



Learning with Kernels

10-715 Fall 2015

Alexander Smola
alex@smola.org

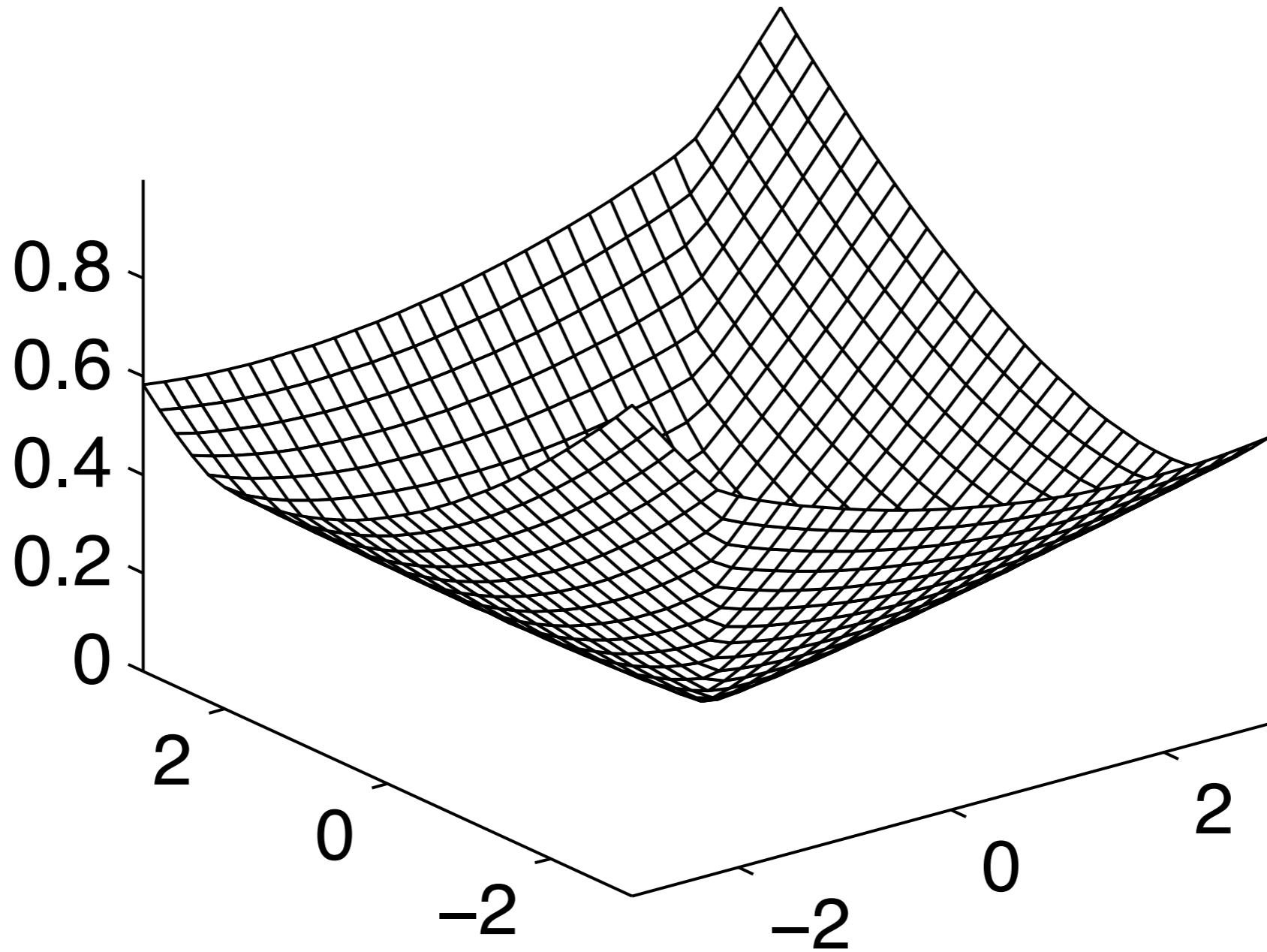
Office hours - after class & my office

Outline

- **Convex Optimization**
 - Unconstrained Optimization
 - Constrained Optimization and Duality
 - Linear and Quadratic programs
- **Support Vector Machines**
 - Classification
 - Regression
 - Novelty Detection
- **Kernels**
 - Feature Space
 - Kernel PCA
 - Kernelized SVM

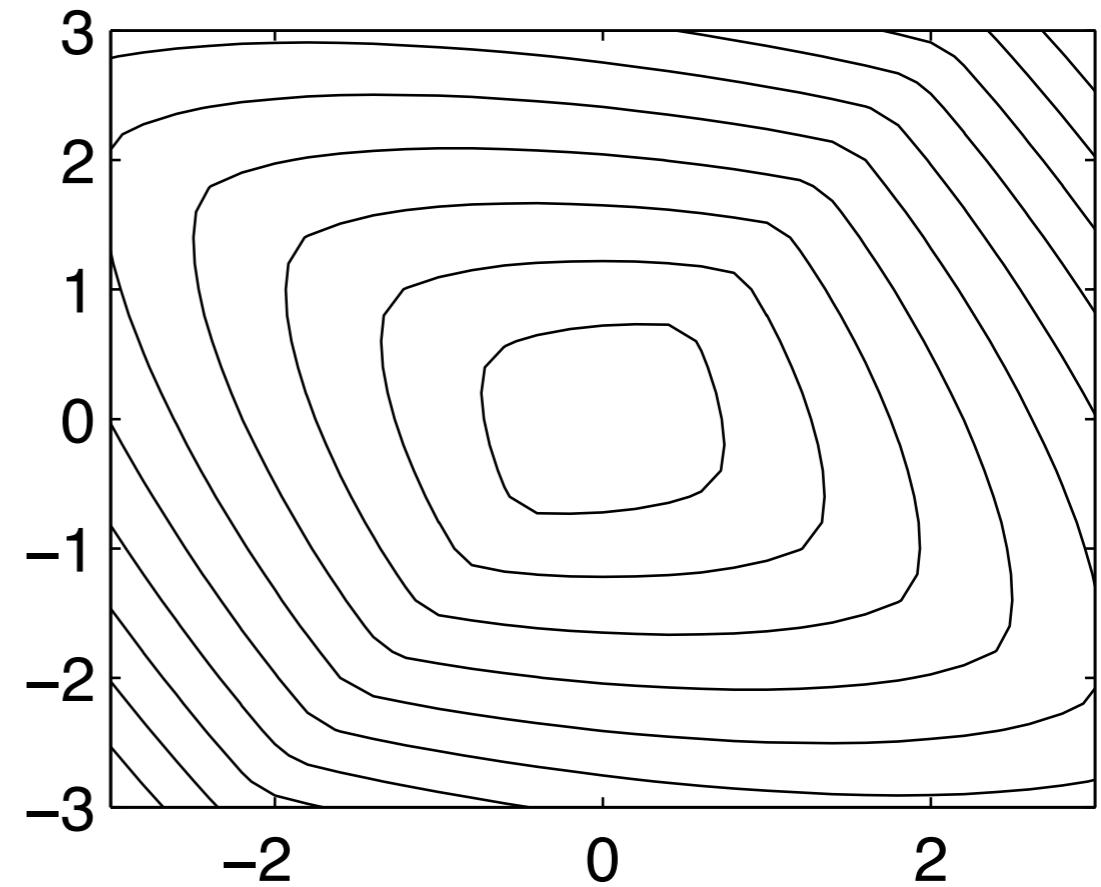
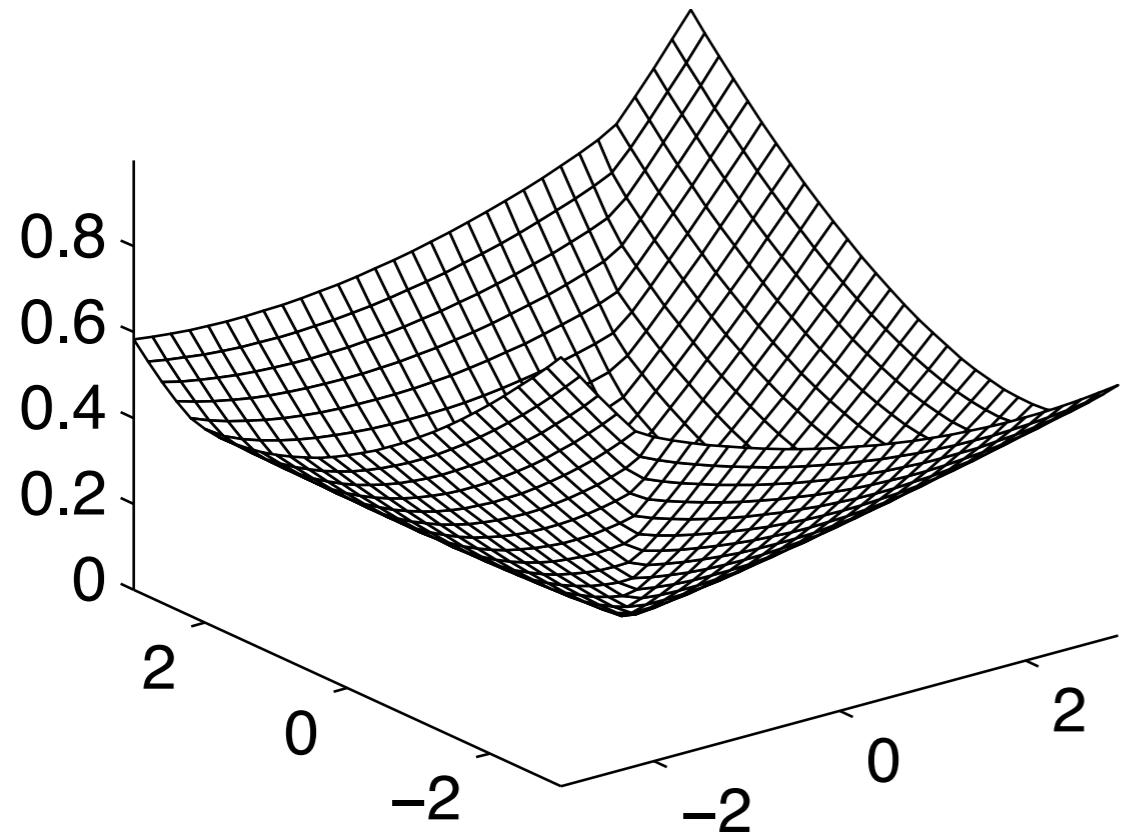
The background image shows a wide, calm lake in the foreground, with a shoreline that curves around a small, sandy peninsula. Beyond the lake, there are rolling hills covered in dry, golden-brown vegetation. In the far distance, a range of mountains is visible against a light blue sky with scattered clouds.

Unconstrained Optimization



Convexity 101

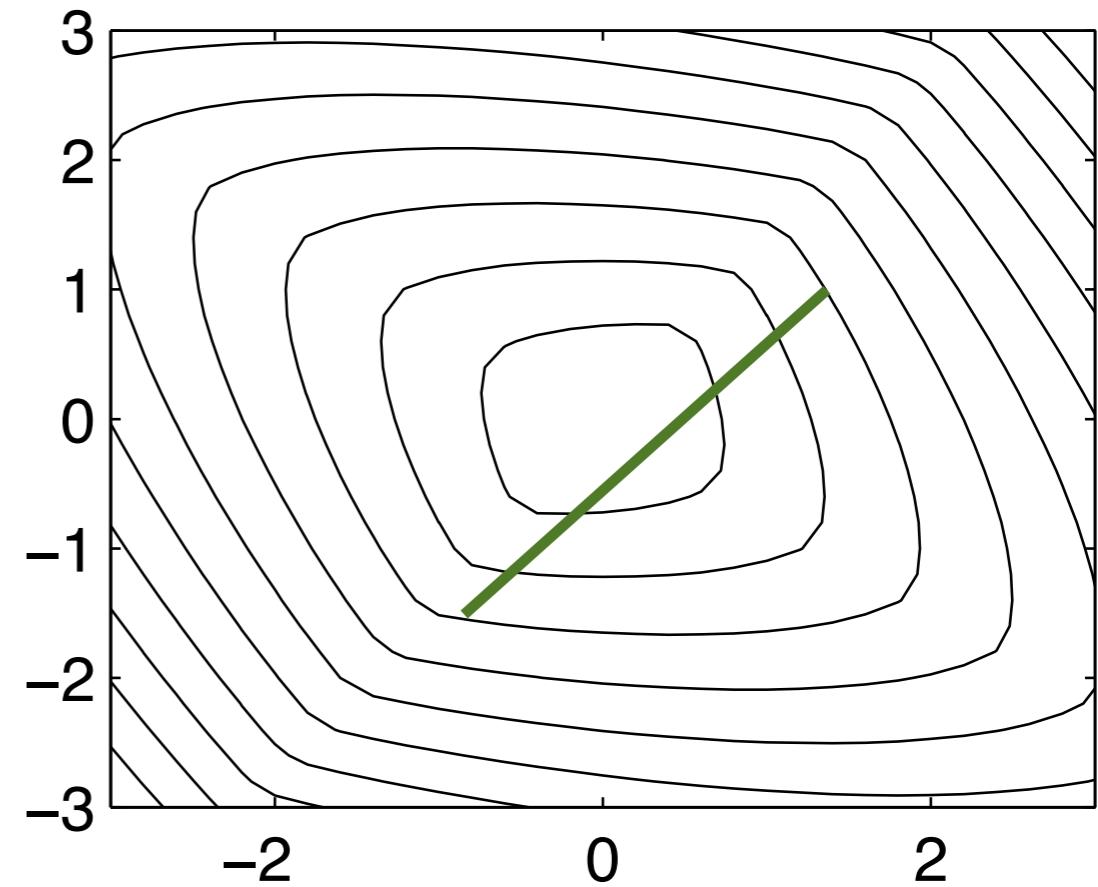
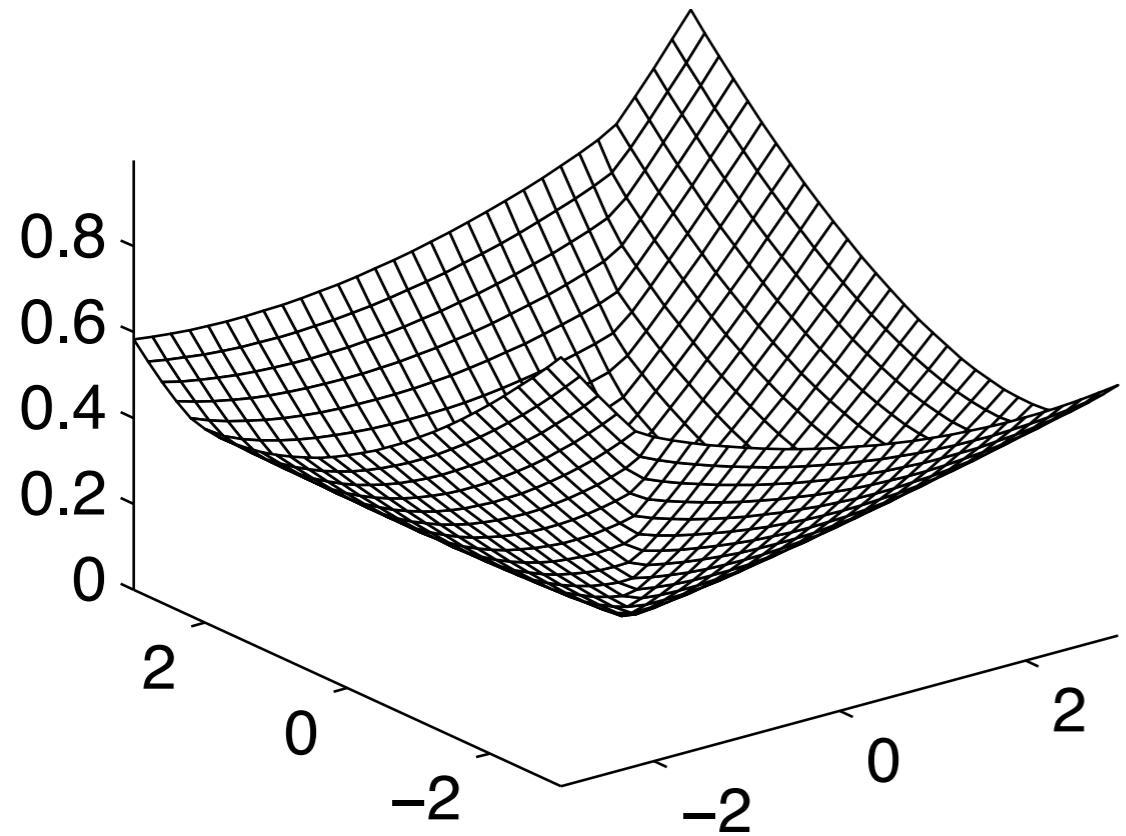
Convexity 101



- **Convex set**
For $x, x' \in X$ it follows that $\lambda x + (1 - \lambda)x' \in X$ for $\lambda \in [0, 1]$
- **Convex function**

$$\lambda f(x) + (1 - \lambda)f(x') \geq f(\lambda x + (1 - \lambda)x') \text{ for } \lambda \in [0, 1]$$

Convexity 101

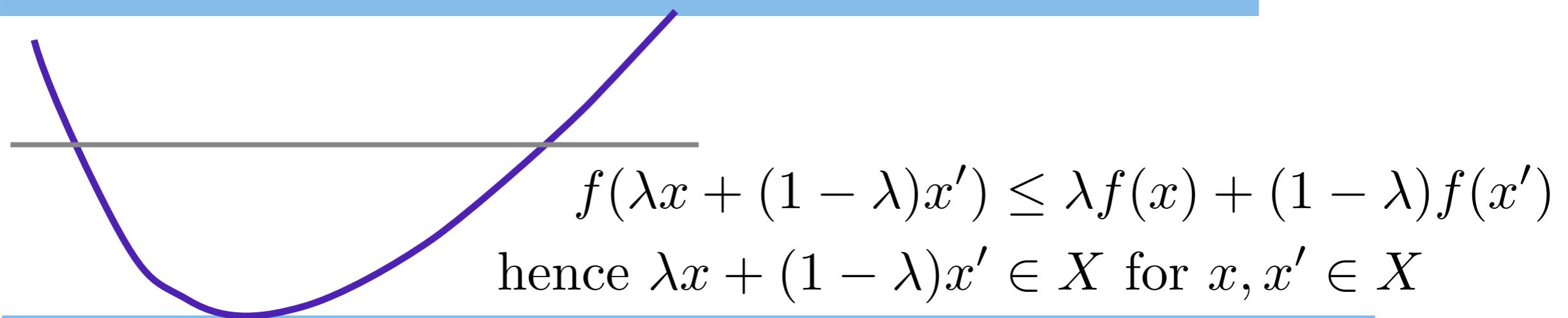


- **Convex set**
For $x, x' \in X$ it follows that $\lambda x + (1 - \lambda)x' \in X$ for $\lambda \in [0, 1]$
- **Convex function**

$$\lambda f(x) + (1 - \lambda)f(x') \geq f(\lambda x + (1 - \lambda)x') \text{ for } \lambda \in [0, 1]$$

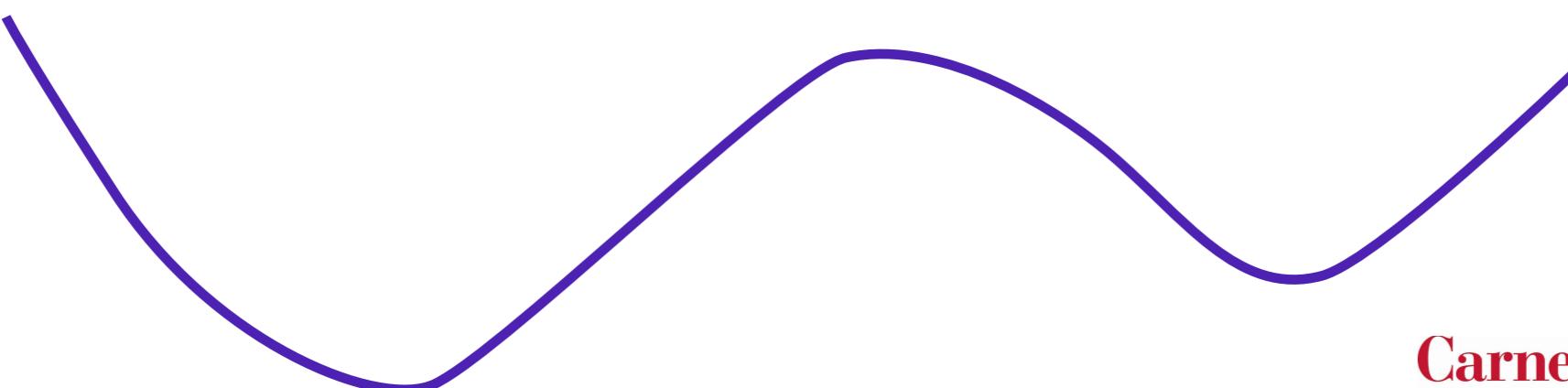
Convexity 101

- Below-set of convex function is convex



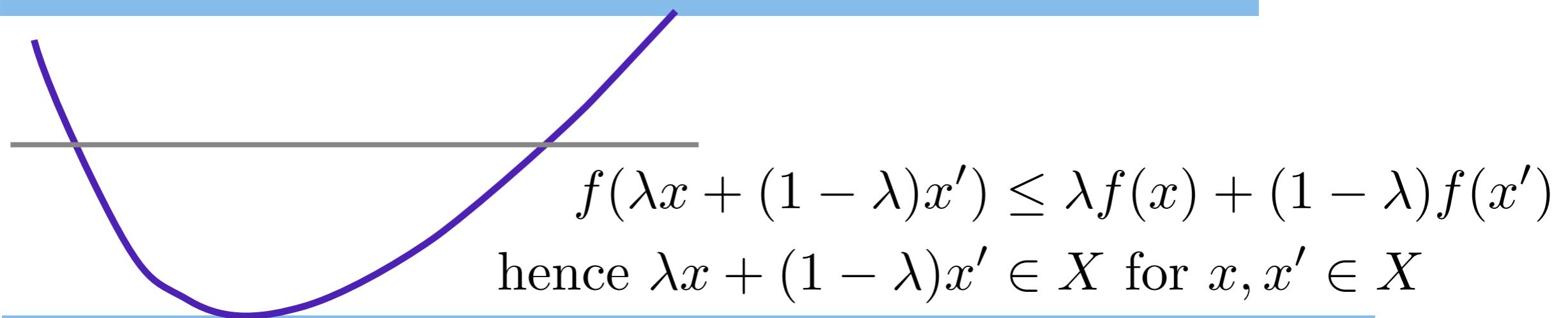
- Convex functions don't have local minima

Proof by contradiction - linear interpolation breaks local minimum condition



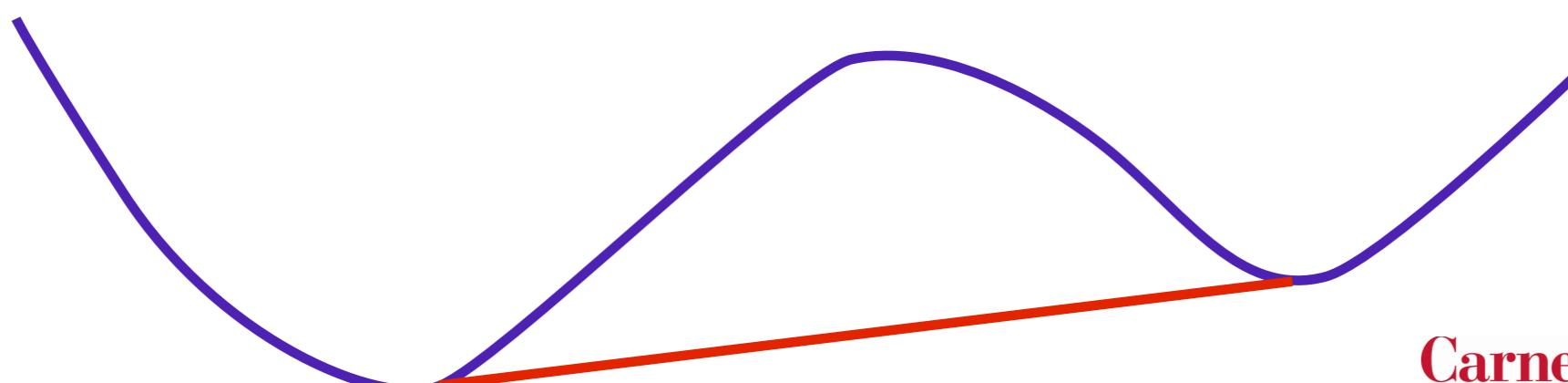
Convexity 101

- Below-set of convex function is convex



- Convex functions don't have local minima

Proof by contradiction - linear interpolation breaks local minimum condition



Convexity 101

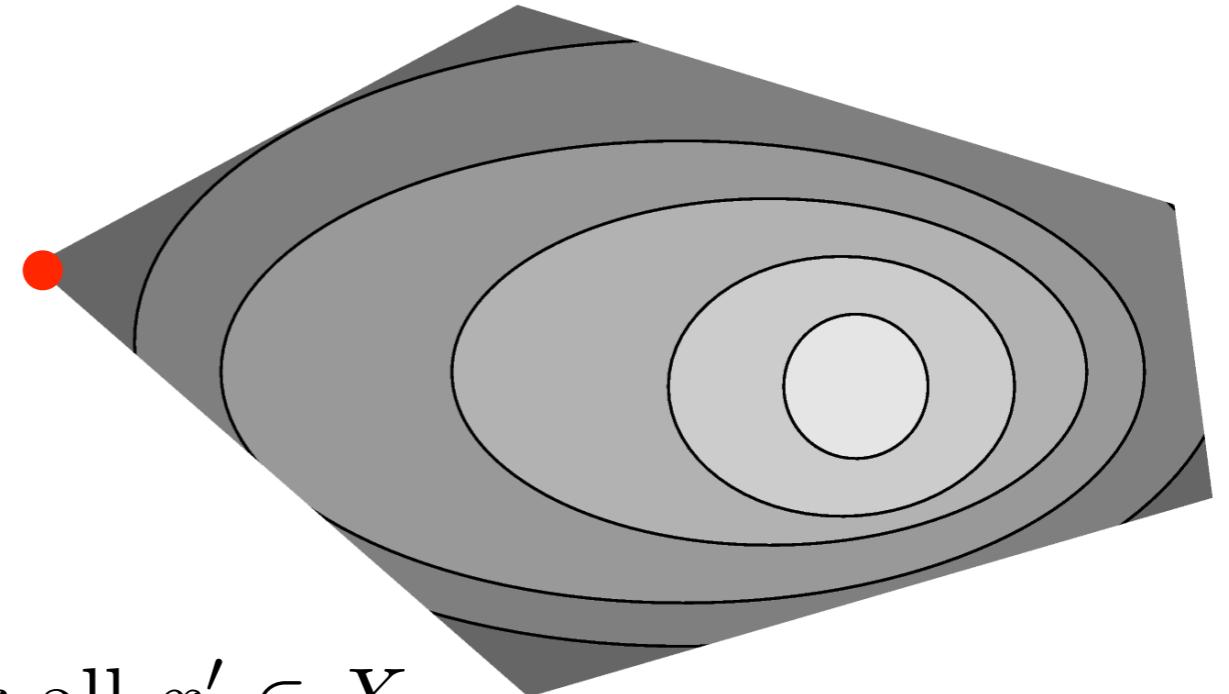
- Vertex of a convex set
Point which cannot
be extrapolated
within convex set

$$\lambda x + (1 - \lambda)x' \notin X \text{ for } \lambda > 1 \text{ for all } x' \in X$$

- Convex hull

$$\text{co } X := \left\{ \bar{x} \left| \bar{x} = \sum_{i=1}^n \alpha_i x_i \text{ where } n \in \mathbb{N}, \alpha_i \geq 0 \text{ and } \sum_{i=1}^n \alpha_i \leq 1 \right. \right\}$$

- Convex hull of set is a convex set



Convexity 101

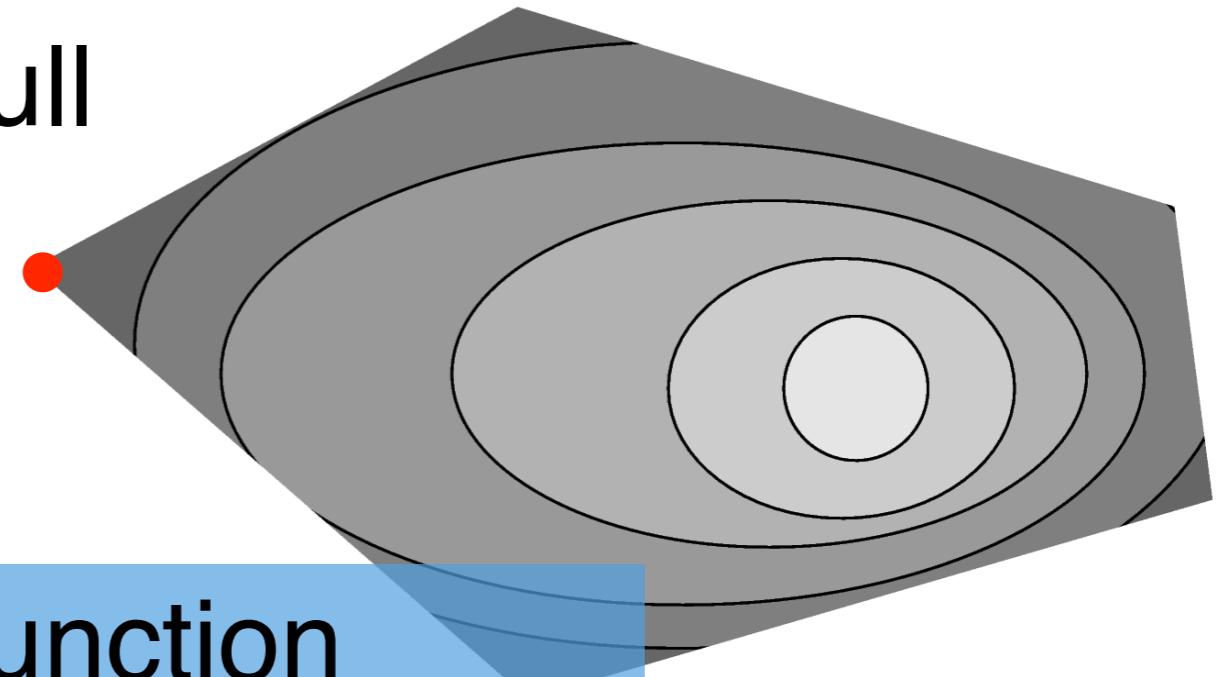
- Supremum on convex hull

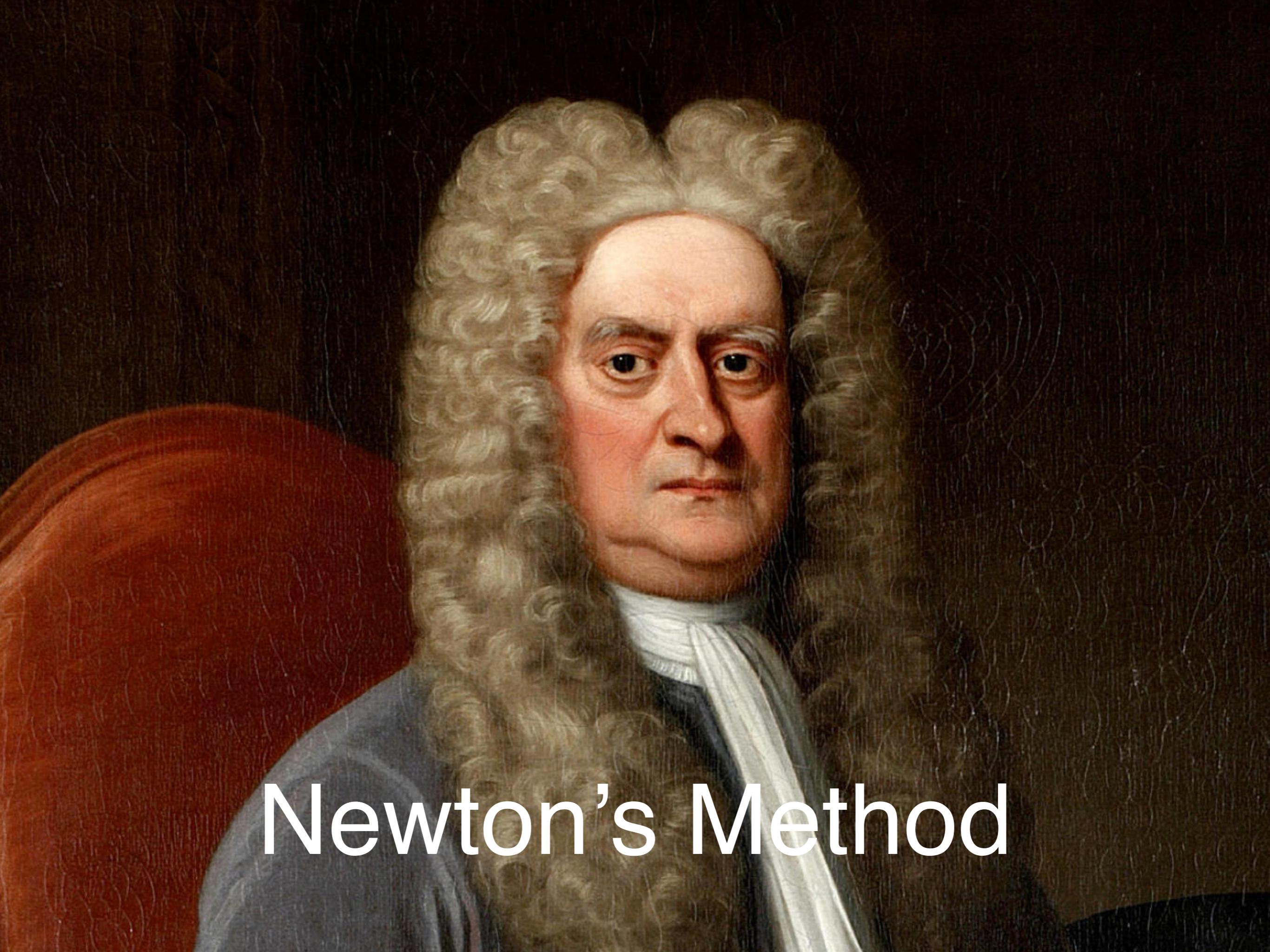
$$\sup_{x \in X} f(x) = \sup_{x \in \text{co}X} f(x)$$

