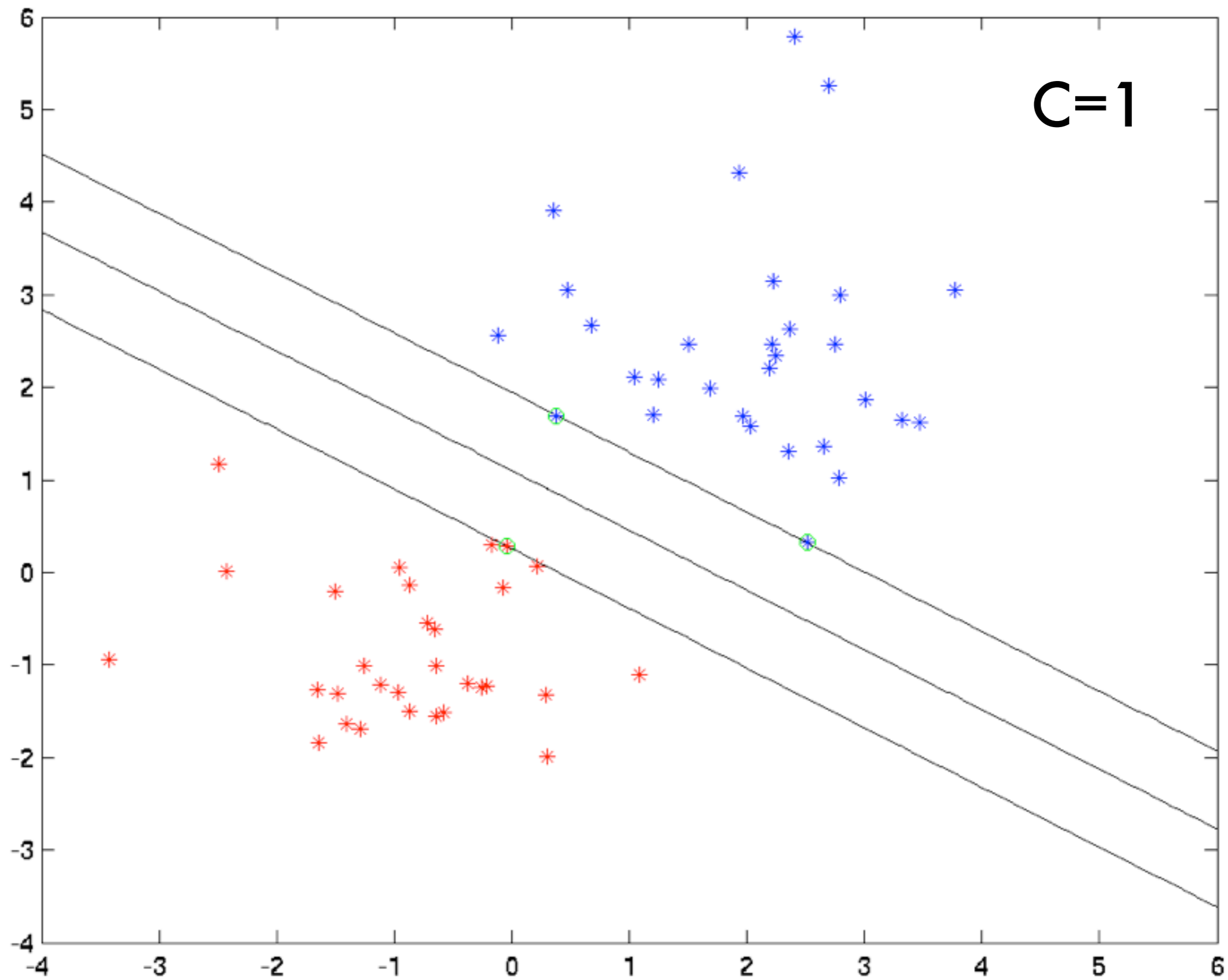


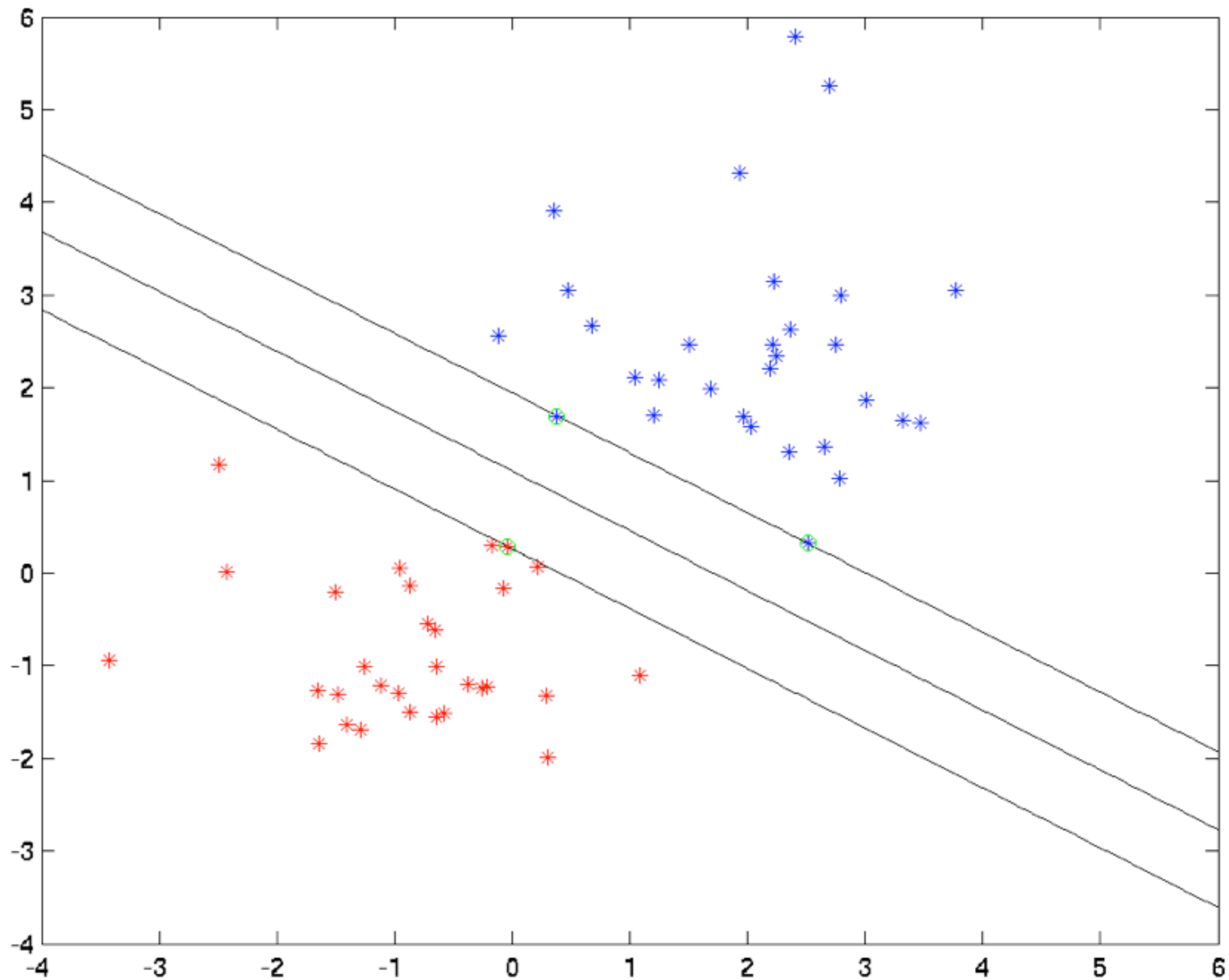


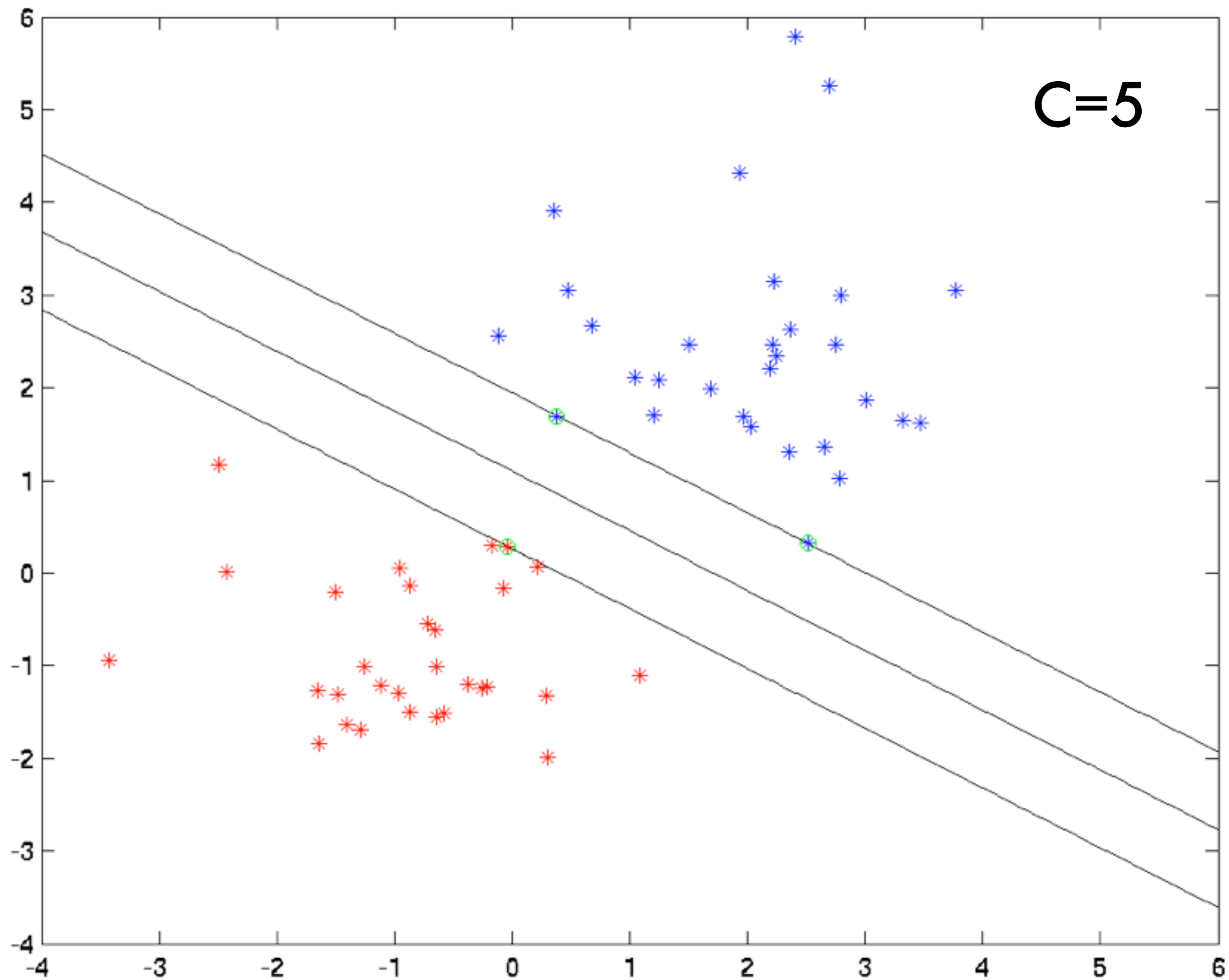


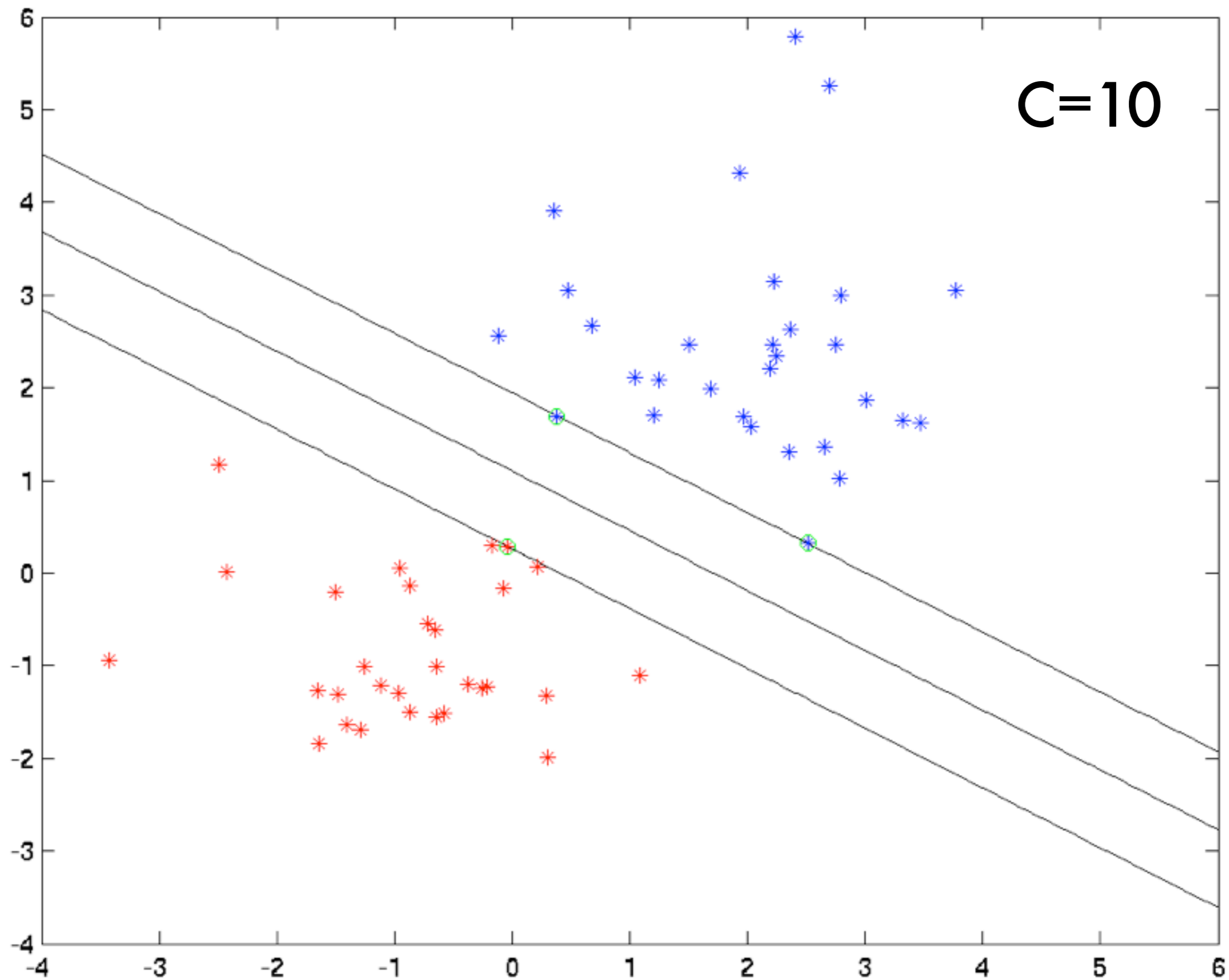


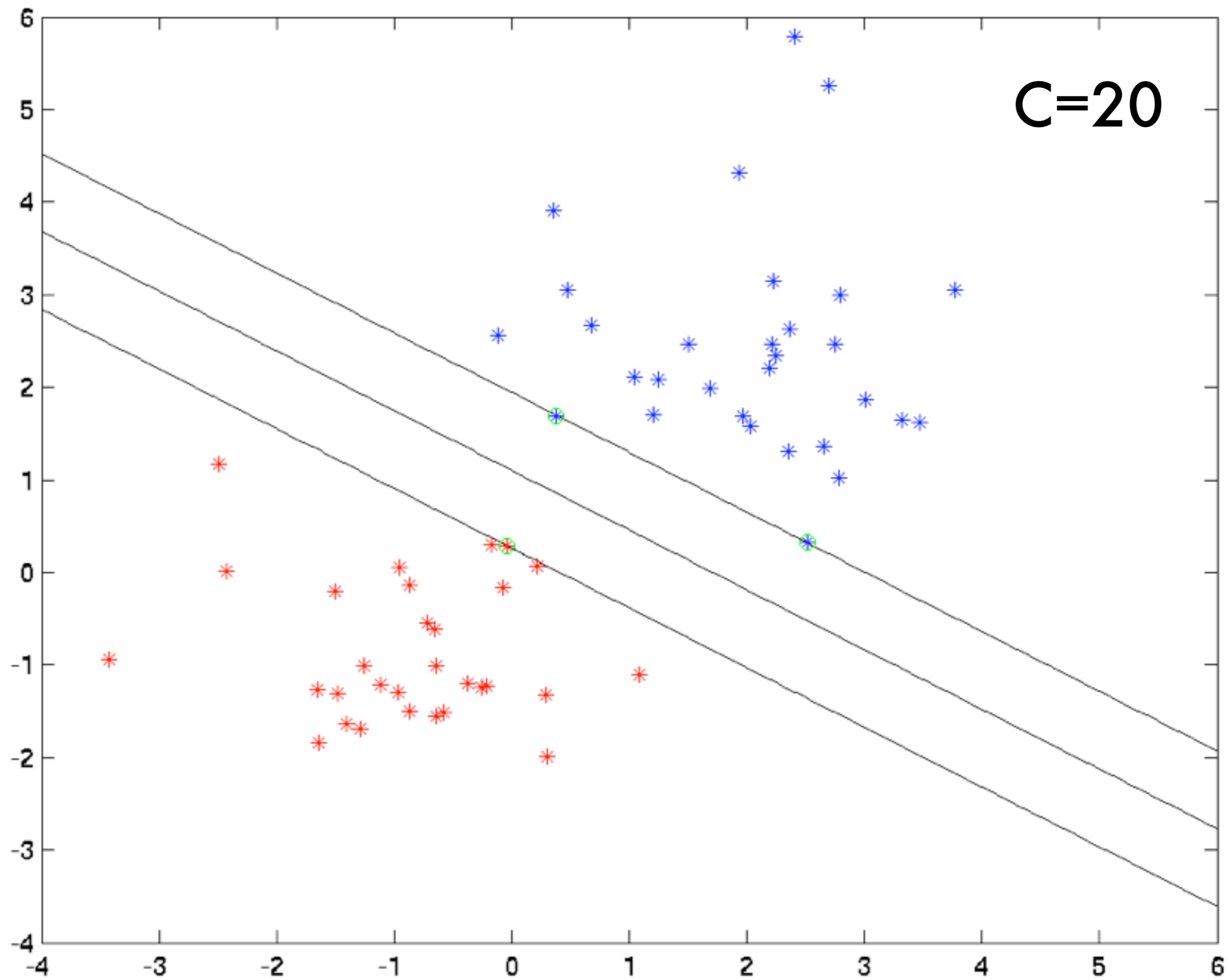


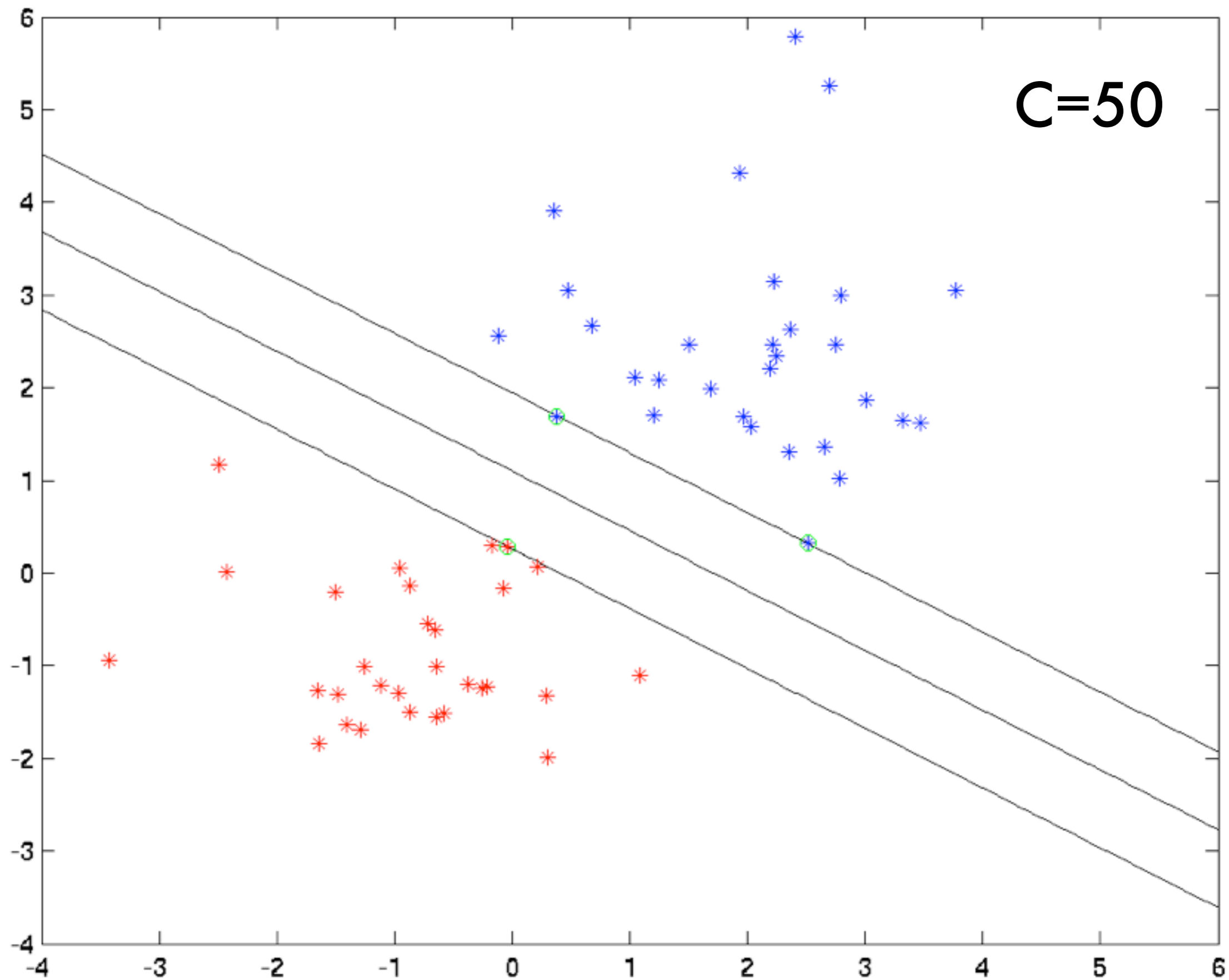


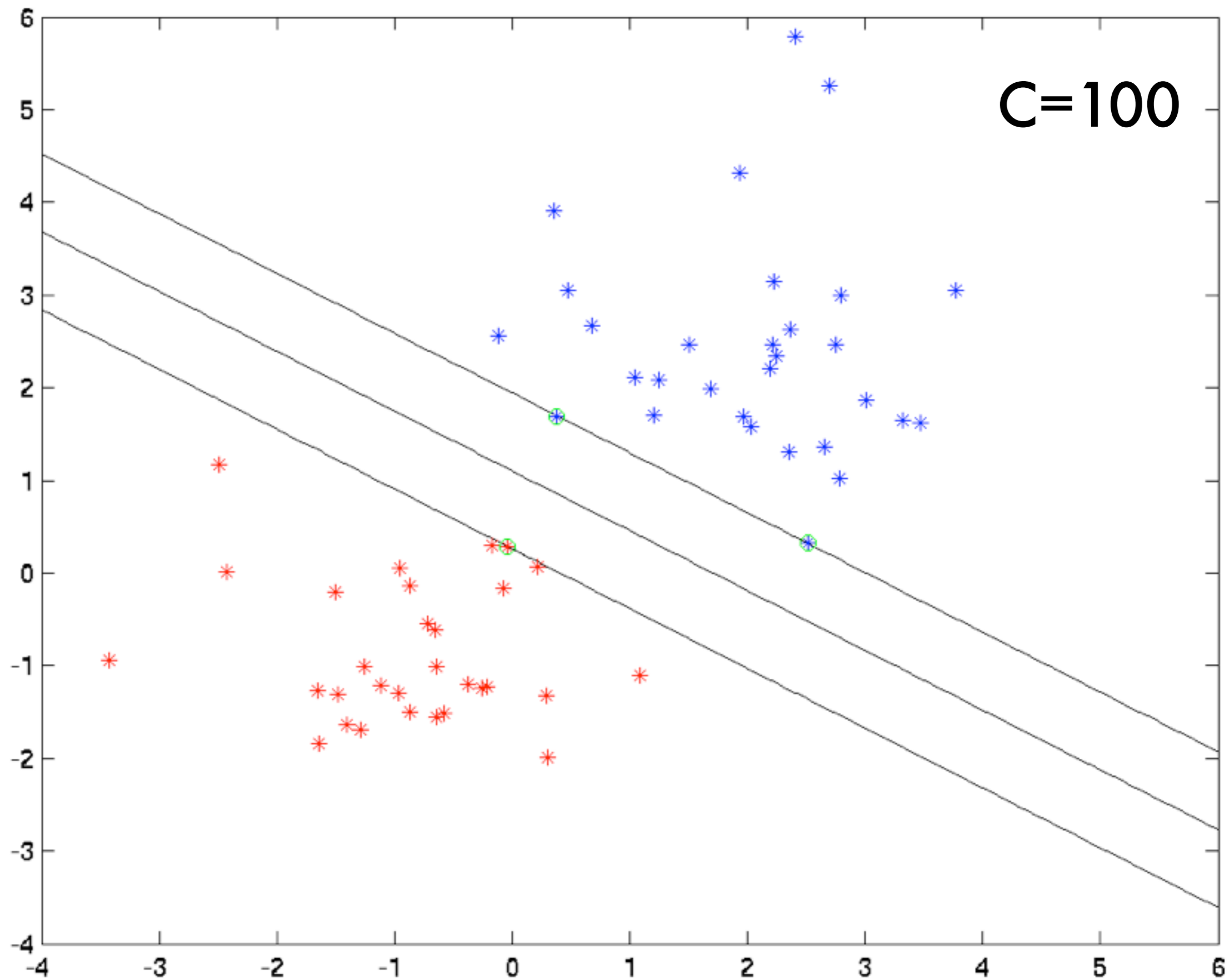


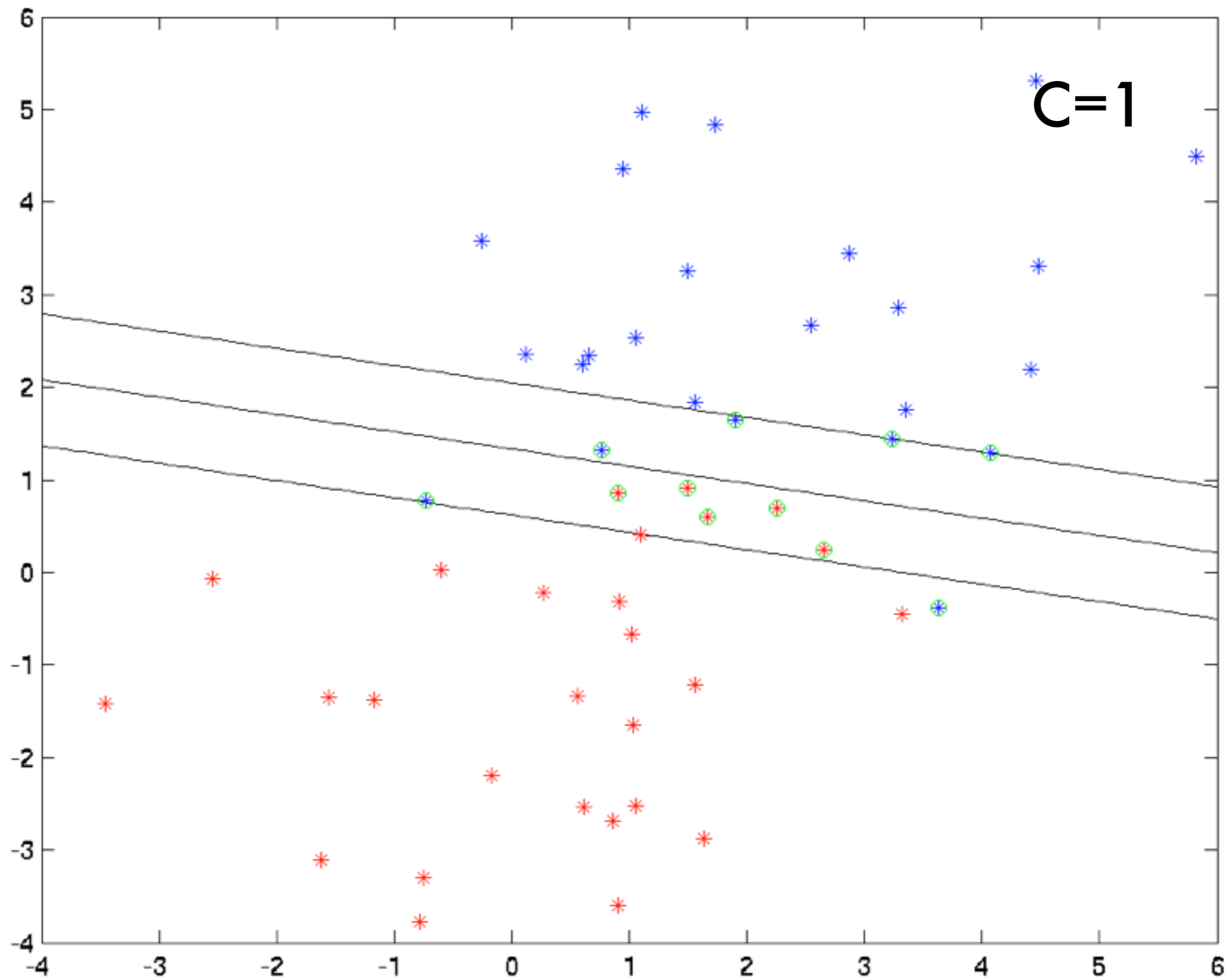


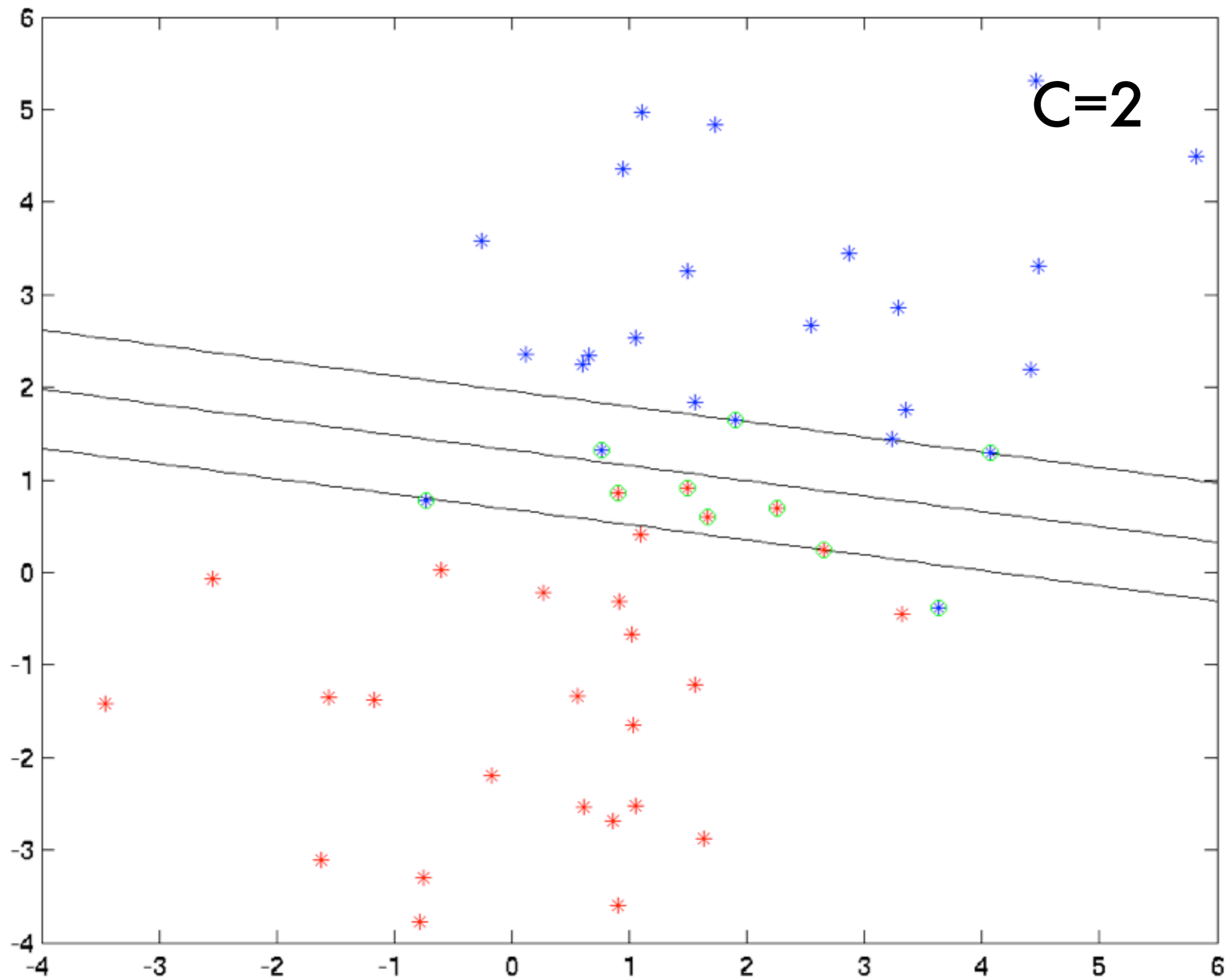


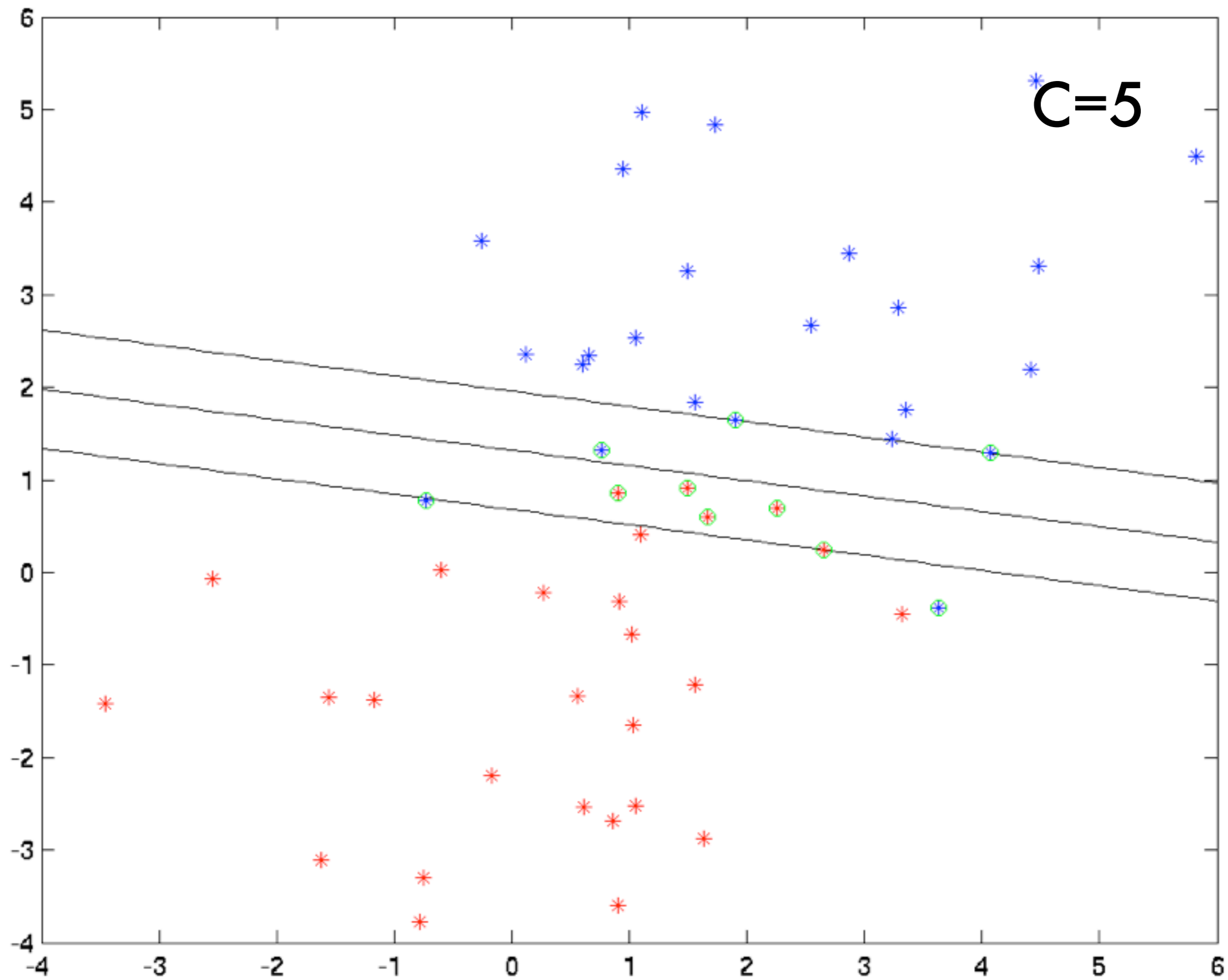


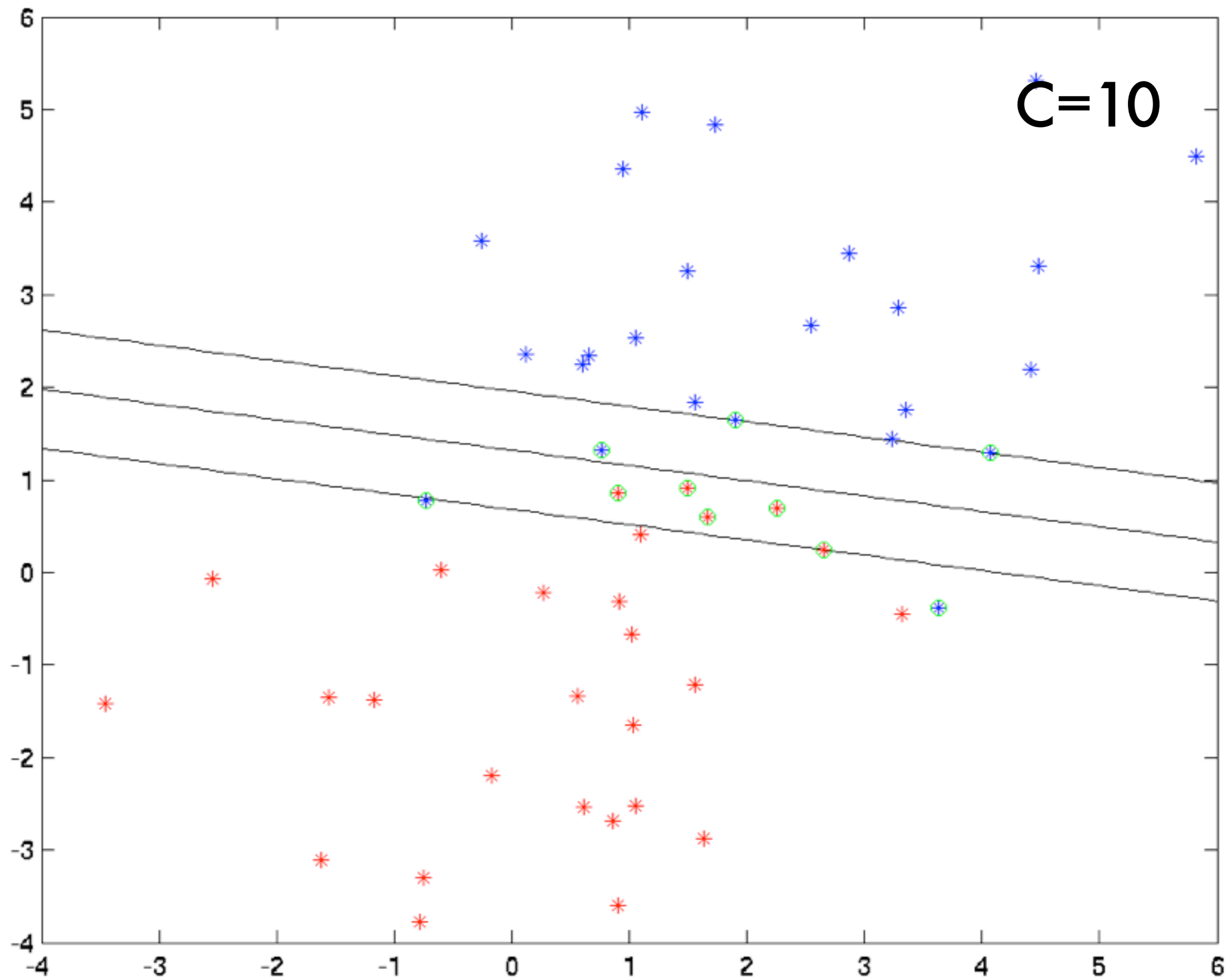


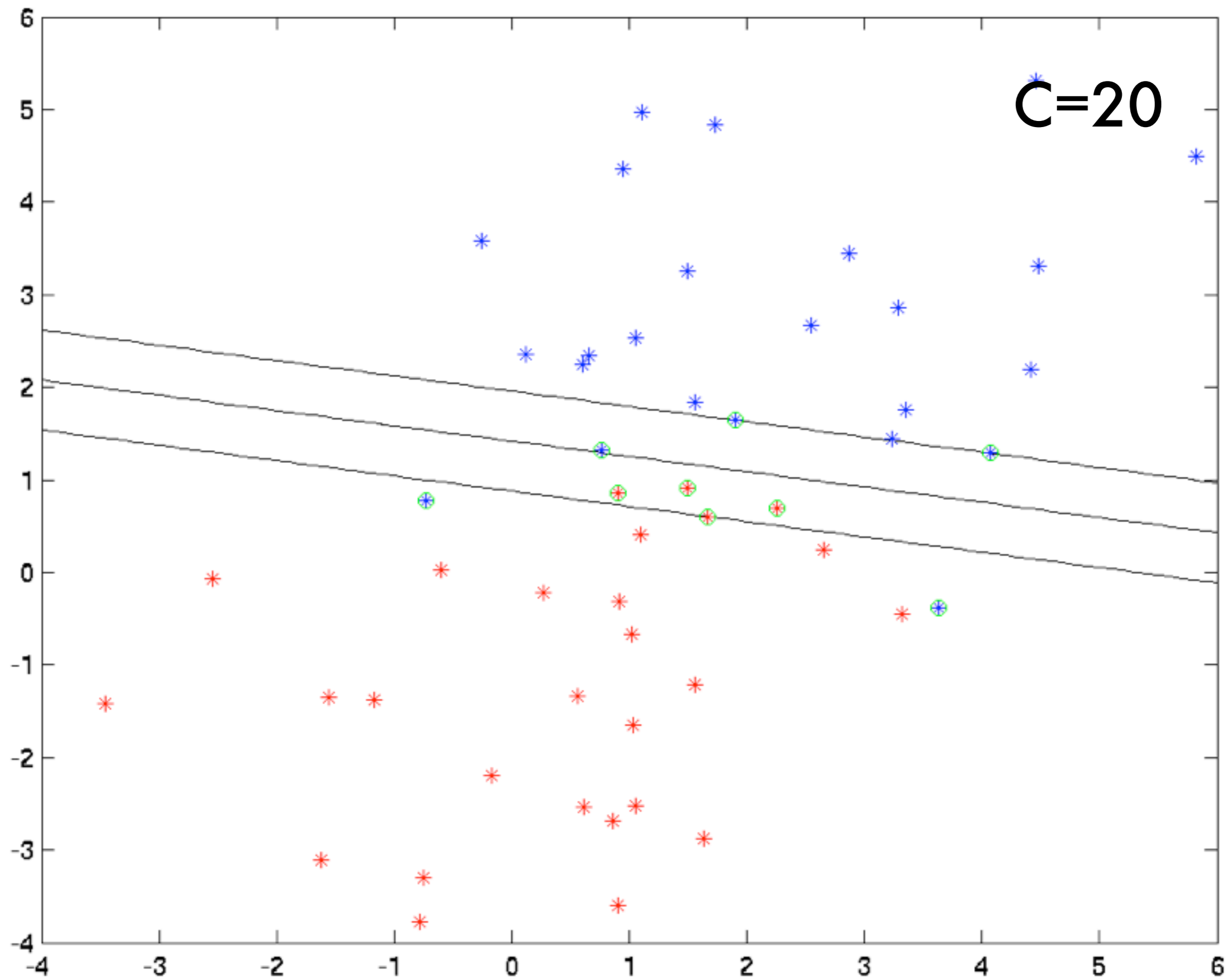


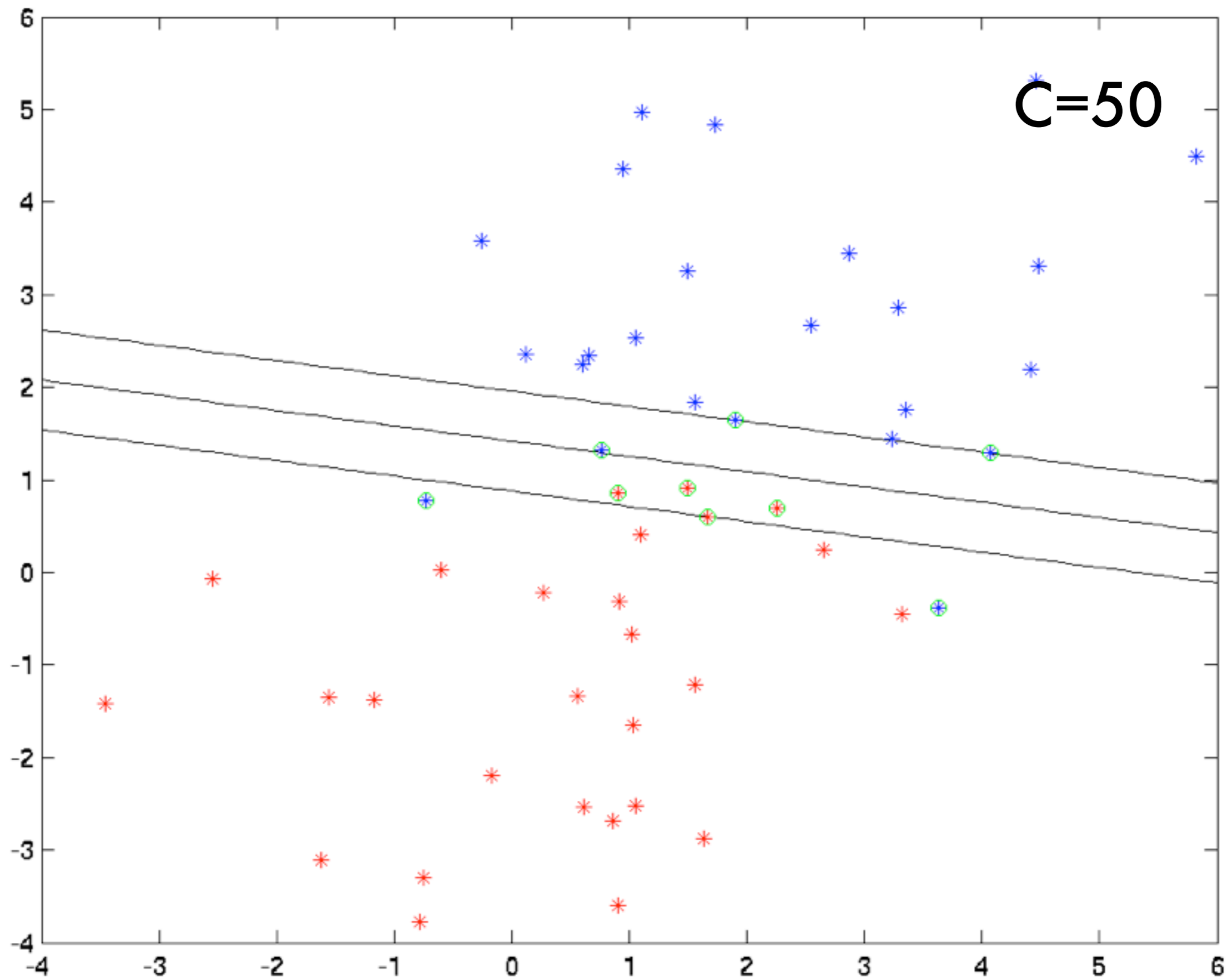


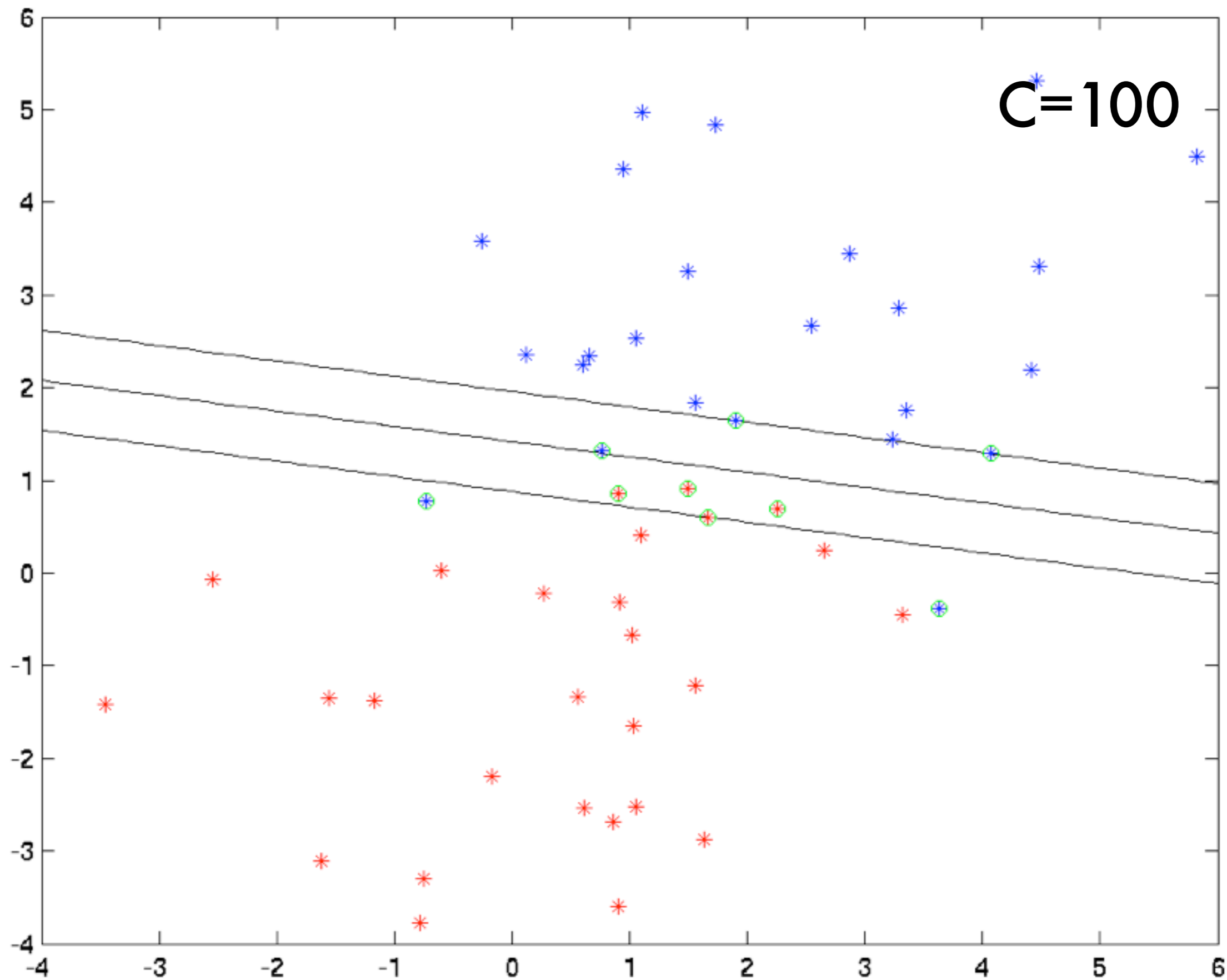












Solving the optimization problem

- Dual problem

$$\text{maximize}_{\alpha} -\frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \langle x_i, x_j \rangle + \sum_i \alpha_i$$