Proof by contradiction

- Maximum over convex function on convex set is obtained on vertex

- Assume that maximum inside line segment
 - Then function cannot be convex
 - Hence it must be on vertex





Newton's Method

Newton Method

- Convex objective function f
- Nonnegative second derivative

$$\partial_x^2 f(x) \succeq 0$$

- Taylor expansion

$$f(x + \delta) = f(x) + \langle \delta, \partial_x f(x) \rangle + \frac{1}{2} \delta^\top \partial_x^2 f(x) \delta + O(\delta^3)$$

gradient

Hessian

- Minimize approximation & iterate until converged

$$x \leftarrow x - [\partial_x^2 f(x)]^{-1} \partial_x f(x)$$

Convergence Analysis

- There exists a region around optimality where Newton's method converges quadratically if f is twice continuously differentiable
- For some region around x^* gradient is well approximated by Taylor expansion

$$\|\partial_x f(x^*) - \partial_x f(x) - \langle x^* - x, \partial_x^2 f(x) \rangle\| \leq \gamma \|x^* - x\|^2$$

- Expand Newton update

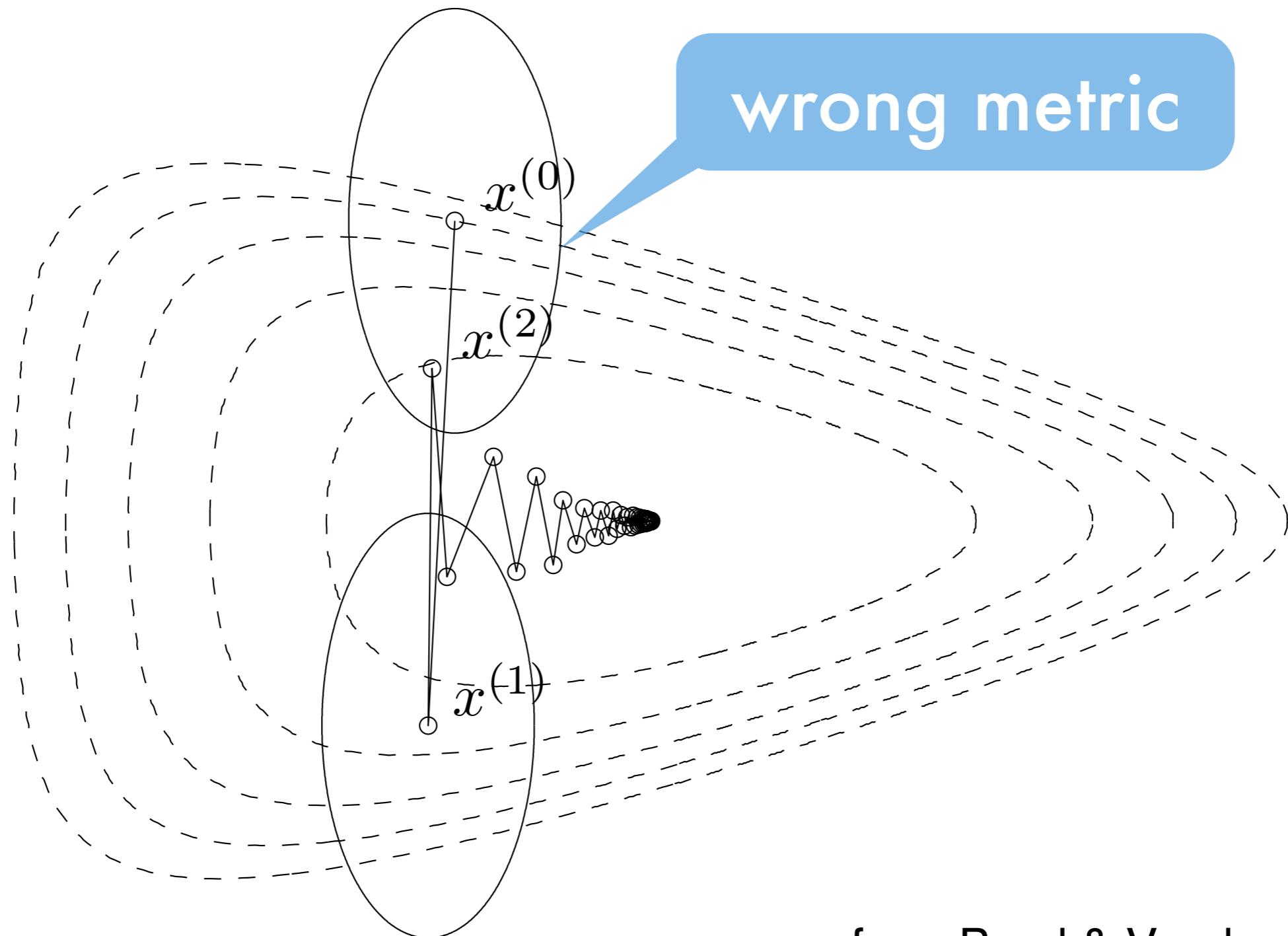
$$\begin{aligned}\|x_{n+1} - x^*\| &= \left\| x_n - x^* - [\partial_x^2 f(x_n)]^{-1} [\partial_x f(x_n) - \partial_x f(x^*)] \right\| \\ &= \left\| [\partial_x^2 f(x_n)]^{-1} [\partial_x^f(x_n)[x_n - x^*] - \partial_x f(x_n) + \partial_x f(x^*)] \right\| \\ &\leq \gamma \left\| [\partial_x^2 f(x_n)]^{-1} \right\| \|x_n - x^*\|^2\end{aligned}$$

Convergence Analysis

- Two convergence regimes
 - As slow as gradient descent outside the region where Taylor expansion is good
$$\|\partial_x f(x^*) - \partial_x f(x) - \langle x^* - x, \partial_x^2 f(x) \rangle\| \leq \gamma \|x^* - x\|^2$$
 - Quadratic convergence once the bound holds
$$\|x_{n+1} - x^*\| \leq \gamma \left\| [\partial_x^2 f(x_n)]^{-1} \right\| \|x_n - x^*\|^2$$
 - Newton method is affine invariant (proof by chain rule)

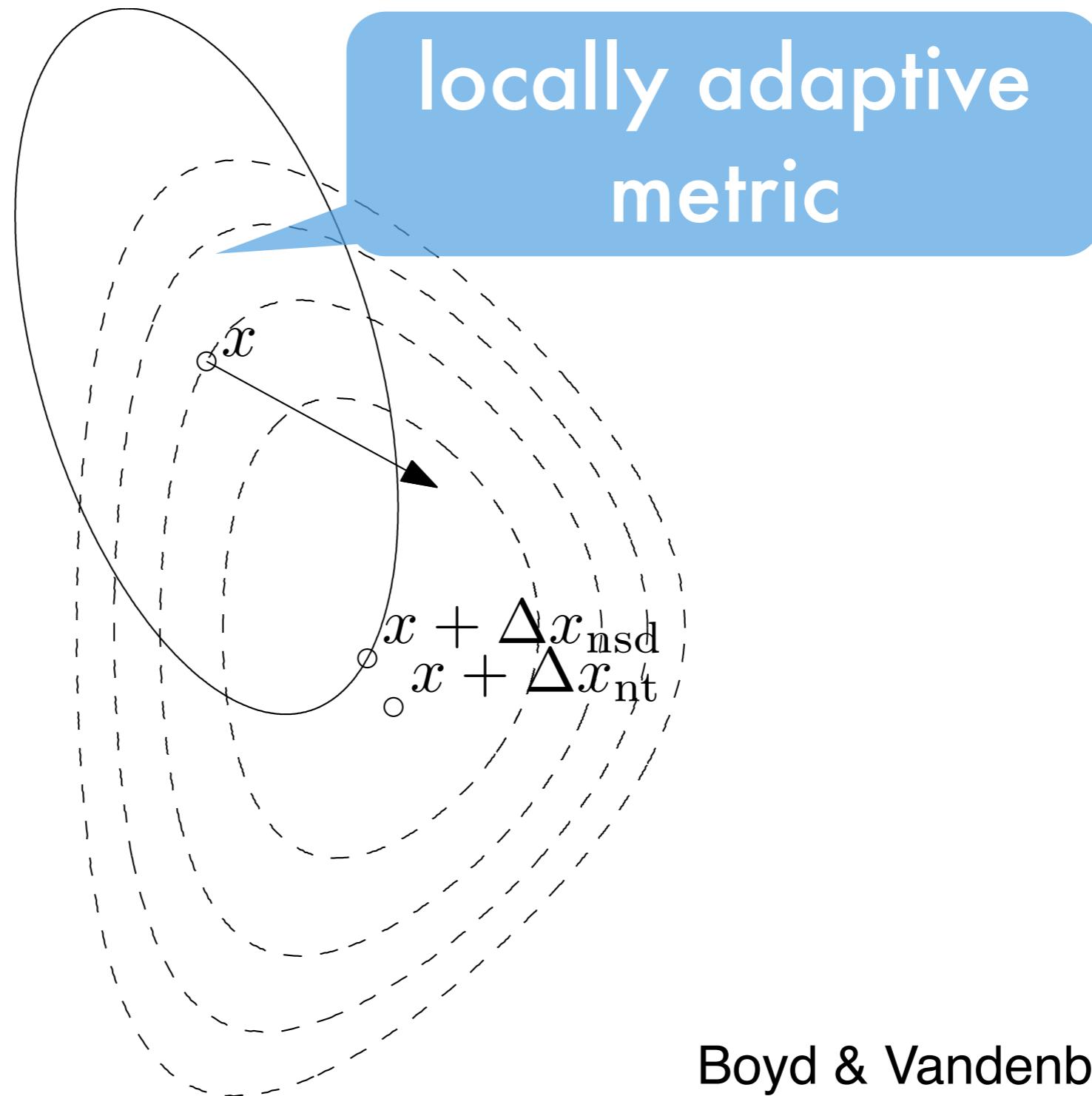
See Boyd and Vandenberghe, Chapter 9.5 for much more

Newton method rescales space



from Boyd & Vandenberghe
Carnegie Mellon University

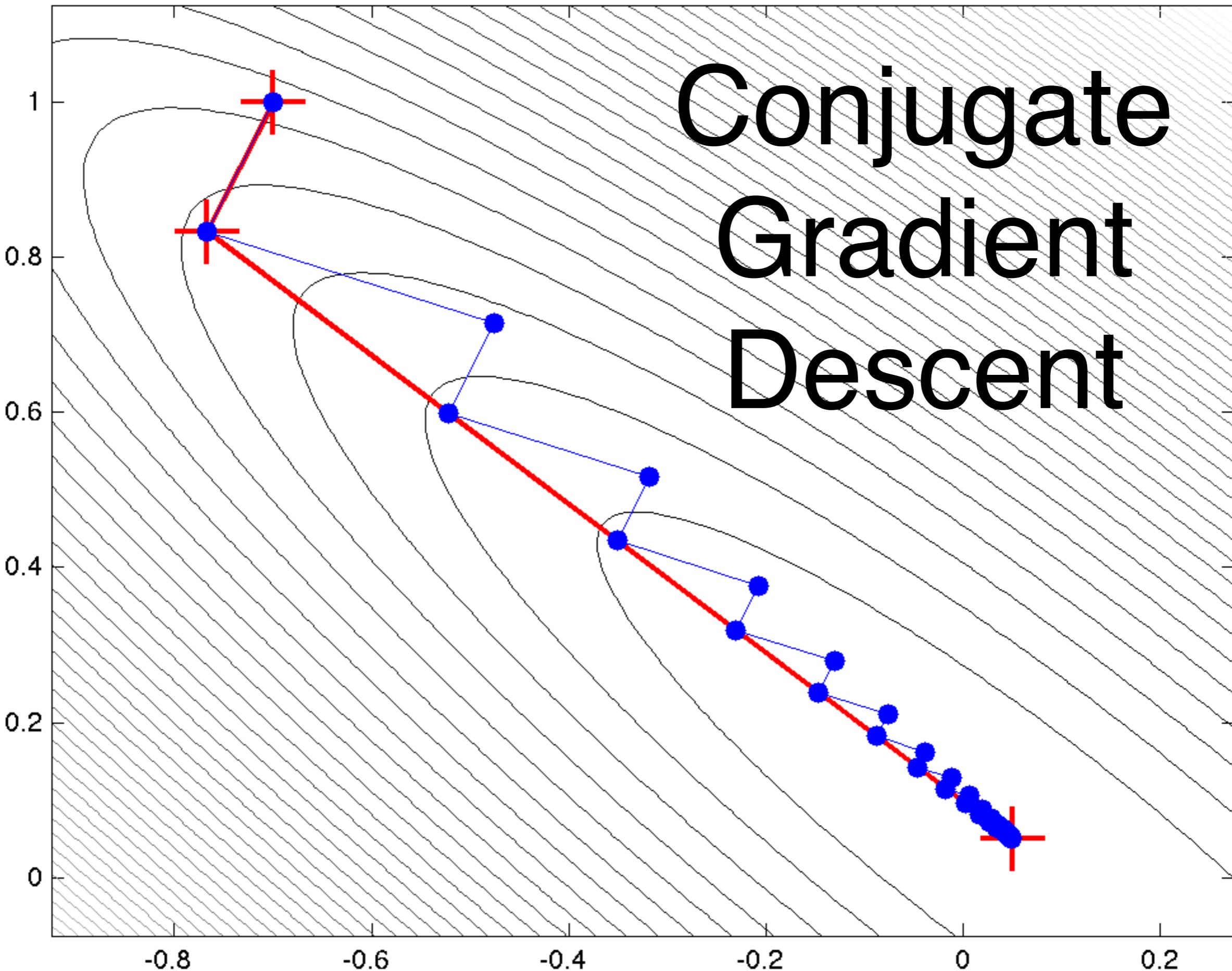
Newton method rescales space



Parallel Newton Method

- Good rate of convergence
- Few passes through data needed
- Parallel aggregation of gradient and Hessian
- Gradient requires $O(d)$ data
- Hessian requires $O(d^2)$ data
- Update step is $O(d^3)$ & nontrivial to parallelize
- Use it only for low dimensional problems

Conjugate Gradient Descent



Key Idea

- Minimizing quadratic function $(K \succeq 0)$

$$f(x) = \frac{1}{2}x^\top Kx - l^\top x + c$$

takes cubic time (e.g. Cholesky factorization)

- Matrix vector products and orthogonalization

- Vectors x, x' are K orthogonal if

$$x^\top K x' = 0$$

- m mutually K orthogonal vectors

$$x_i \in \mathbb{R}^m$$

- form a basis

- allow expansion

- solve linear system

$$z = \sum_{i=1}^m x_i \frac{x_i^\top K z}{x_i^\top K x_i}$$

$$z = \sum_{i=1}^m x_i \frac{x_i^\top y}{x_i^\top K x_i} \text{ for } y = K z$$

Proof

- m mutually K orthogonal vectors $x_i \in \mathbb{R}^m$
 - form a basis
 - allow expansion
 - solve linear system
- $$z = \sum_{i=1}^m x_i \frac{x_i^\top K z}{x_i^\top K x_i}$$
- $$z = \sum_{i=1}^m x_i \frac{x_i^\top y}{x_i^\top K x_i} \text{ for } y = K z$$

- Show linear independence by contradiction

$$\sum_i \alpha_i x_i = 0 \text{ hence } 0 = x_j^\top K \sum_i \alpha_i x_i = x_j^\top K x_j \alpha_j$$

- Reconstruction - expand z into basis

$$z = \sum_i \alpha_i x_i \text{ hence } x_j^\top K z = x_j^\top K \sum_i \alpha_i x_i = x_j^\top K x_j \alpha_j$$

- For linear system plug in $y = K z$

Conjugate Gradient Descent

- Gradient computation

$$f(x) = \frac{1}{2}x^\top Kx - l^\top x + c \text{ hence } g(x) = Kx - l$$

- Algorithm

initialize x_0 and $v_0 = g_0 = Kx_0 - l$ and $i = 0$

repeat

$$x_{i+1} = x_i - v_i \frac{g_i^\top v_i}{v_i^\top K v_i}$$

$$g_{i+1} = Kx_{i+1} - l$$

$$v_{i+1} = -g_{i+1} + v_i \frac{g_{i+1}^\top K v_i}{v_i^\top K v_i}$$

$$i \leftarrow i + 1$$

until $g_i = 0$

deflation step

K orthogonal

Proof - Deflation property

$$x_{i+1} = x_i - v_i \frac{g_i^\top v_i}{v_i^\top K v_i}$$

$$g_{i+1} = Kx_{i+1} - l$$

$$v_{i+1} = -g_{i+1} + v_i \frac{g_{i+1}^\top K v_i}{v_i^\top K v_i}$$

- First assume that the v_i are K orthogonal and show that x_{i+1} is optimal in span of $\{v_1 \dots v_i\}$
- Enough if we show that $v_j^\top g_i = 0$ for all $j < i$
 - For $j=i$ expand
$$\begin{aligned} v_i^\top g_{i+1} &= v_i^\top \left[Kx_i - l - K v_i \frac{g_i^\top v_i}{v_i^\top K v_i} \right] \\ &= v_i^\top g_i - v_i^\top K v_i \frac{g_i^\top v_i}{v_i^\top K v_i} = 0 \end{aligned}$$
 - For smaller j a consequence of K orthogonality

Proof - K orthogonality

$$x_{i+1} = x_i - v_i \frac{g_i^\top v_i}{v_i^\top K v_i}$$

$$g_{i+1} = Kx_{i+1} - l$$

$$v_{i+1} = -g_{i+1} + v_i \frac{g_{i+1}^\top K v_i}{v_i^\top K v_i}$$

- Need to check that v_{i+1} is K orthogonal to all v_j (rest automatically true by construction)

$$v_j^\top K v_{i+1} = -v_j^\top K g_{i+1} + v_j^\top K v_i \frac{g_{i+1}^\top K v_i}{v_i^\top K v_i}$$

0 by deflation

0 by K orthogonality

Properties

- Subspace expansion method for optimality
(g, Kg, K^2g, K^3g, \dots)
- Focuses on leading eigenvalues
- Often sufficient to take only a few steps
(whenever the eigenvalues decay rapidly)

Extensions

Generic Method

Compute Hessian $K_i := f''(x_i)$ and update α_i, β_i with

$$\alpha_i = -\frac{g_i^\top v_i}{v_i^\top K_i v_i}$$

$$\beta_i = \frac{g_{i+1}^\top K_i v_i}{v_i^\top K_i v_i}$$

x and v updates

This requires calculation of the Hessian at each iteration.

Fletcher–Reeves [163]

Find α_i via a line search and use Theorem 6.20 (iii) for β_i

$$\alpha_i = \operatorname{argmin}_\alpha f(x_i + \alpha v_i)$$

$$\beta_i = \frac{g_{i+1}^\top g_{i+1}}{g_i^\top g_i}$$

Polak–Ribiere [398]

Find α_i via a line search

$$\alpha_i = \operatorname{argmin}_\alpha f(x_i + \alpha v_i)$$

$$\beta_i = \frac{(g_{i+1} - g_i)^\top g_{i+1}}{g_i^\top g_i}$$

Experimentally, Polak–Ribiere tends to be better than Fletcher–Reeves.

Broyden-Fletcher-Goldfarb-Shanno



Basic Idea

- Newton-like method to compute descent direction

$$\delta_i = B_i^{-1} \partial_x f(x_{i-1})$$

- Line search on f in direction

$$x_{i+1} = x_i - \alpha_i \delta_i$$

- Update B with rank 2 matrix

$$B_{i+1} = B_i + u_i u_i^\top + v_i v_i^\top$$

- Require that Quasi-Newton condition holds

$$B_{i+1}(x_{i+1} - x_i) = \partial_x f(x_{i+1}) - \partial_x f(x_i)$$

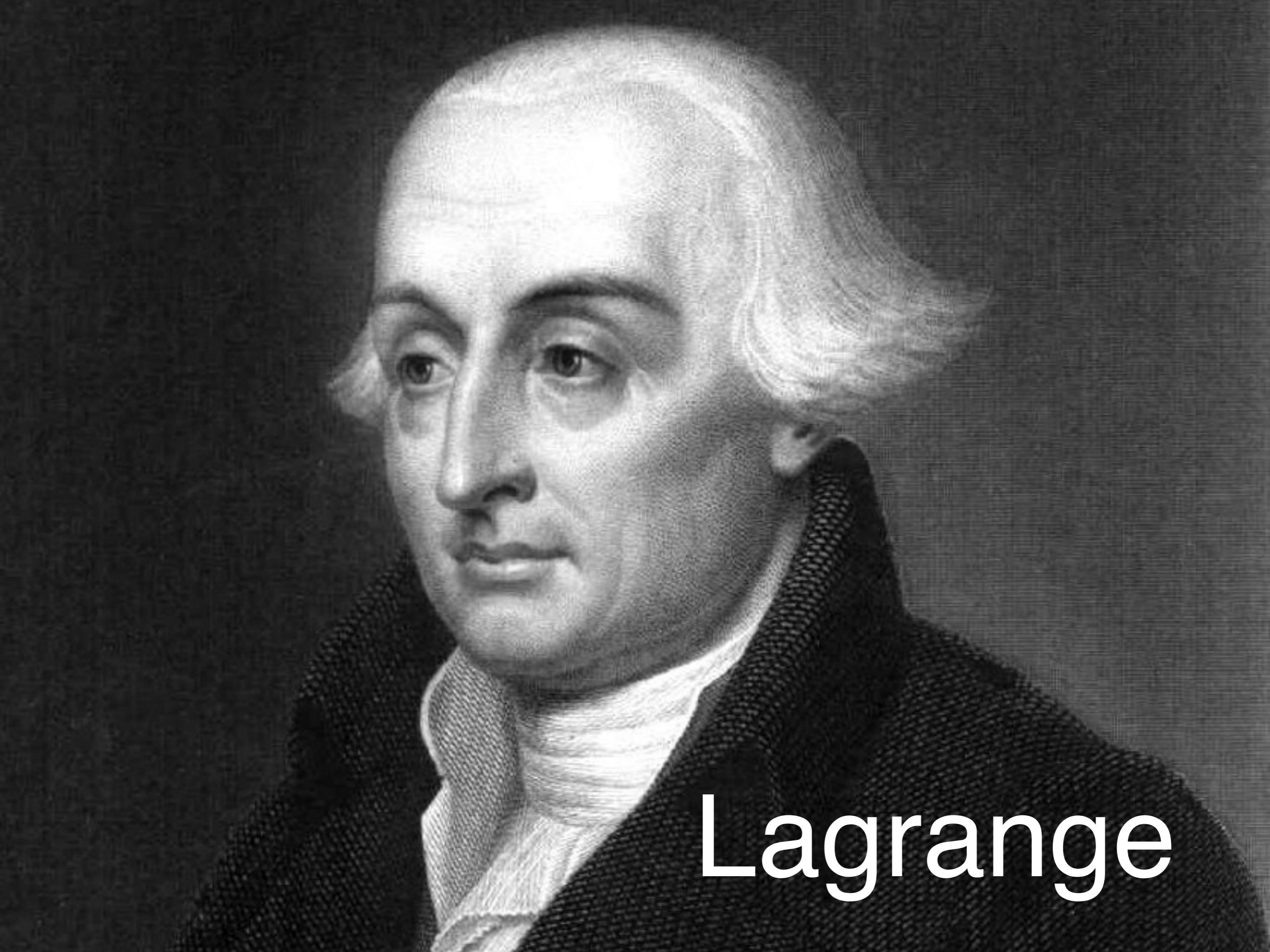
$$B_{i+1} = B_i + \frac{g_i g_i^\top}{\alpha_i \delta_i^\top g_i} - \frac{B_i \delta_i \delta_i^\top B_i}{\delta_i^\top B_i \delta_i}$$

Properties

- Simple rank 2 update for B
- Use matrix inversion lemma to update inverse
- Memory-limited versions L-BFGS
- Use toolbox if possible (TAO, MATLAB)
(typically slower if you implement it yourself)
- Works well for nonlinear nonconvex objectives
(often even for nonsmooth objectives)

An aerial photograph of the Palio di Siena, a traditional horse race held twice yearly in the Italian town of Siena. The image shows a massive crowd of spectators filling the Piazza del Campo, which is roughly circular in shape. The crowd is densest in the center and along the perimeter, where it spills onto the surrounding city streets. In the background, the historic buildings of the city, including the Palazzo Pubblico, are visible. Several horses and jockeys are seen on the track, which is a dirt surface. The overall scene is one of a major public event.

Constrained Convex Optimization



Lagrange

Constrained Convex Minimization

- Optimization problem

$$\underset{x}{\text{minimize}} \quad f(x)$$

subject to $c_i(x) \leq 0$ for all i

- Common constraints

- linear inequality constraints

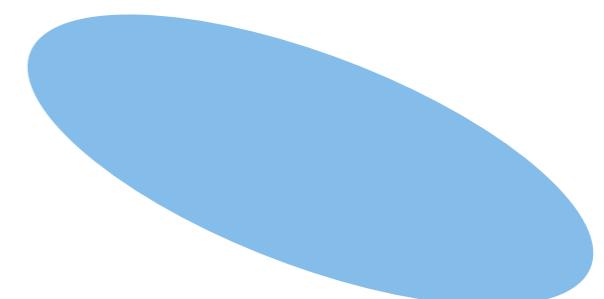
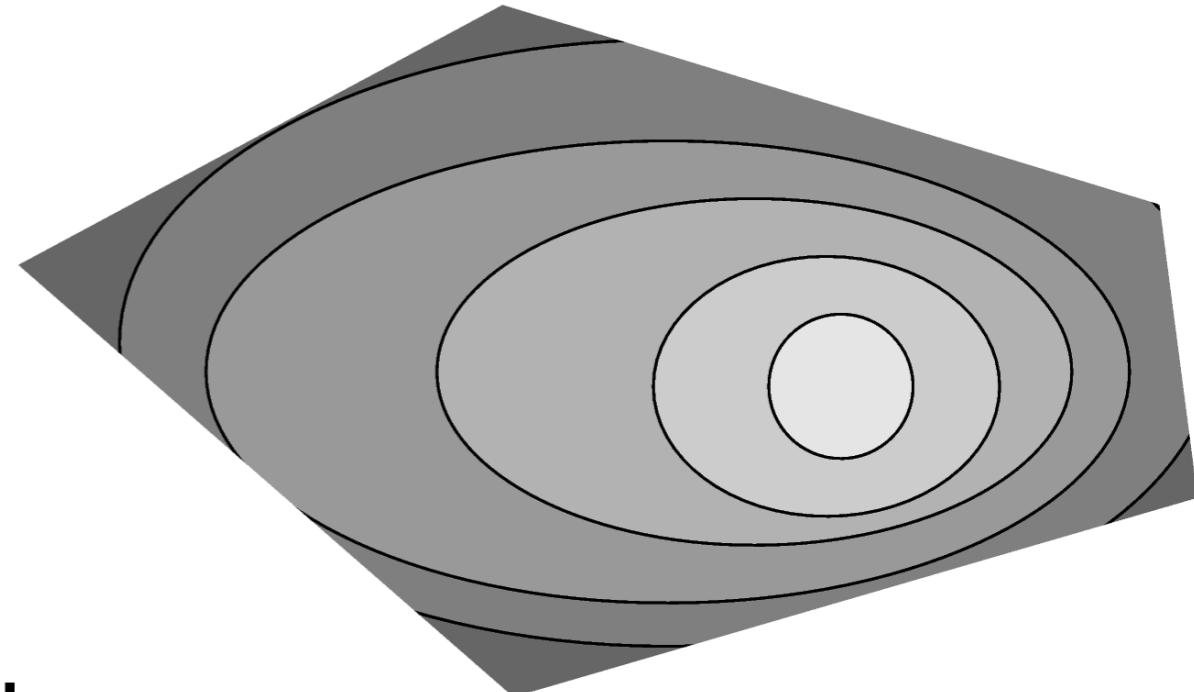
$$\langle w_i, x \rangle + b_i \leq 0$$

- quadratic cone constraints

$$x^\top Q x + b^\top x \leq c \text{ with } Q \succeq 0$$

- semidefinite constraints

$$M \succeq 0 \text{ or } M_0 + \sum_i x_i M_i \succeq 0$$



Constrained Convex Minimization

- Optimization problem

$$\underset{x}{\text{minimize}} \quad f(x)$$

subject to $c_i(x) \leq 0$ for

- Common constraints

- linear inequality constraints

$$\langle w_i, x \rangle + b_i \leq 0$$

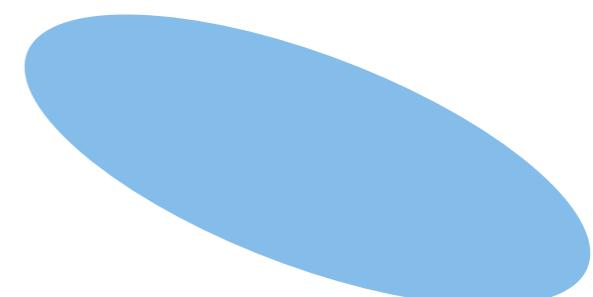
- quadratic cone constraints

$$x^\top Q x + b^\top x \leq c \text{ with } Q \succeq 0$$

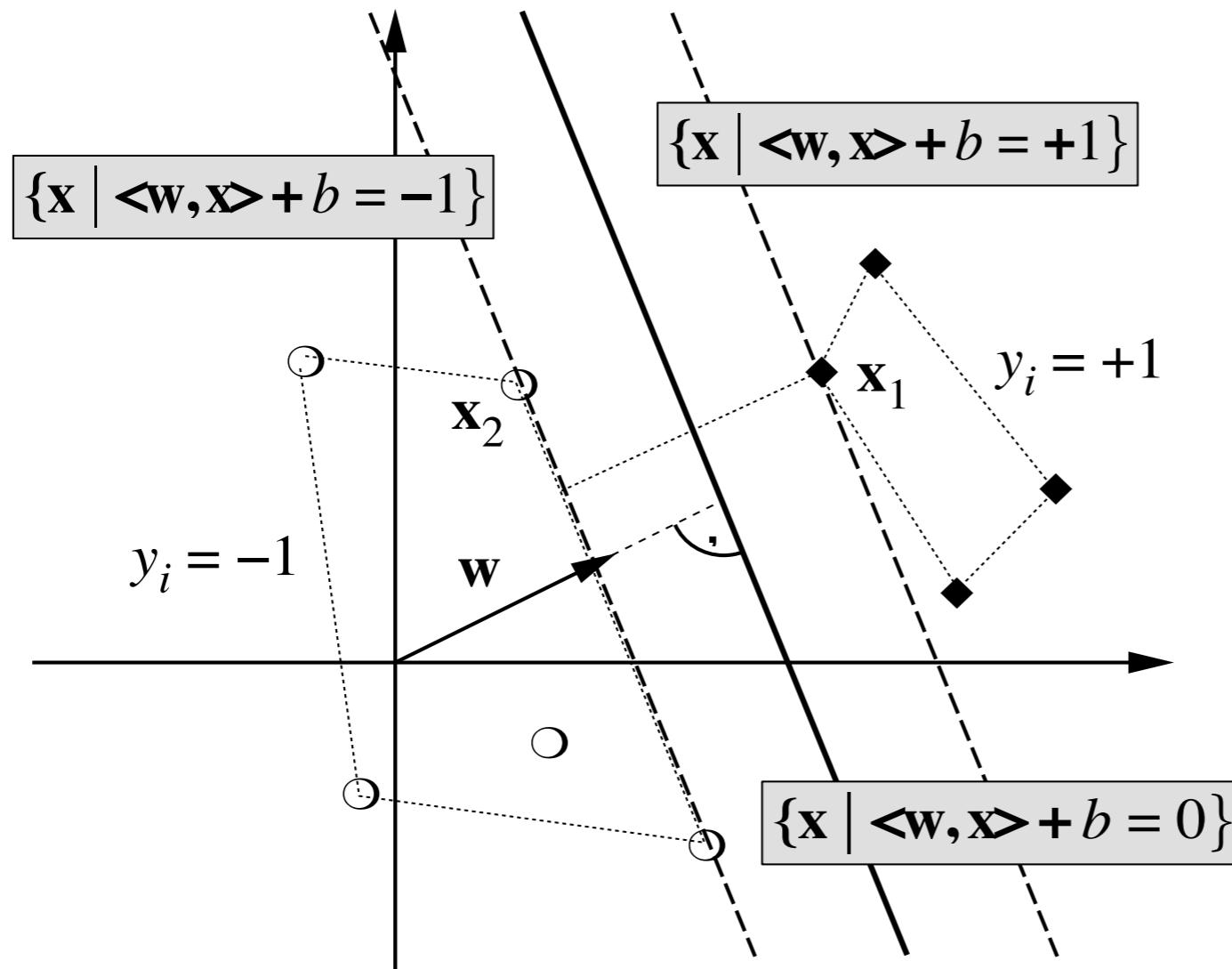
- semidefinite constraints

$$M \succeq 0 \text{ or } M_0 + \sum_i x_i M_i \succeq 0$$

Equality is special case
Why?



Example - Support Vectors



$$\langle \mathbf{w}, \mathbf{x}_1 \rangle + b = 1$$

$$\langle \mathbf{w}, \mathbf{x}_2 \rangle + b = -1$$

$$\text{hence } \langle \mathbf{w}, \mathbf{x}_1 - \mathbf{x}_2 \rangle = 2$$

$$\text{hence } \left\langle \frac{\mathbf{w}}{\|\mathbf{w}\|}, \mathbf{x}_1 - \mathbf{x}_2 \right\rangle = \frac{2}{\|\mathbf{w}\|}$$

margin

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

Lagrange Multipliers

- Lagrange function

$$L(x, \alpha) := f(x) + \sum_{i=1}^n \alpha_i c_i(x) \text{ where } \alpha_i \geq 0$$

- Saddlepoint Condition

If there are x^* and nonnegative α^* such that

$$L(x^*, \alpha) \leq L(x^*, \alpha^*) \leq L(x, \alpha^*)$$

then x^* is an optimal solution to the constrained optimization problem

Proof

$$L(x^*, \alpha) \leq L(x^*, \alpha^*) \leq L(x, \alpha^*)$$

- From first inequality we see that x^* is feasible
$$(\alpha_i - \alpha_i^*)c_i(x^*) \leq 0 \text{ for all } \alpha_i \geq 0$$
- Setting some $\alpha_i = 0$ yields KKT conditions
$$\alpha_i^* c_i(x^*) = 0$$
- Consequently we have

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

This proves optimality

Constraint gymnastics (all three conditions are equivalent)

- **Slater's condition**

There exists some x such that for all i

$$c_i(x) < 0$$

- **Karlin's condition**

For all nonnegative a there exists some x such that

$$\sum_i \alpha_i c_i(x) \leq 0$$

- **Strict constraint qualification**

The feasible region contains at least two distinct elements and there exists an x in X such that all $c_i(x)$ are strictly convex at x with respect to X

Necessary Kuhn-Tucker Conditions

- Assume optimization problem
 - satisfies the constraint qualifications
 - has convex differentiable objective + constraints
- Then the KKT conditions are necessary & sufficient

$$\partial_x L(x^*, \alpha^*) = \partial_x f(x^*) + \sum_i \alpha_i^* \partial_x c_i(x^*) = 0 \text{ (Saddlepoint in } x^*)$$

$$\partial_{\alpha_i} L(x^*, \alpha^*) = c_i(x^*) \leq 0 \text{ (Saddlepoint in } \alpha^*)$$

$$\sum_i \alpha_i^* c_i(x^*) = 0 \text{ (Vanishing KKT-gap)}$$

Yields algorithm for solving optimization problems
Solve for saddlepoint and KKT conditions

Proof

$$\begin{aligned} f(x) - f(x^*) &\geq [\partial_x f(x^*)]^\top (x - x^*) && \text{(by convexity)} \\ &= - \sum_i \alpha_i^* [\partial_x c_i(x^*)]^\top (x - x^*) && \text{(by Saddlepoint in } x^*) \\ &\geq - \sum_i \alpha_i^* (c_i(x) - c_i(x^*)) && \text{(by convexity)} \\ &= \sum_i \alpha_i^* c_i(x) && \text{(by vanishing KKT gap)} \\ &\geq 0 \end{aligned}$$

Linear and Quadratic Programs



Linear Programs

- **Objective**

$$\underset{x}{\text{minimize}} \quad c^\top x \text{ subject to } Ax + d \leq 0$$

- **Lagrange function**

$$L(x, \alpha) = c^\top x + \alpha^\top (Ax + d)$$

- **Optimality conditions**

$$\partial_x L(x, \alpha) = A^\top \alpha + c = 0$$

$$\partial_\alpha L(x, \alpha) = Ax + d \leq 0$$

$$0 = \alpha^\top (Ax + d)$$

$$0 \leq \alpha$$

- **Dual problem**

$$\underset{i}{\text{maximize}} \quad d^\top \alpha \text{ subject to } A^\top \alpha + c = 0 \text{ and } \alpha \geq 0$$

Linear Programs

- Objective

$$\underset{x}{\text{minimize}} \quad c^T x \text{ subject to } Ax + d \leq 0$$

- Lagrange function

$$L(x, \alpha) = c^T x + \alpha^T (Ax + d)$$

- Optimality conditions

$$\partial_x L(x, \alpha) = A^T \alpha + c = 0$$

$$\partial_\alpha L(x, \alpha) = Ax + d \leq 0$$

$$0 = \alpha^T (Ax + d)$$

$$0 \leq \alpha$$

- Dual problem

$$\underset{i}{\text{maximize}} \quad d^T \alpha \text{ subject to } A^T \alpha + c = 0 \text{ and } \alpha \geq 0$$

Linear Programs

- Objective

$$\underset{x}{\text{minimize}} \quad c^\top x \text{ subject to } Ax + d \leq 0$$

- Lagrange function

$$L(x, \alpha) = c^\top x + \alpha^\top (Ax + d)$$

- Optimality conditions

$$\partial_x L(x, \alpha) = A^\top \alpha + c = 0$$

$$\partial_\alpha L(x, \alpha) = Ax + d \leq 0$$

$$0 = \alpha^\top (Ax + d)$$

$$0 \leq \alpha$$

- Dual problem

$$\underset{i}{\text{maximize}} \quad d^\top \alpha \text{ subject to } A^\top \alpha + c = 0 \text{ and } \alpha \geq 0$$

Linear Programs

- Primal

$$\underset{x}{\text{minimize}} \quad c^\top x \text{ subject to } Ax + d \leq 0$$

- Dual

$$\underset{i}{\text{maximize}} \quad d^\top \alpha \text{ subject to } A^\top \alpha + c = 0 \text{ and } \alpha \geq 0$$

- Free variables become equality constraints
- Equality constraints become free variables
- Inequalities become inequalities
- Dual of dual is primal

Quadratic Programs

- Objective

$$\underset{x}{\text{minimize}} \frac{1}{2}x^\top Qx + c^\top x \text{ subject to } Ax + d \leq 0$$

- Lagrange function

$$L(x, \alpha) = \frac{1}{2}x^\top Qx + c^\top x + \alpha^\top (Ax + d)$$

- Optimality conditions

$$\partial_x L(x, \alpha) = Qx + A^\top \alpha + c = 0$$

$$\partial_\alpha L(x, \alpha) = Ax + d \leq 0$$

$$0 = \alpha^\top (Ax + d)$$

$$0 \leq \alpha$$

plug into L

Quadratic Program

- Eliminating x from the Lagrangian via

$$Qx + A^\top \alpha + c = 0$$

- Lagrange function

$$\begin{aligned} L(x, \alpha) &= \frac{1}{2}x^\top Qx + c^\top x + \alpha^\top (Ax + d) \\ &= -\frac{1}{2}x^\top Qx + \alpha^\top d \\ &= -\frac{1}{2}(A^\top \alpha + c)^\top Q^{-1}(A^\top \alpha + c) + \alpha^\top d \\ &= -\frac{1}{2}\alpha^\top A Q^{-1} A^\top \alpha + \alpha^\top [d - A Q^{-1} c] - \frac{1}{2}c^\top Q^{-1} c \end{aligned}$$

subject to $\alpha \geq 0$

Quadratic Program

- Eliminating x from the Lagrangian via

$$Qx + A^\top \alpha + c = 0$$

- Lagrange function

$$\begin{aligned} L(x, \alpha) &= \frac{1}{2}x^\top Qx + c^\top x + \alpha^\top (Ax + d) \\ &= -\frac{1}{2}x^\top Qx + \alpha^\top d \\ &= -\frac{1}{2}(A^\top \alpha + c)^\top Q^{-1}(A^\top \alpha + c) + \alpha^\top d \\ &= -\frac{1}{2}\alpha^\top A Q^{-1} A^\top \alpha + \alpha^\top [d - A Q^{-1} c] - \frac{1}{2}c^\top Q^{-1} c \end{aligned}$$

dual

subject to $\alpha \geq 0$

Quadratic Programs

- Primal

$$\underset{x}{\text{minimize}} \frac{1}{2} x^\top Q x + c^\top x \text{ subject to } Ax + d \leq 0$$

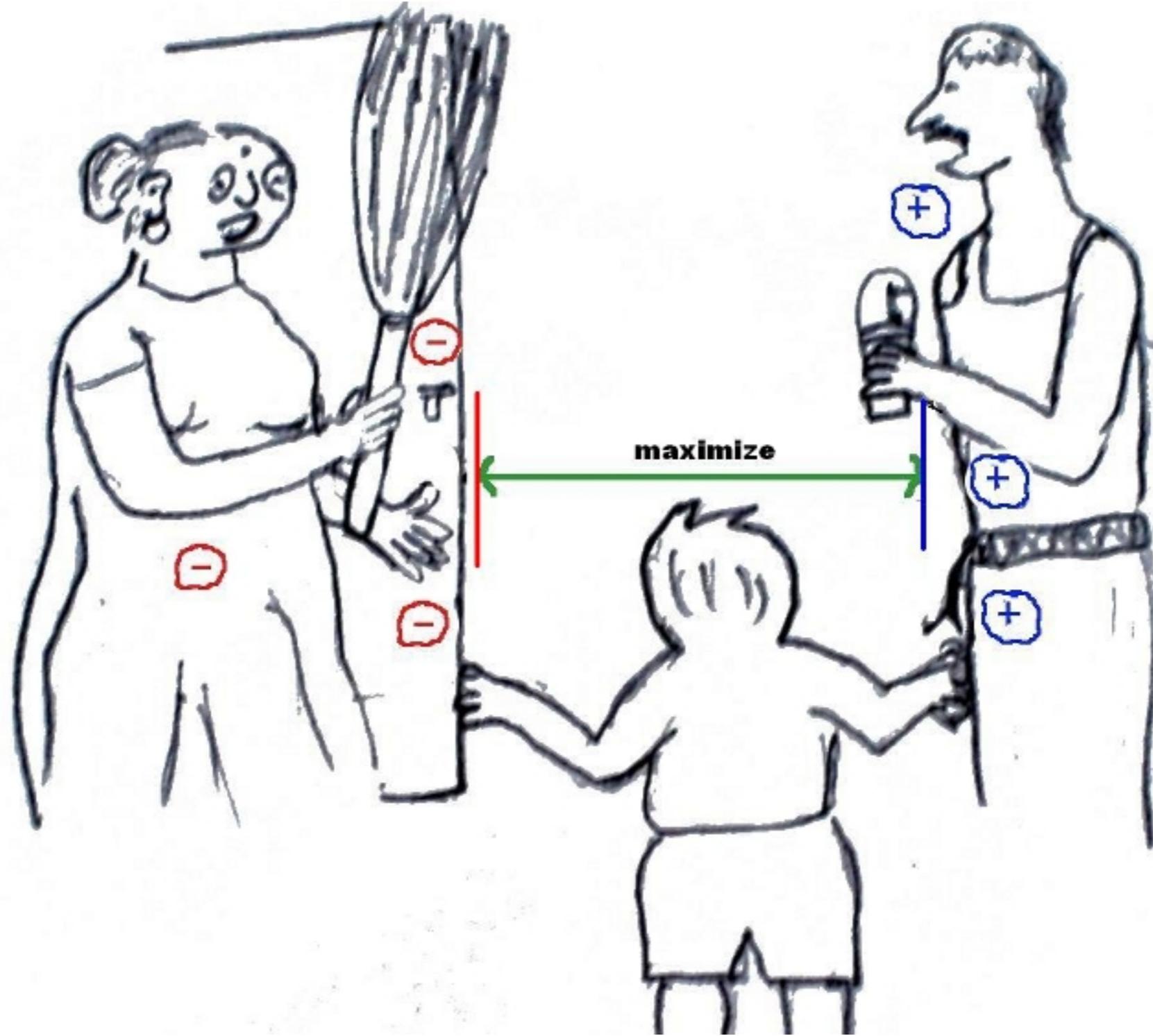
- Dual

$$\underset{\alpha}{\text{minimize}} \frac{1}{2} \alpha^\top A Q^{-1} A^\top \alpha + \alpha^\top [A Q^{-1} c - d] \text{ subject to } \alpha \geq 0$$

- Dual constraints are simpler
- Possibly many fewer variables
- Dual of dual is not (always) primal
(e.g. in SVMs x is in a Hilbert Space)

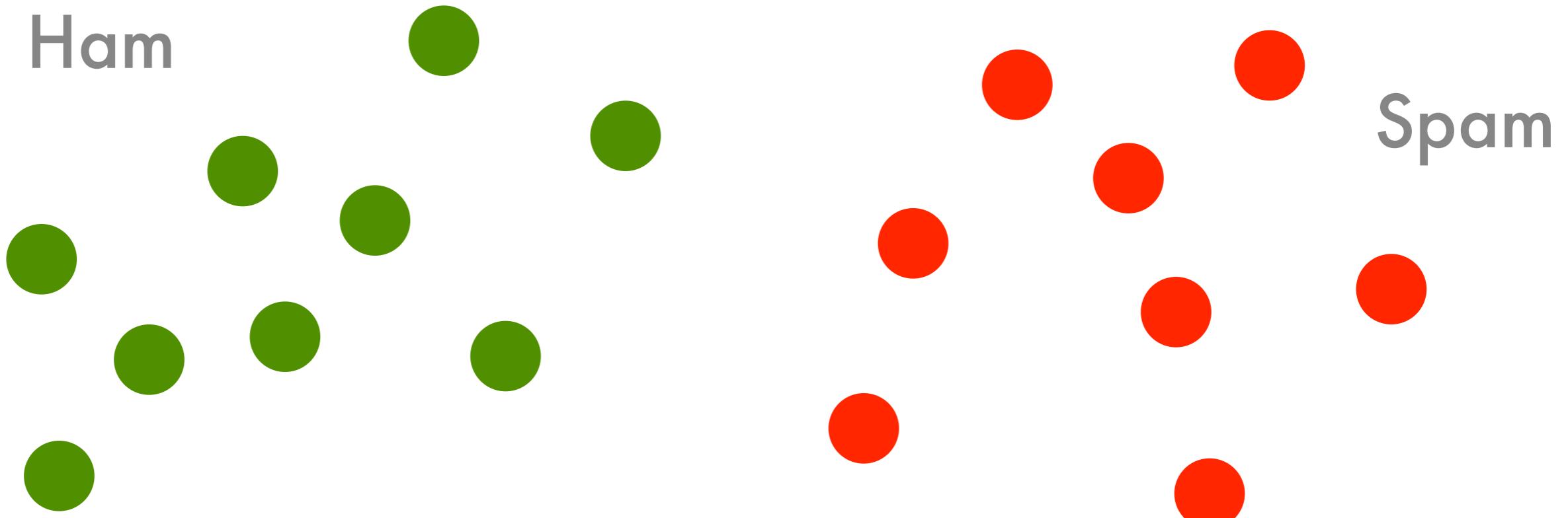
Outline

- **Convex Optimization**
 - Unconstrained Optimization
 - Constrained Optimization and Duality
 - Linear and Quadratic programs
- **Support Vector Machines**
 - Classification
 - Regression
 - Novelty Detection
- **Kernels**
 - Feature Space
 - Kernel PCA
 - Kernelized SVM