$$\text{subject to } \sum_i \alpha_i y_i = 0 \text{ and } \alpha_i \in [0, C]$$

- If problem is small enough (1000s of variables) we can use off-the-shelf solver (CVXOPT, CPLEX, OOQP, LOQO)
- For larger problem use fact that only SVs matter and solve in blocks (active set method).



MAGIC Etch A Sketch® SCREEN

Nonlinear
Separation

Horizontal
Dial

OHIO ART *The World of Toys*®

Vertical
Dial

MAGIC SCREEN IS GLASS SET IN DURABLE PLASTIC FRAME
USE WITH CARE

The Kernel Trick

- **Linear soft margin problem**

$$\text{minimize}_{w,b} \frac{1}{2} \|w\|^2 + C \sum_i \xi_i$$

subject to $y_i [\langle w, x_i \rangle + b] \geq 1 - \xi_i$ and $\xi_i \geq 0$

- **Dual problem**

$$\text{maximize}_{\alpha} - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \langle x_i, x_j \rangle + \sum_i \alpha_i$$

subject to $\sum_i \alpha_i y_i = 0$ and $\alpha_i \in [0, C]$

- **Support vector expansion**

$$f(x) = \sum_i \alpha_i y_i \langle x_i, x \rangle + b$$

The Kernel Trick

- **Linear soft margin problem**

$$\underset{w, b}{\text{minimize}} \quad \frac{1}{2} \|w\|^2 + C \sum_i \xi_i$$

subject to $y_i [\langle w, \phi(x_i) \rangle + b] \geq 1 - \xi_i$ and $\xi_i \geq 0$

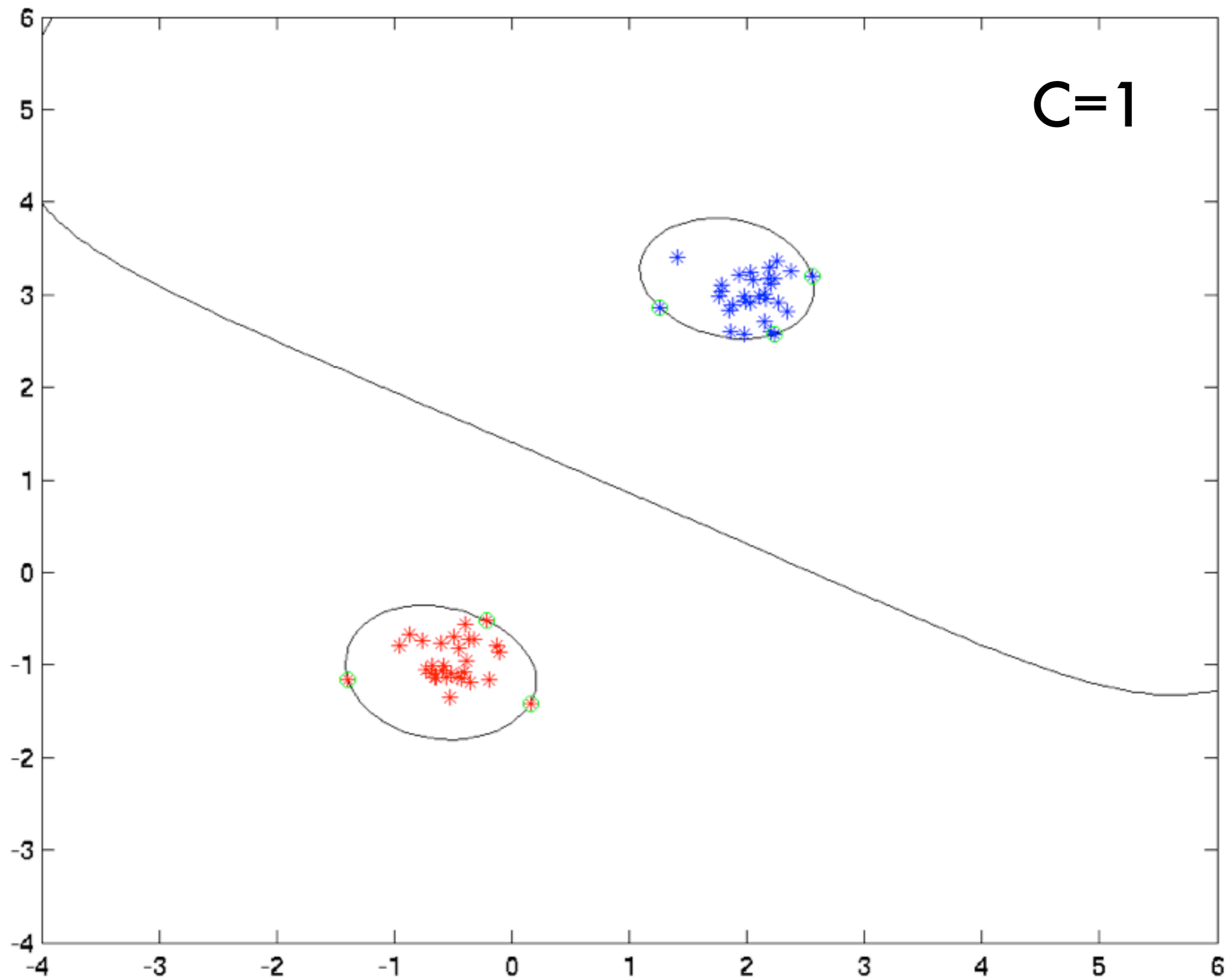
- **Dual problem**

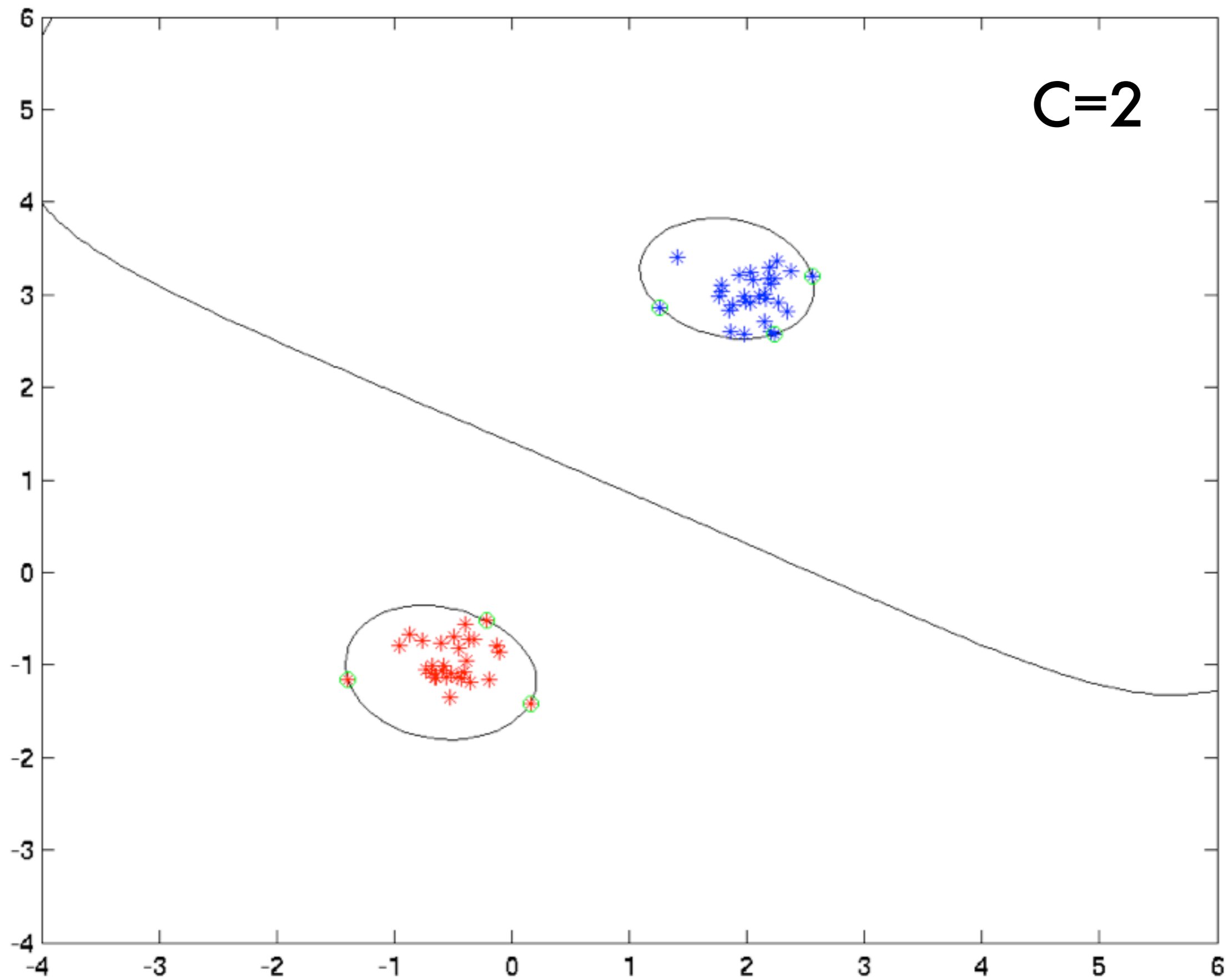
$$\underset{\alpha}{\text{maximize}} \quad -\frac{1}{2} \sum_{i, j} \alpha_i \alpha_j y_i y_j k(x_i, x_j) + \sum_i \alpha_i$$

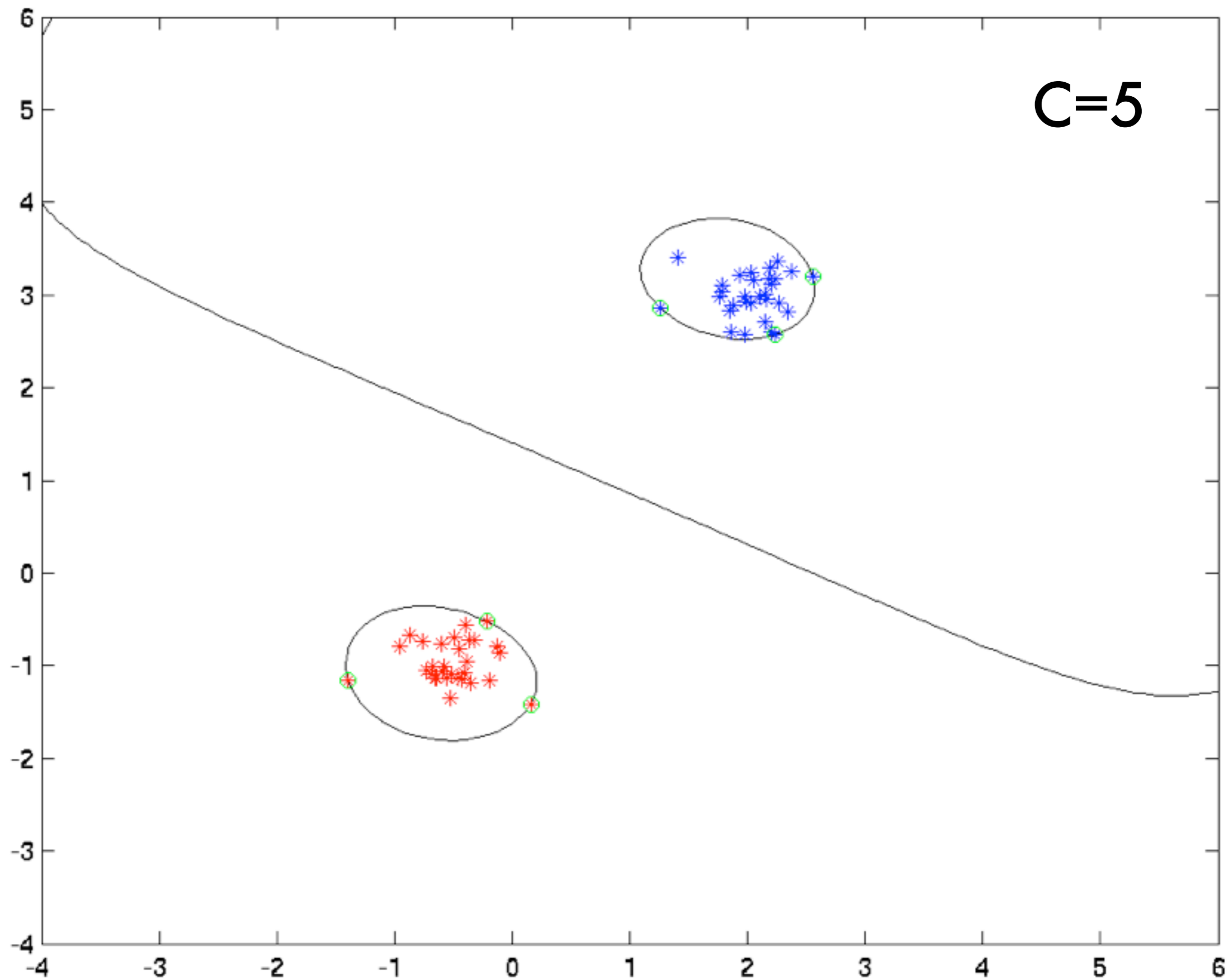
subject to $\sum_i \alpha_i y_i = 0$ and $\alpha_i \in [0, C]$

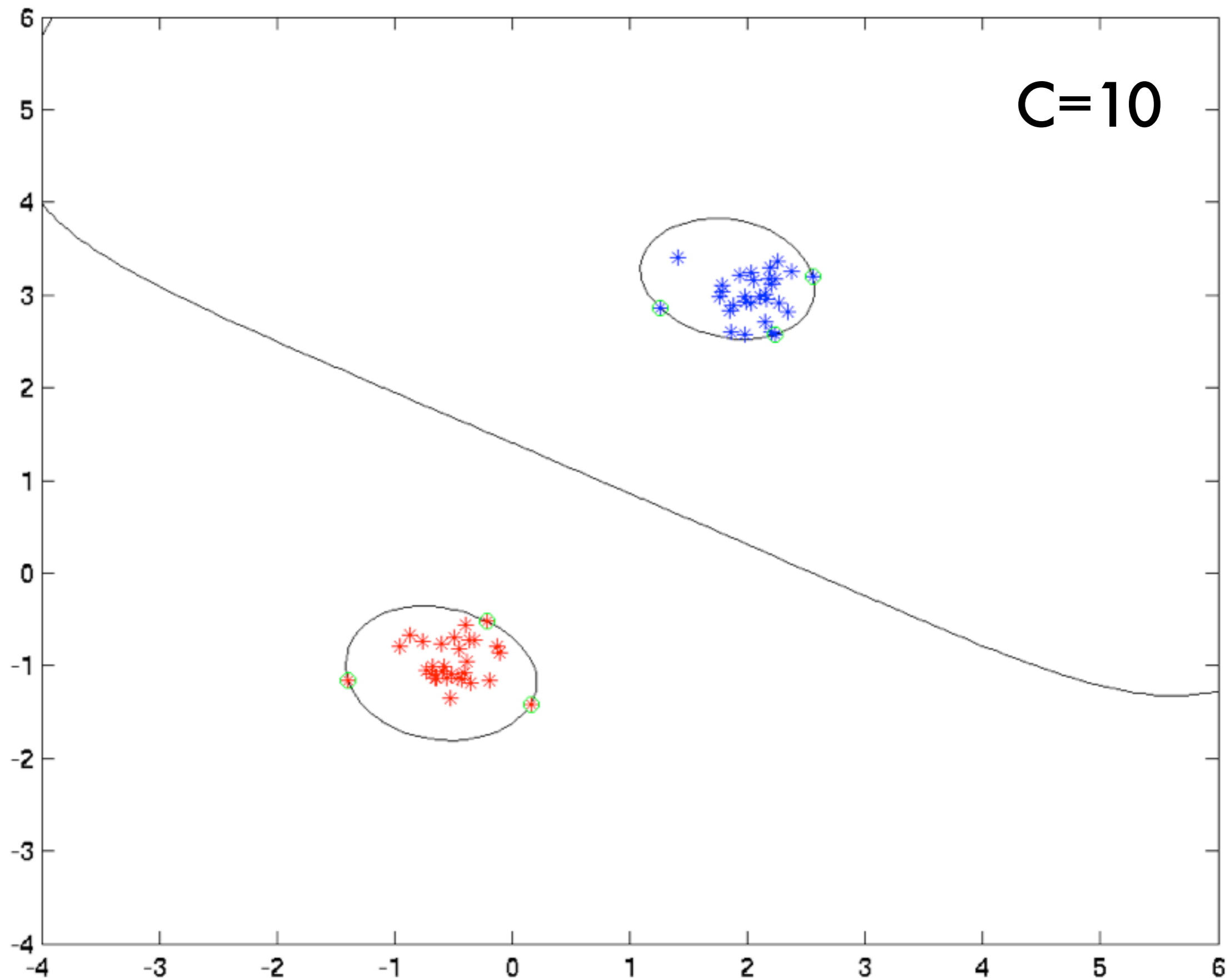
- **Support vector expansion**

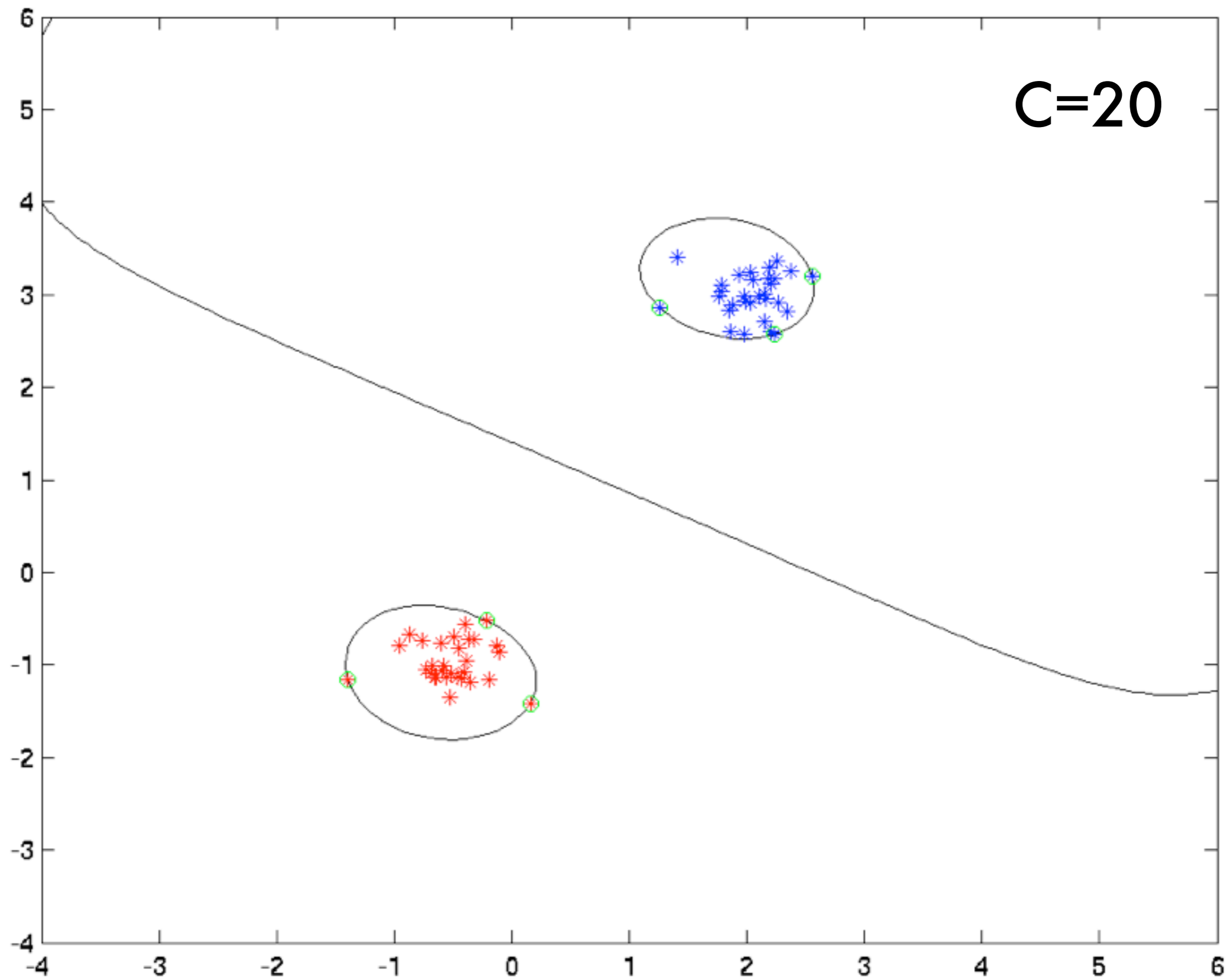
$$f(x) = \sum_i \alpha_i y_i k(x_i, x) + b$$

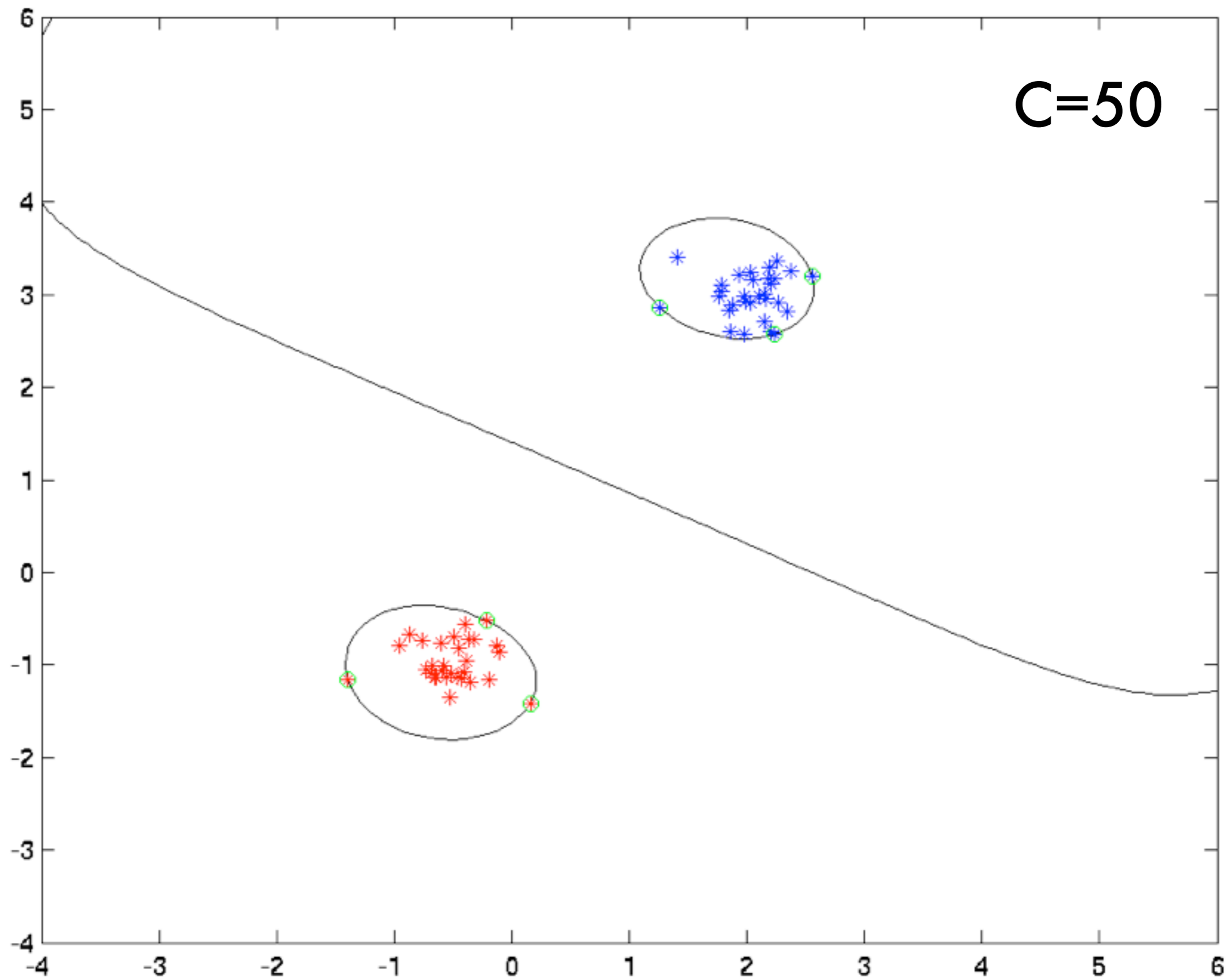


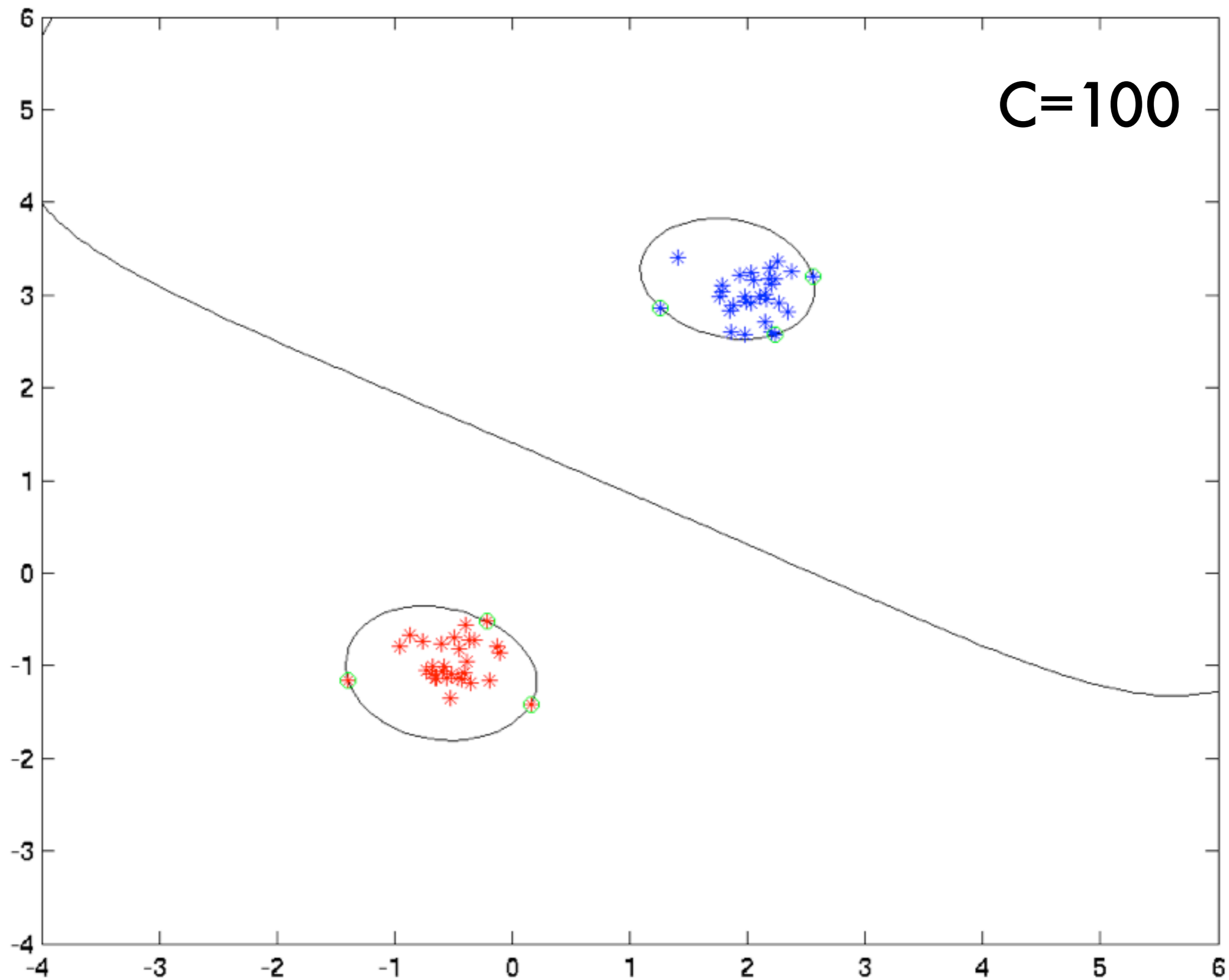


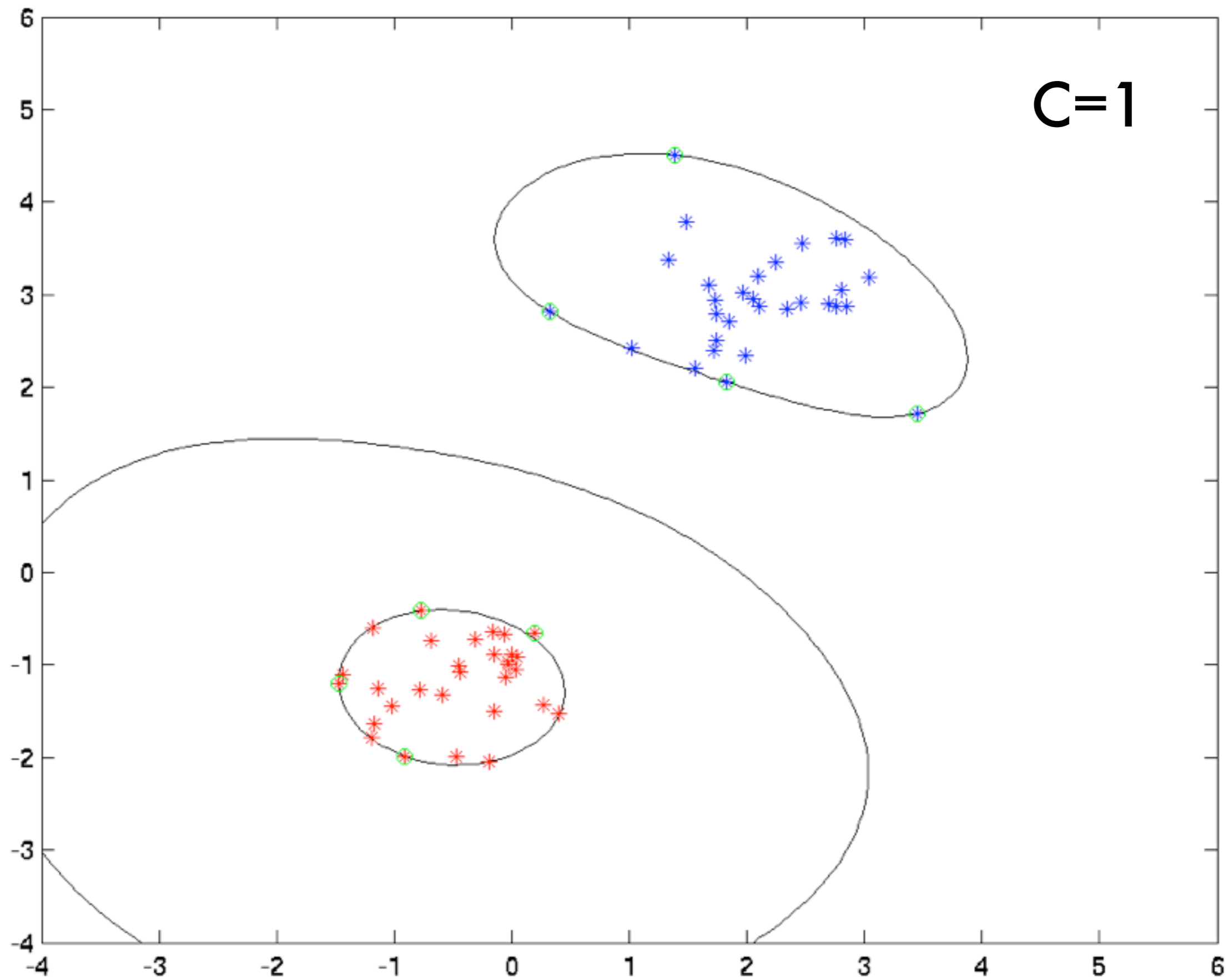


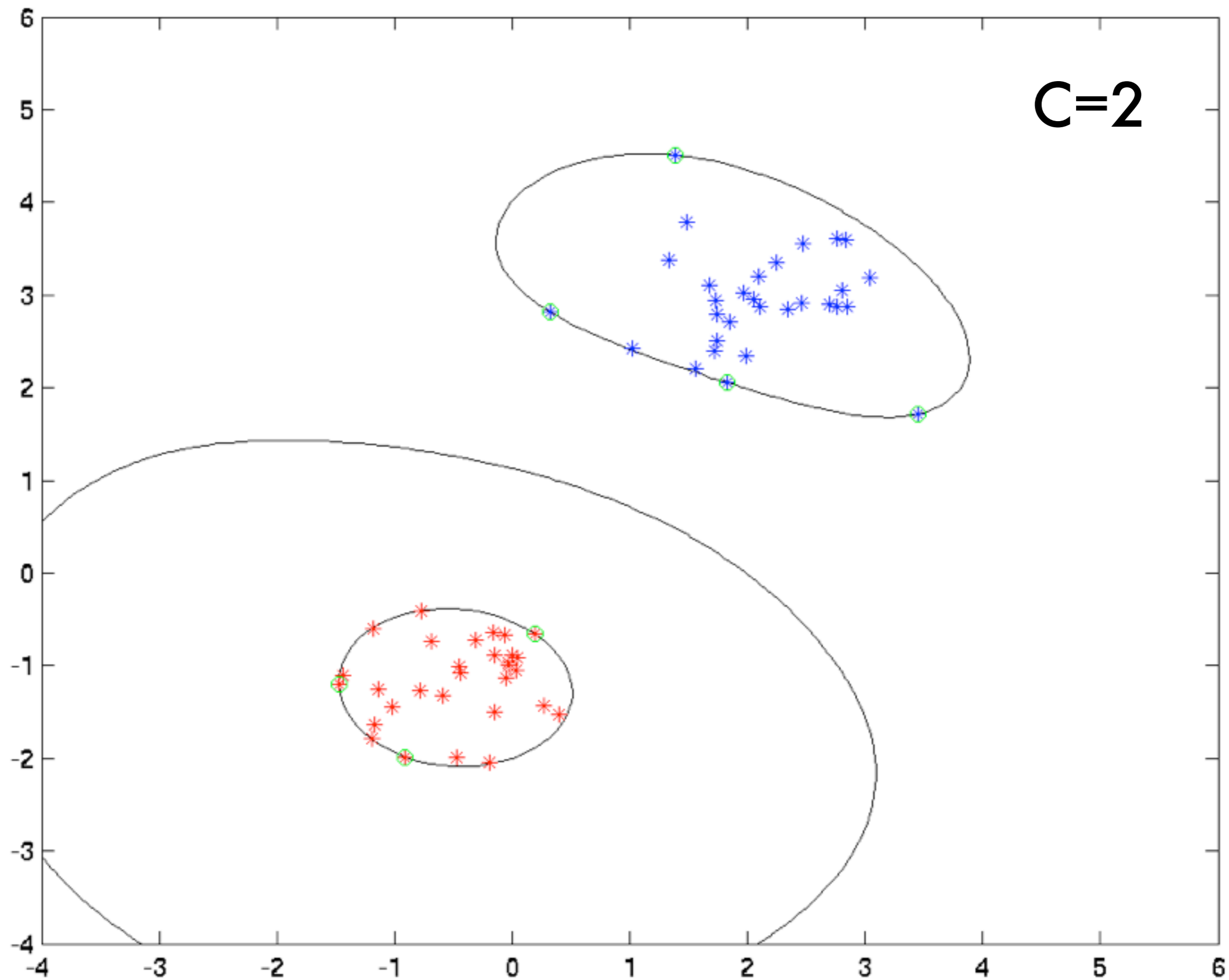


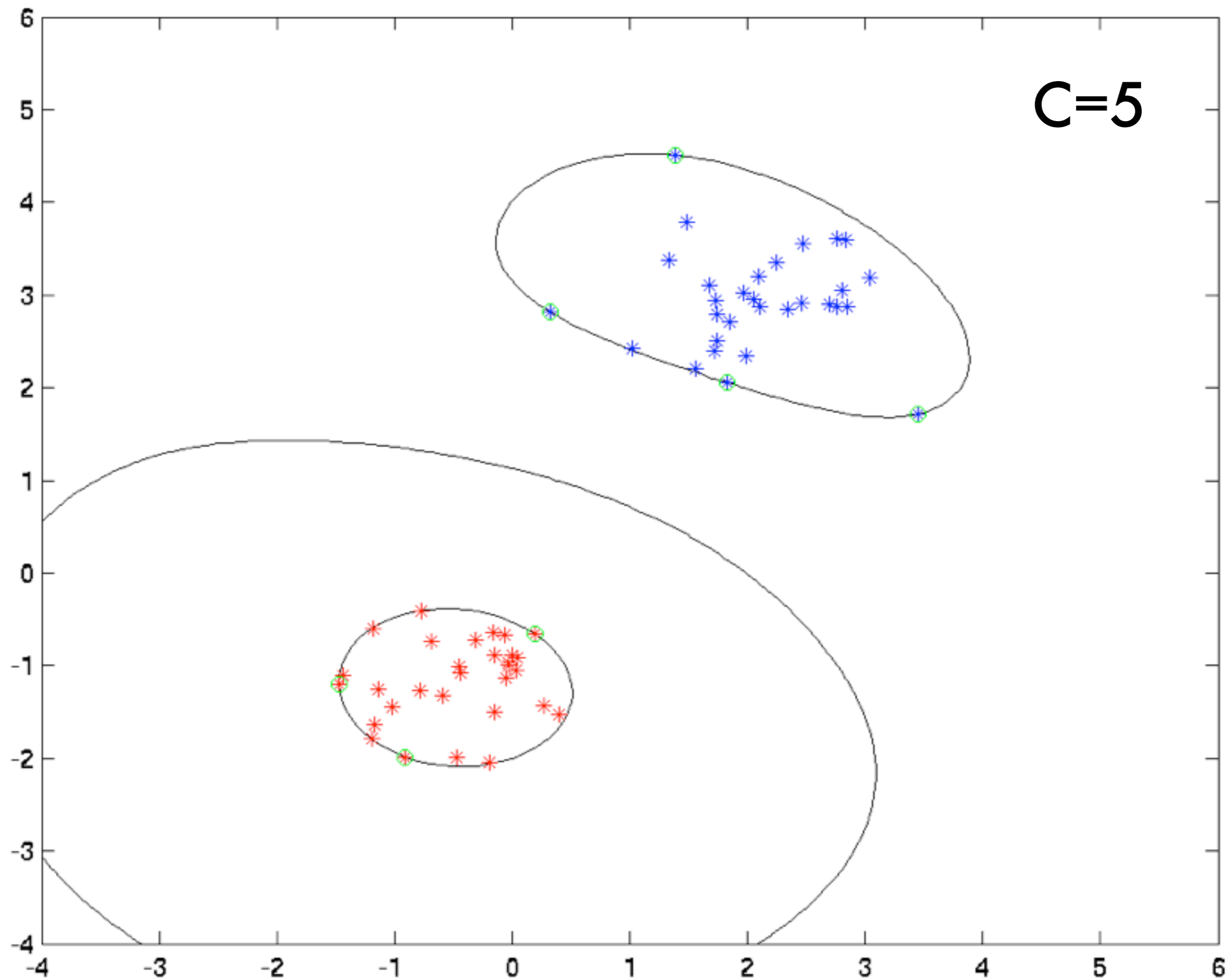


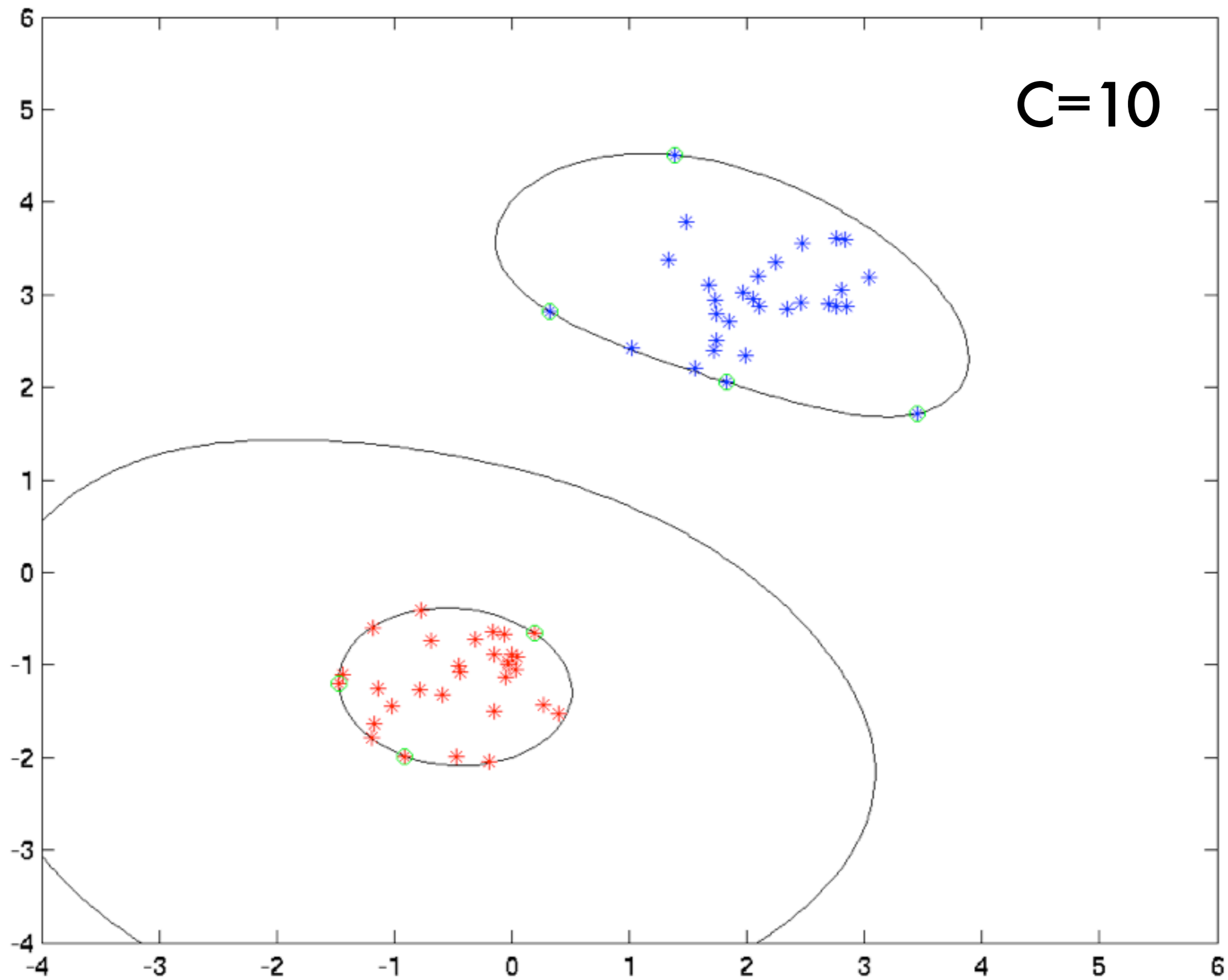


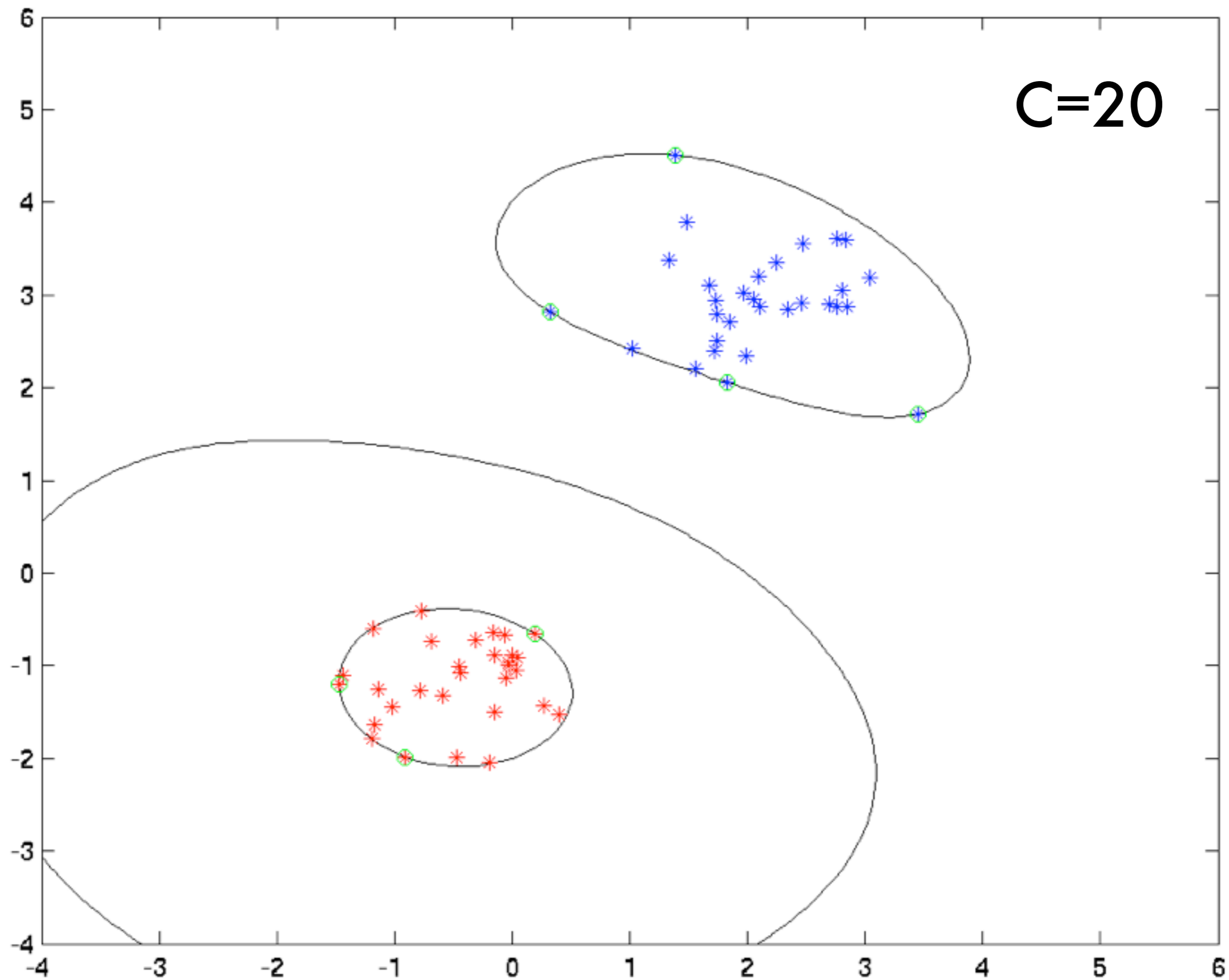


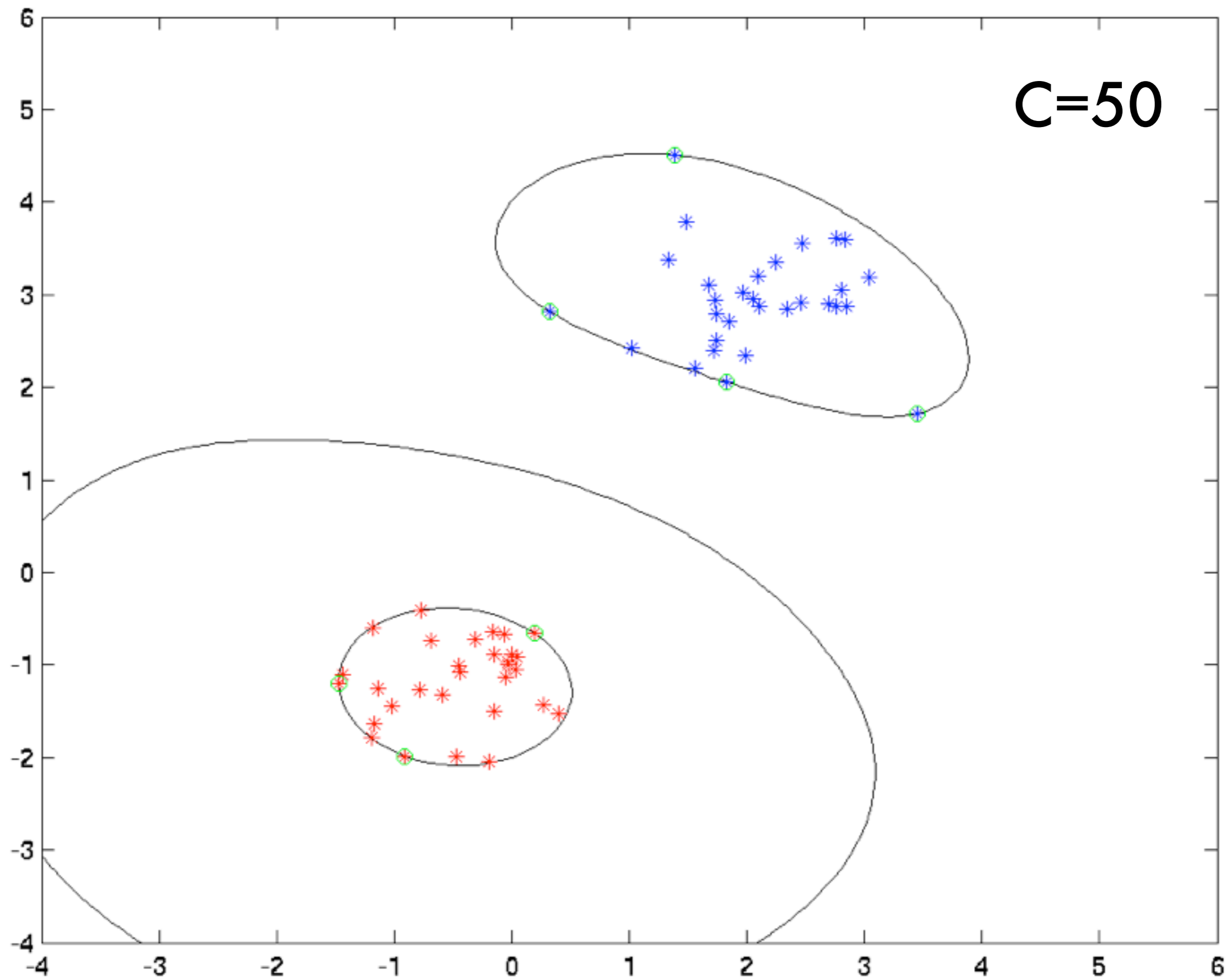


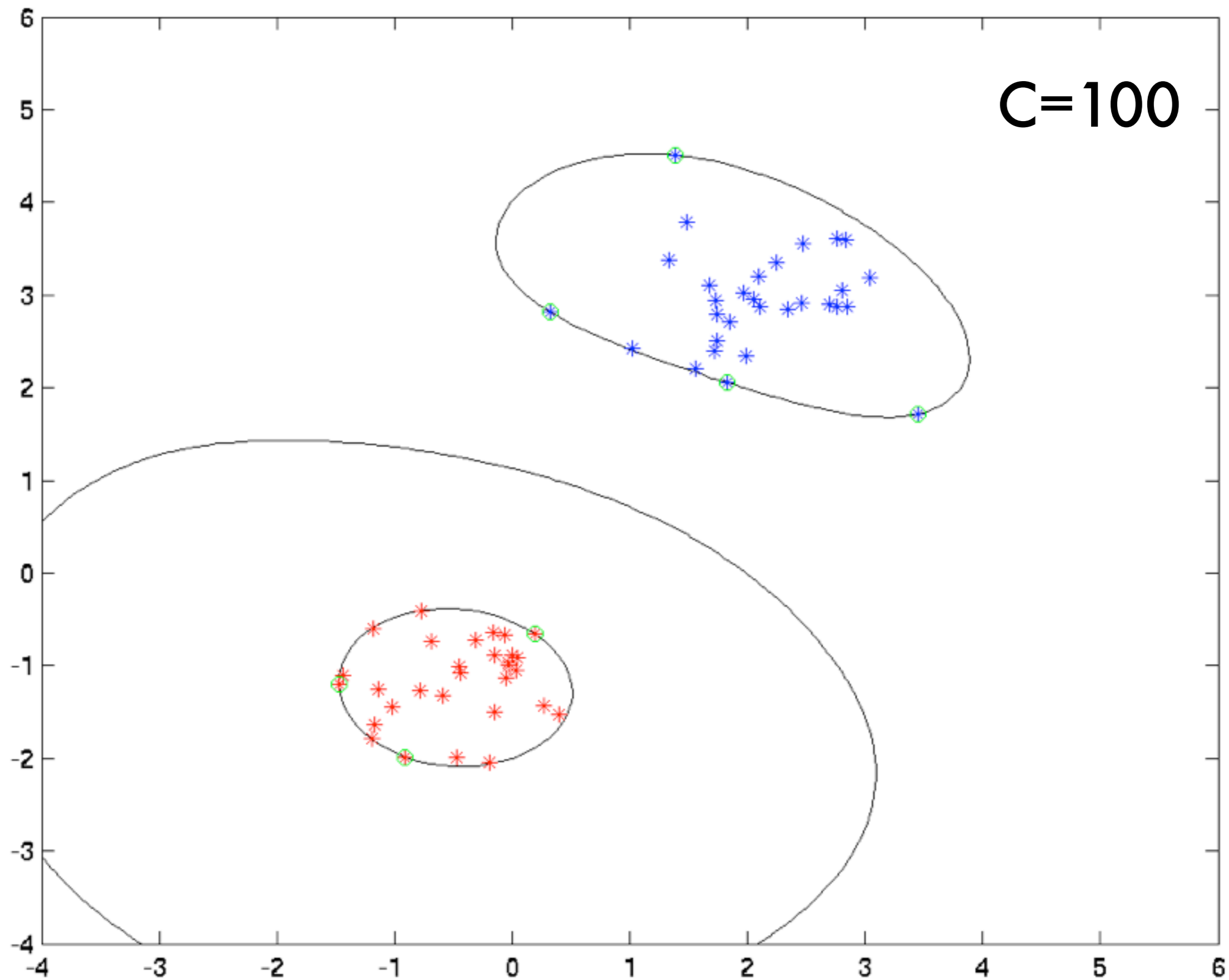


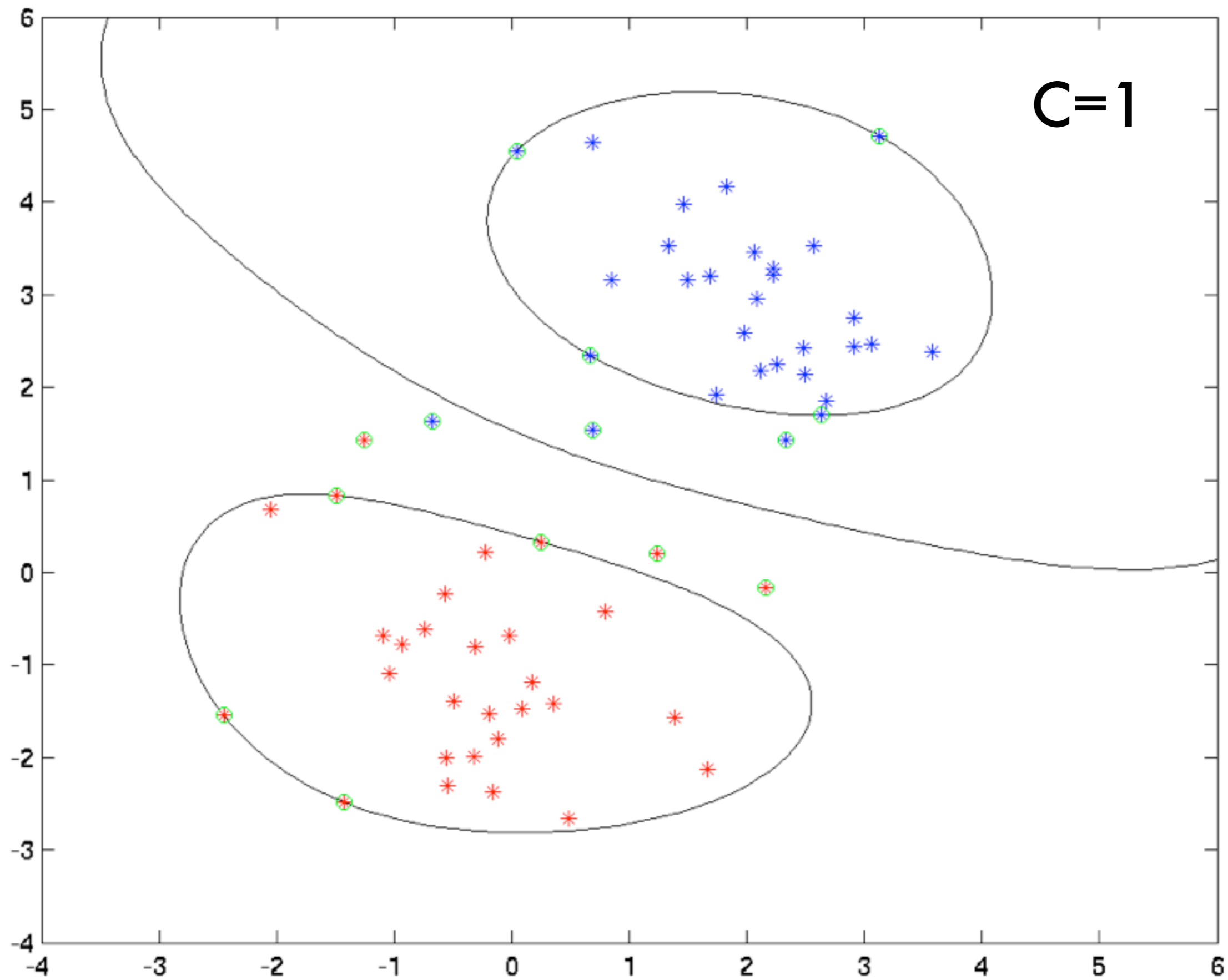


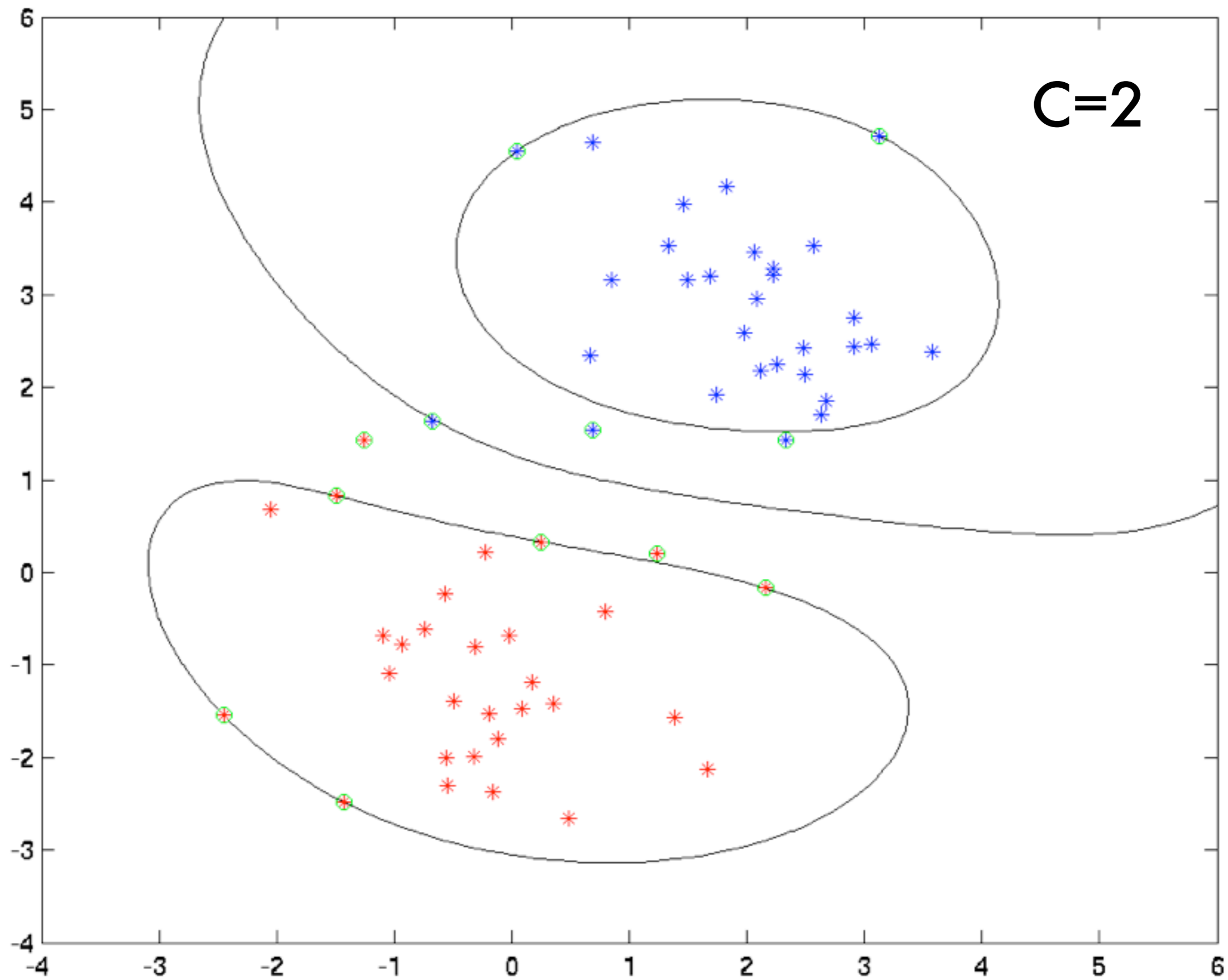


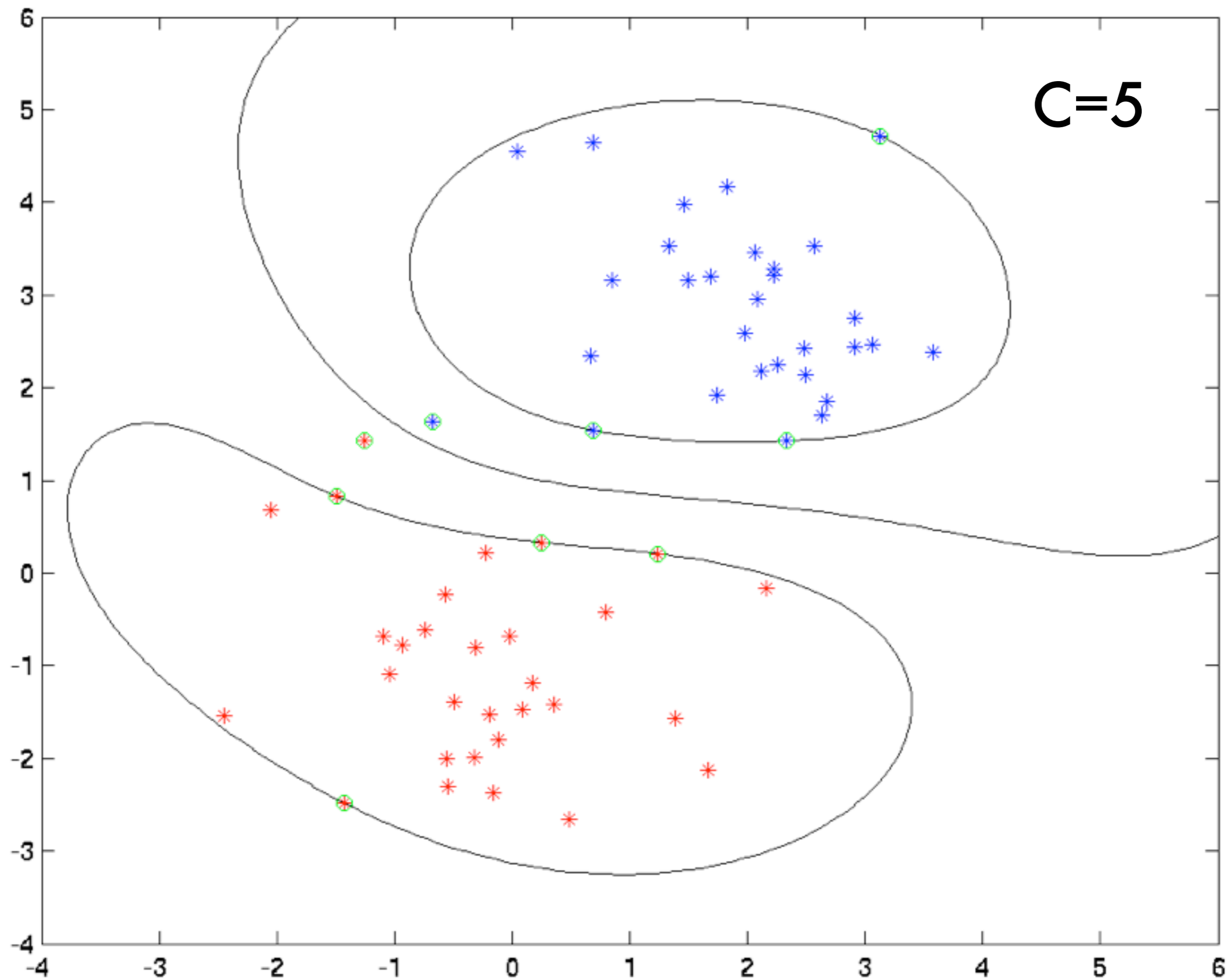


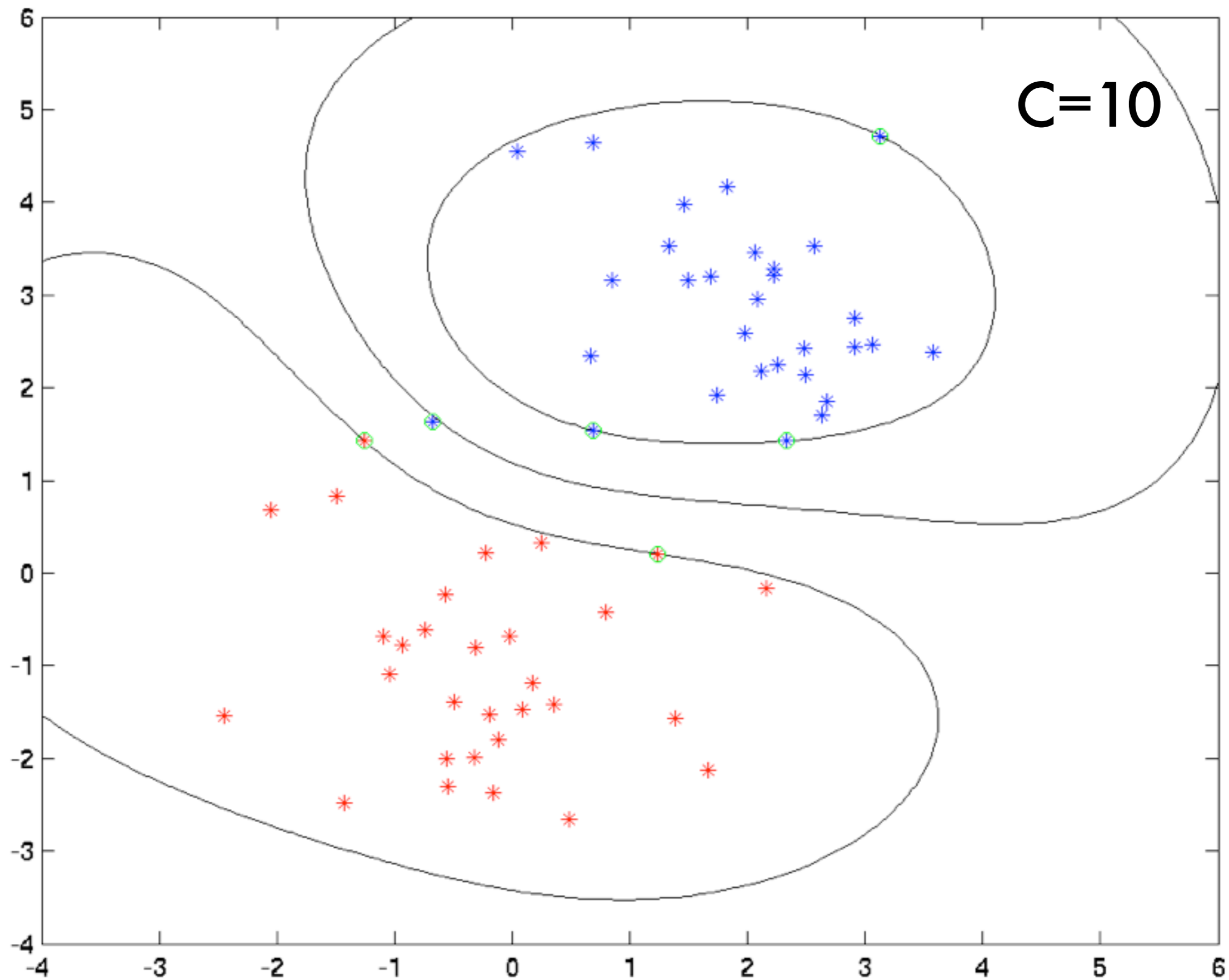


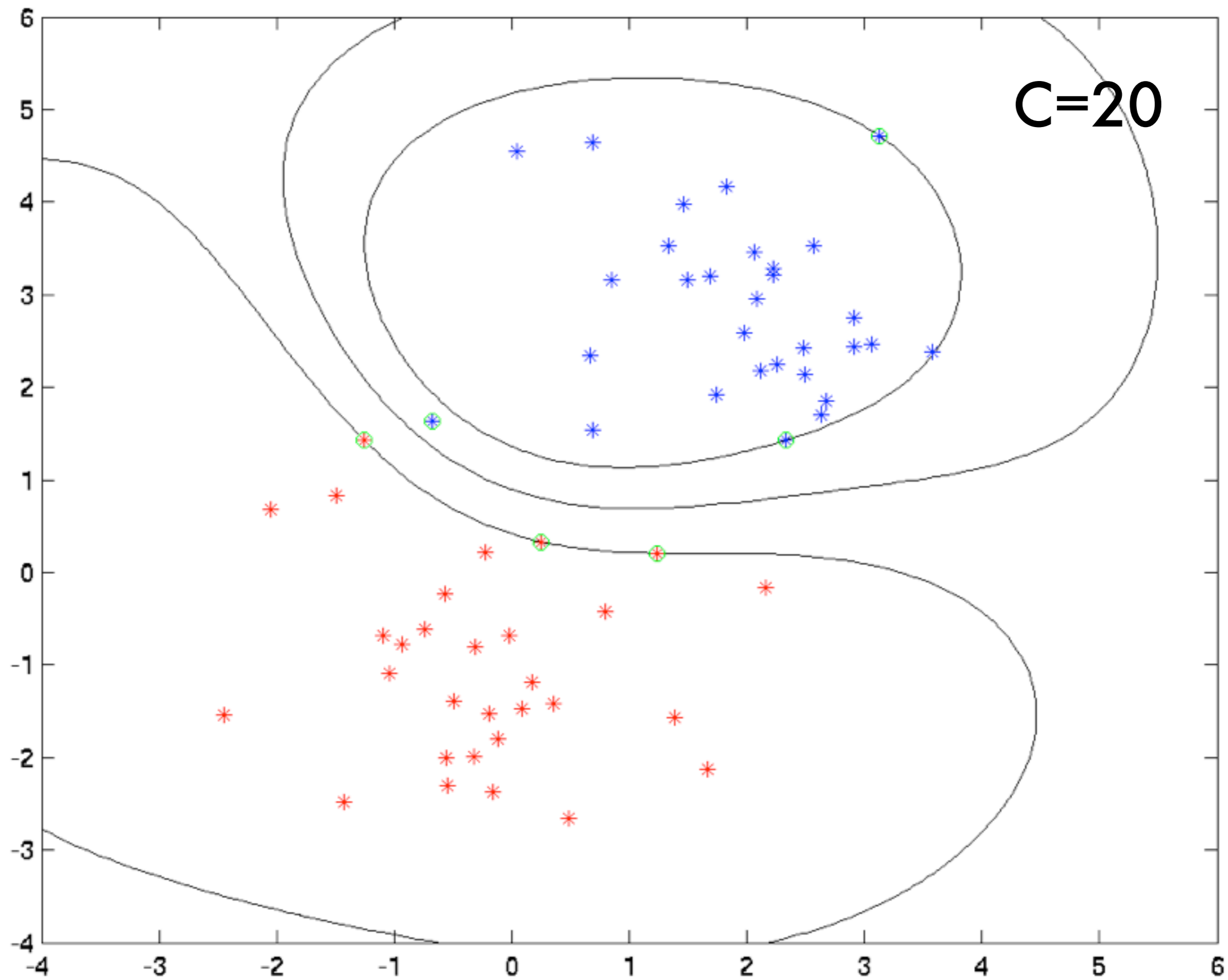