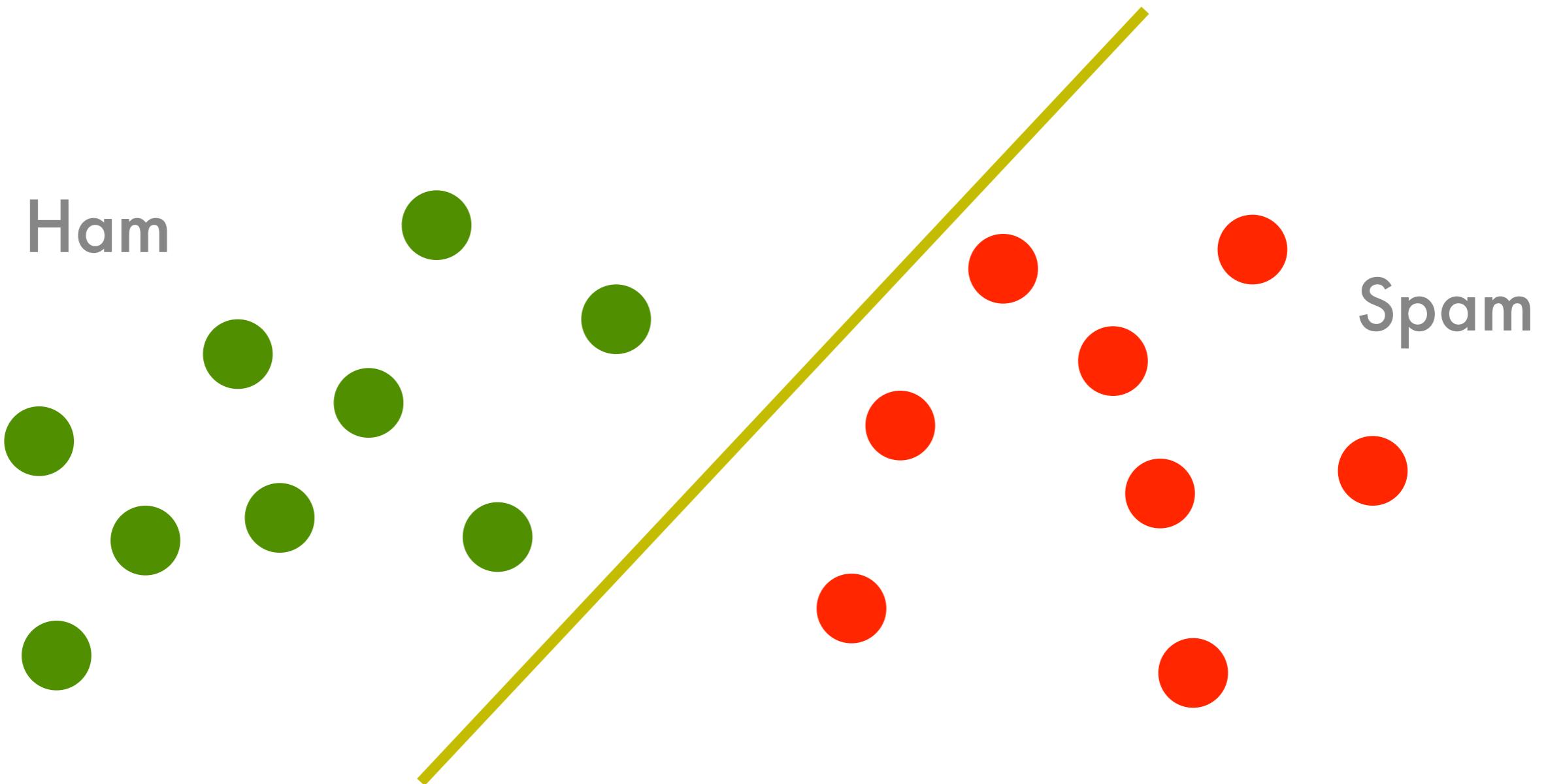


Support Vector Machines

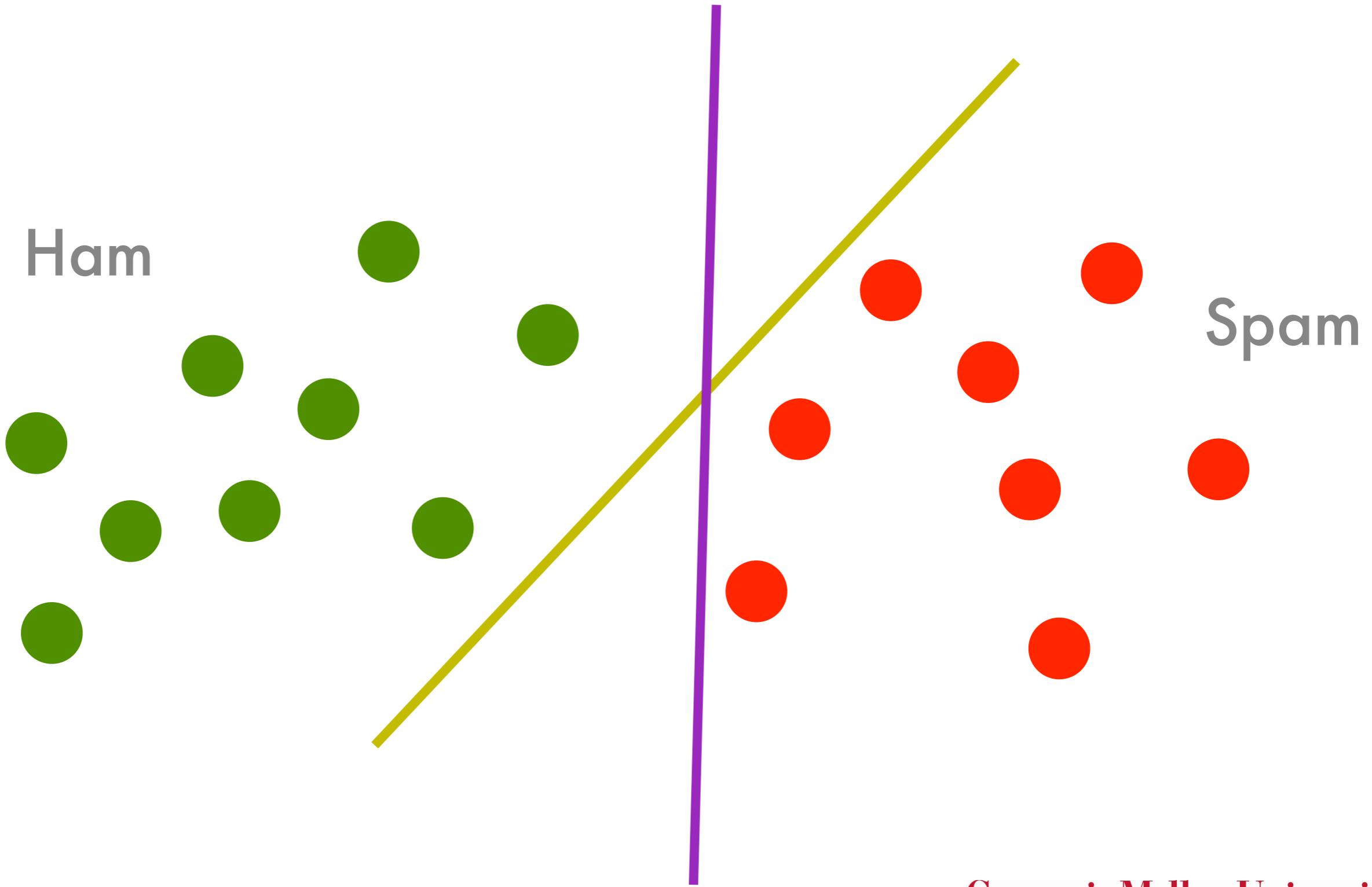
Linear Separator



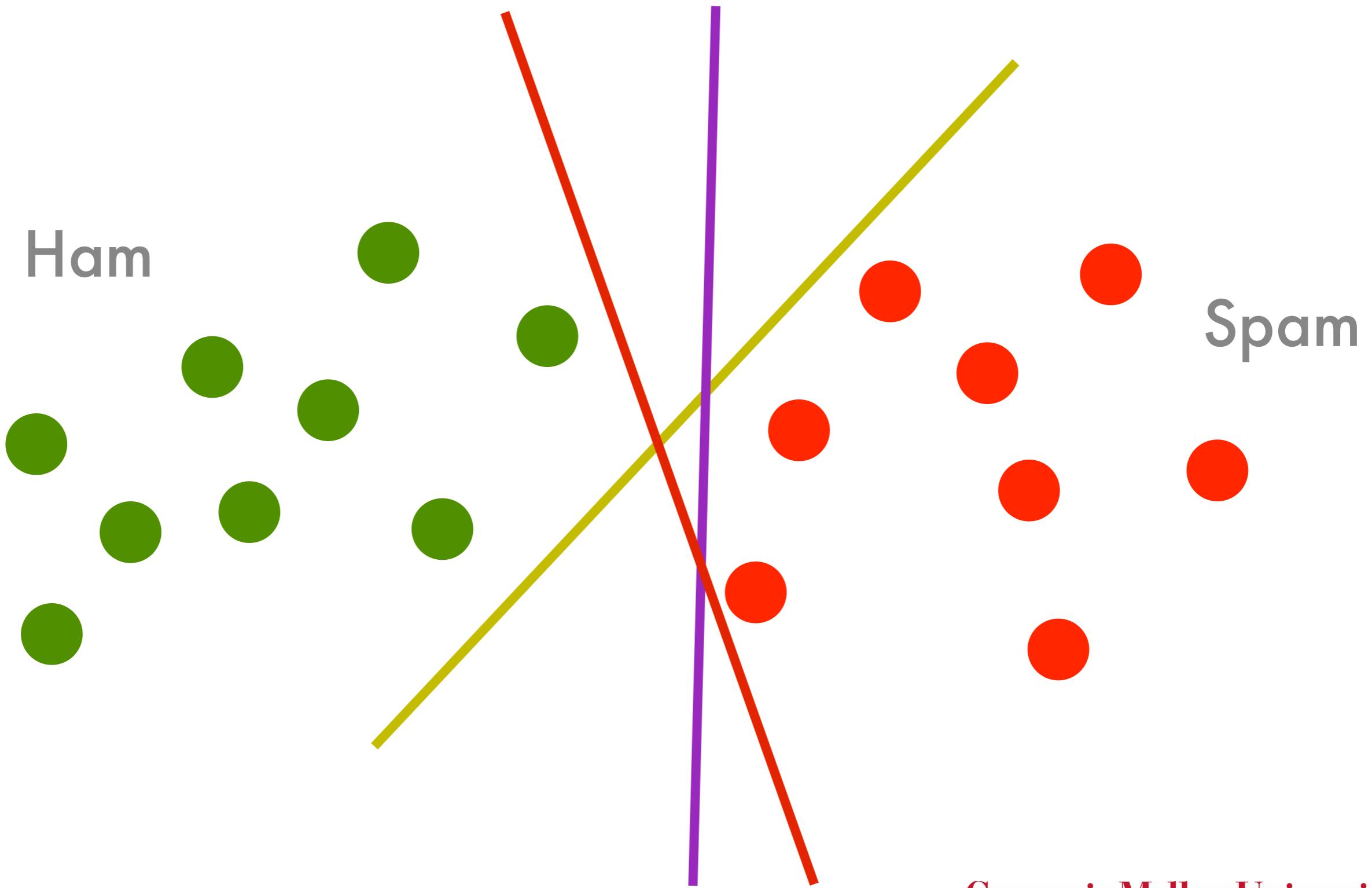
Linear Separator



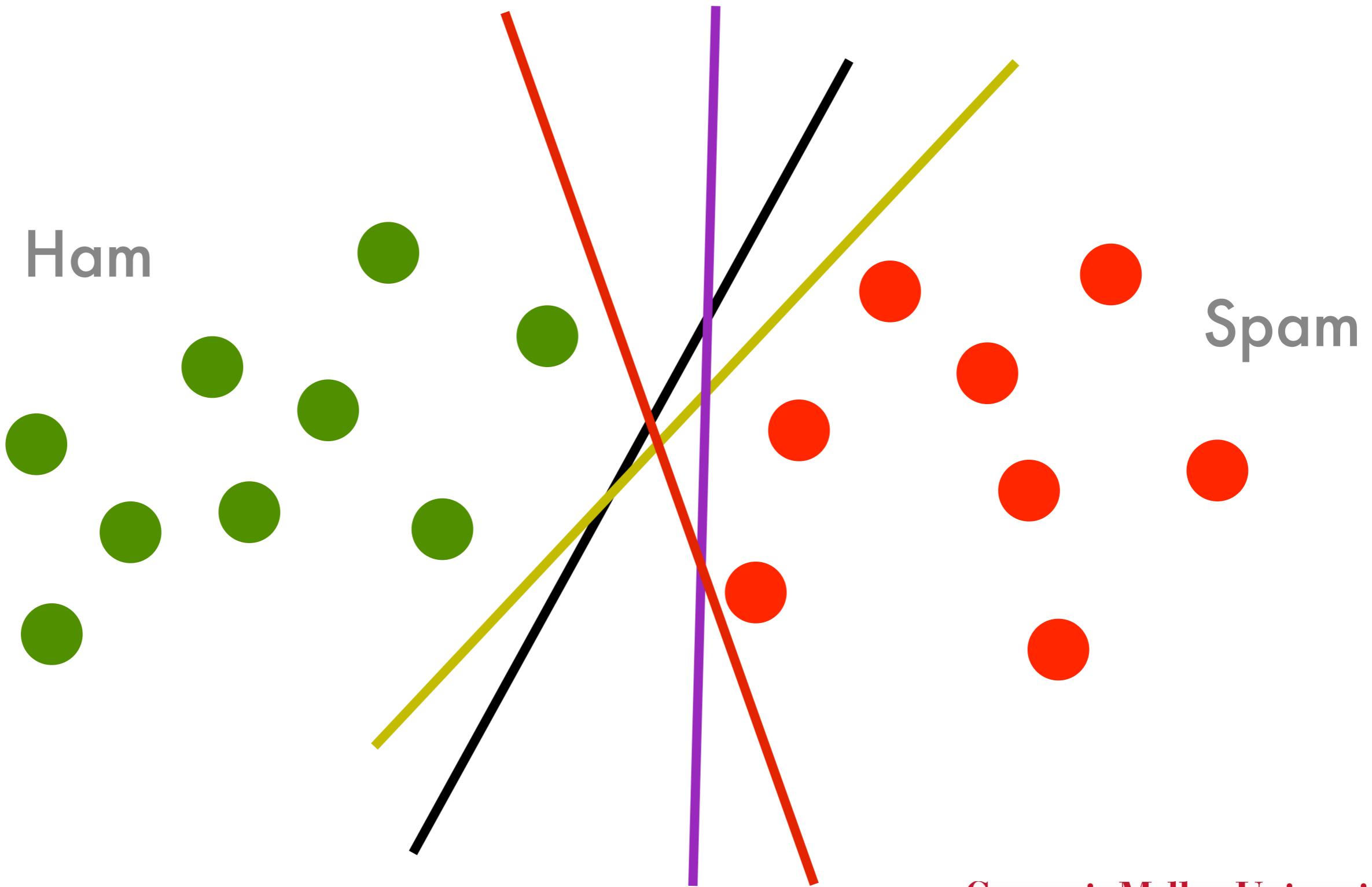
Linear Separator



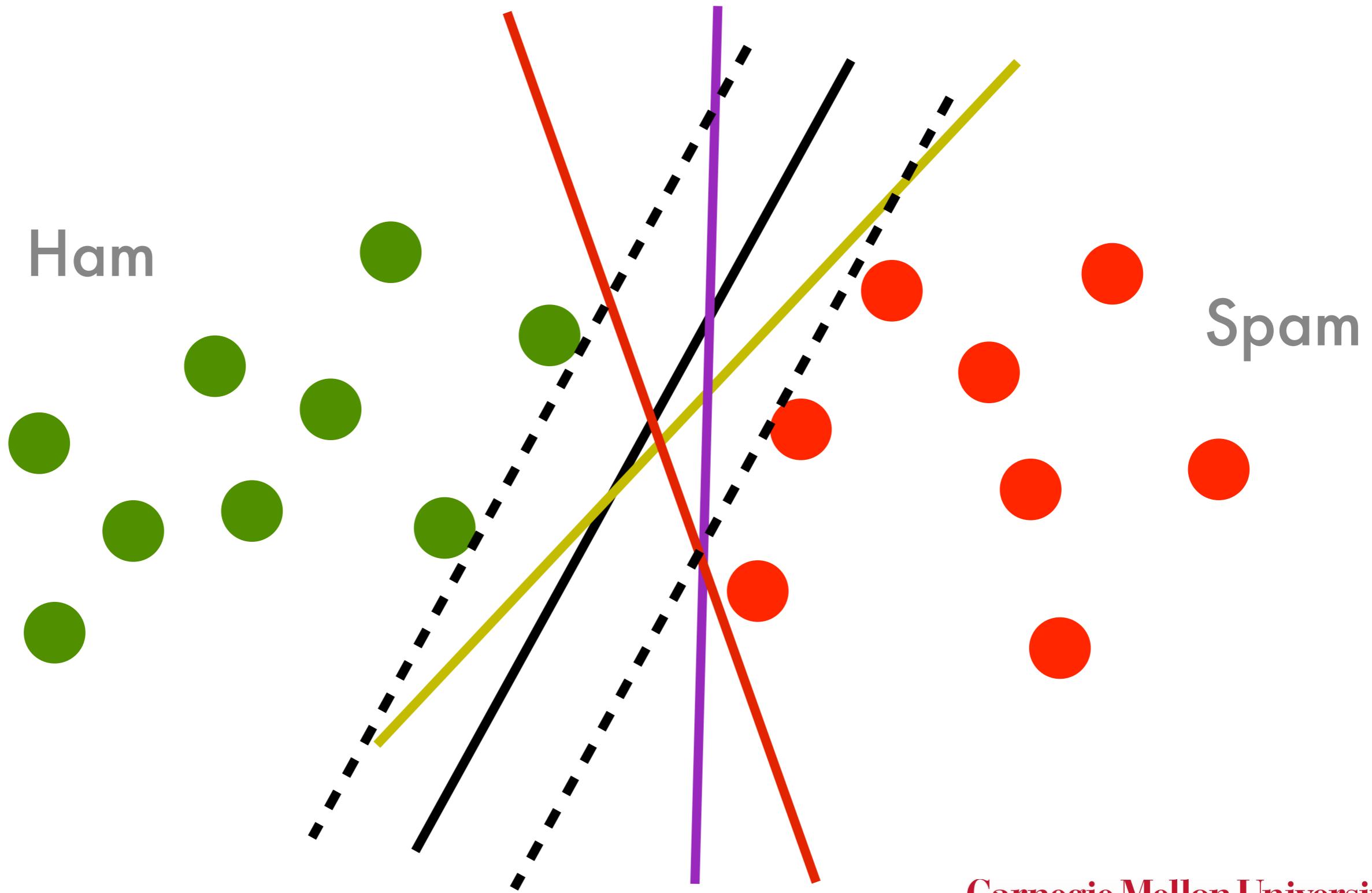
Linear Separator



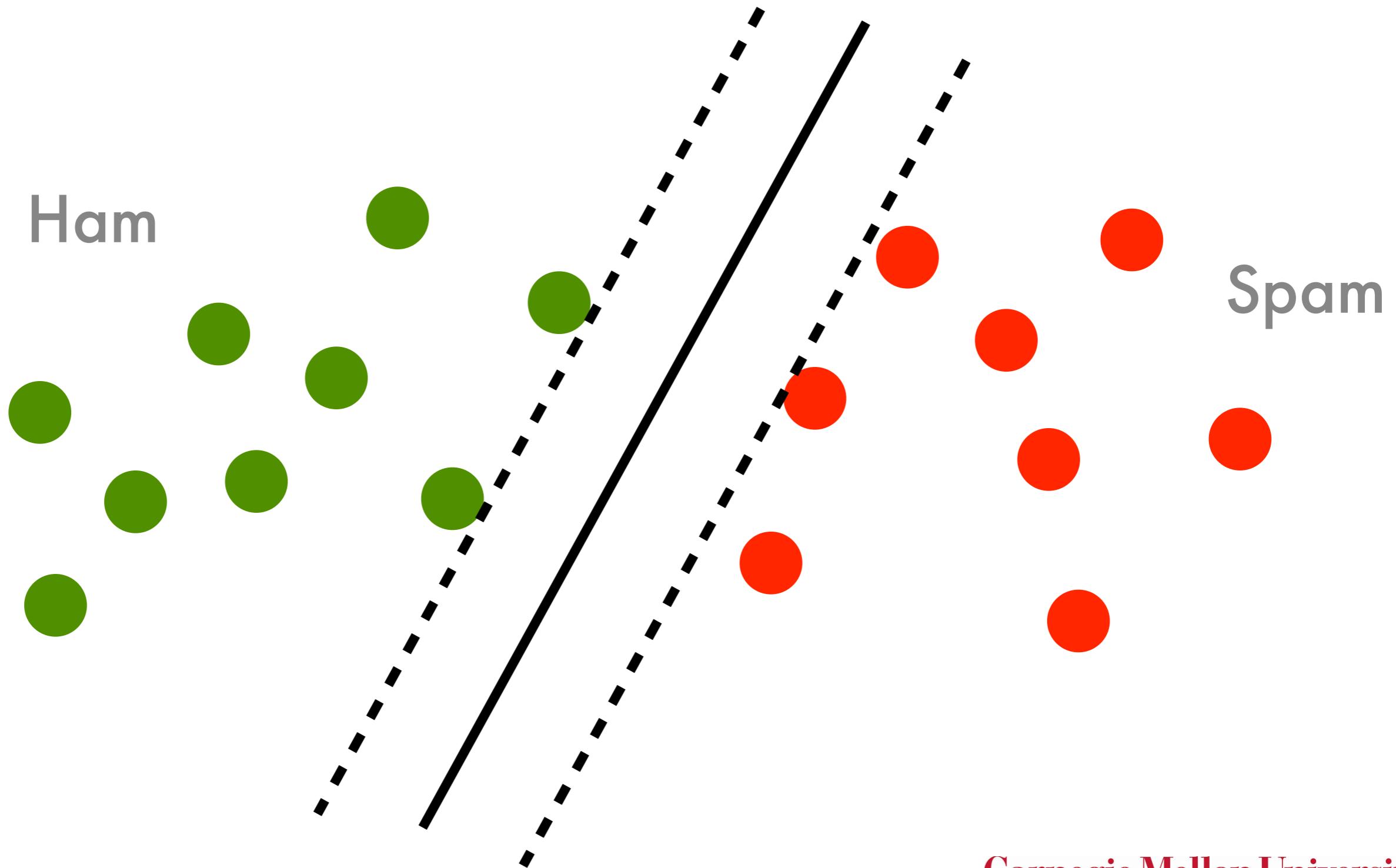
Linear Separator



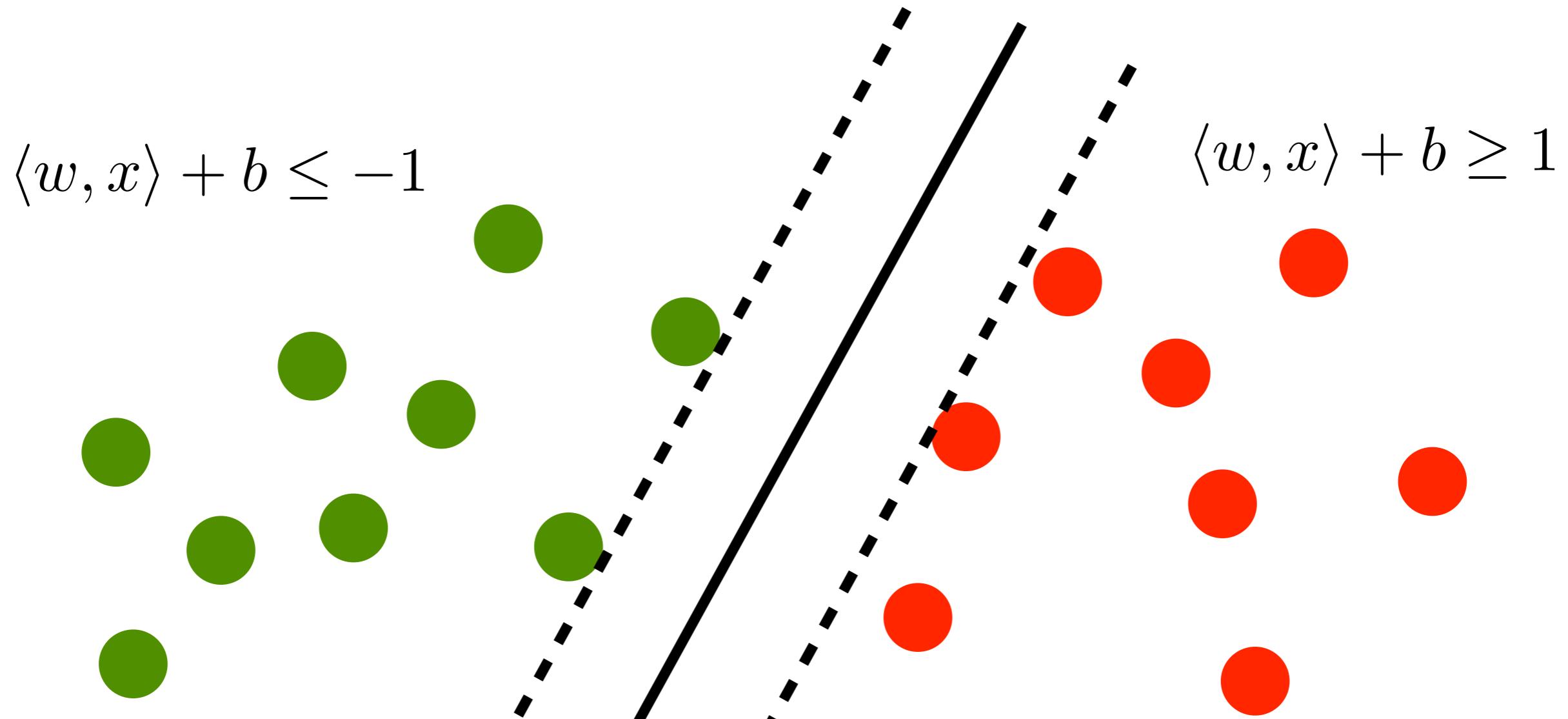
Linear Separator



Linear Separator

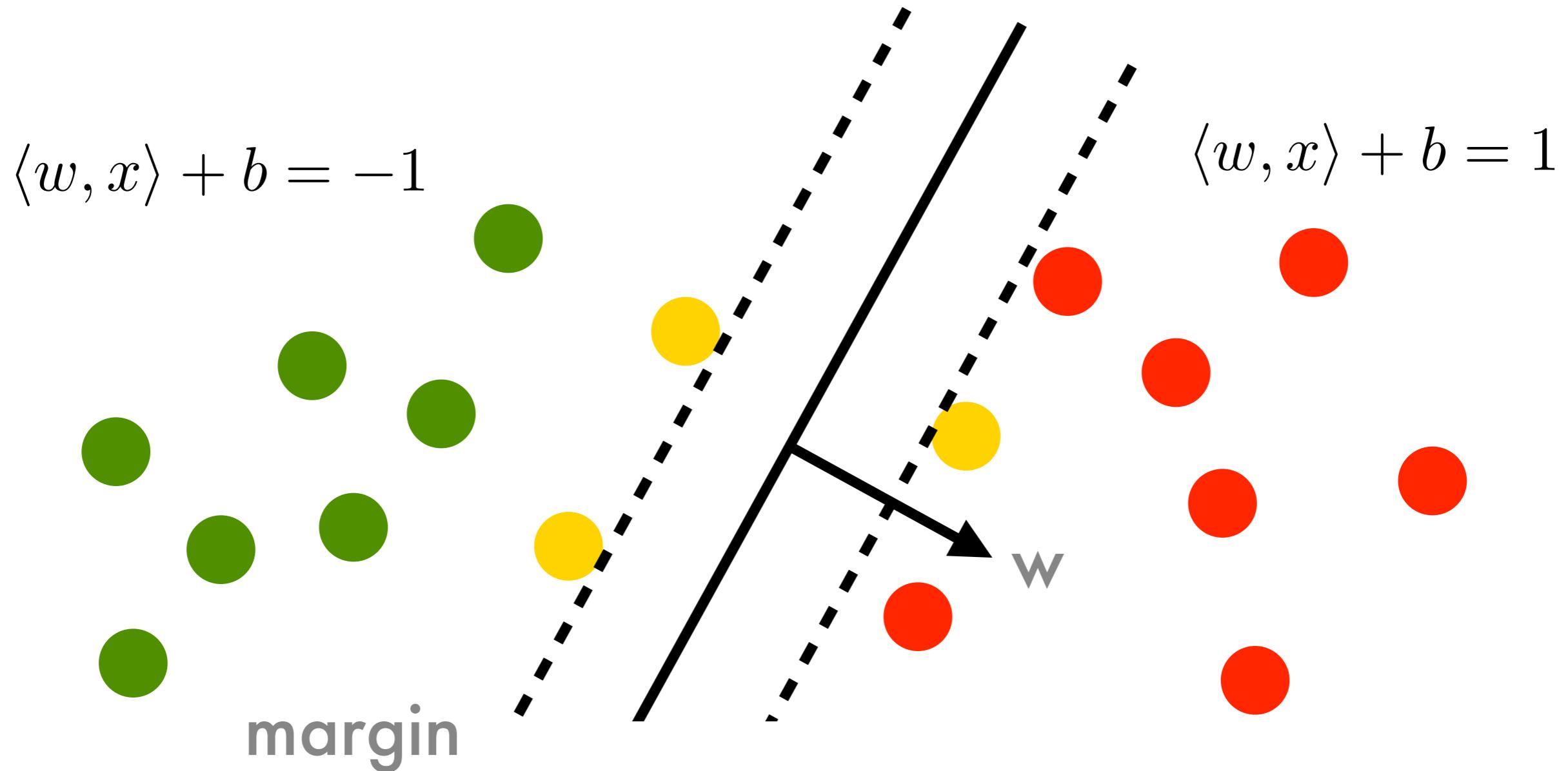


Large Margin Classifier



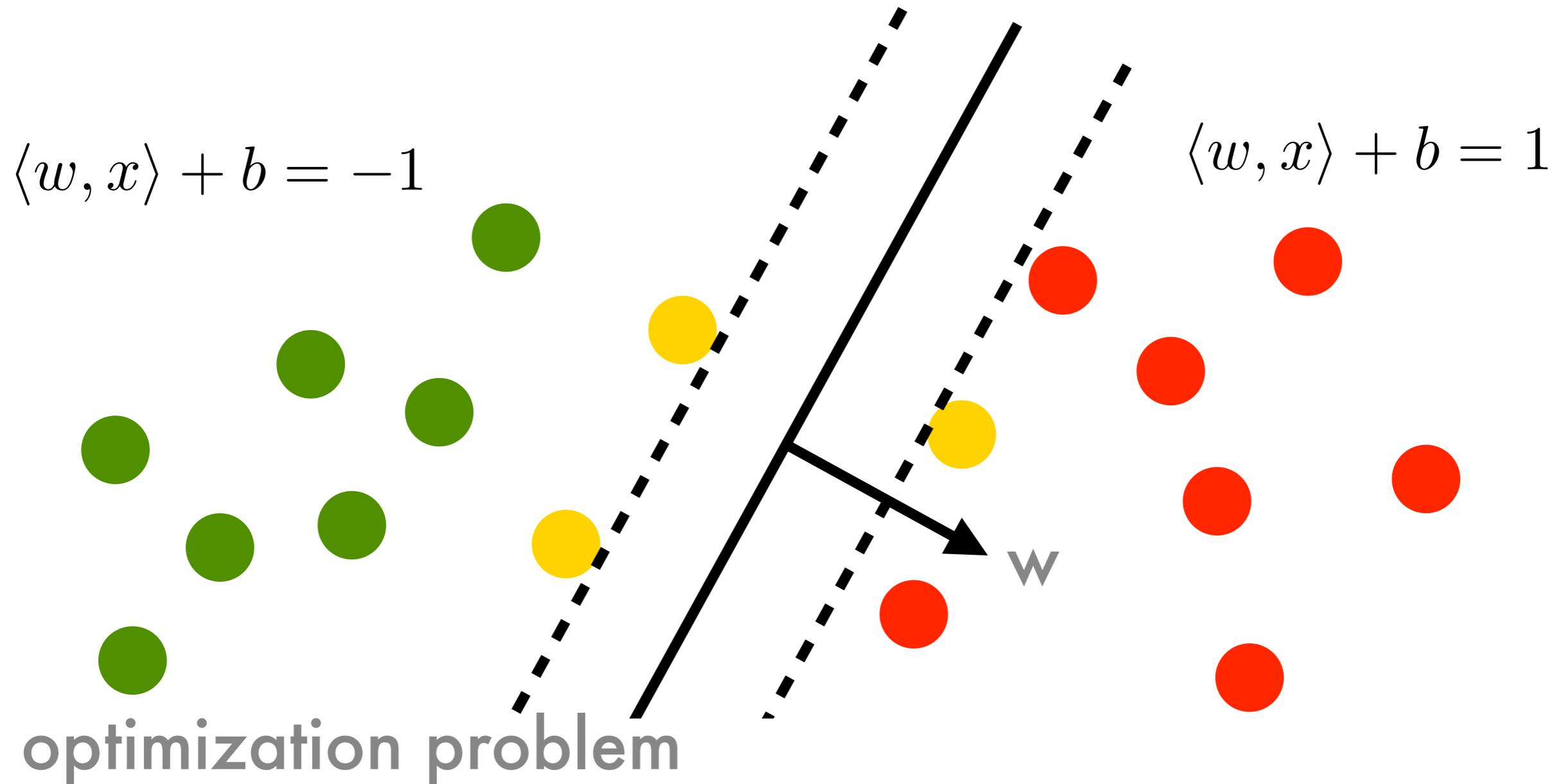
linear function
 $f(x) = \langle w, x \rangle + b$

Large Margin Classifier



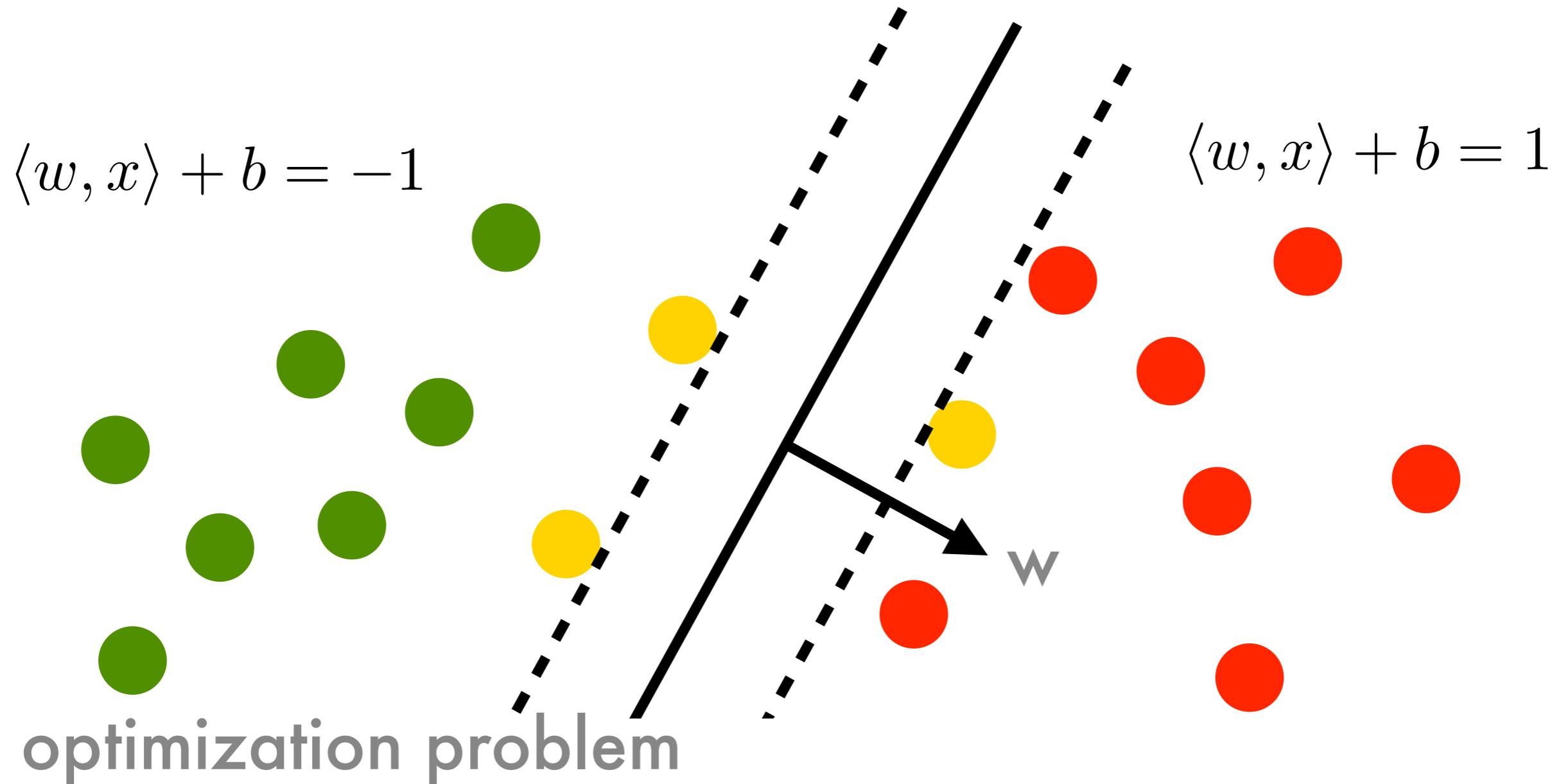
$$\frac{\langle x_+ - x_-, w \rangle}{2\|w\|} = \frac{1}{2\|w\|} [\langle x_+, w \rangle + b - (\langle x_-, w \rangle + b)] = \frac{1}{\|w\|}$$

Large Margin Classifier



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

Large Margin Classifier



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

Dual Problem

- Primal optimization problem

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

- Lagrange function

constraint

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

Optimality in w, b is at saddle point with a

- 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

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

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

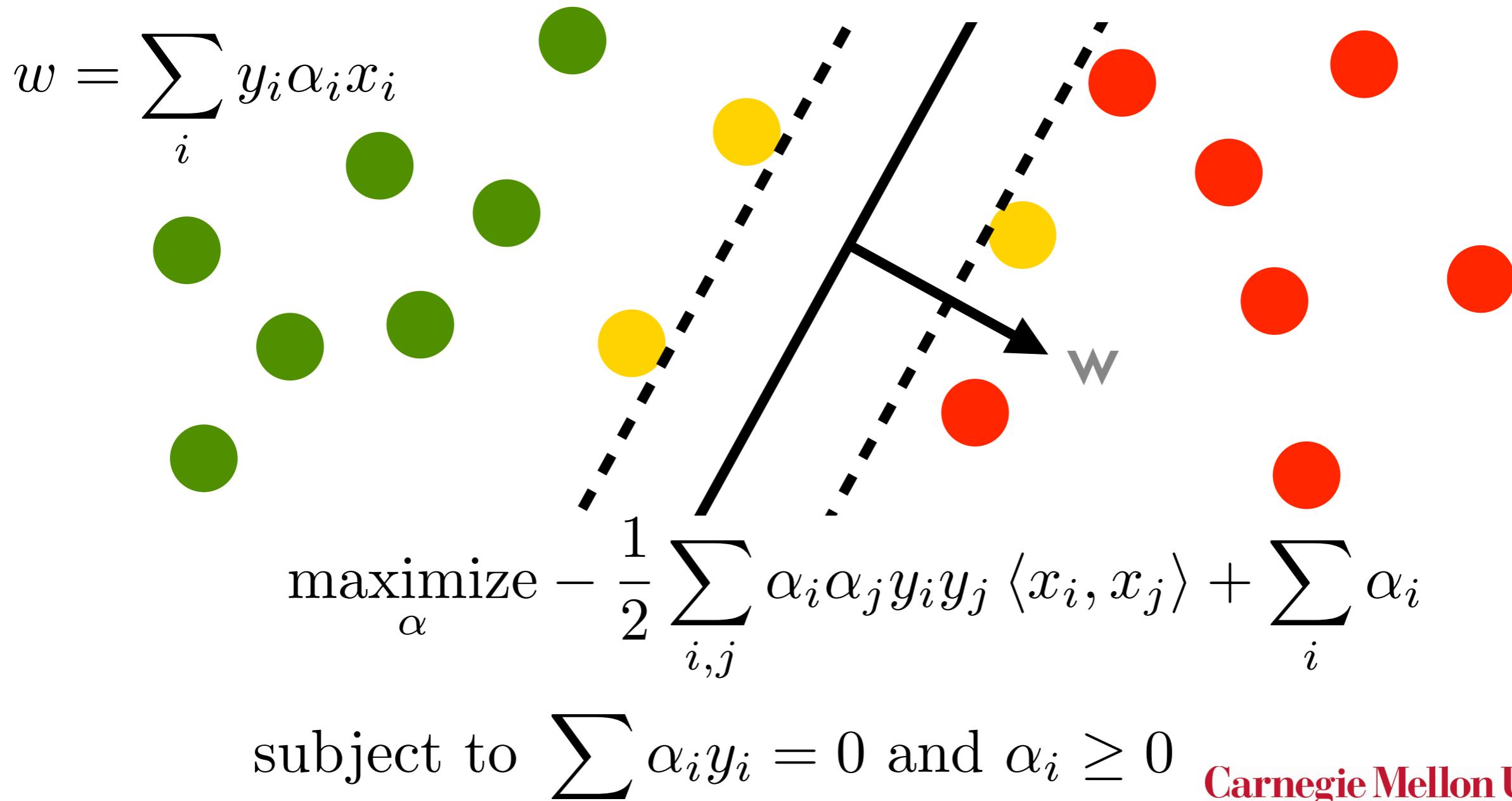
- Plugging terms back into L yields

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

subject to $\sum_i \alpha_i y_i = 0$ and $\alpha_i \geq 0$

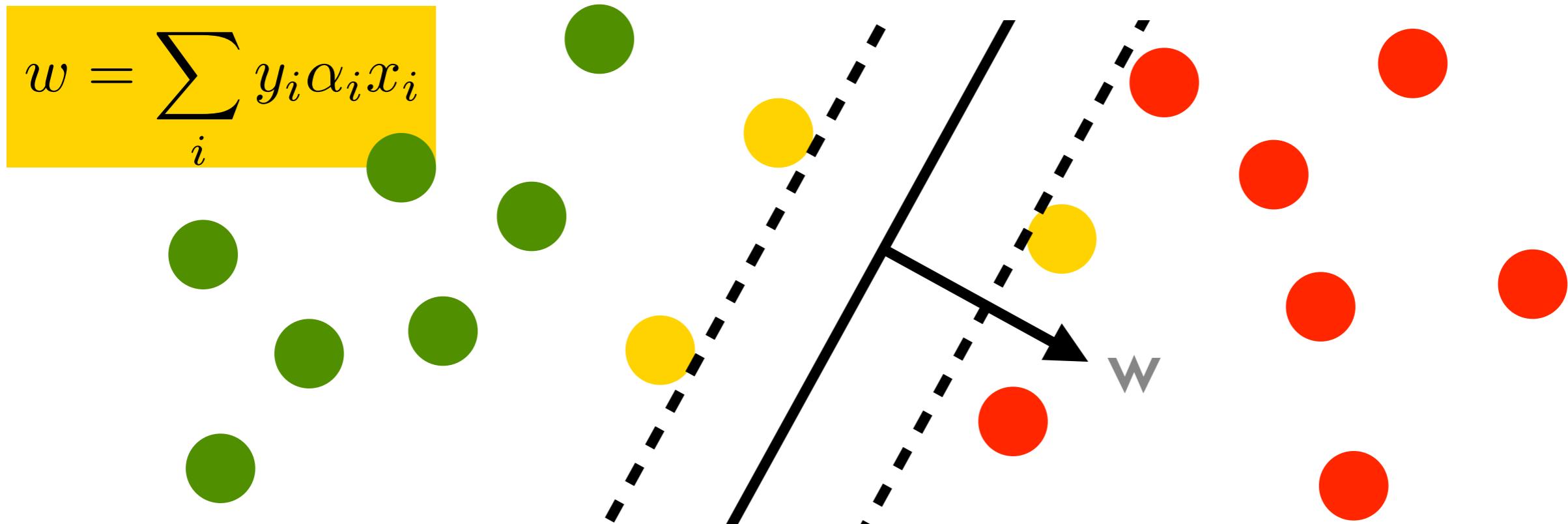
Support Vector Machines

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



Support Vectors

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



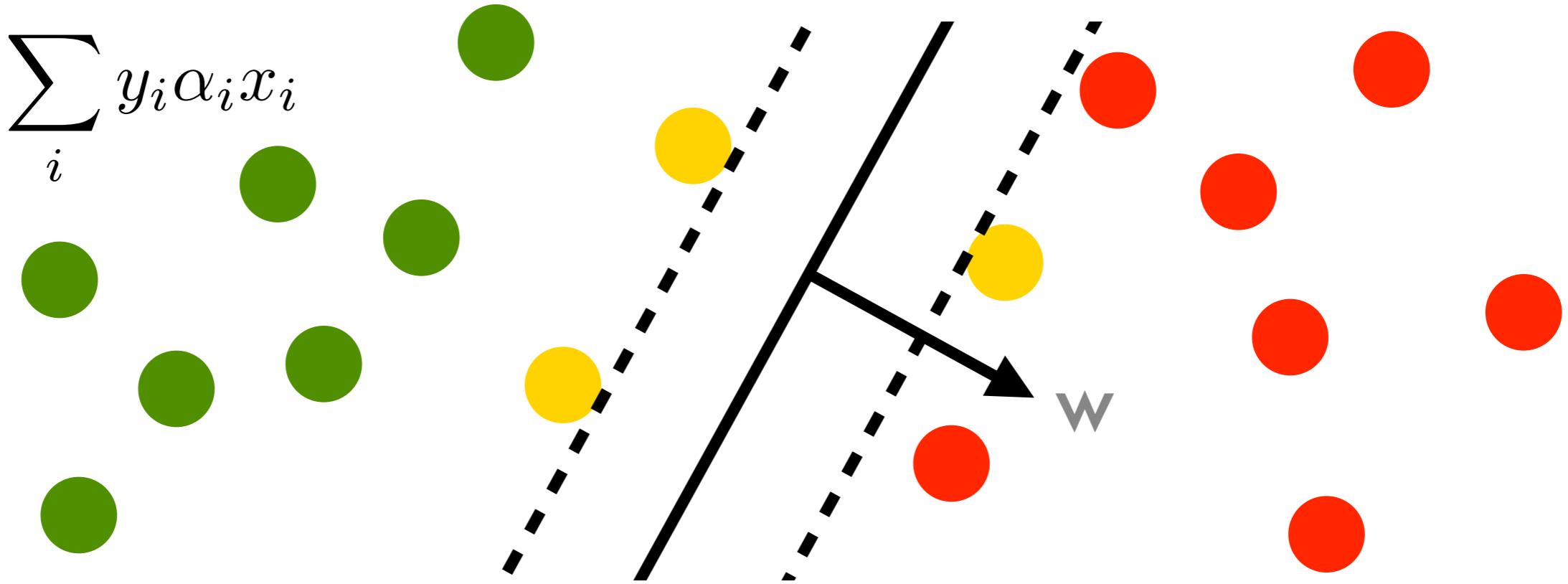
Karush Kuhn Tucker
Optimality condition
 $\alpha_i [y_i [\langle w, x_i \rangle + b] - 1] = 0$



$$\begin{aligned}\alpha_i &= 0 \\ \alpha_i > 0 &\Rightarrow y_i [\langle w, x_i \rangle + b] = 1\end{aligned}$$