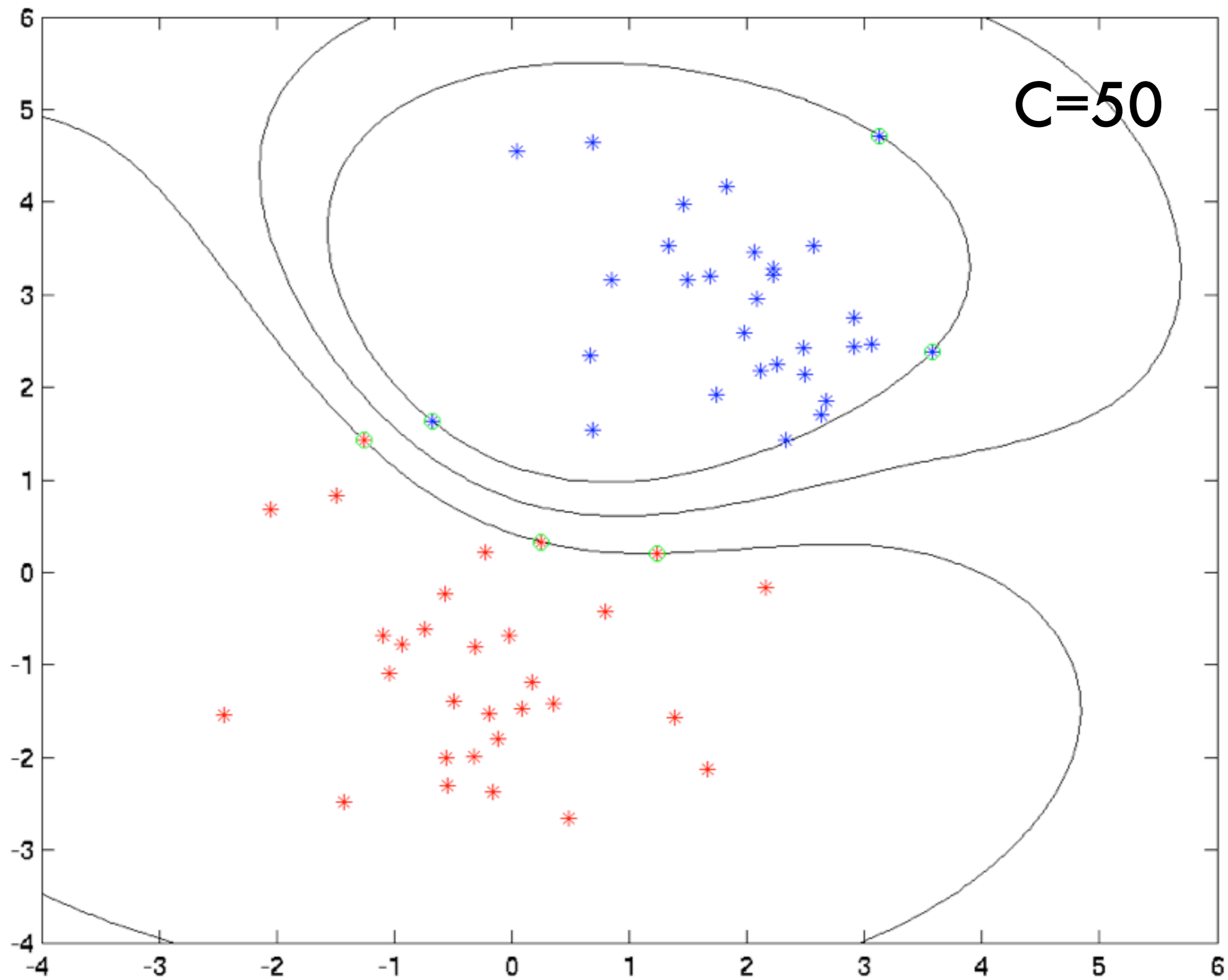


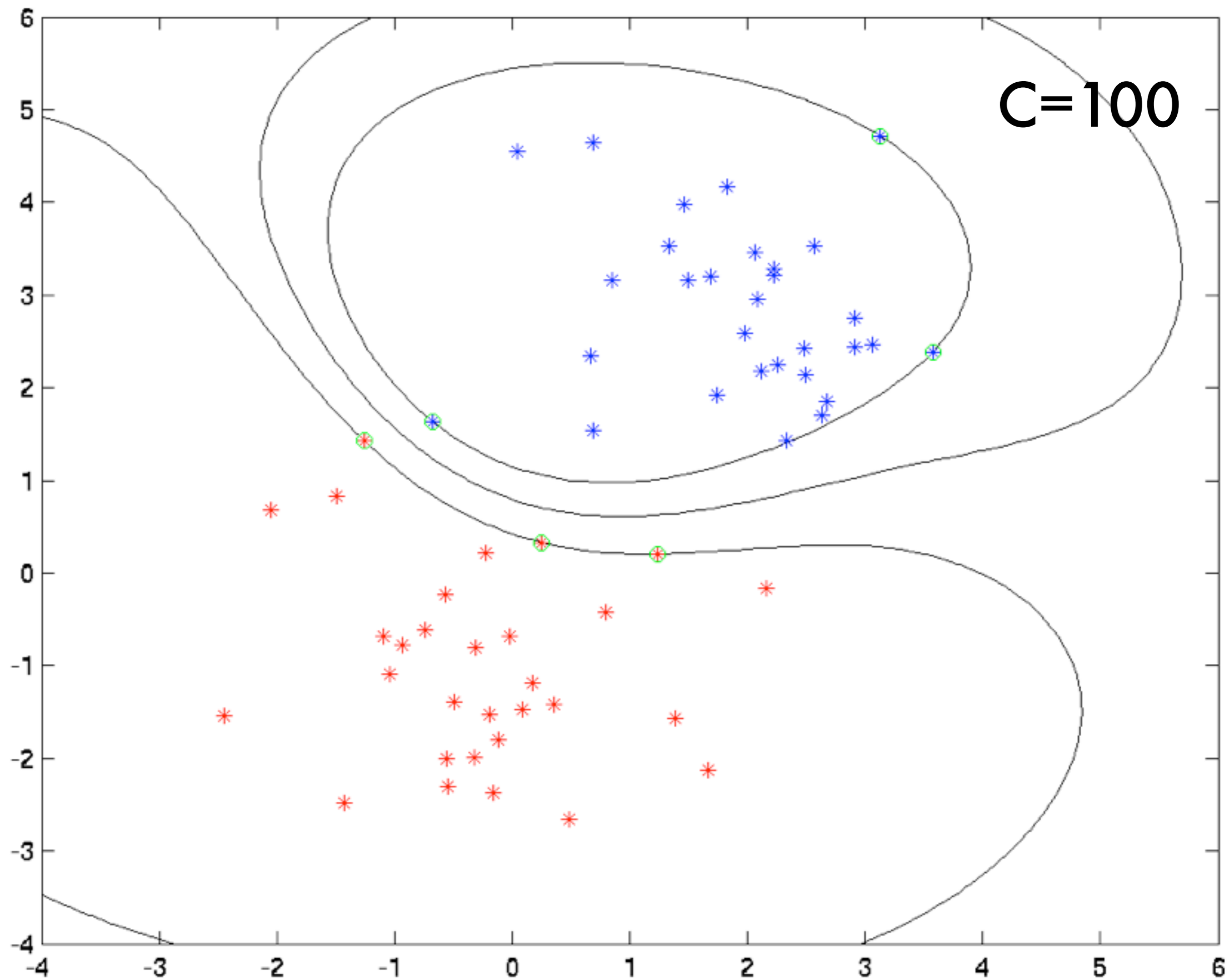


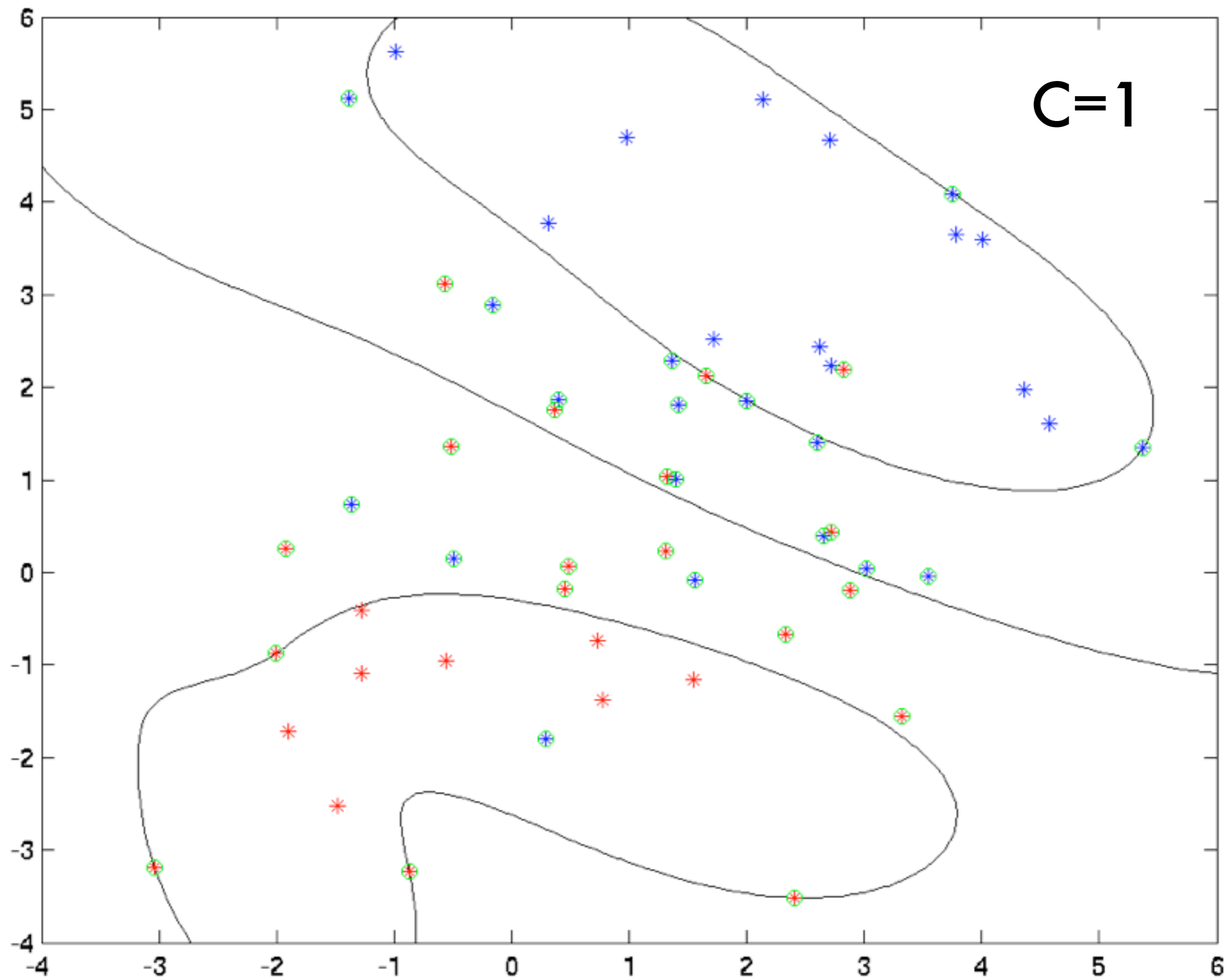


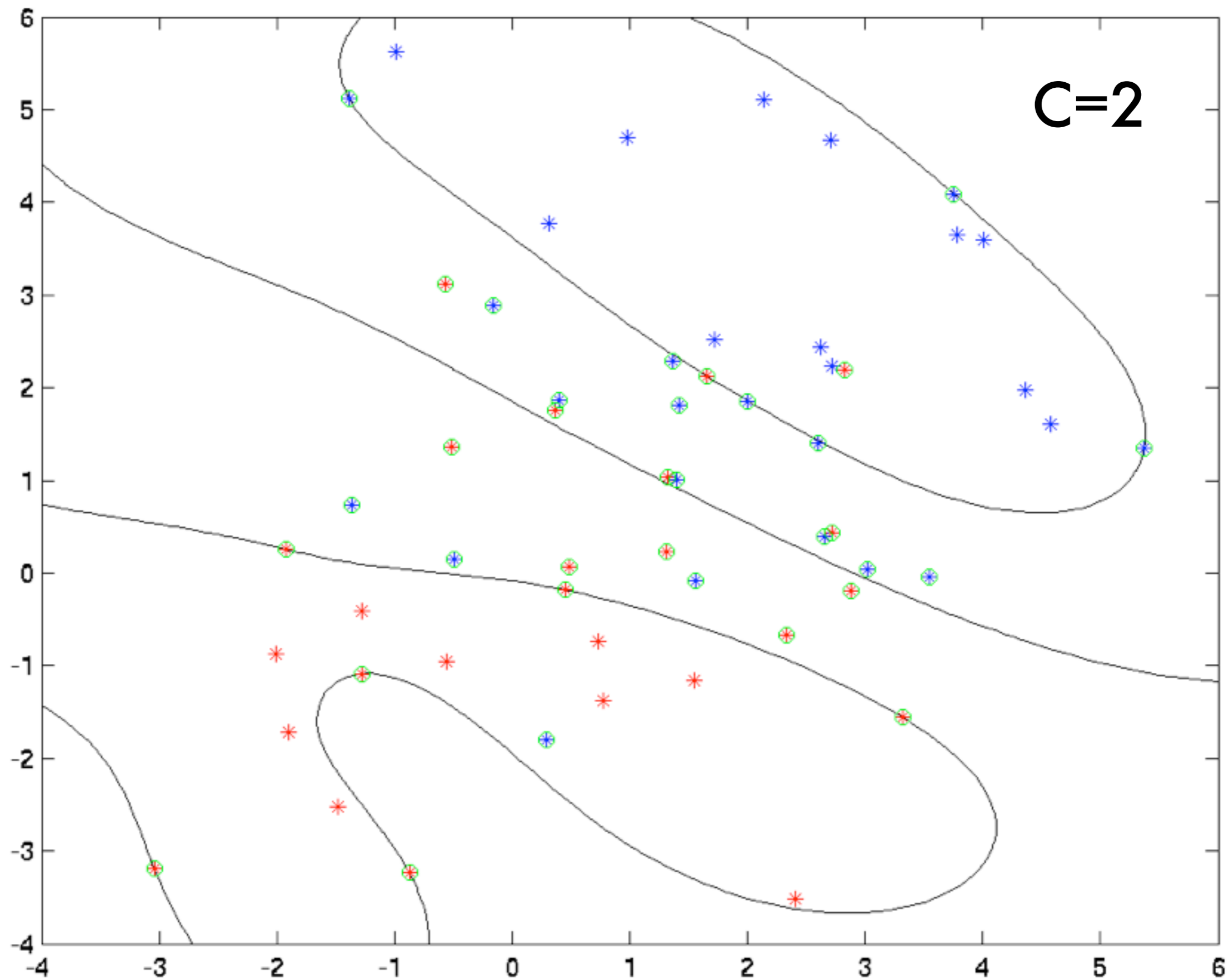


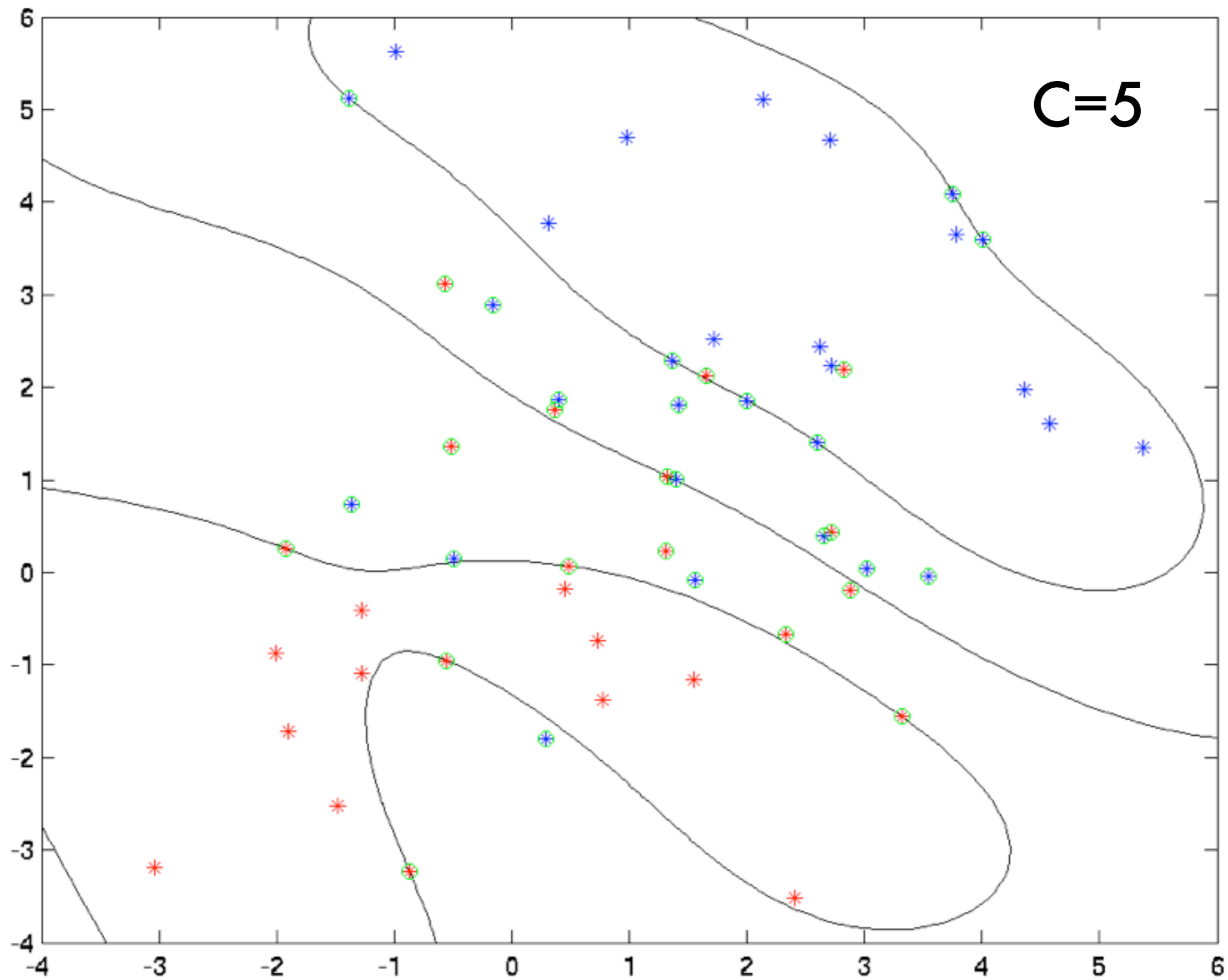


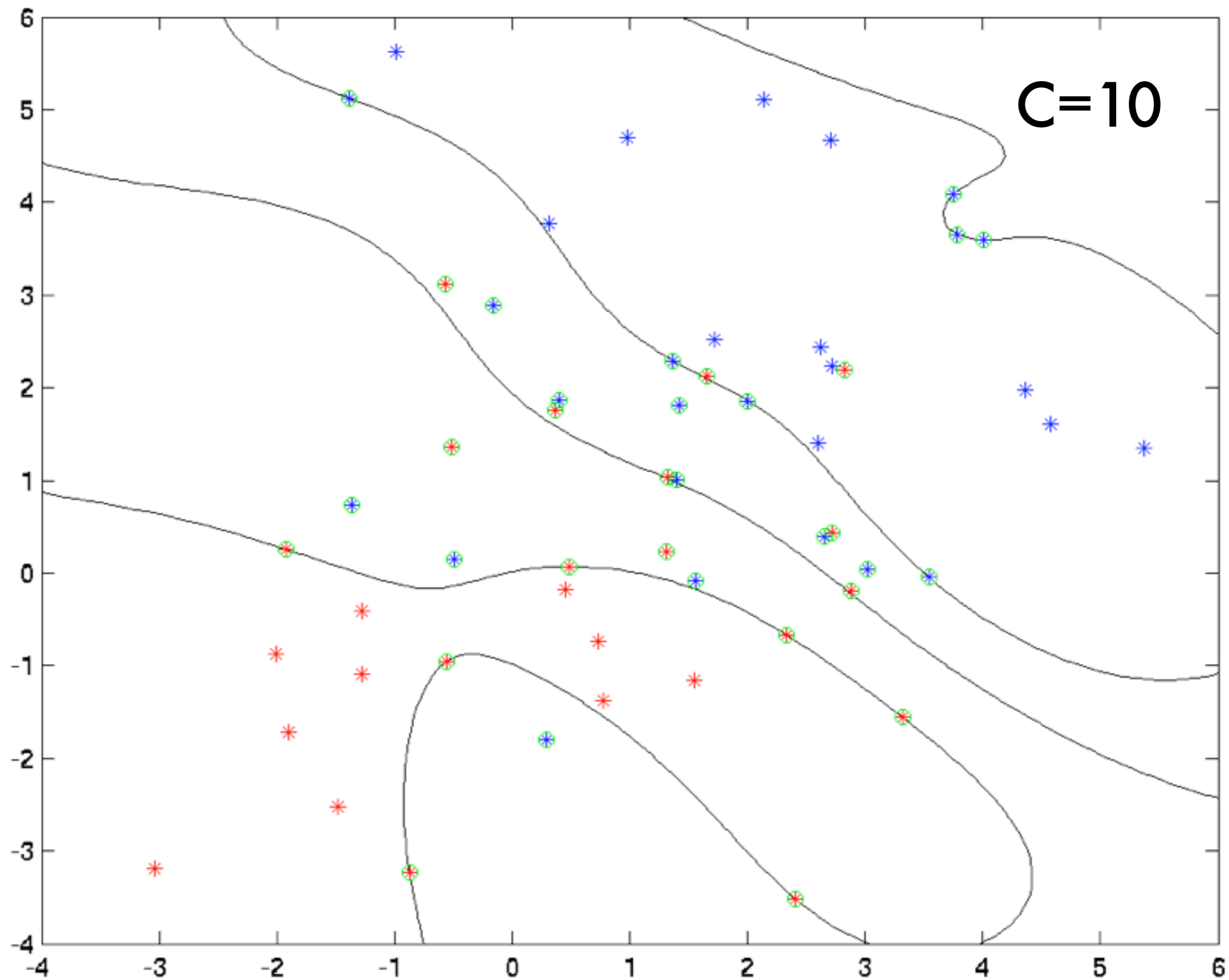


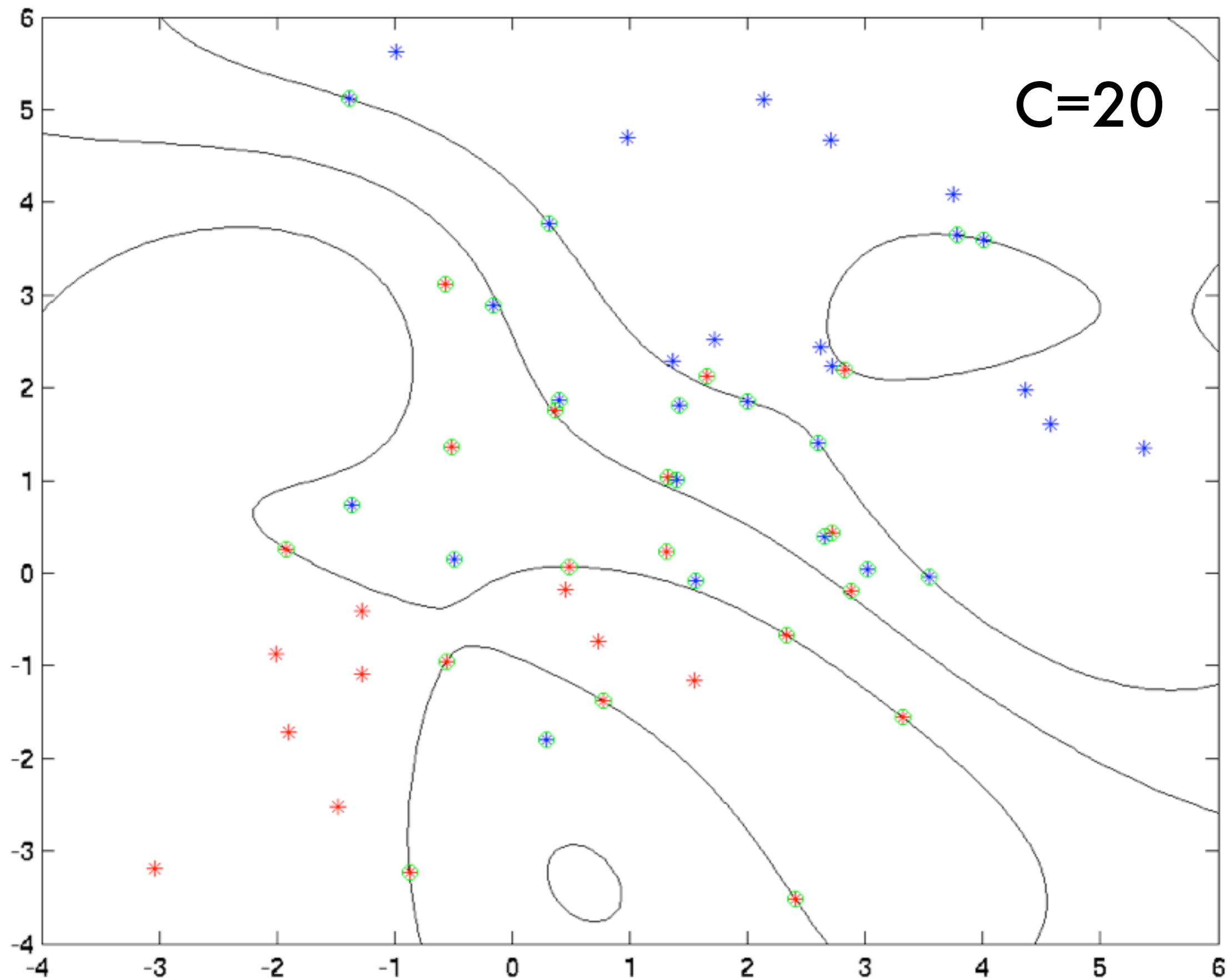


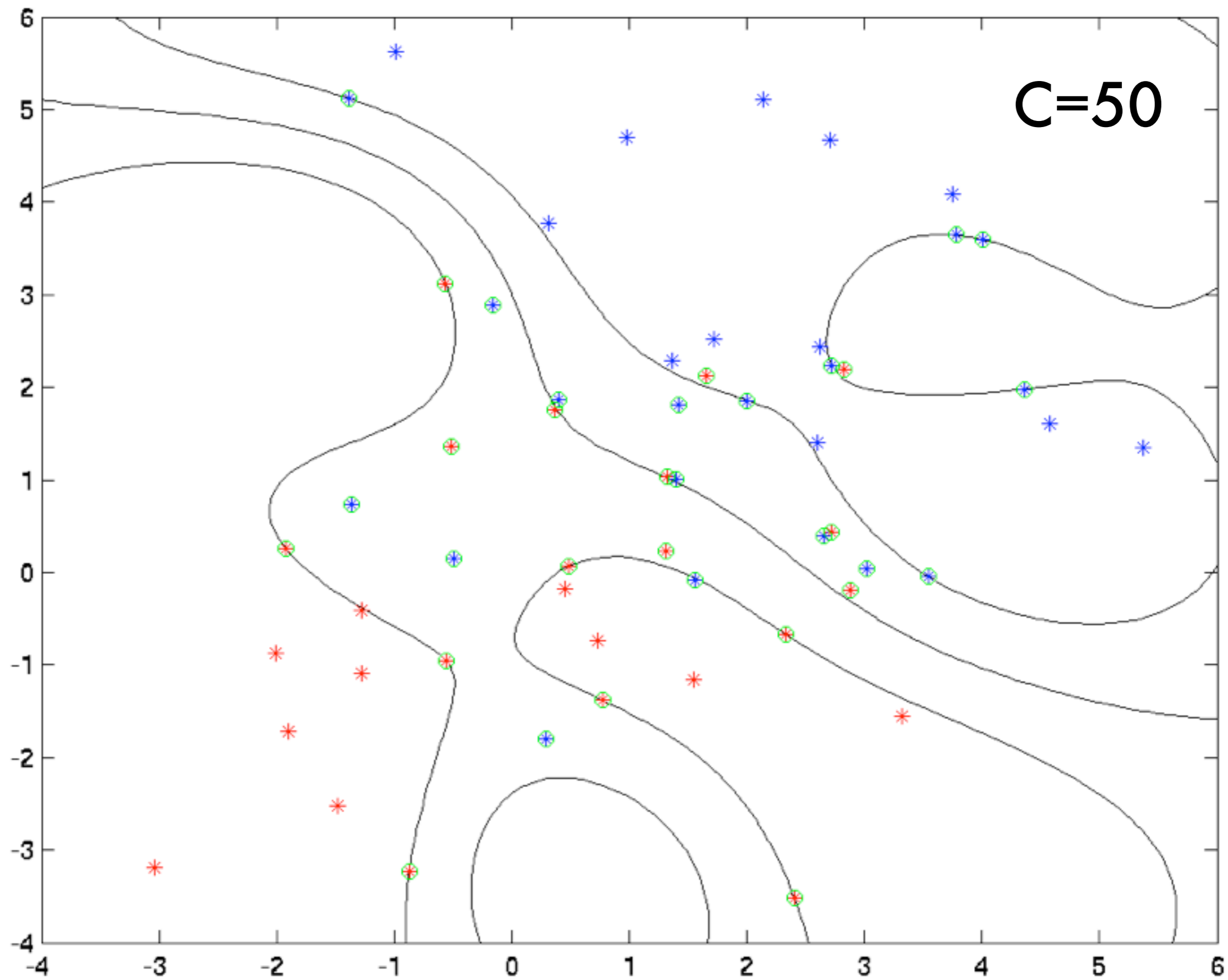


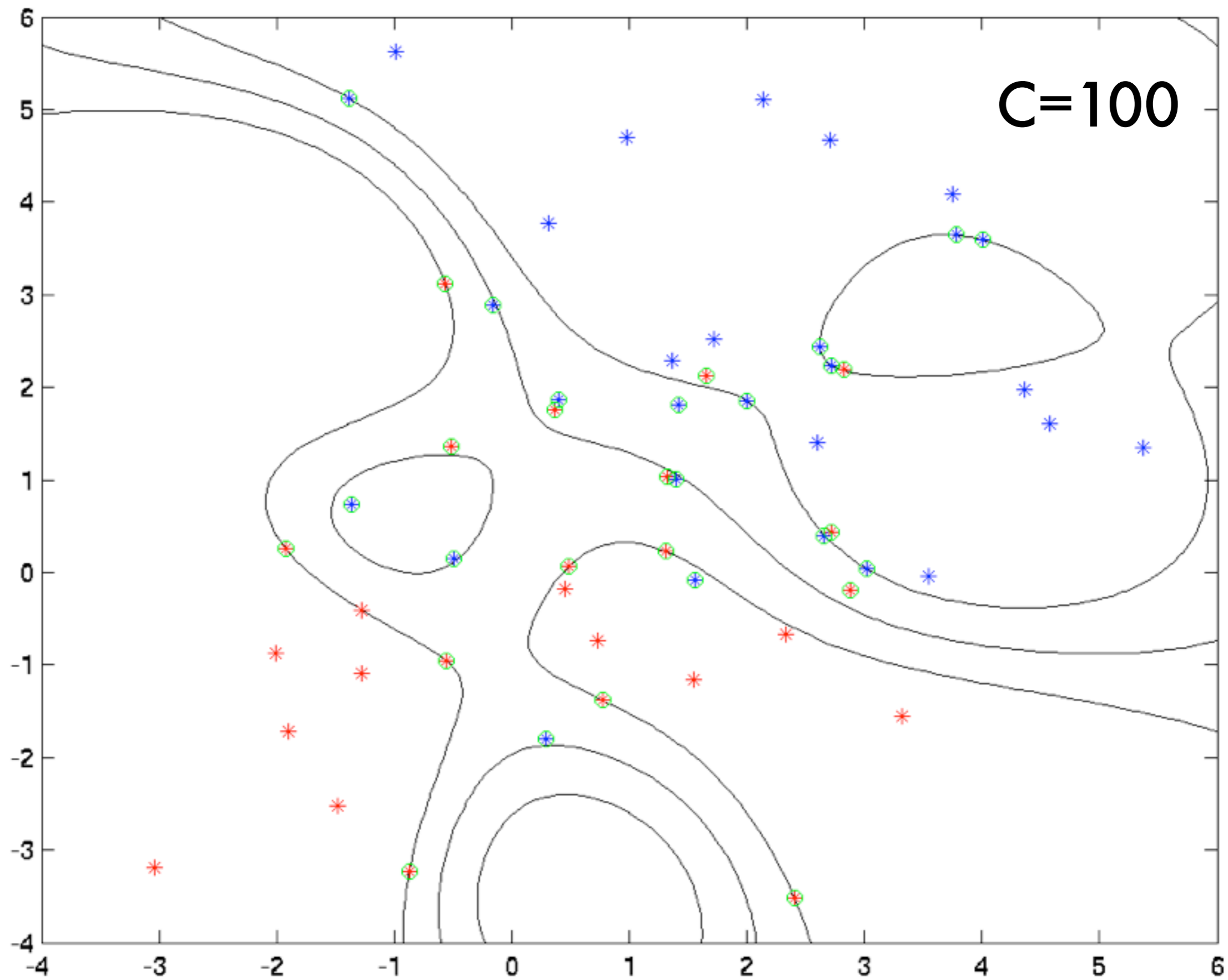




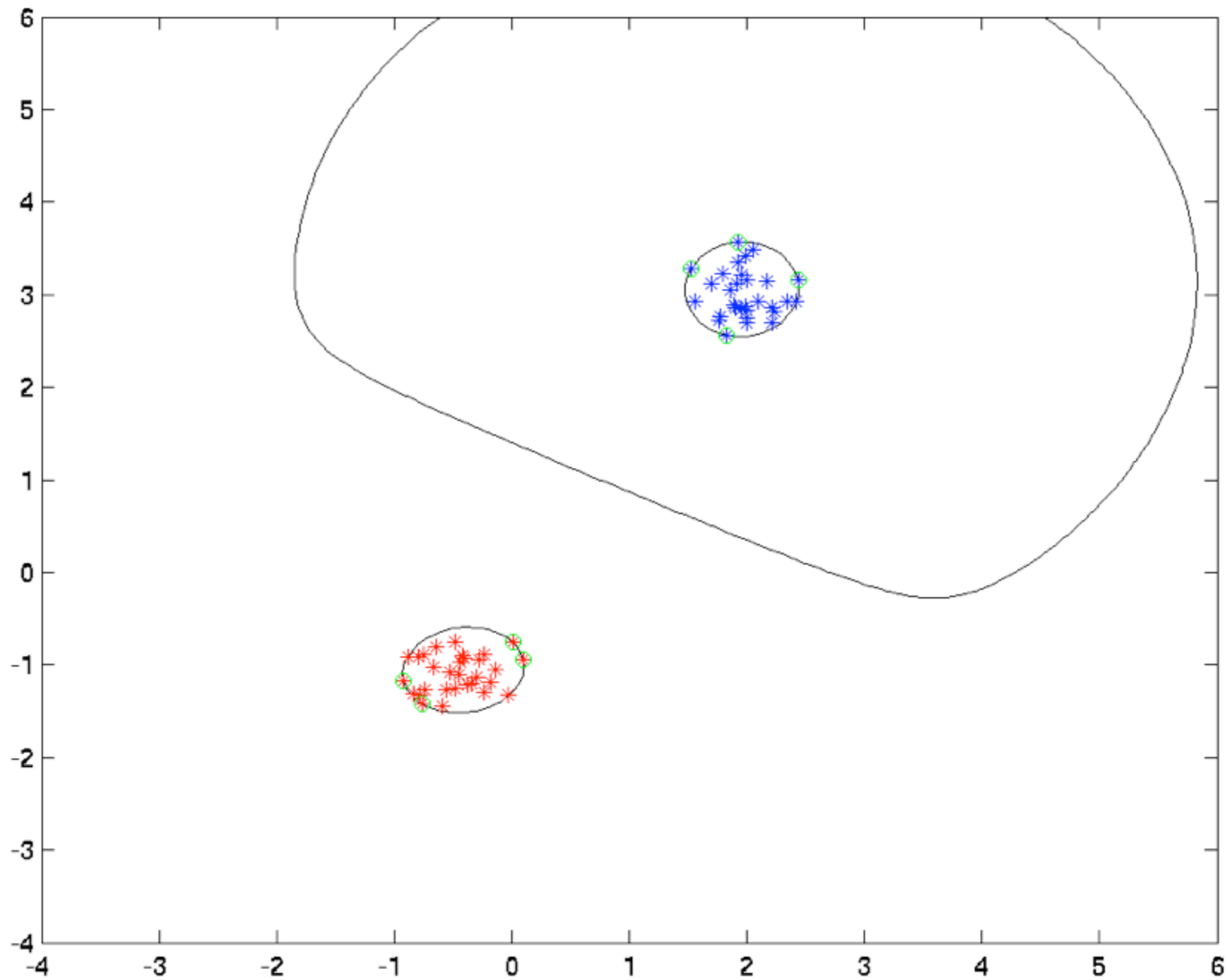


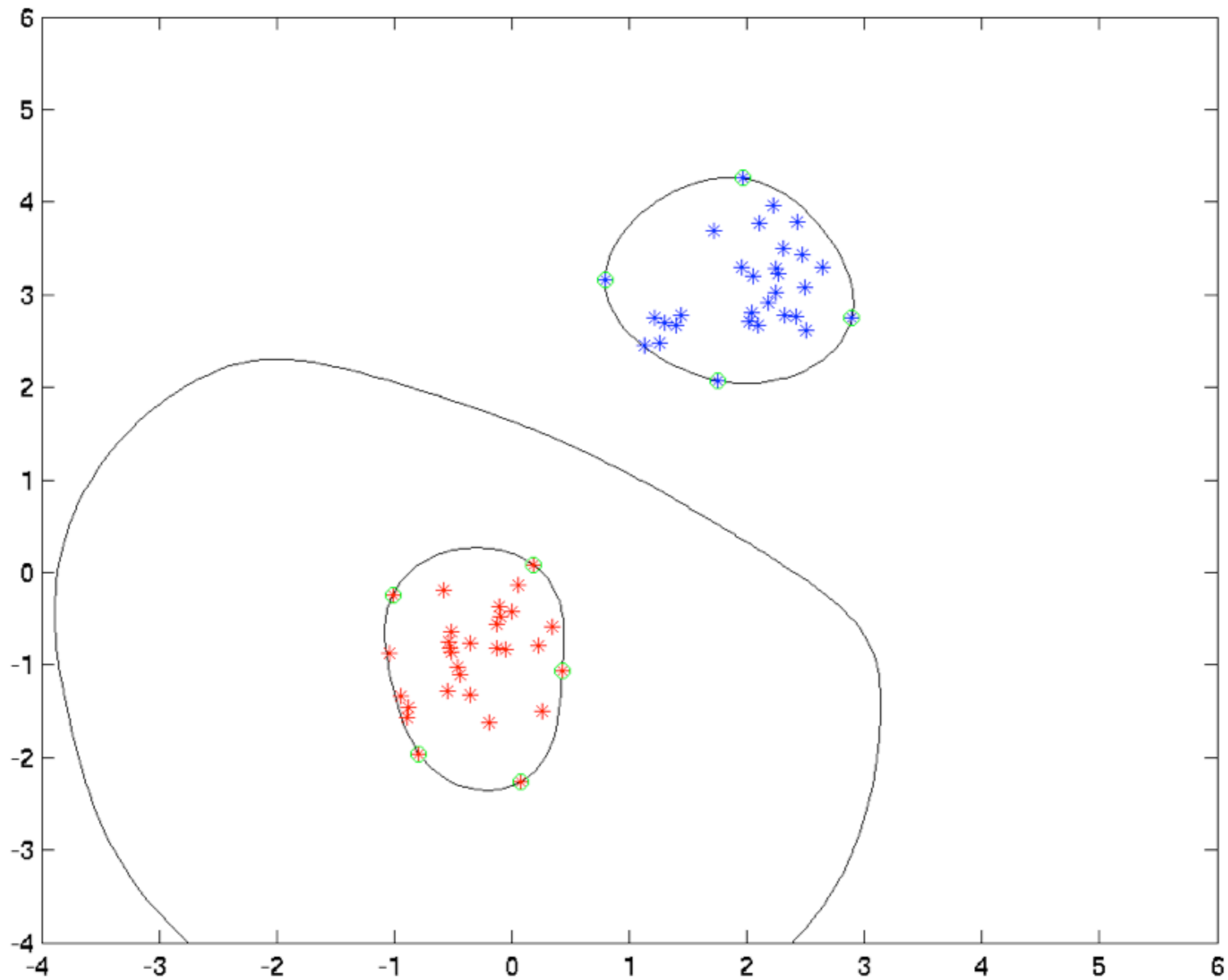


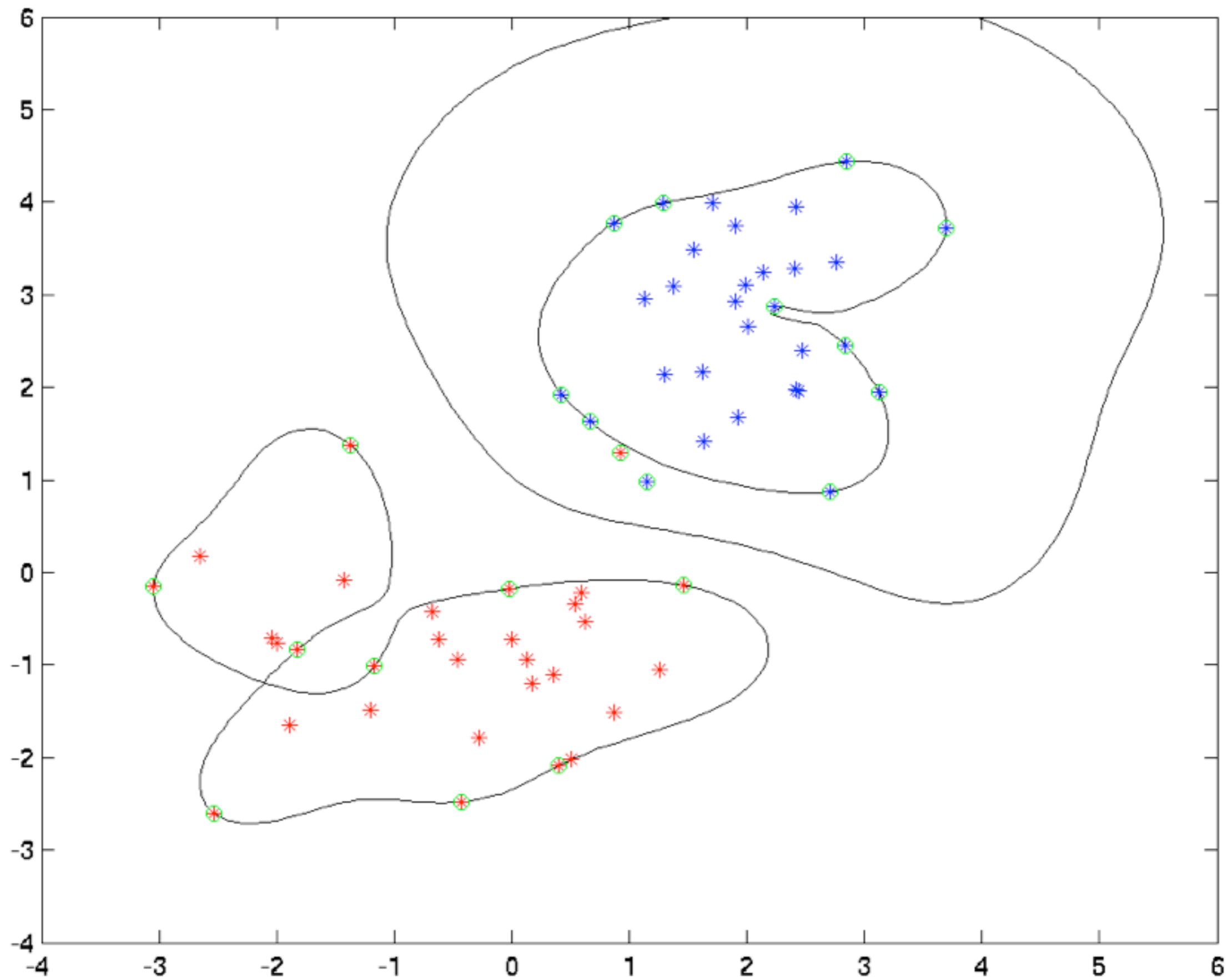


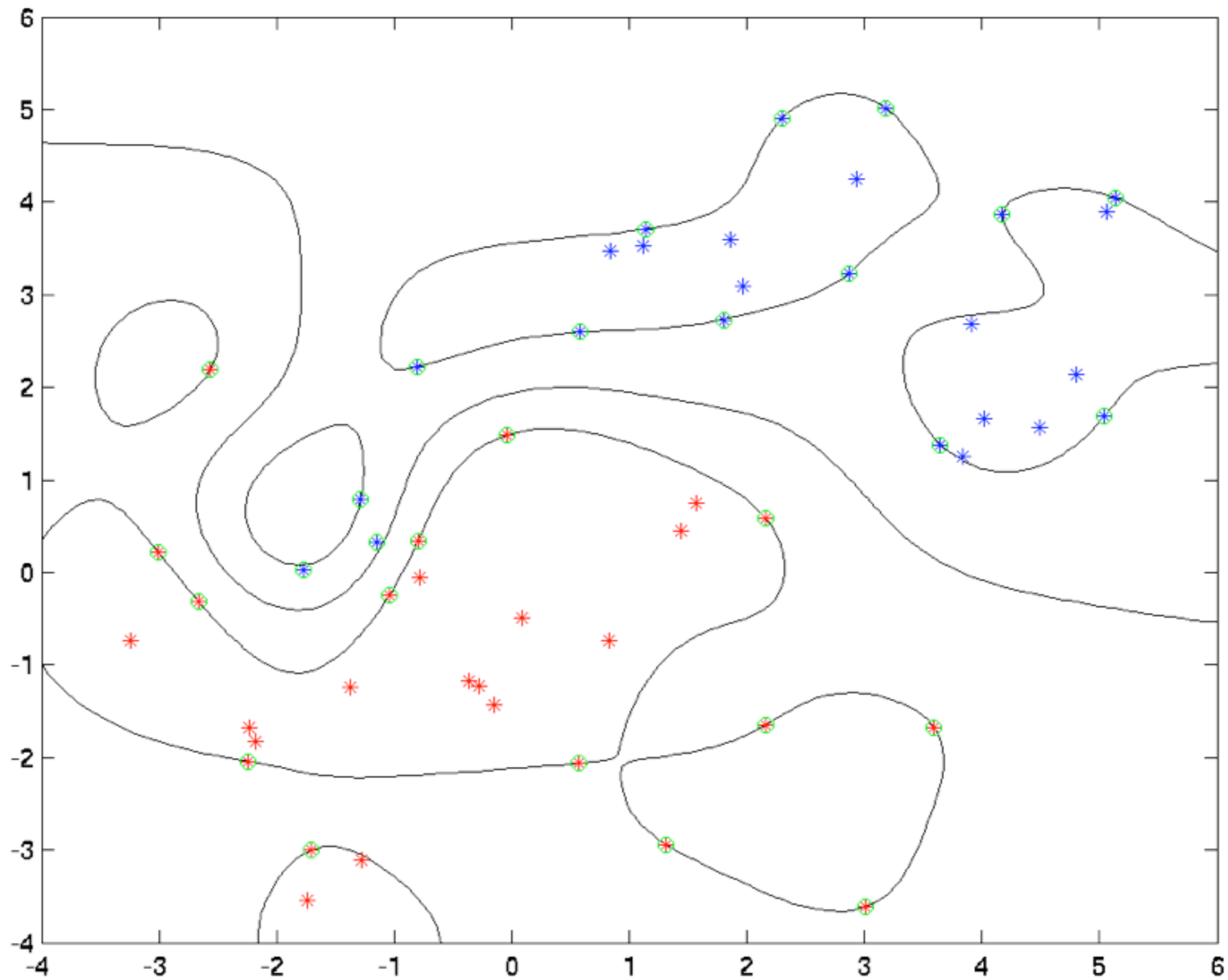


And now with a narrower kernel

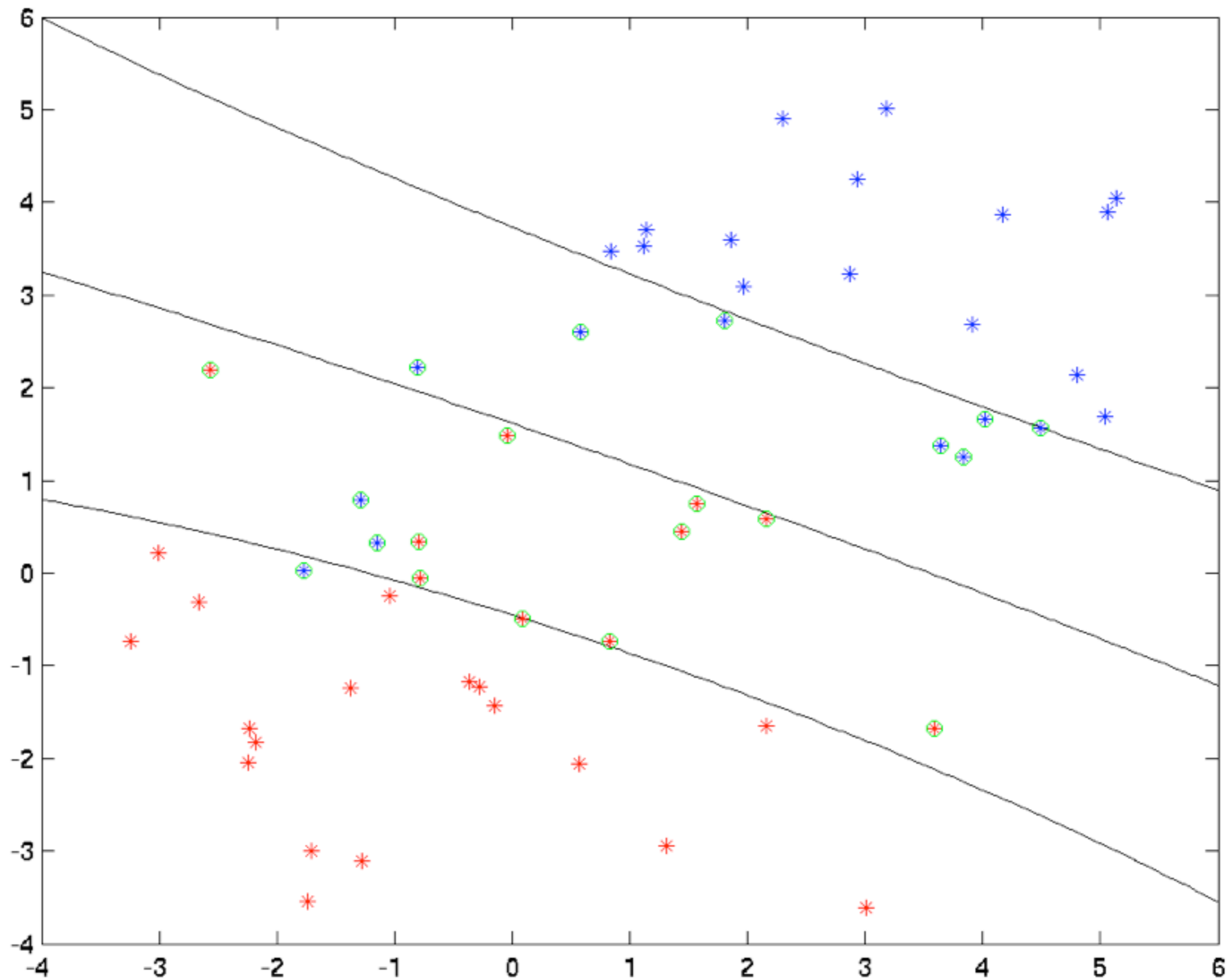