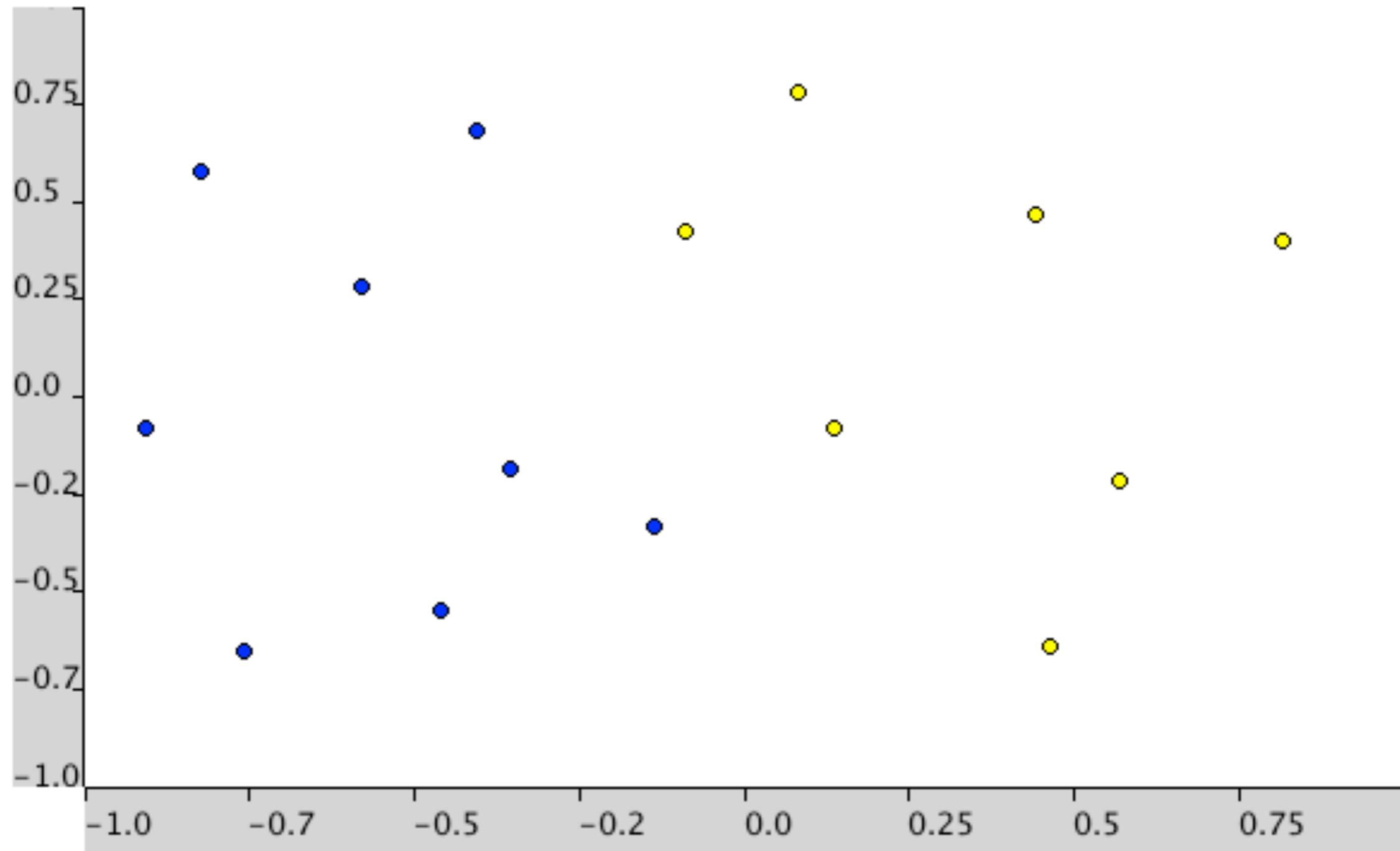
Properties

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



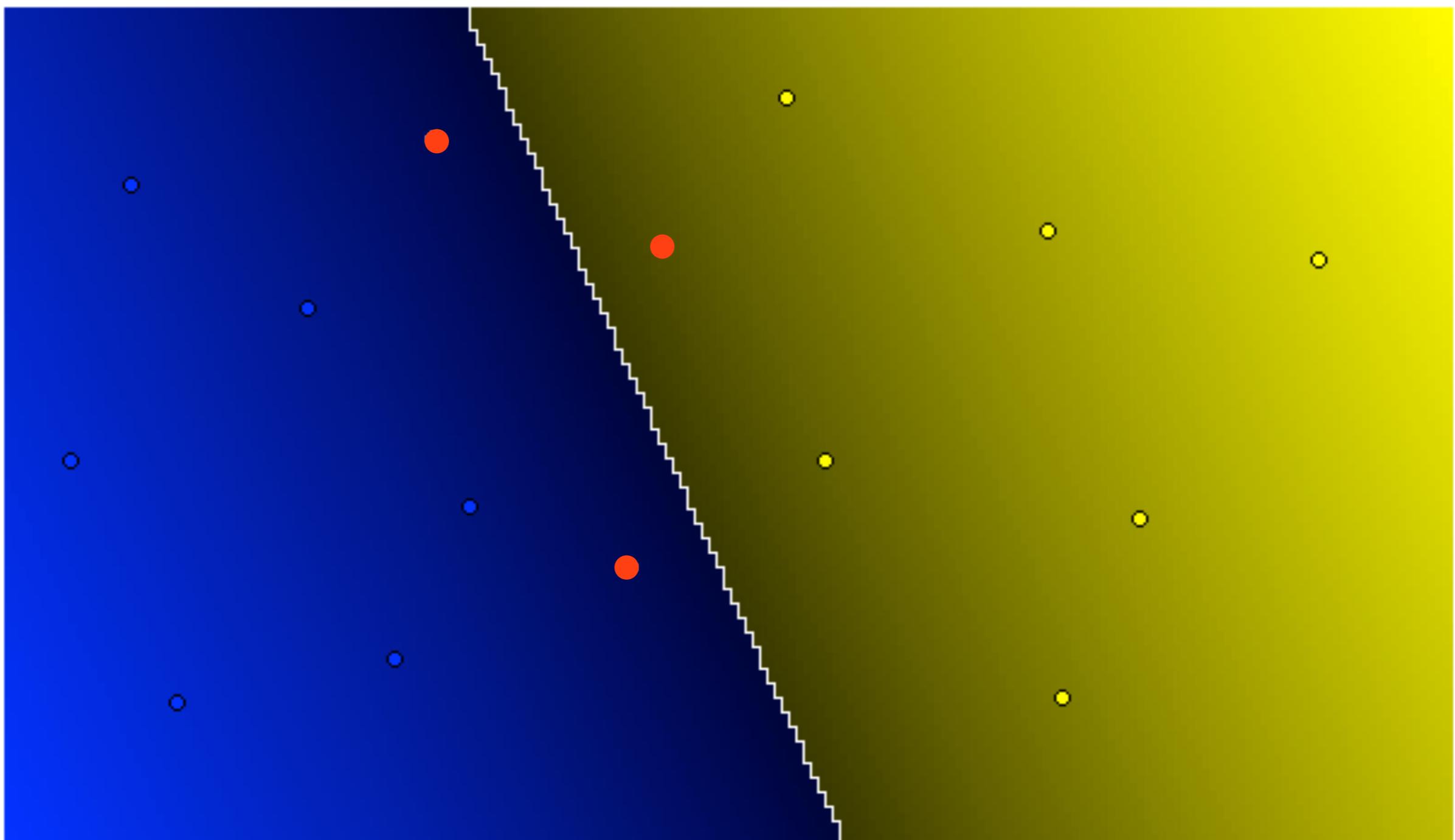
- Weight vector w as weighted linear combination of instances
- Only points on margin matter (ignore the rest and get same solution)
- Only inner products matter
 - Quadratic program
 - We can replace the inner product by a kernel
- Keeps instances away from the margin

Example

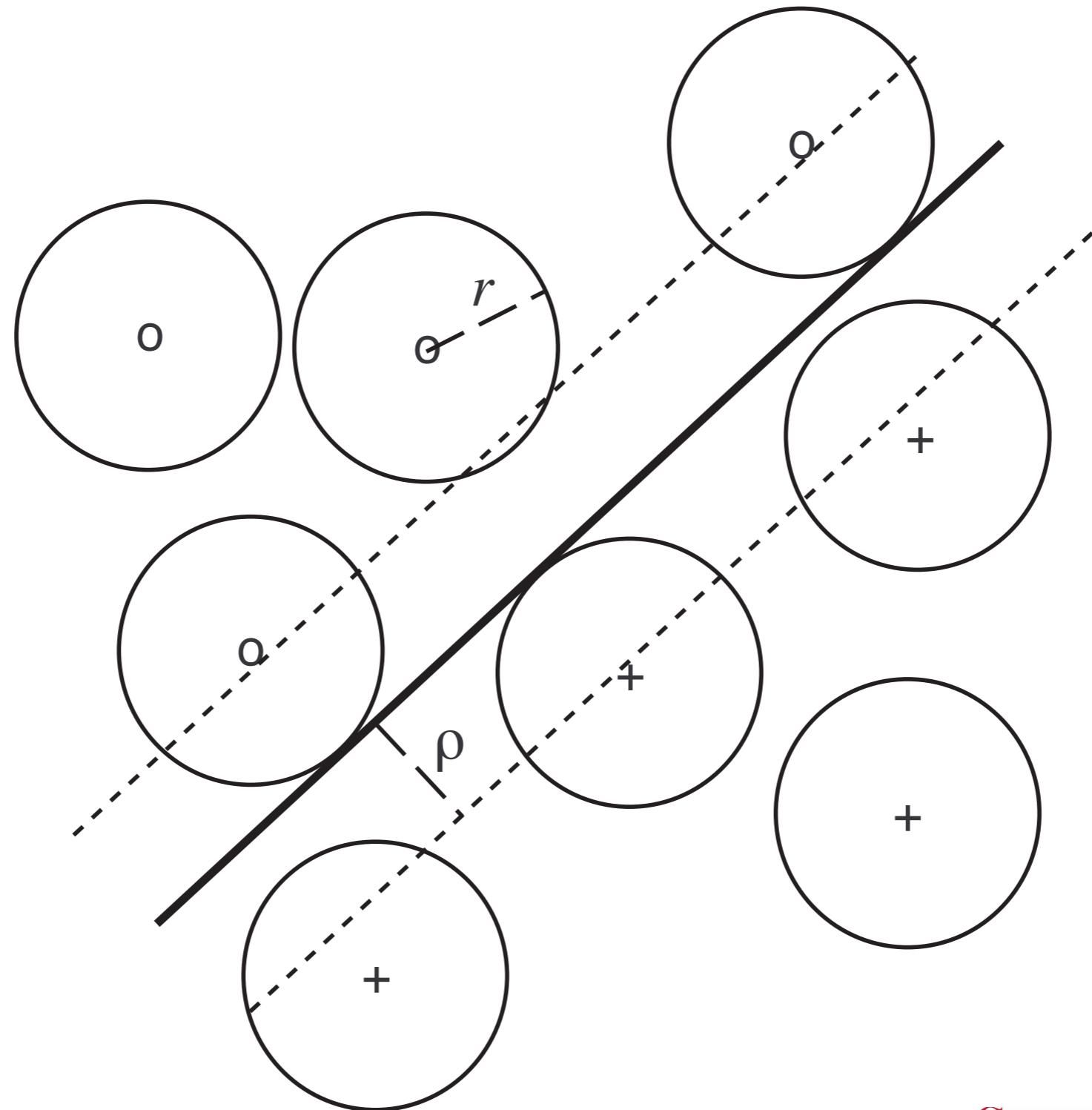


Example

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



Why large margins?

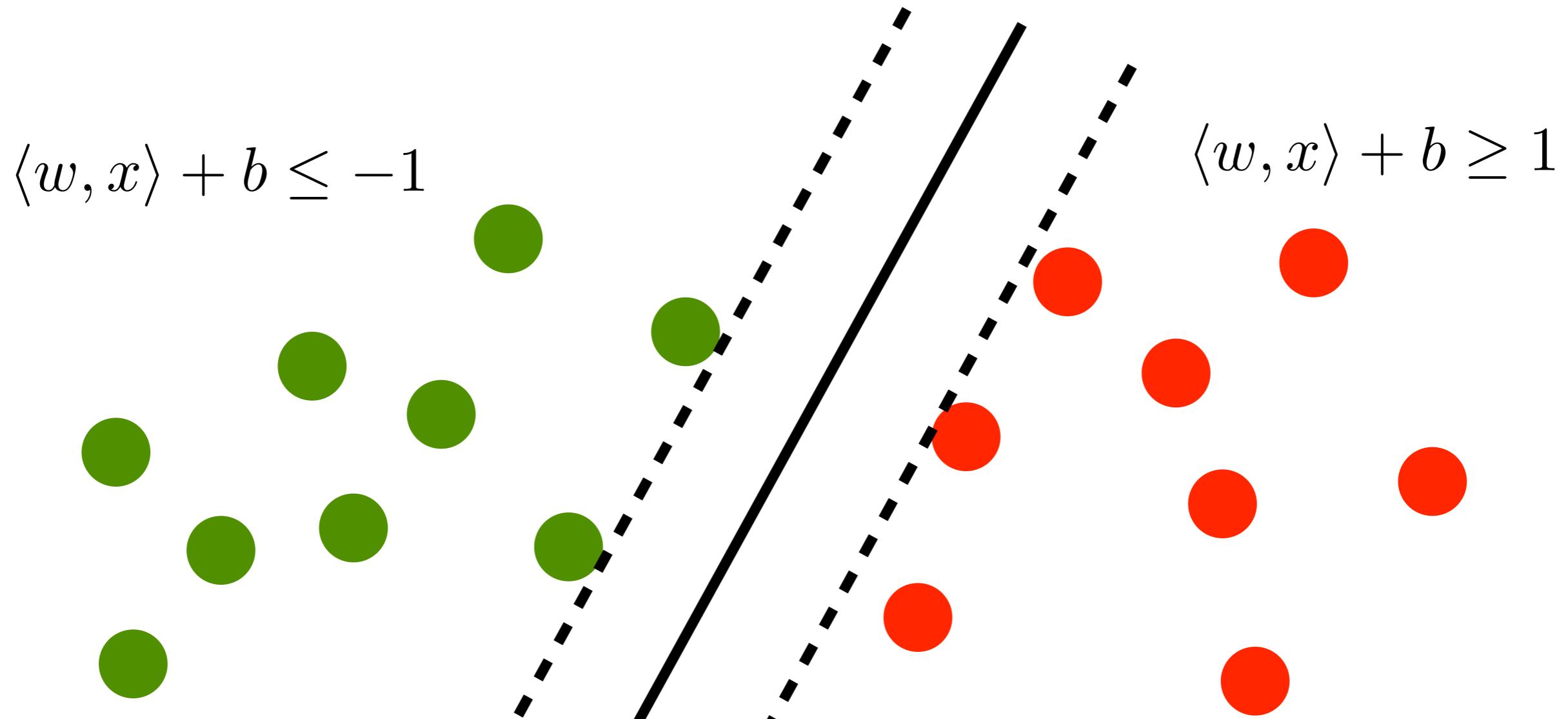




**SOFT
MARCH**

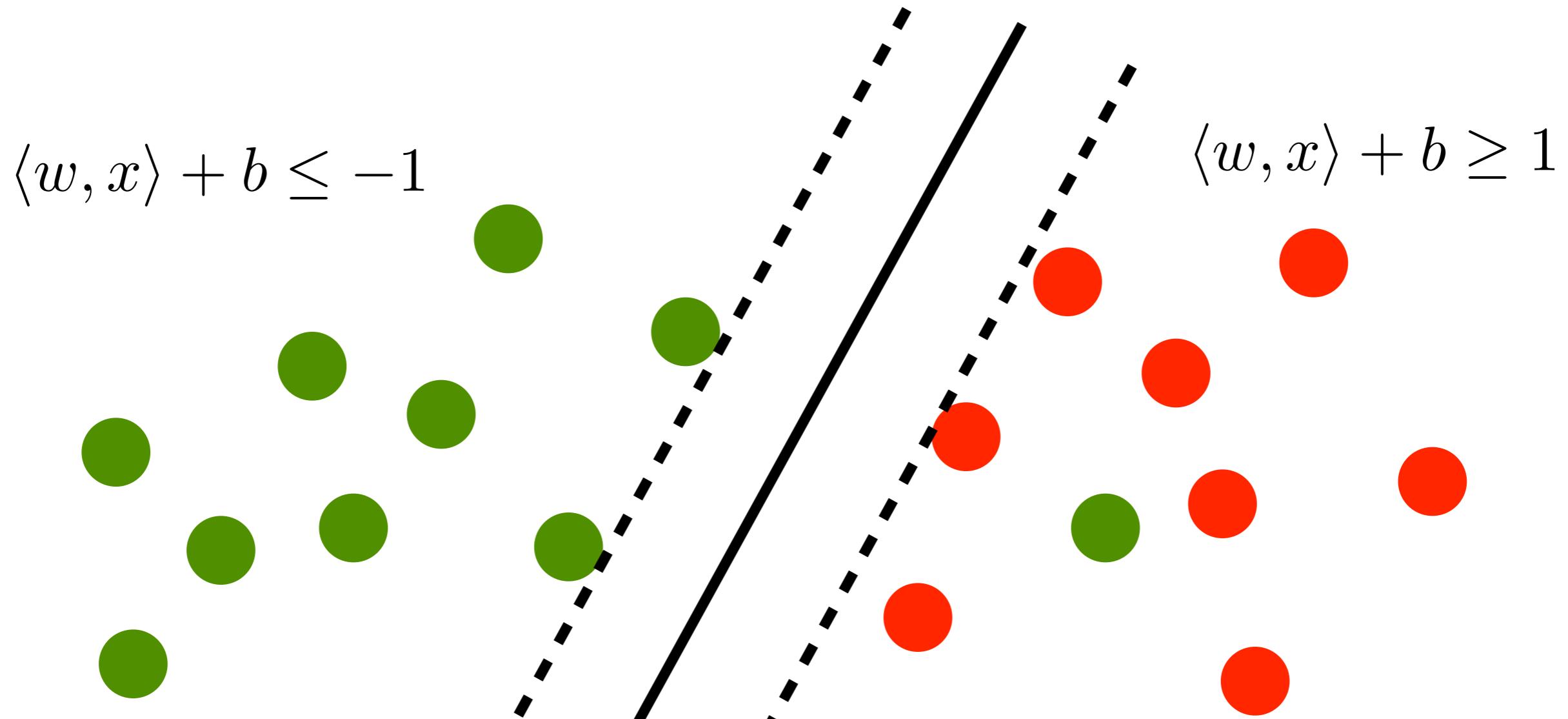
CLASSIFIERS

Large Margin Classifier

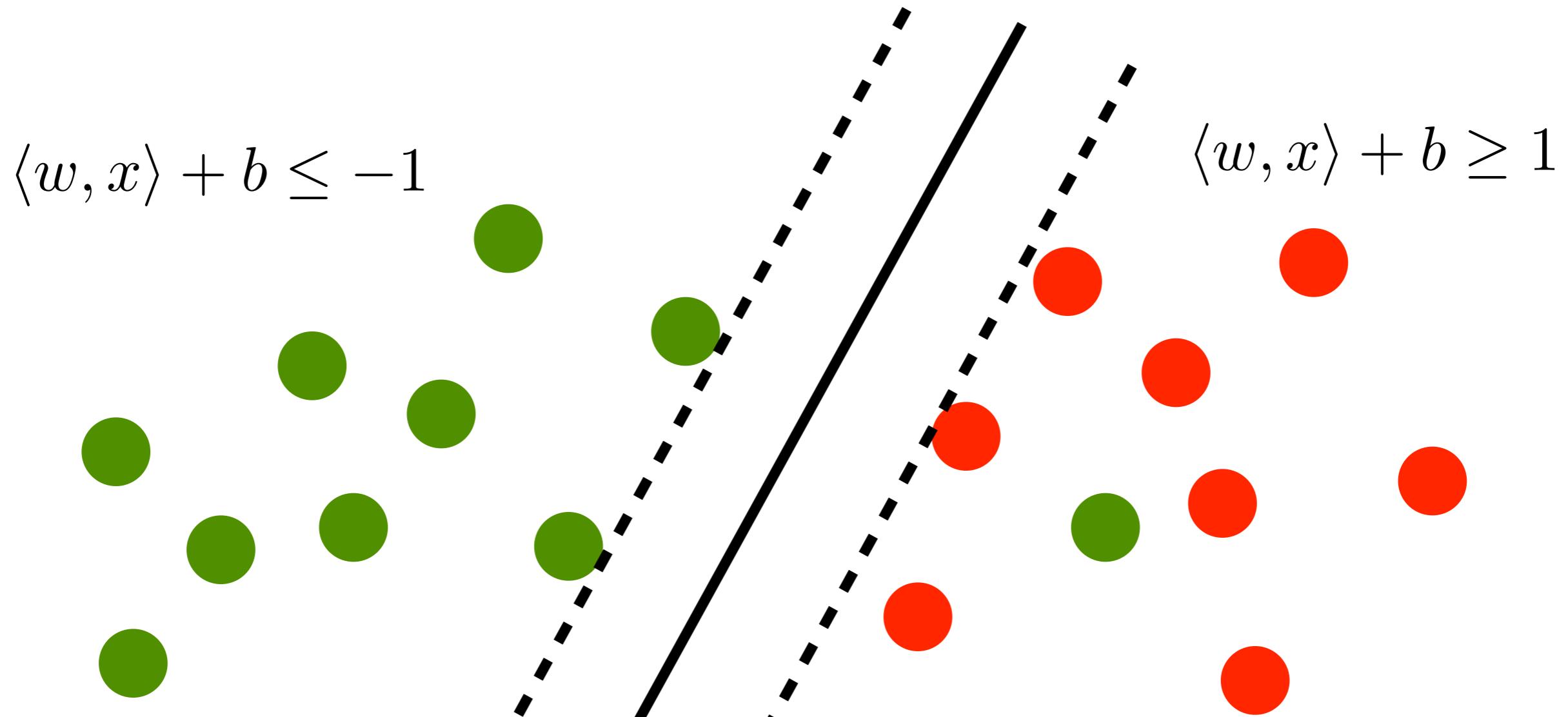


linear function
 $f(x) = \langle w, x \rangle + b$

Large Margin Classifier



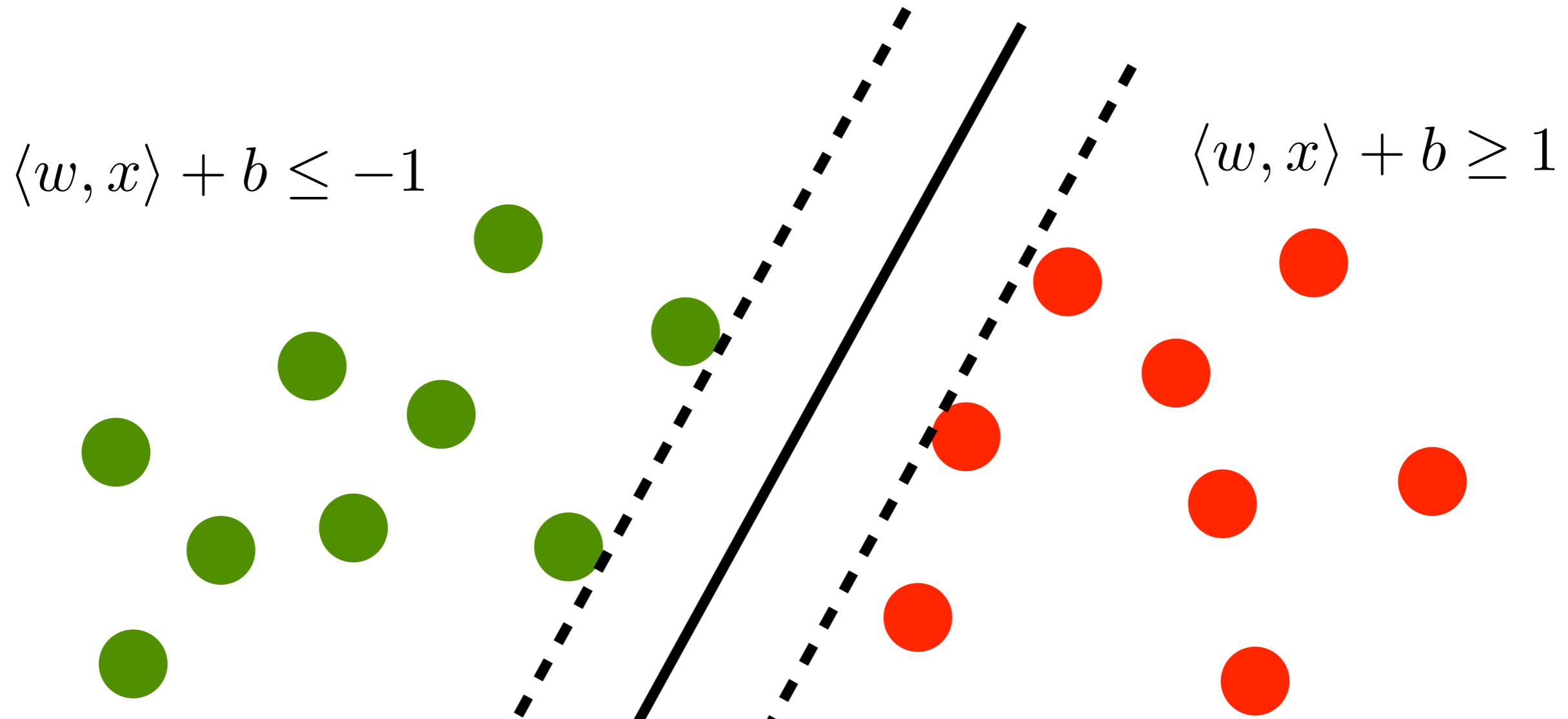
Large Margin Classifier



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

linear separator
is impossible

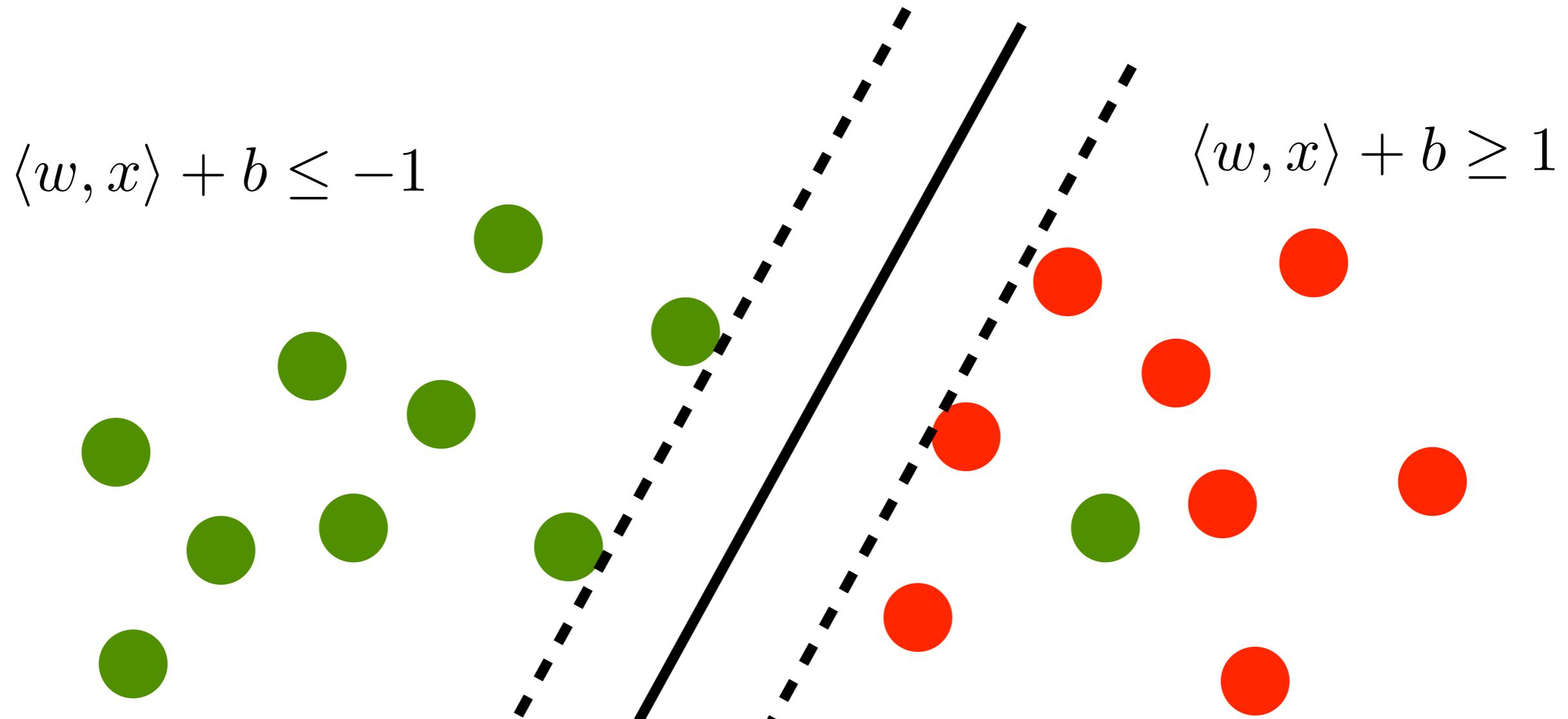
Large Margin Classifier



Theorem (Minsky & Papert)

Finding the minimum error separating hyperplane is NP hard

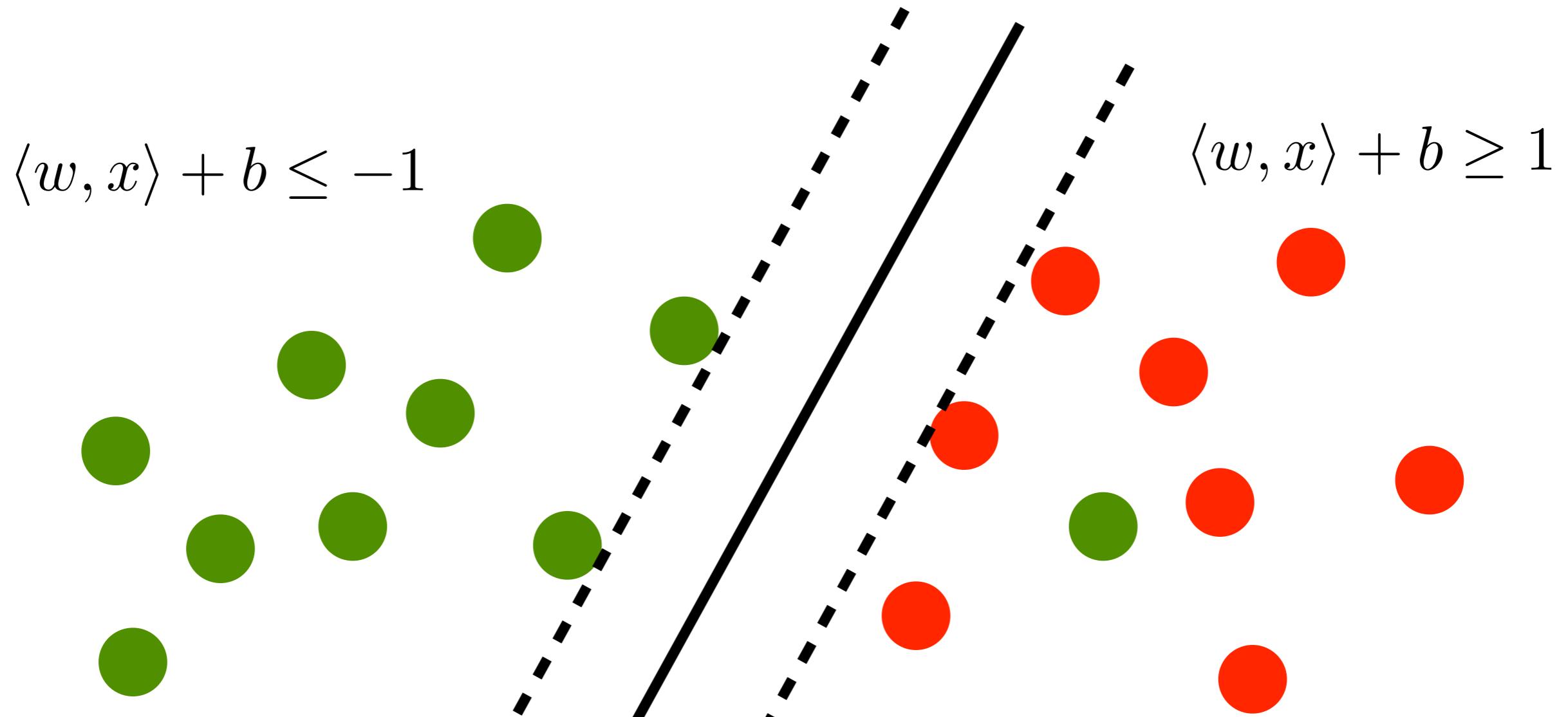
Large Margin Classifier



Theorem (Minsky & Papert)

Finding the minimum error separating hyperplane is NP hard

Large Margin Classifier

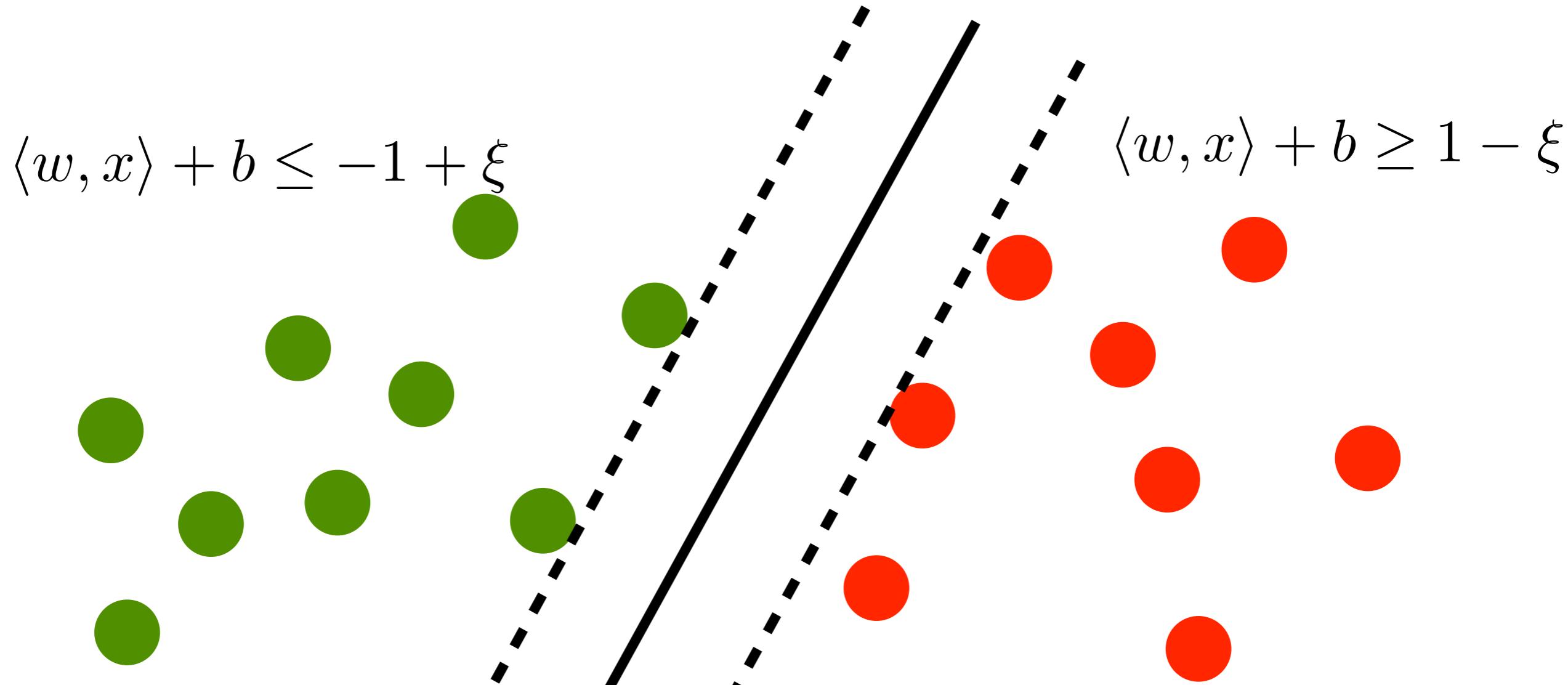


**minimum error separator
is impossible**

Theorem (Minsky & Papert)

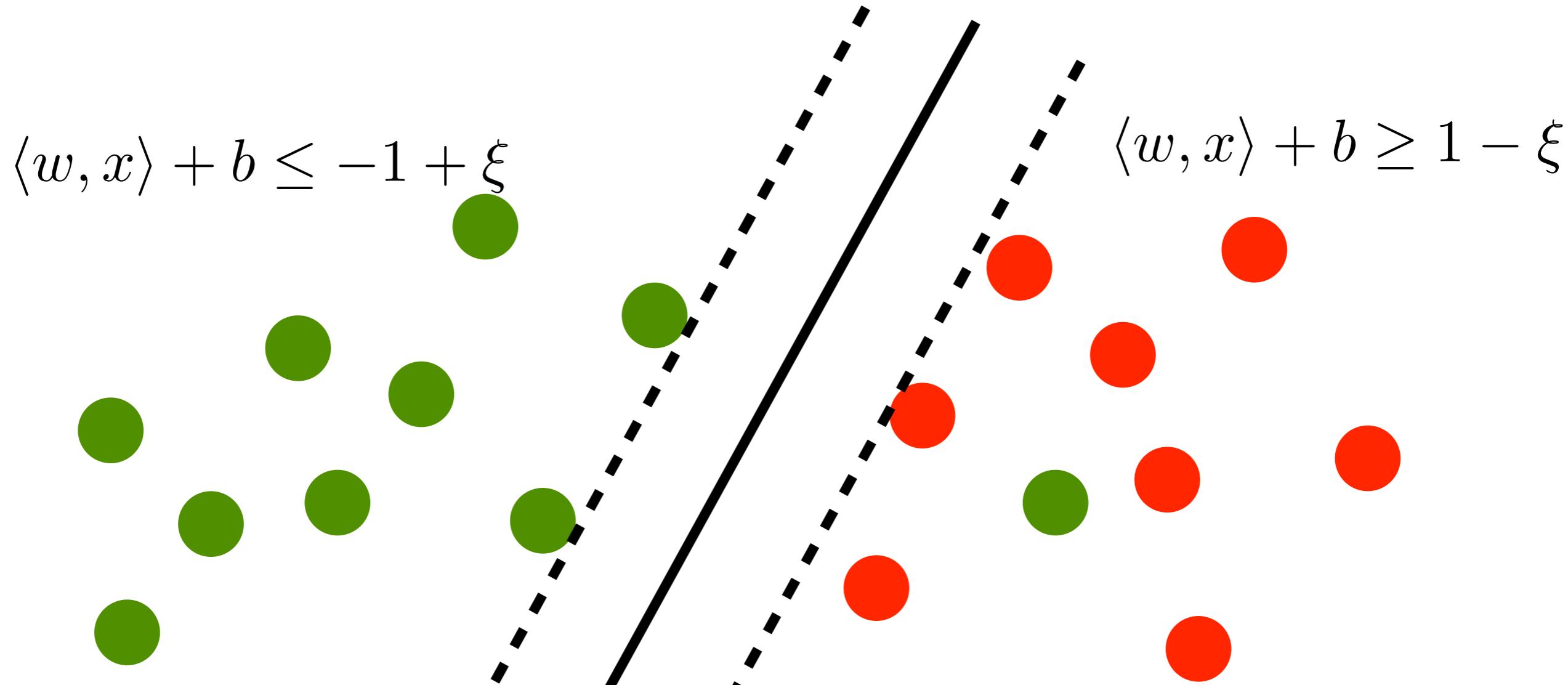
Finding the minimum error separating hyperplane is NP hard

Adding slack variables



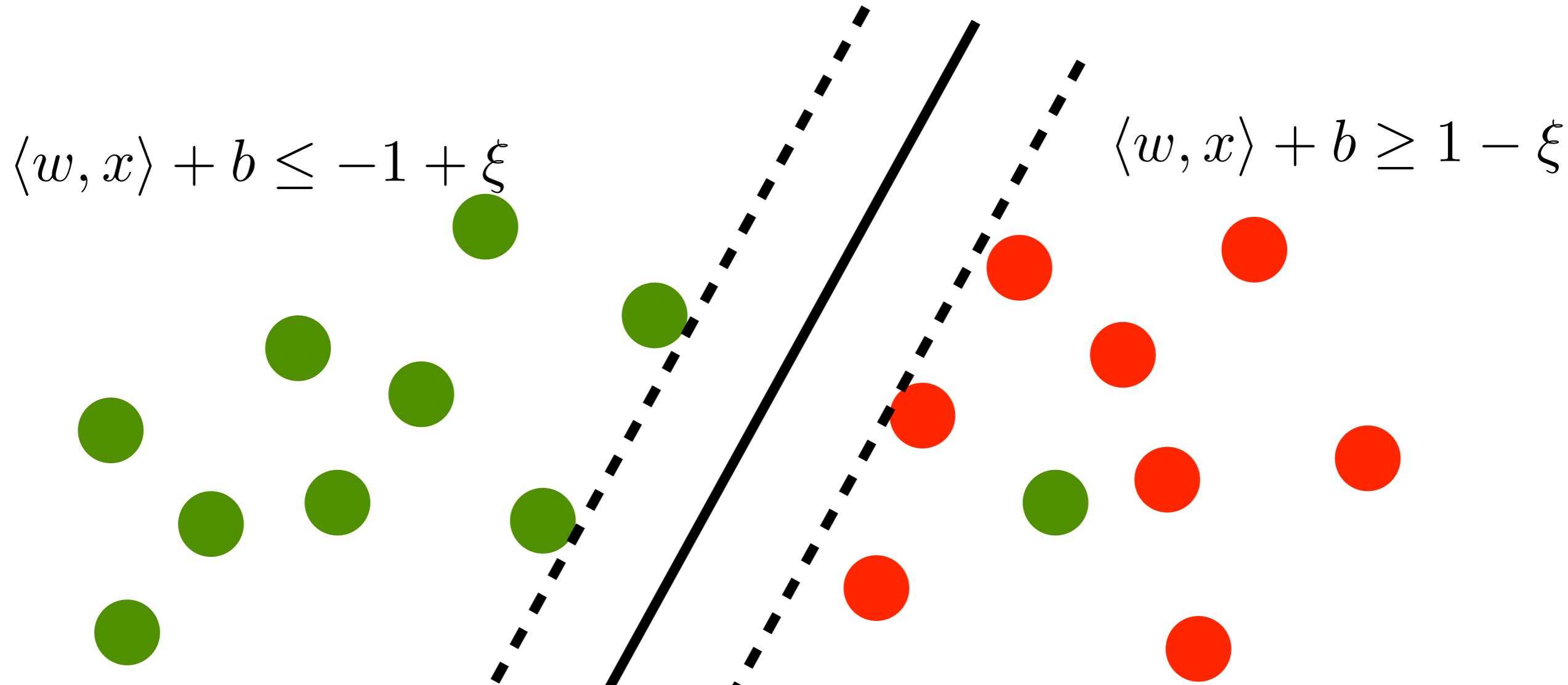
Convex optimization problem

Adding slack variables



Convex optimization problem

Adding slack variables



Convex optimization problem

minimize amount
of slack

Adding slack variables

- Hard margin problem

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

- With slack variables

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

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

Problem is always feasible. $w = 0$ and $b = 0$ and $\xi_i = 1$

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

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

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

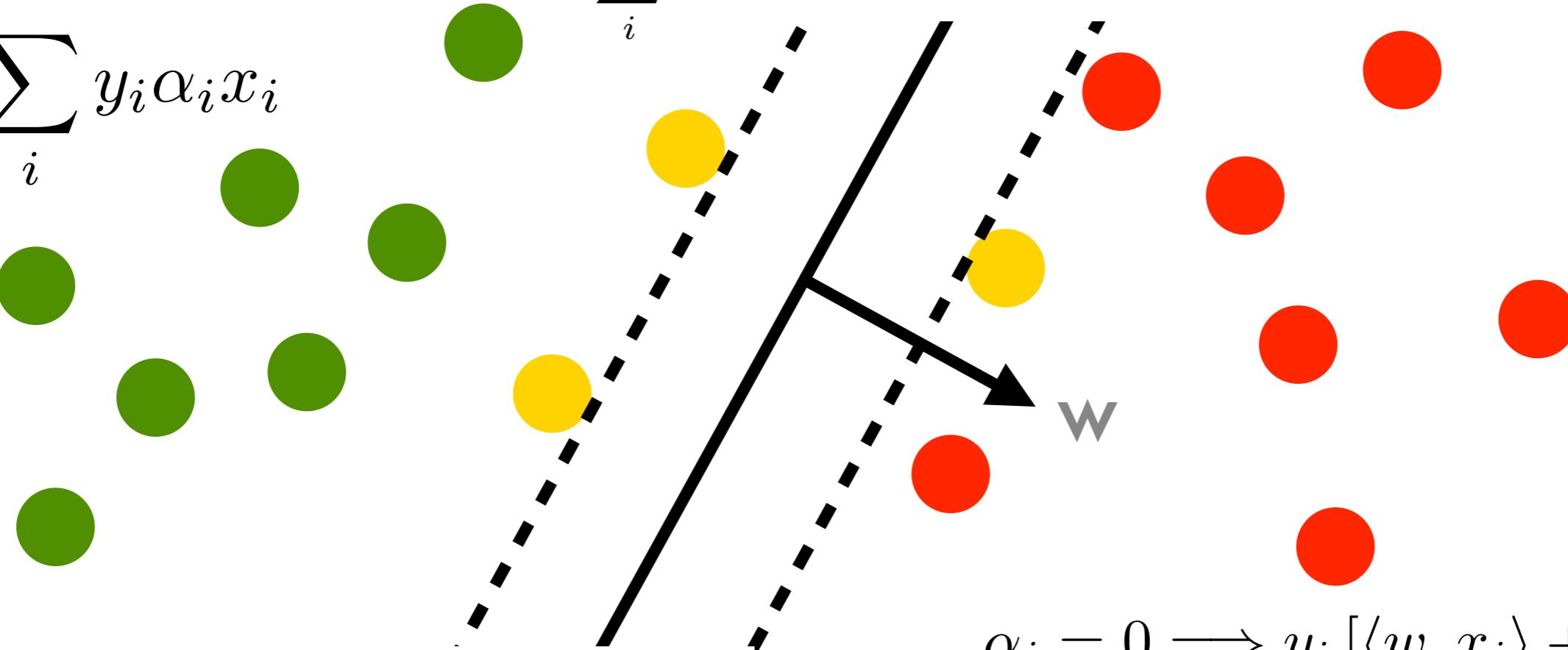
bound
influence

Karush Kuhn Tucker Conditions

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

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

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



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

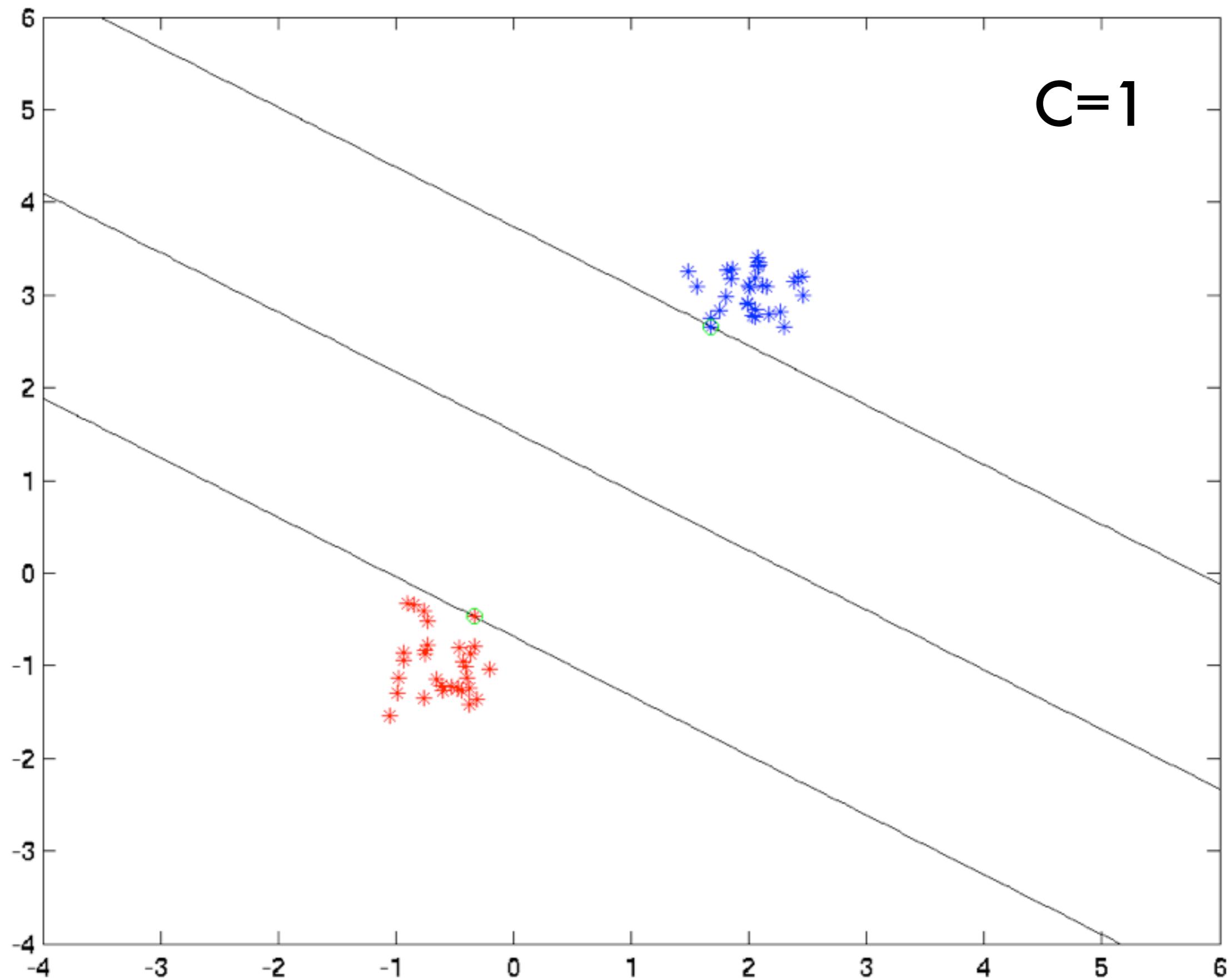
$$0 < \alpha_i < C \implies y_i [\langle w, x_i \rangle + b] = 1$$

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

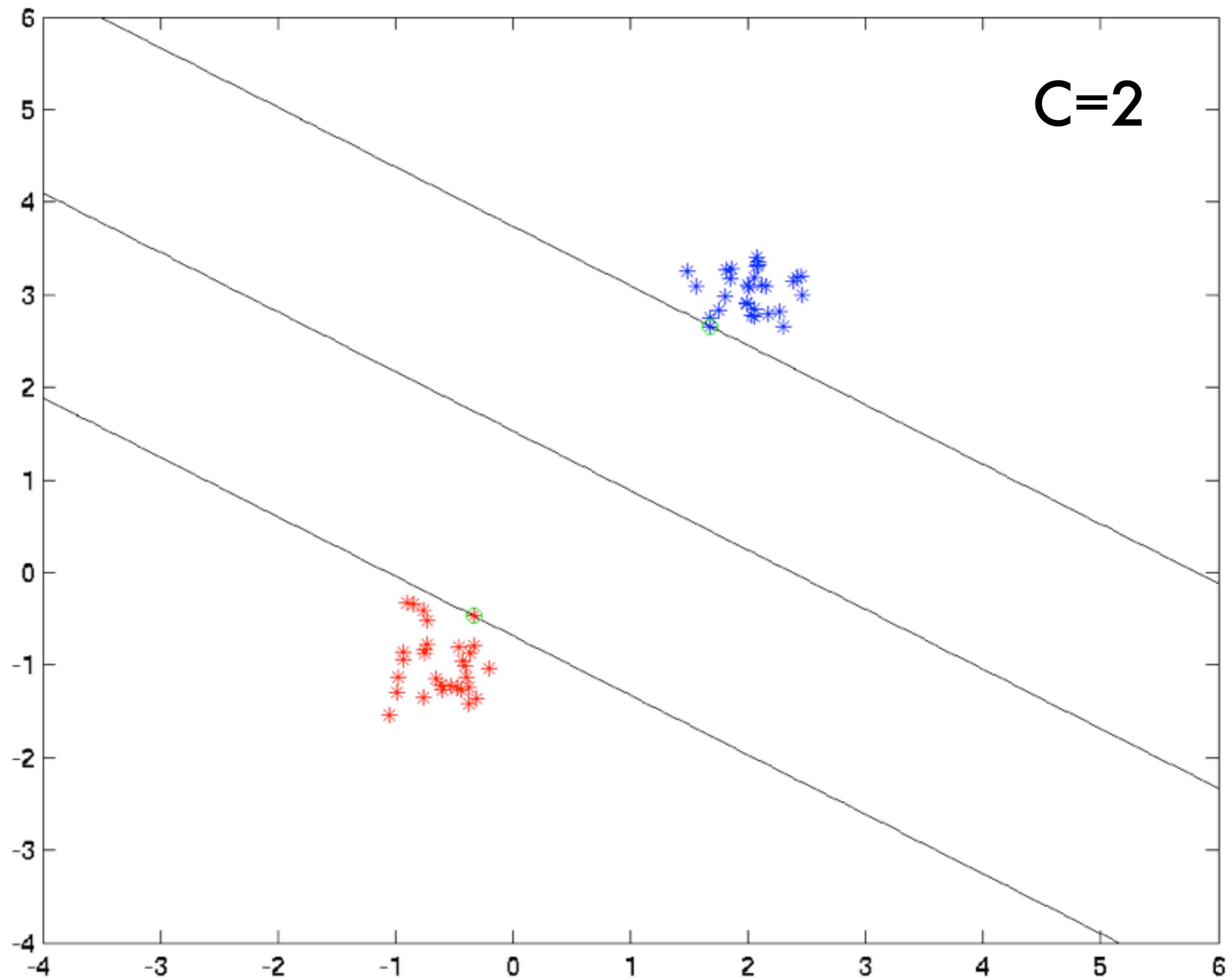
$$\alpha_i [y_i [\langle w, x_i \rangle + b] + \xi_i - 1] = 0$$