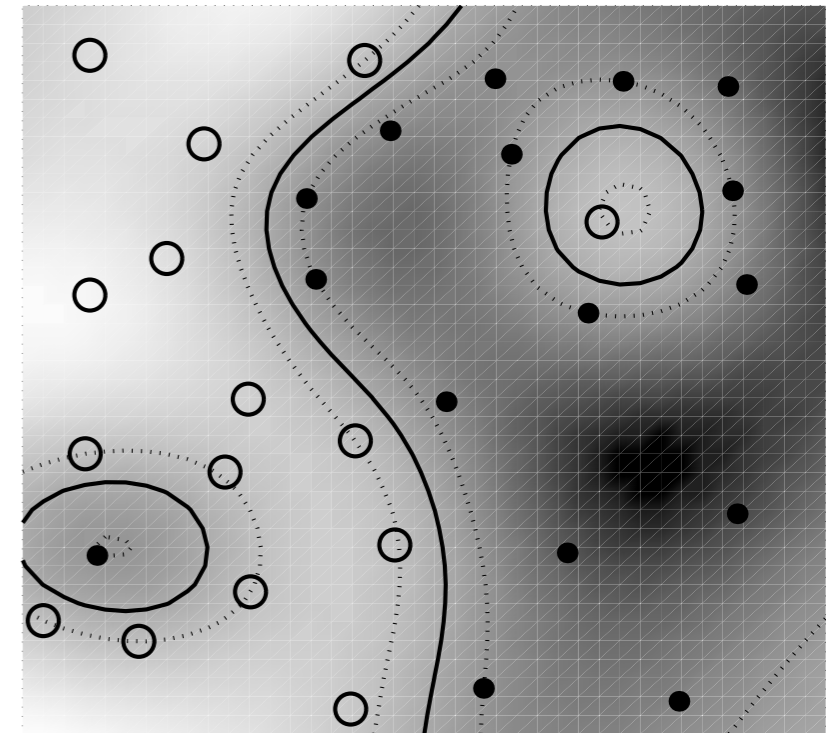
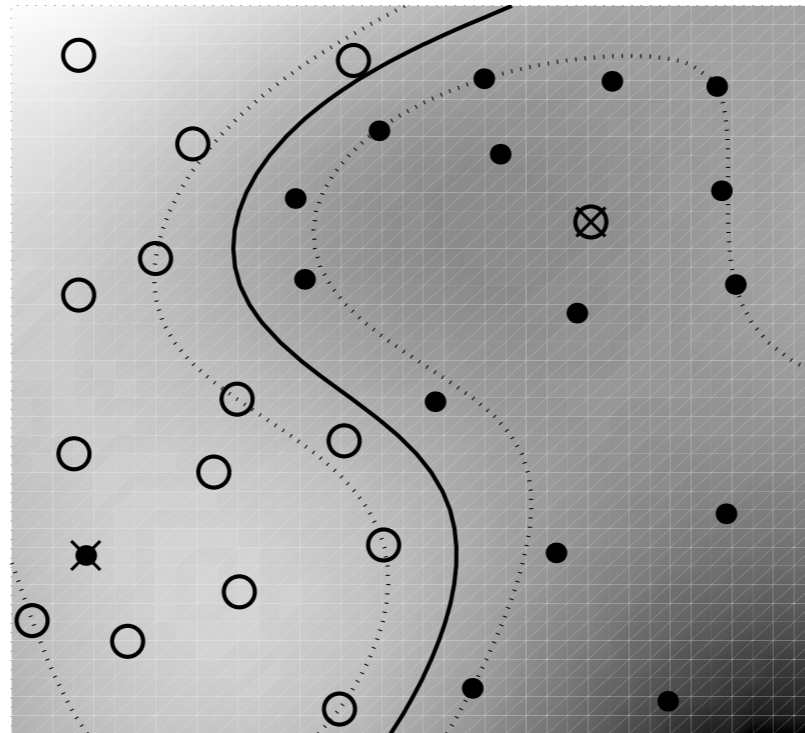
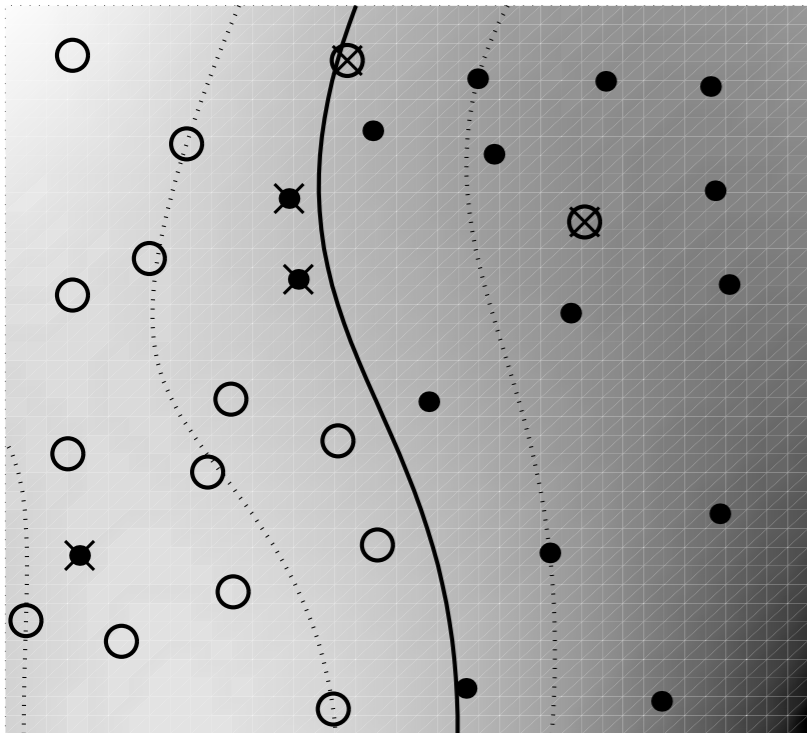




And now with a very wide kernel



Nonlinear separation



- Increasing C allows for more nonlinearities
- Decreases number of errors
- SV boundary need not be contiguous
- Kernel width adjusts function class



MAGIC Etch A Sketch® SCREEN



Risk and Loss

OHIO ART 'The World of Toys'
MAGIC SCREEN IS GLASS SET IN DURABLE PLASTIC FRAME
USE WITH CARE

Loss function point of view

- **Constrained quadratic program**

$$\underset{w,b}{\text{minimize}} \quad \frac{1}{2} \|w\|^2 + C \sum_i \xi_i$$