$$\eta_i \xi_i = 0$$

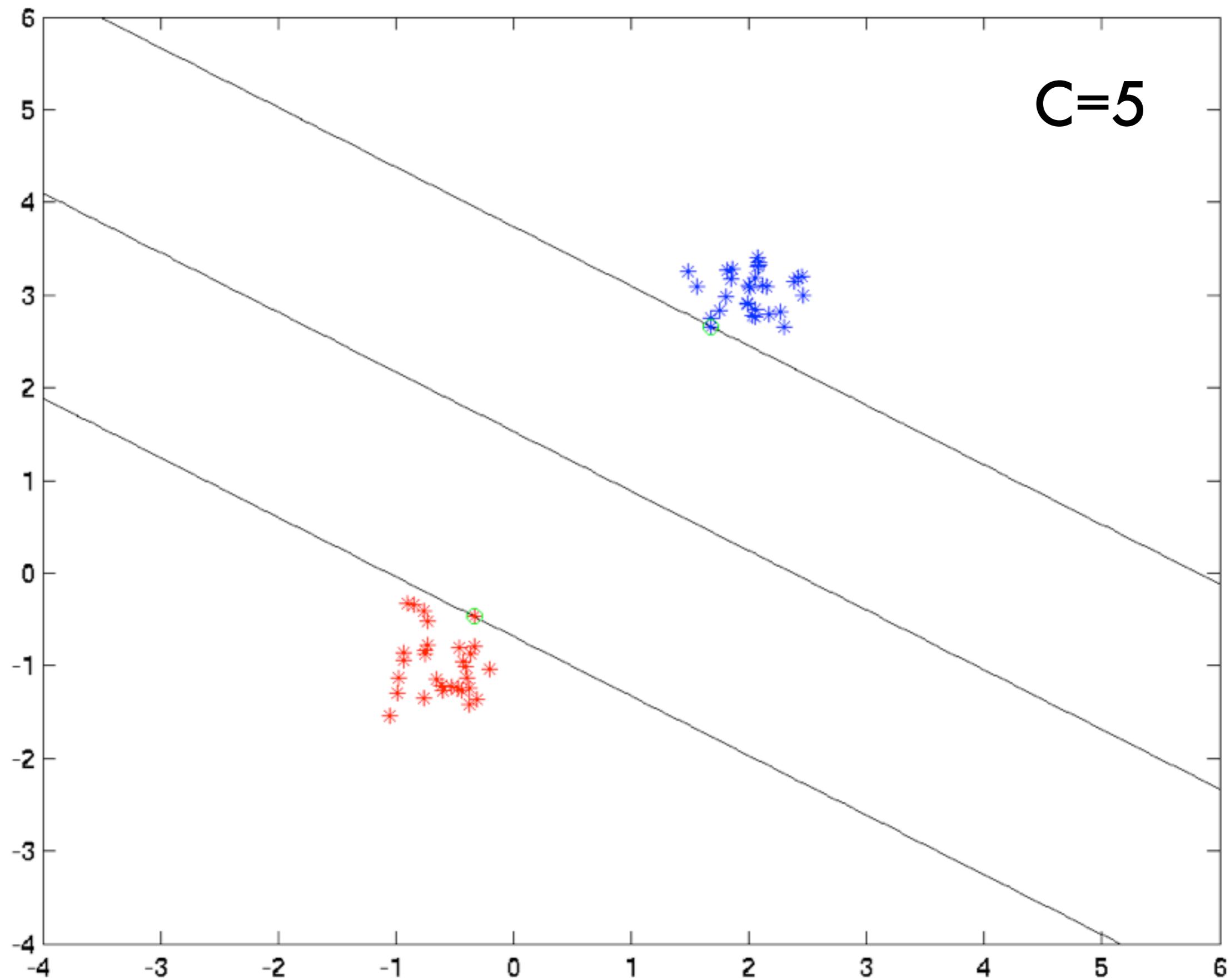
C=1



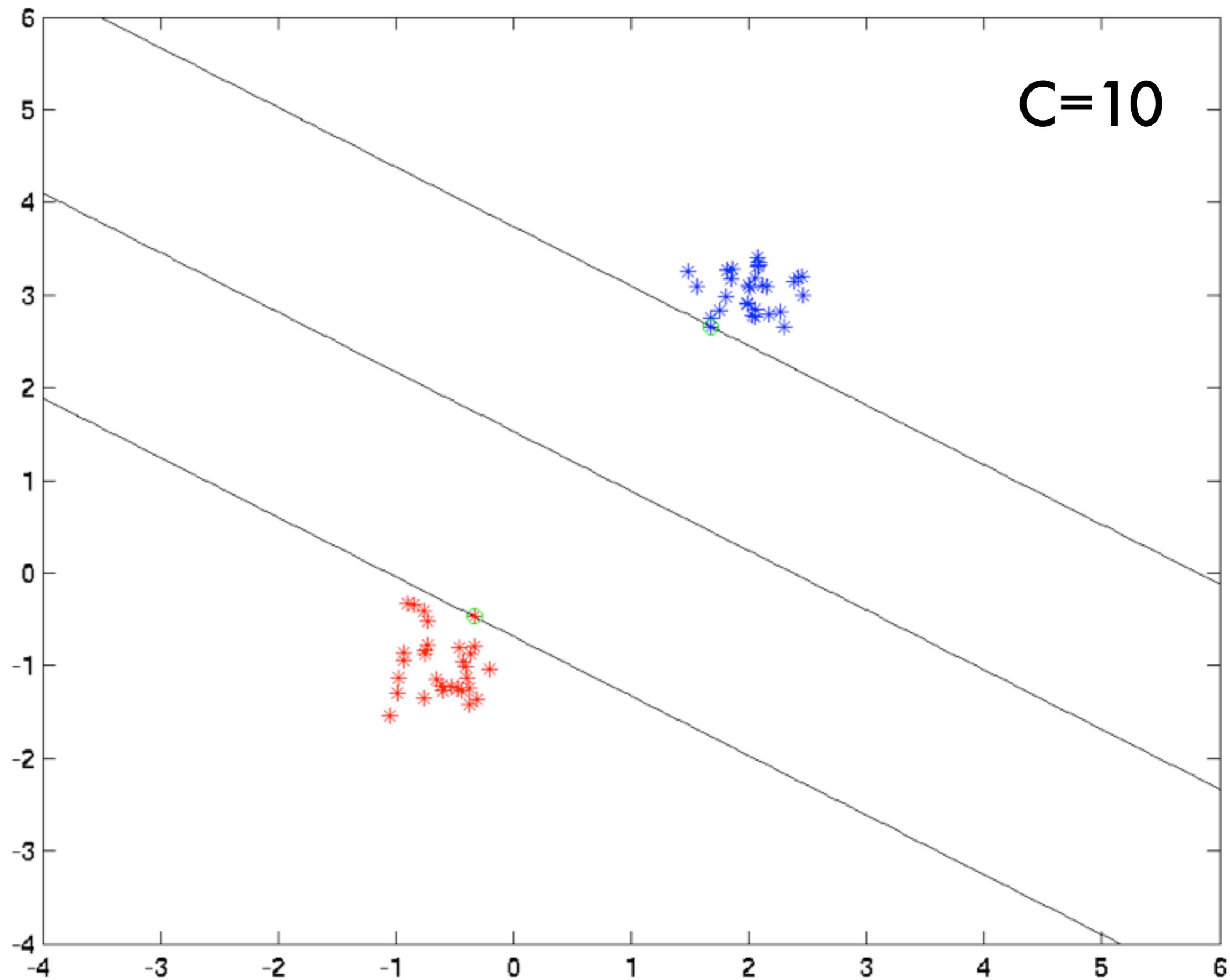
C=2



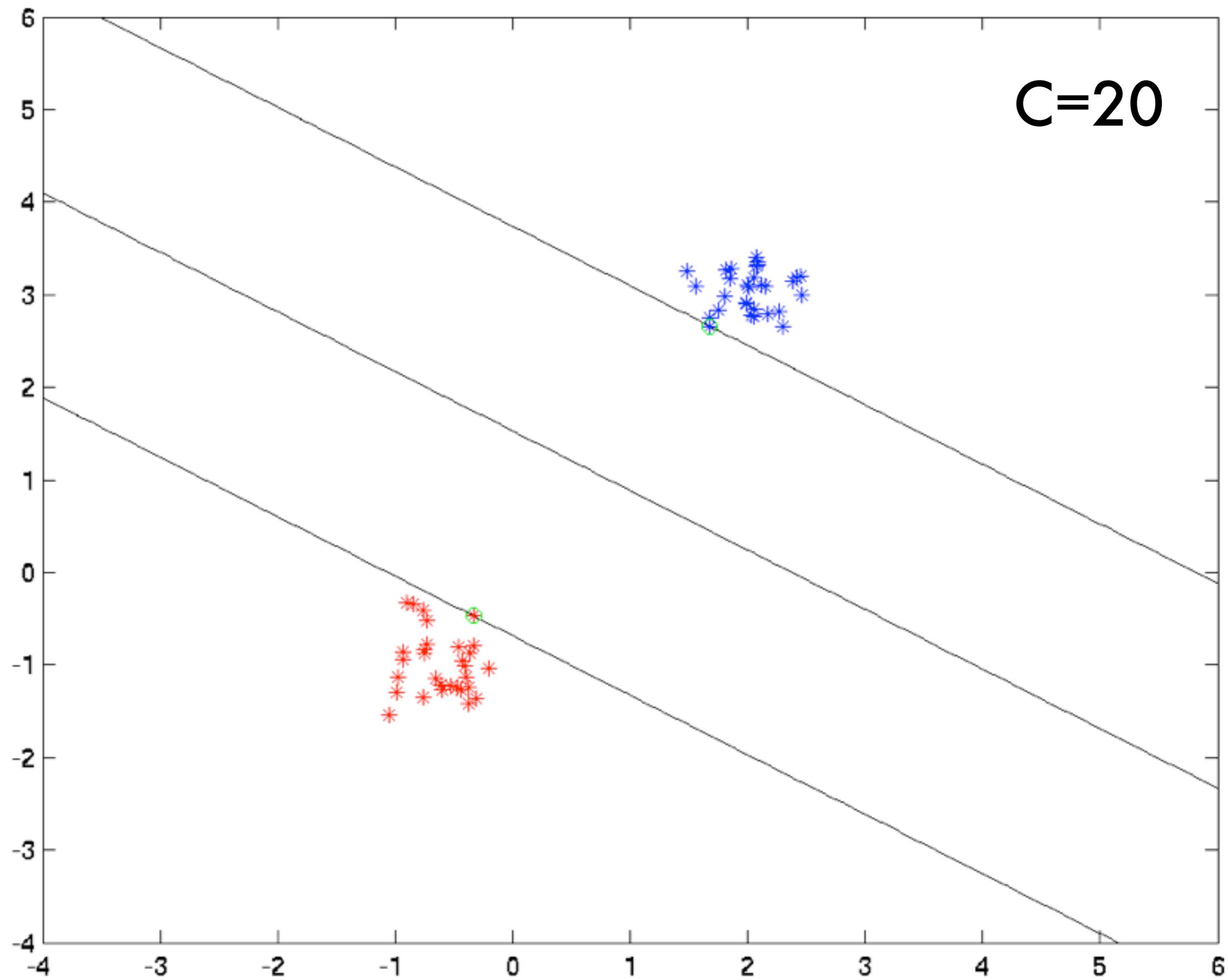
C=5



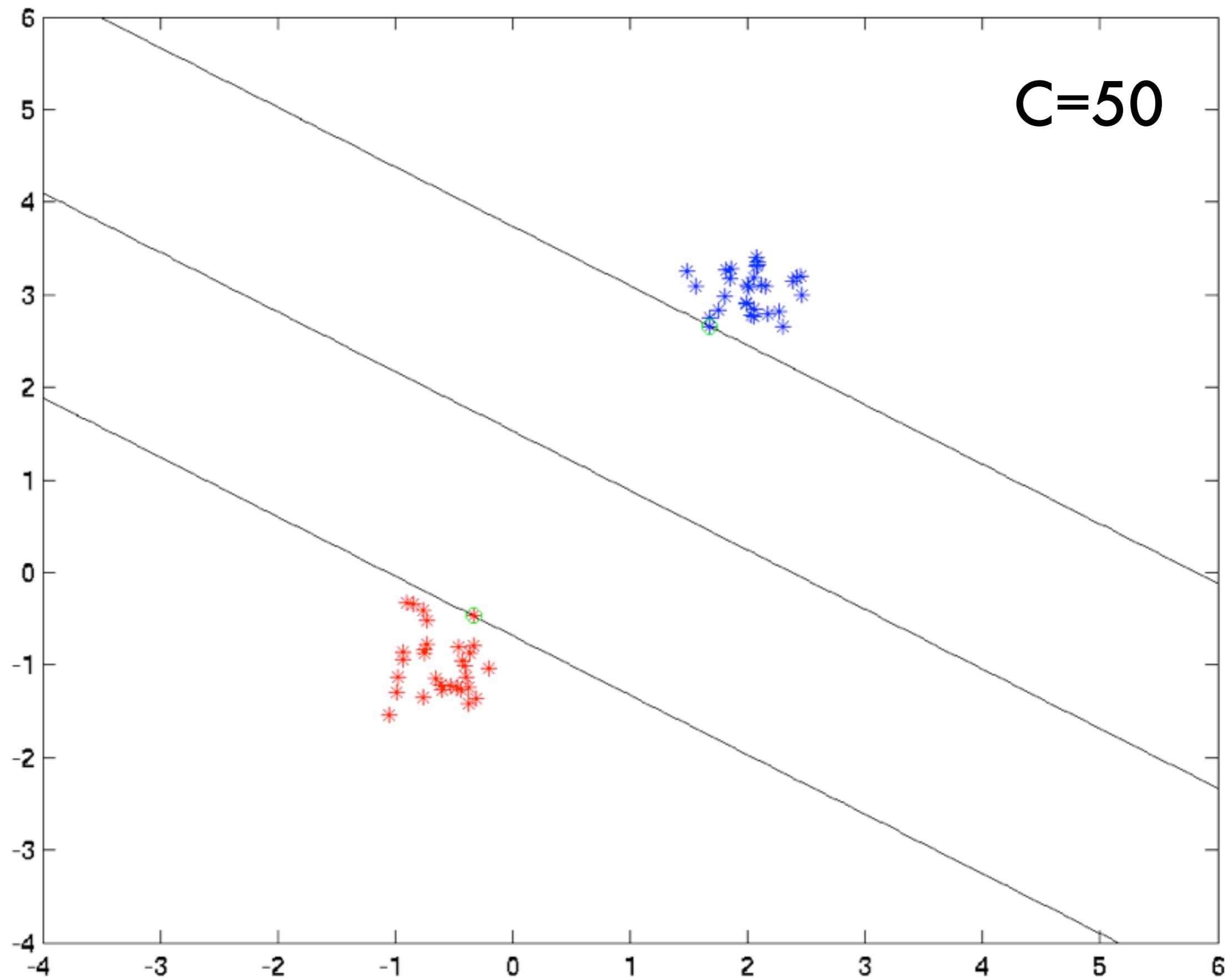
C=10



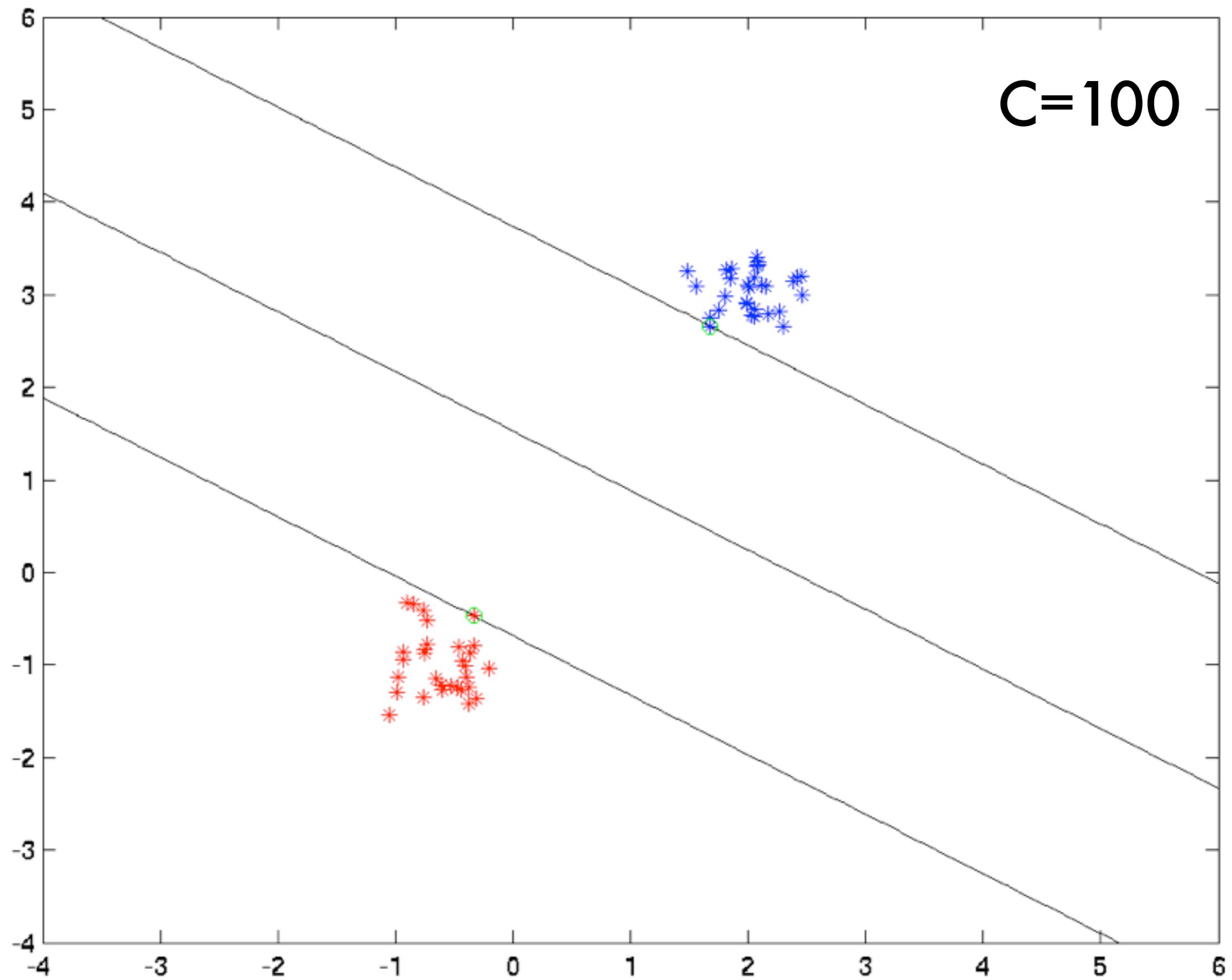
C=20



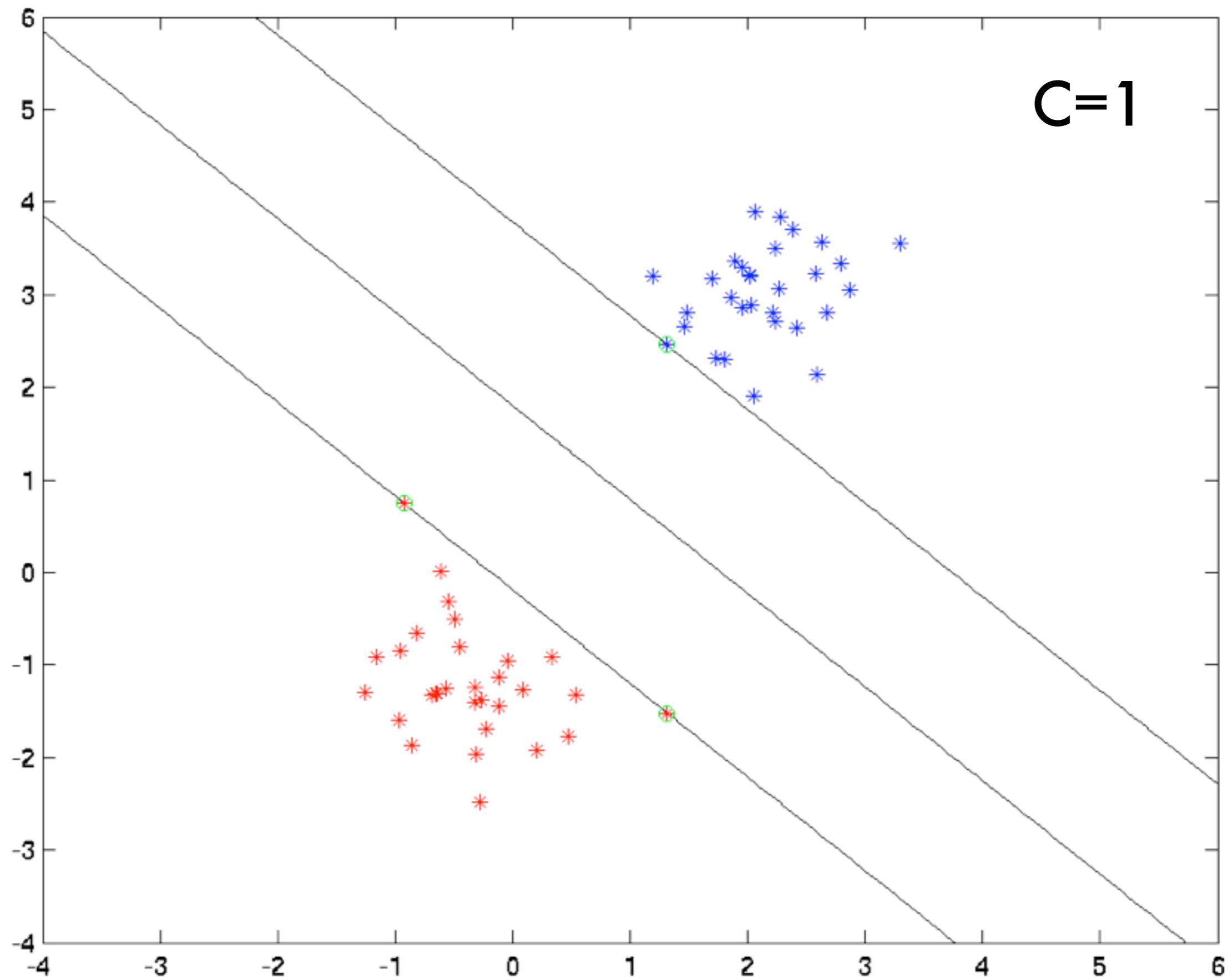
C=50



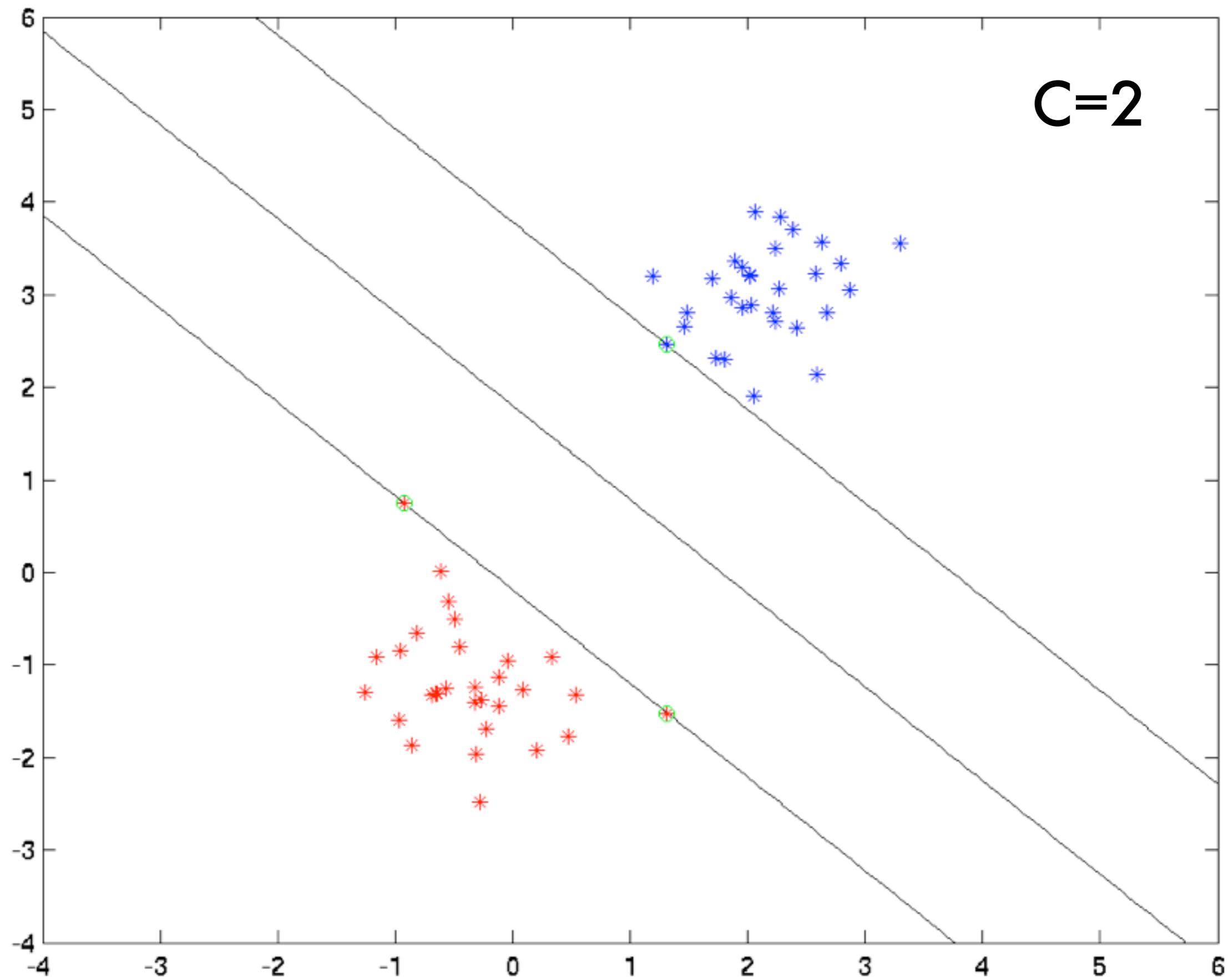
C=100



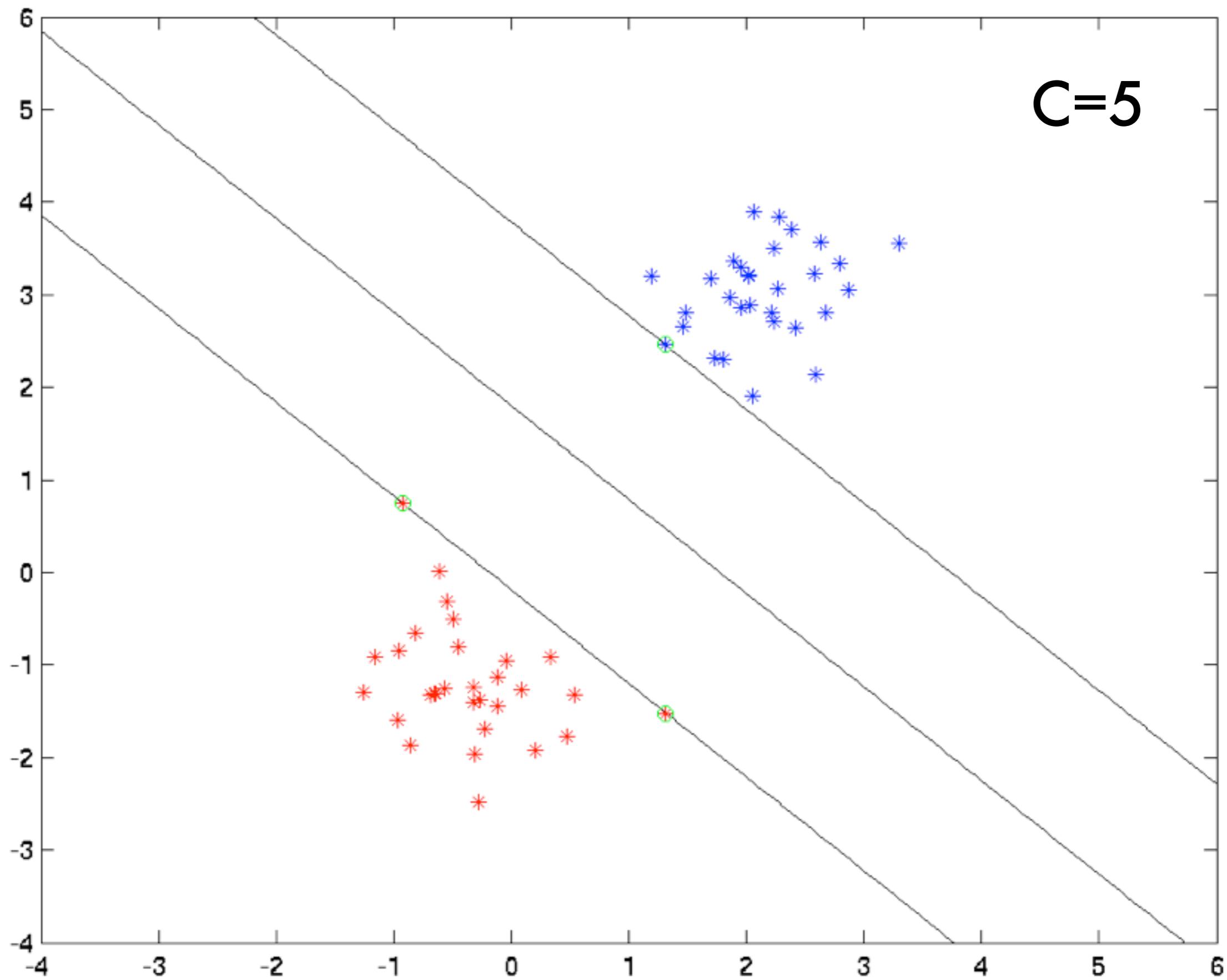
C=1



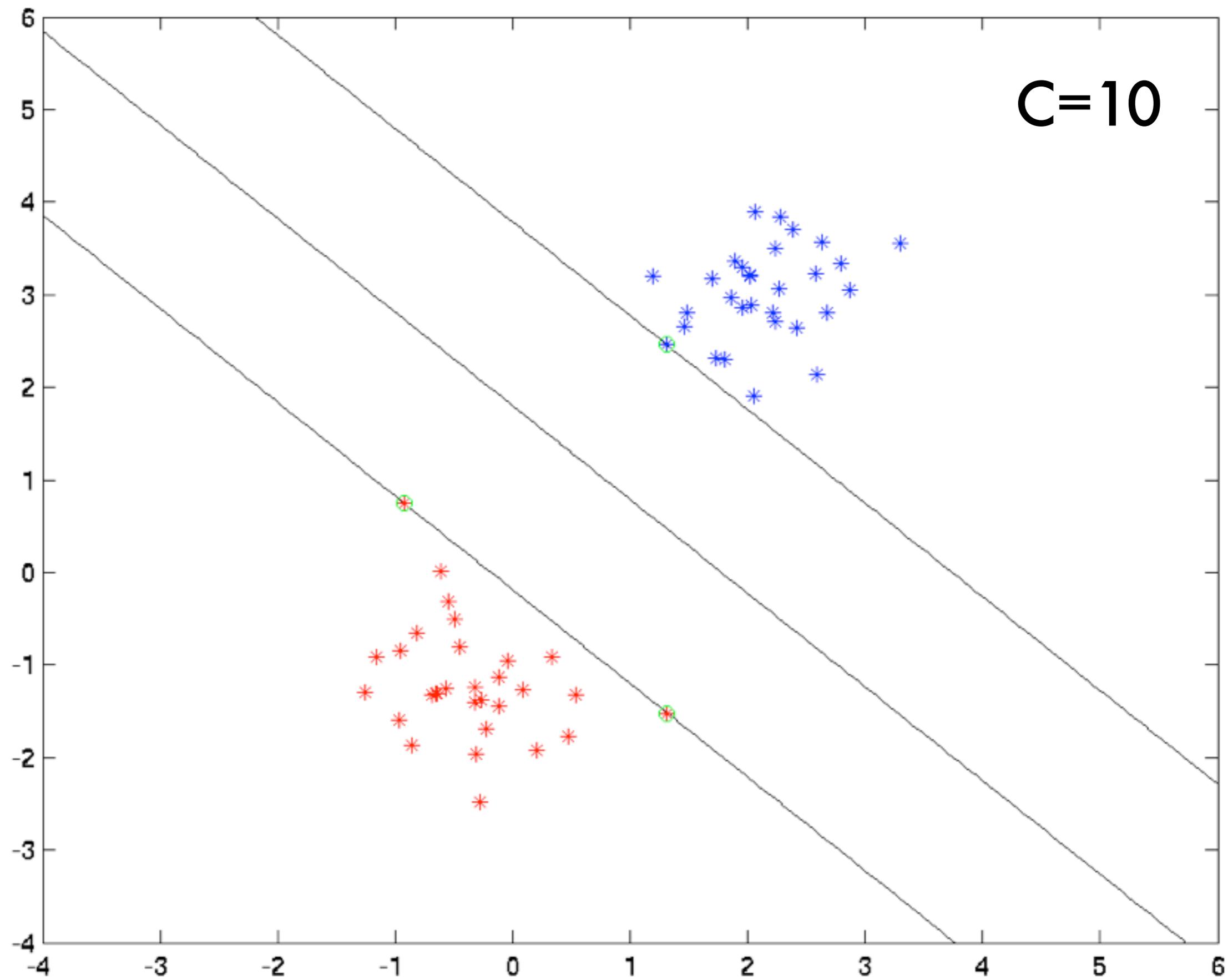
C=2



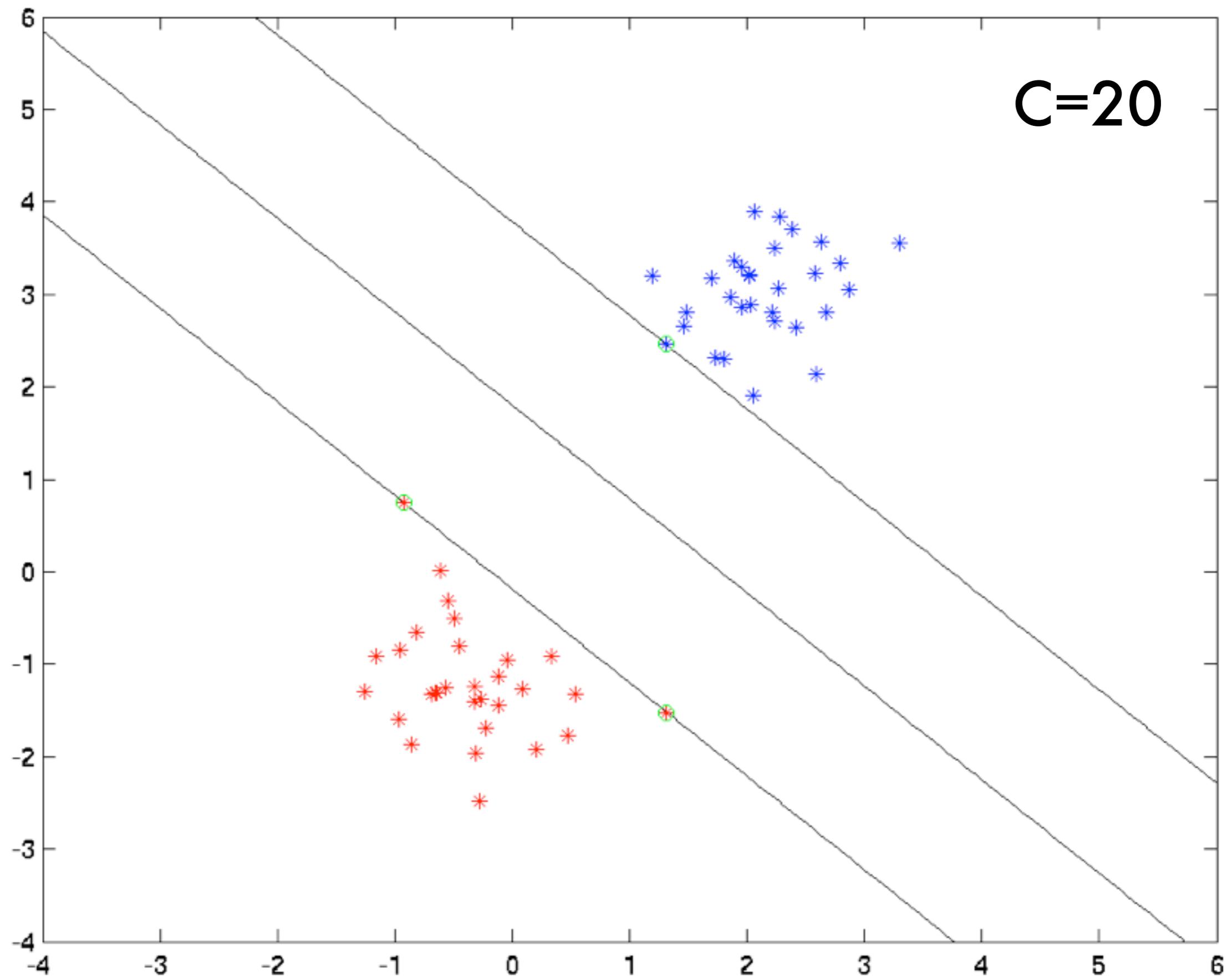
C=5