subject to $y_i [\langle w, x_i \rangle + b] \geq 1 - \xi_i$ and $\xi_i \geq 0$

- **Risk minimization setting**

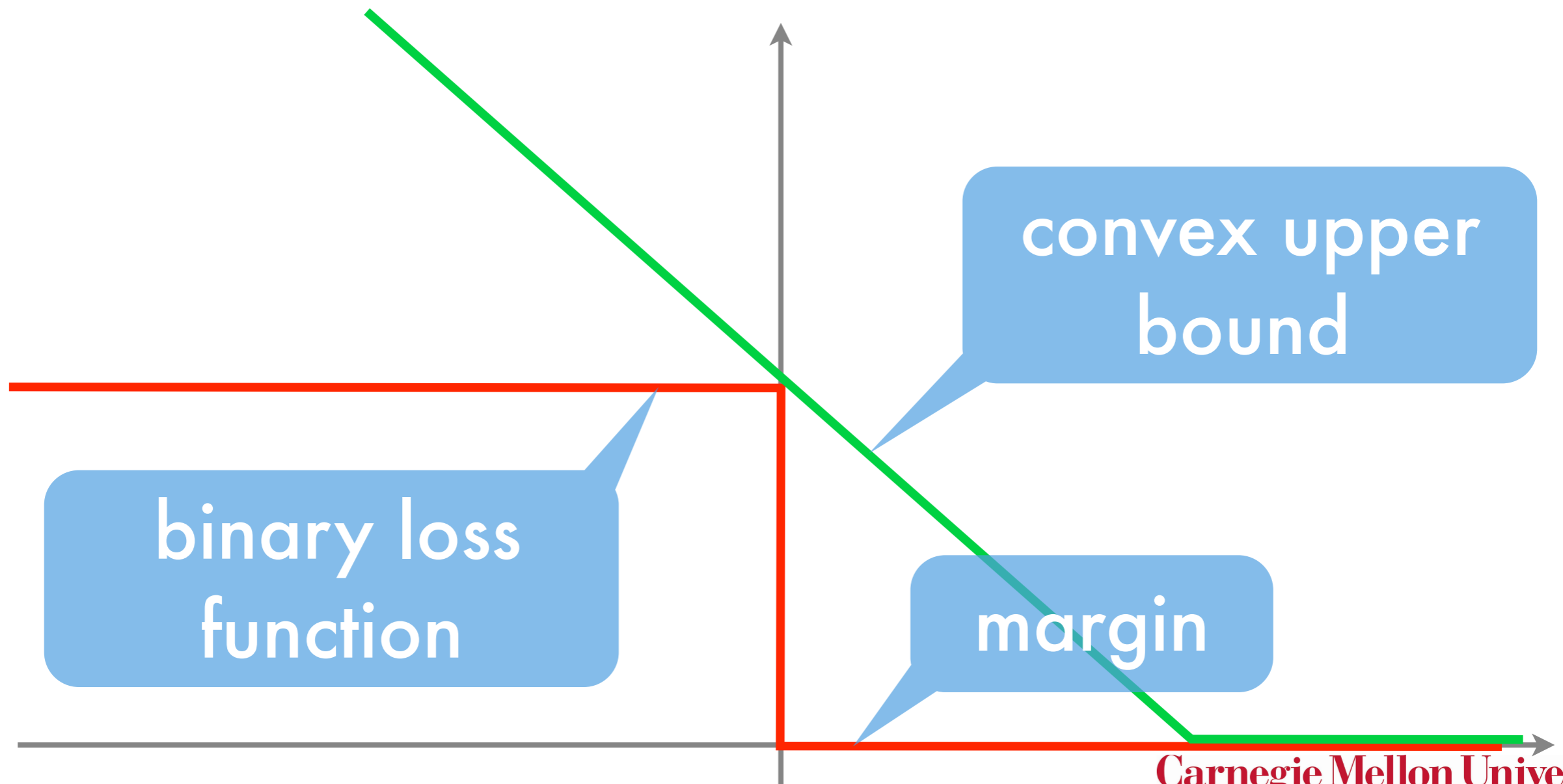
$$\underset{w,b}{\text{minimize}} \quad \frac{1}{2} \|w\|^2 + C \sum_i \max [0, 1 - y_i [\langle w, x_i \rangle + b]]$$

empirical risk

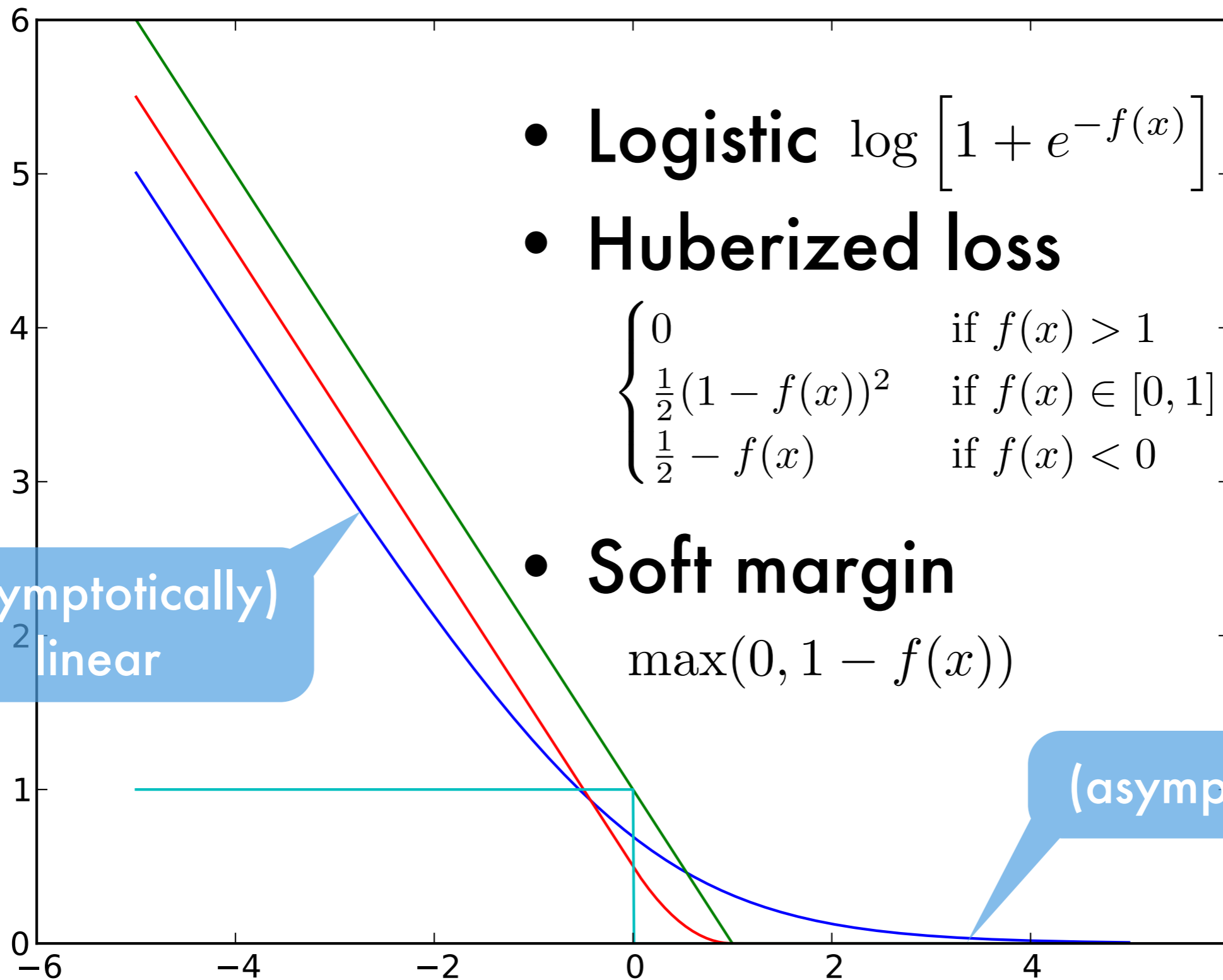
Follows from finding minimal slack variable for given (w,b) pair.

Soft margin as proxy for binary

- **Soft margin loss** $\max(0, 1 - yf(x))$
- **Binary loss** $\{yf(x) < 0\}$



More loss functions



Risk minimization view

- Find function f minimizing classification error

$$R[f] := \mathbf{E}_{x,y \sim p(x,y)} [\{y f(x) > 0\}]$$

- Compute empirical average

$$R_{\text{emp}}[f] := \frac{1}{m} \sum_{i=1}^m \{y_i f(x_i) > 0\}$$

- Minimization is nonconvex
- Overfitting as we minimize empirical error
- Compute convex upper bound on the loss
- Add regularization for capacity control

$$R_{\text{reg}}[f] := \frac{1}{m} \sum_{i=1}^m \max(0, 1 - y_i f(x_i)) + \lambda \Omega[f]$$

regularization

how to control λ

Summary

- **Support Vector Classification**
Large Margin Separation, optimization problem
- **Properties**
Support Vectors, kernel expansion
- **Soft margin classifier**
Dual problem, robustness