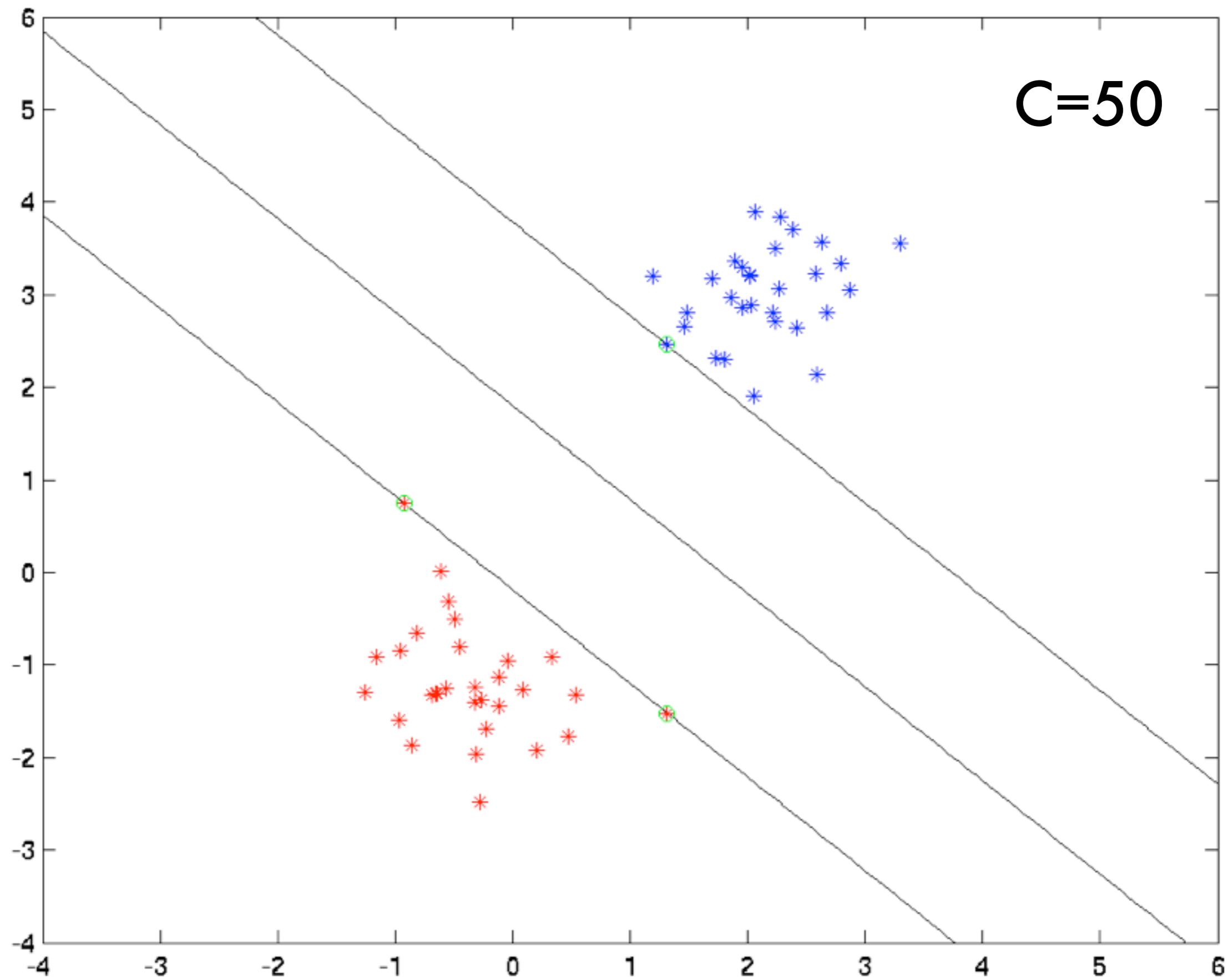
C=10



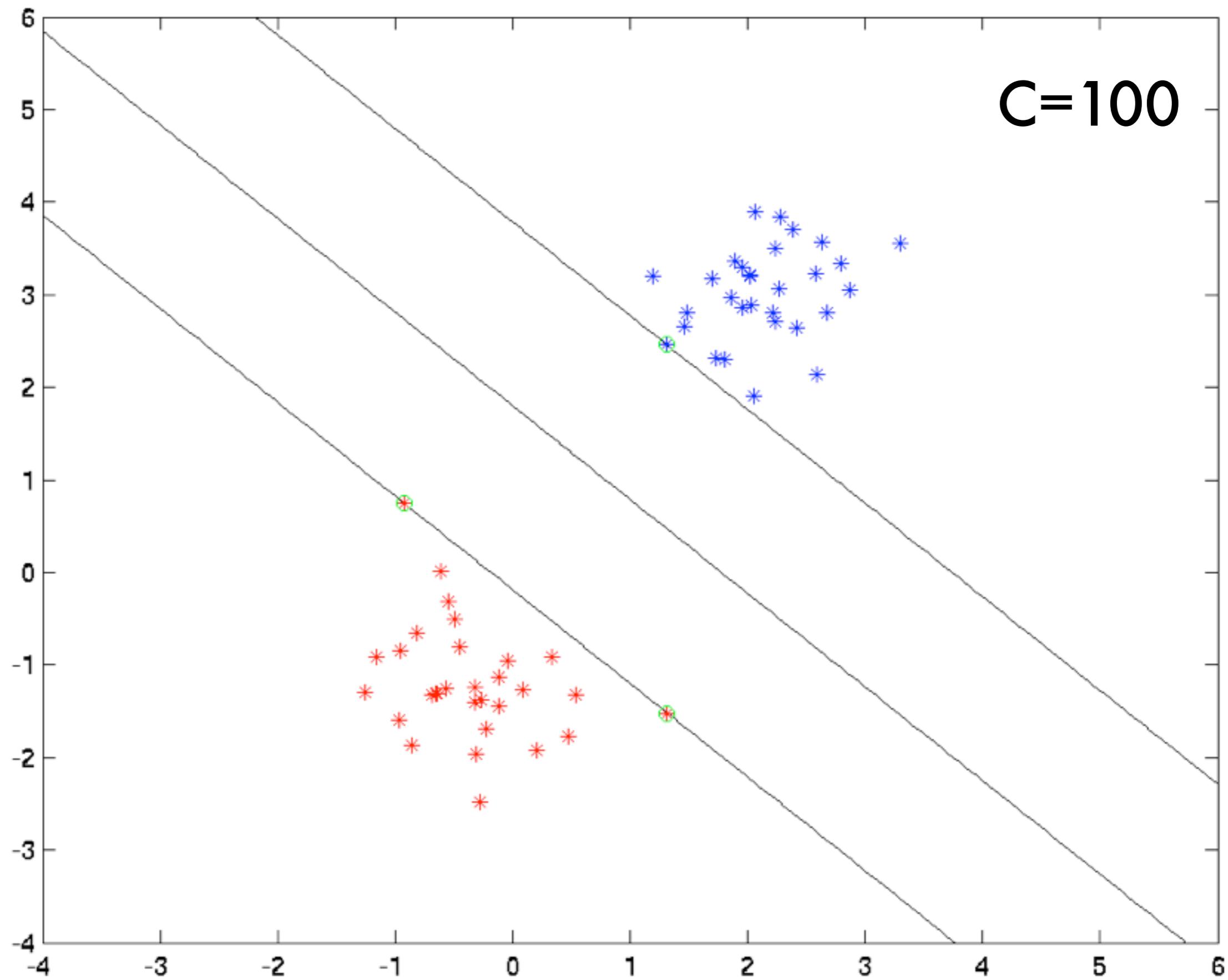
C=20



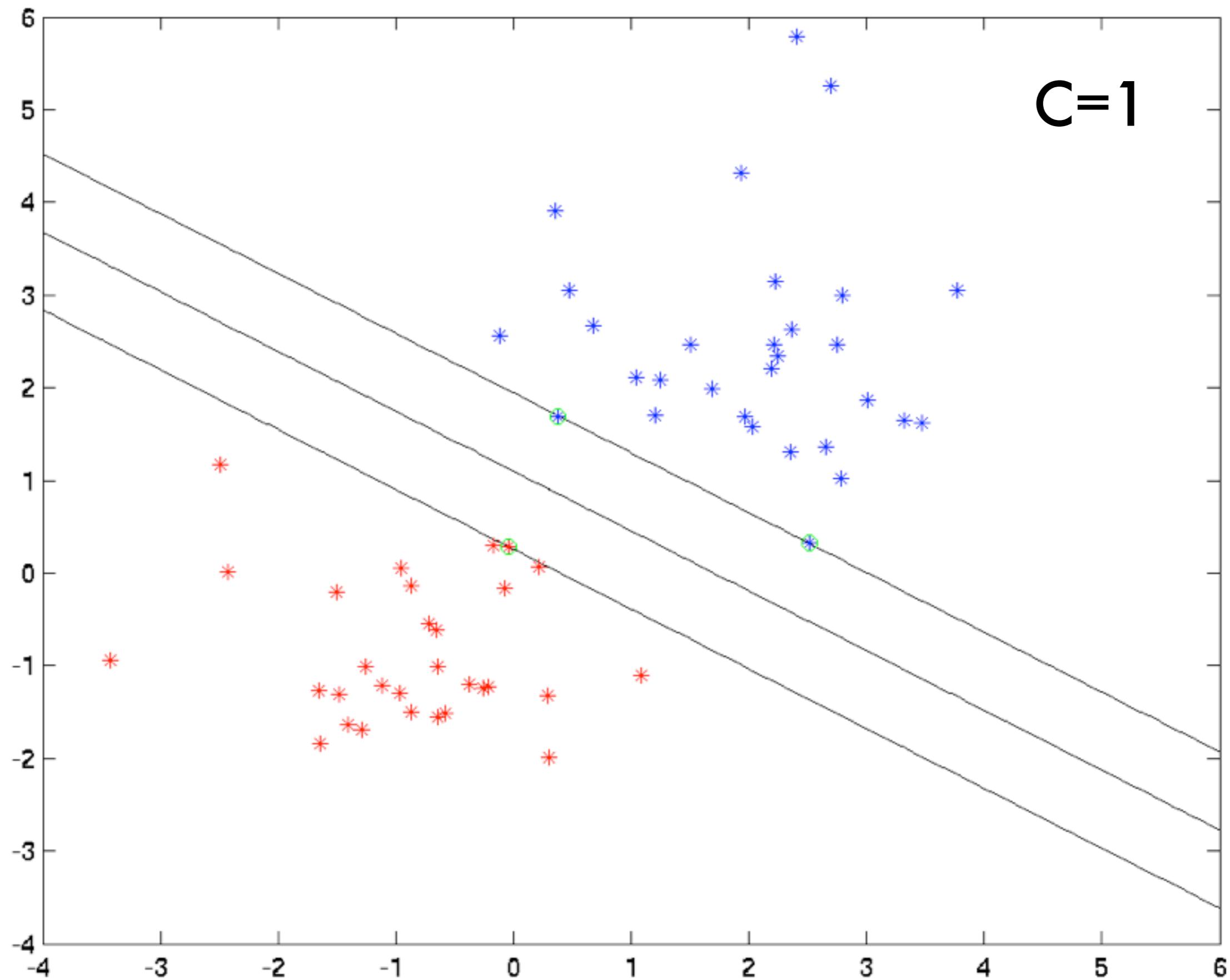
C=50

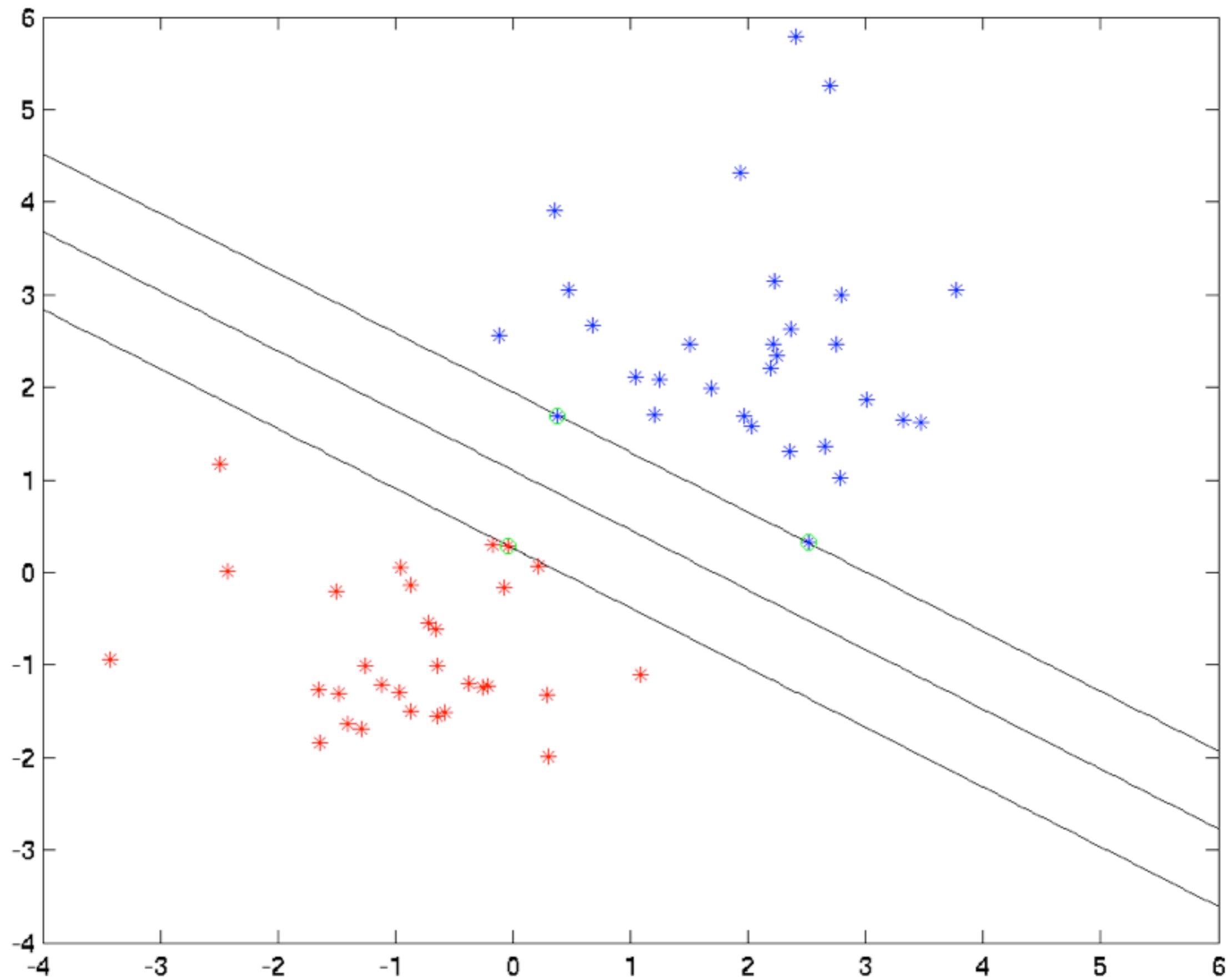


C=100

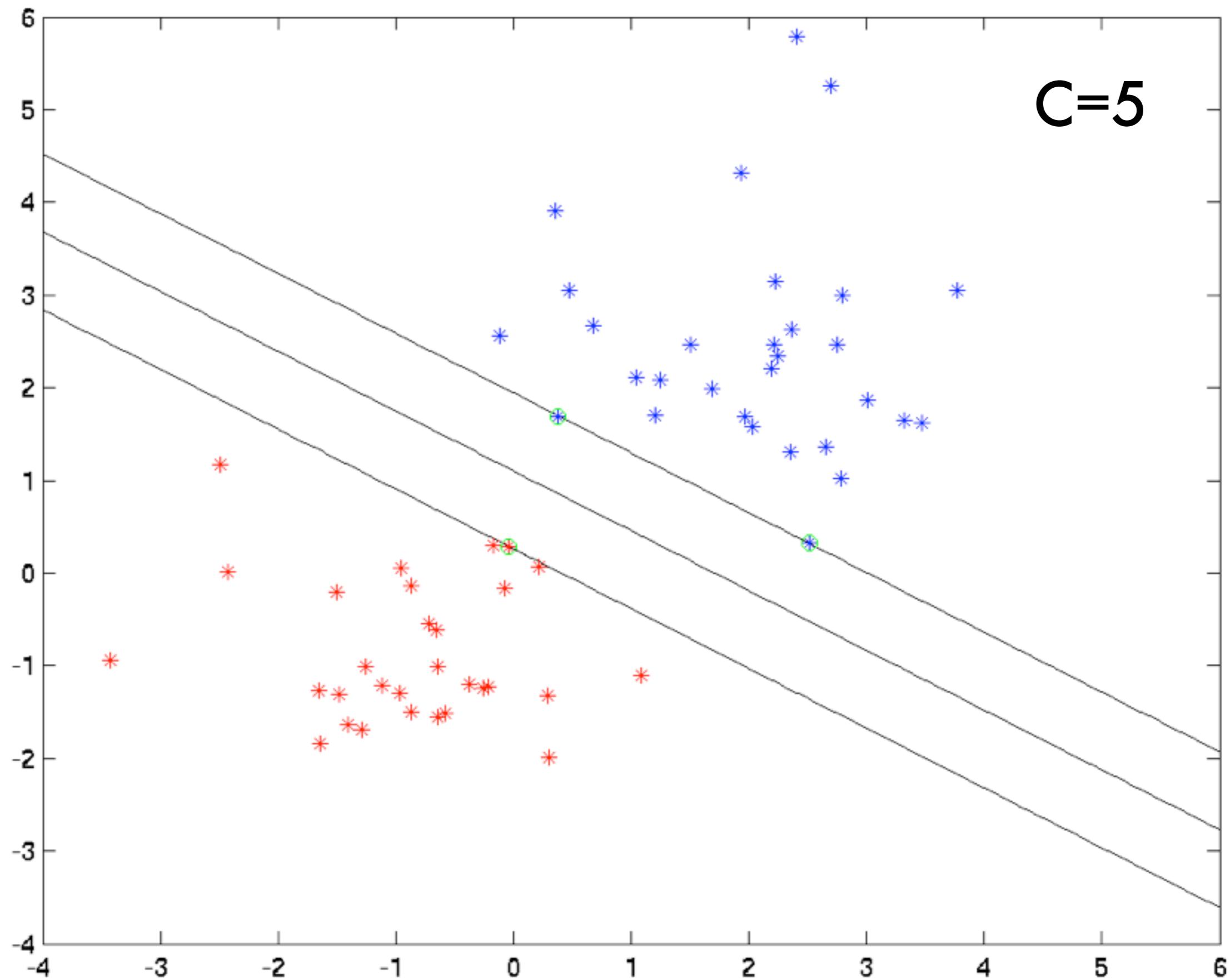


C=1

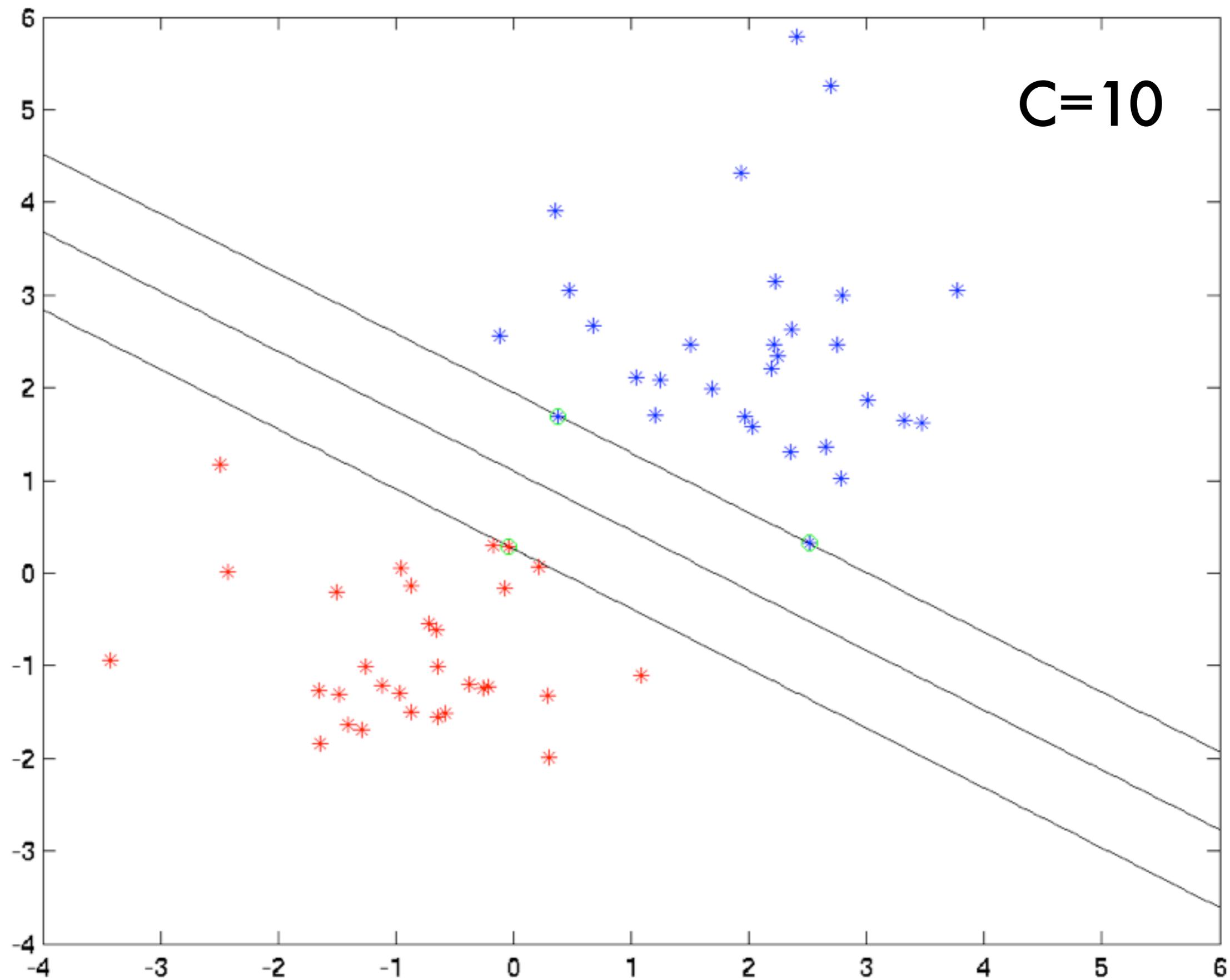




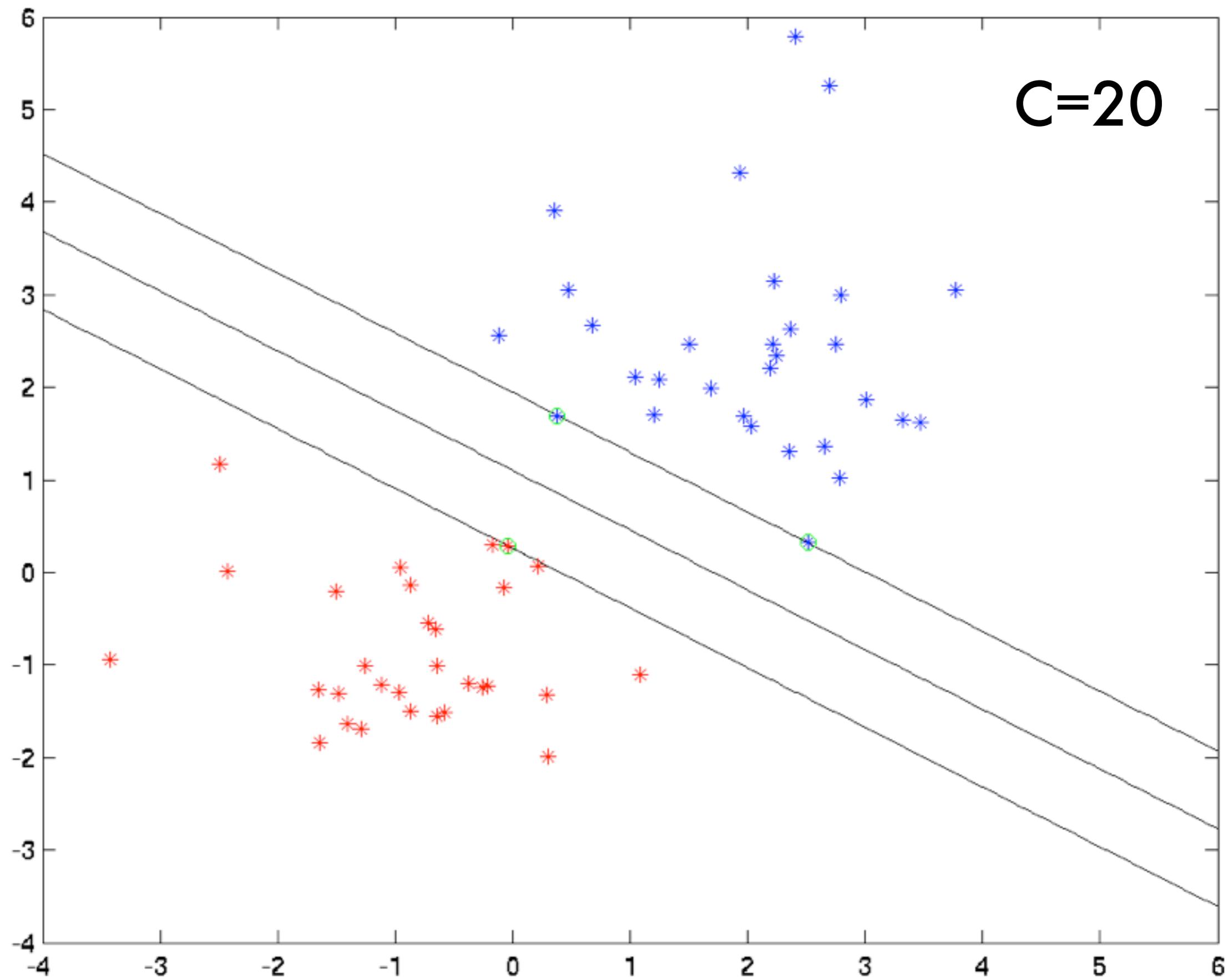
C=5



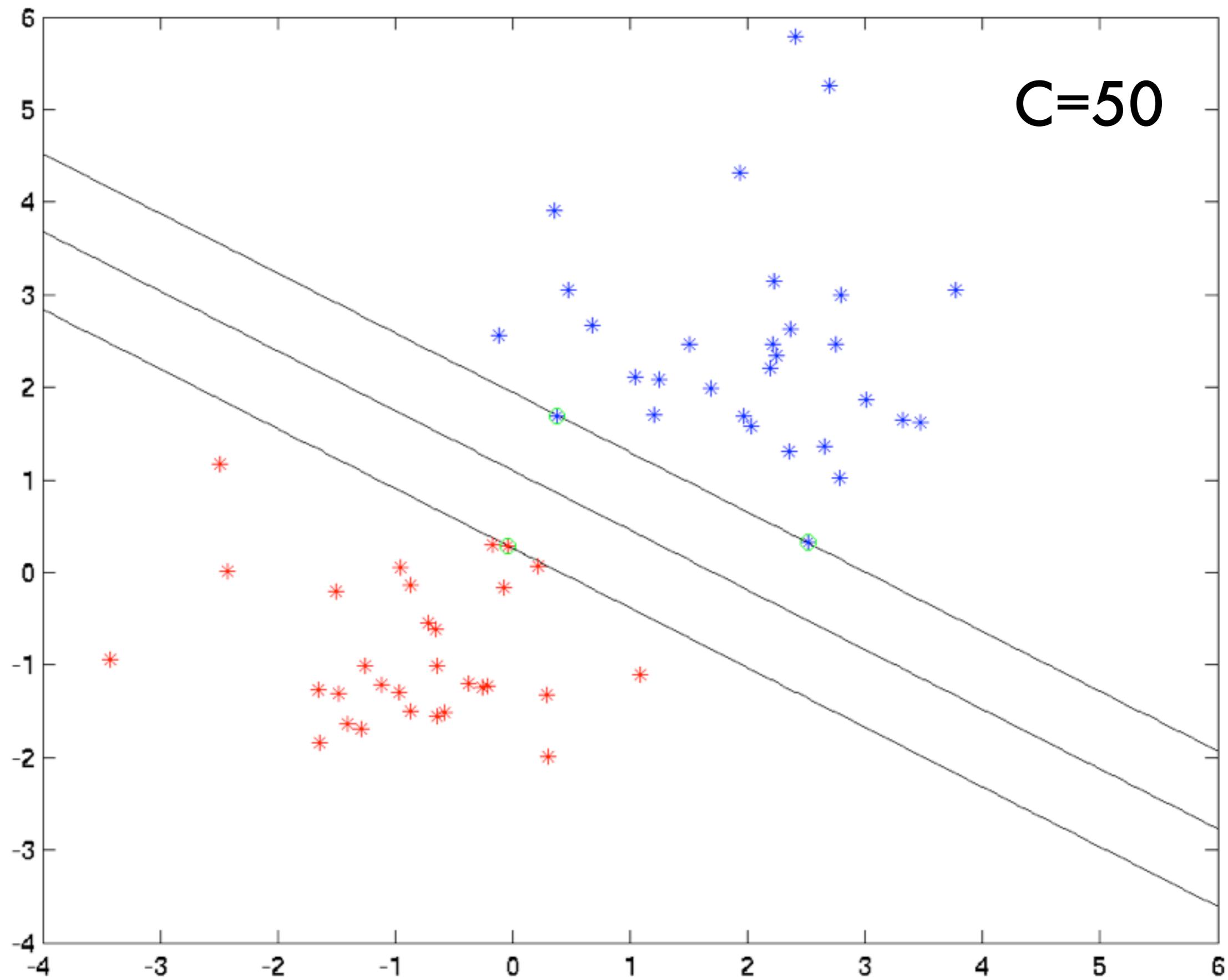
C=10



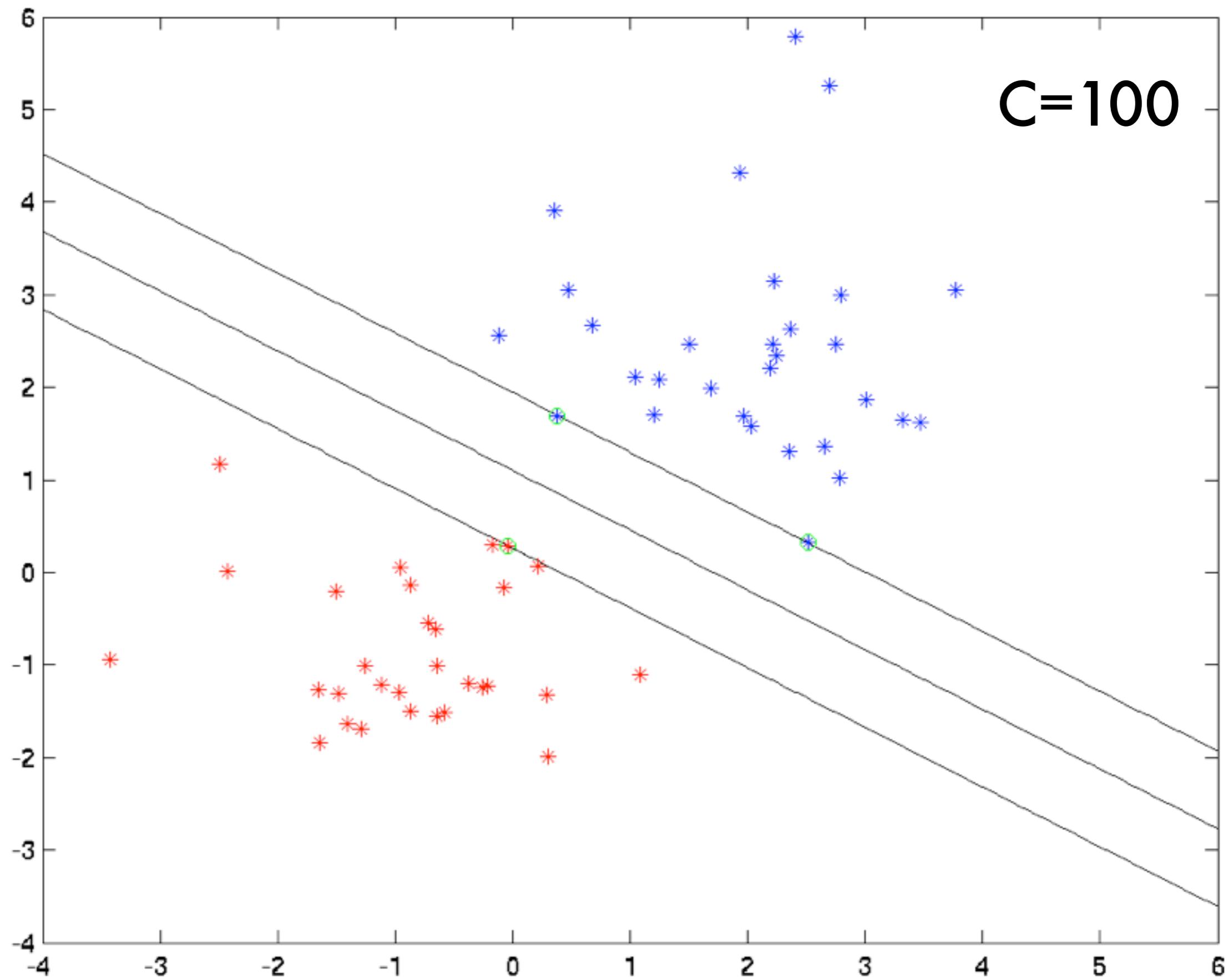
C=20

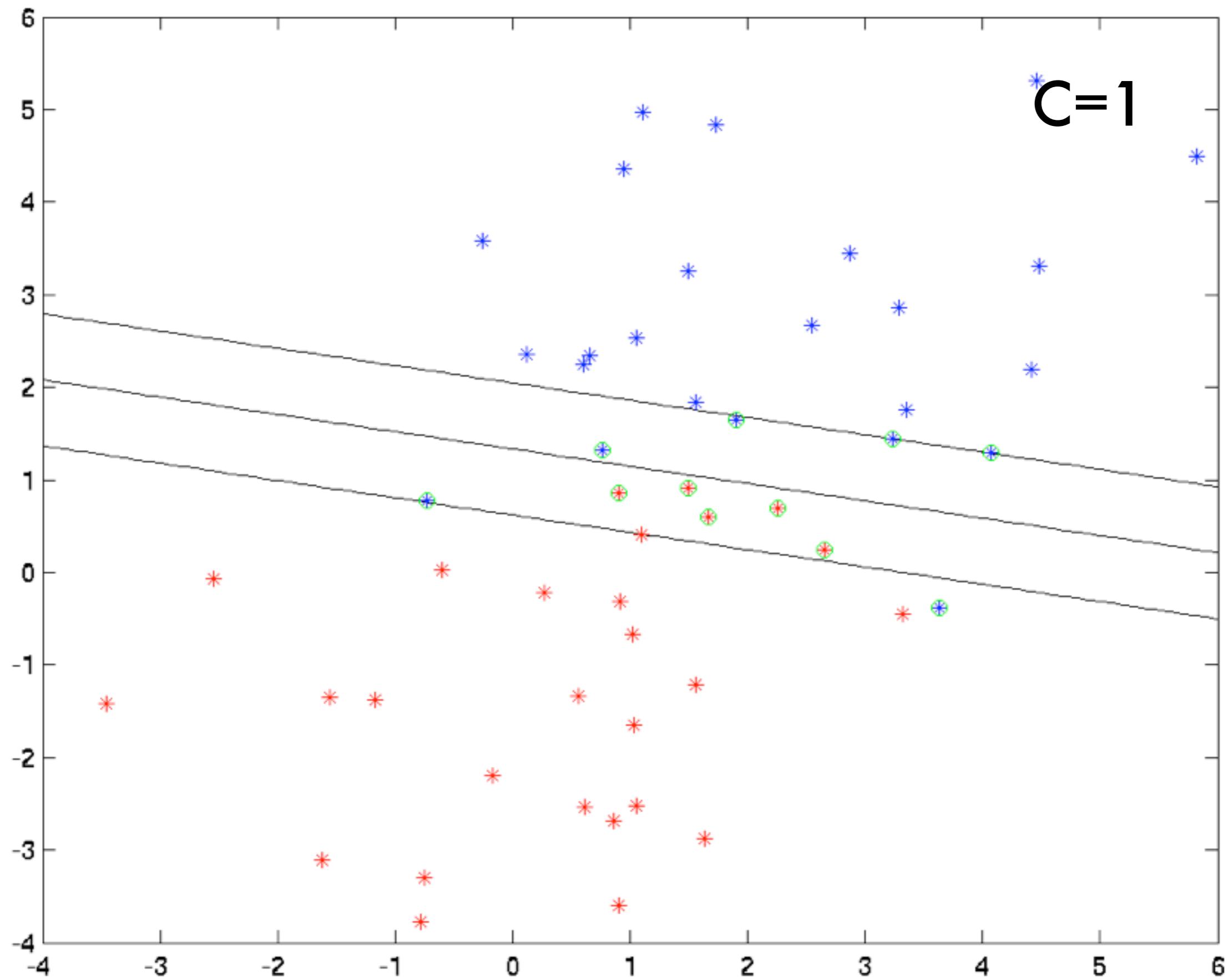


C=50

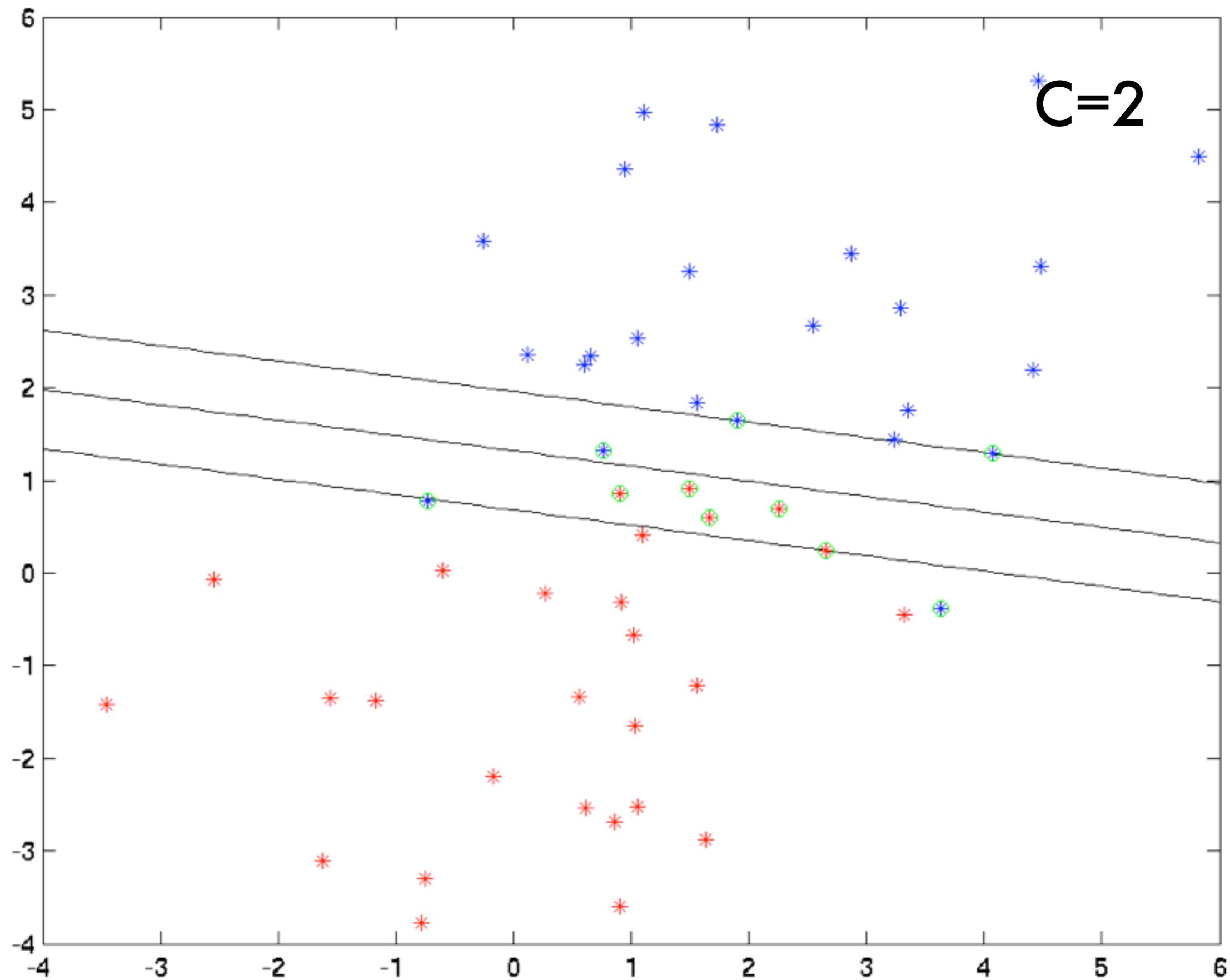


C=100

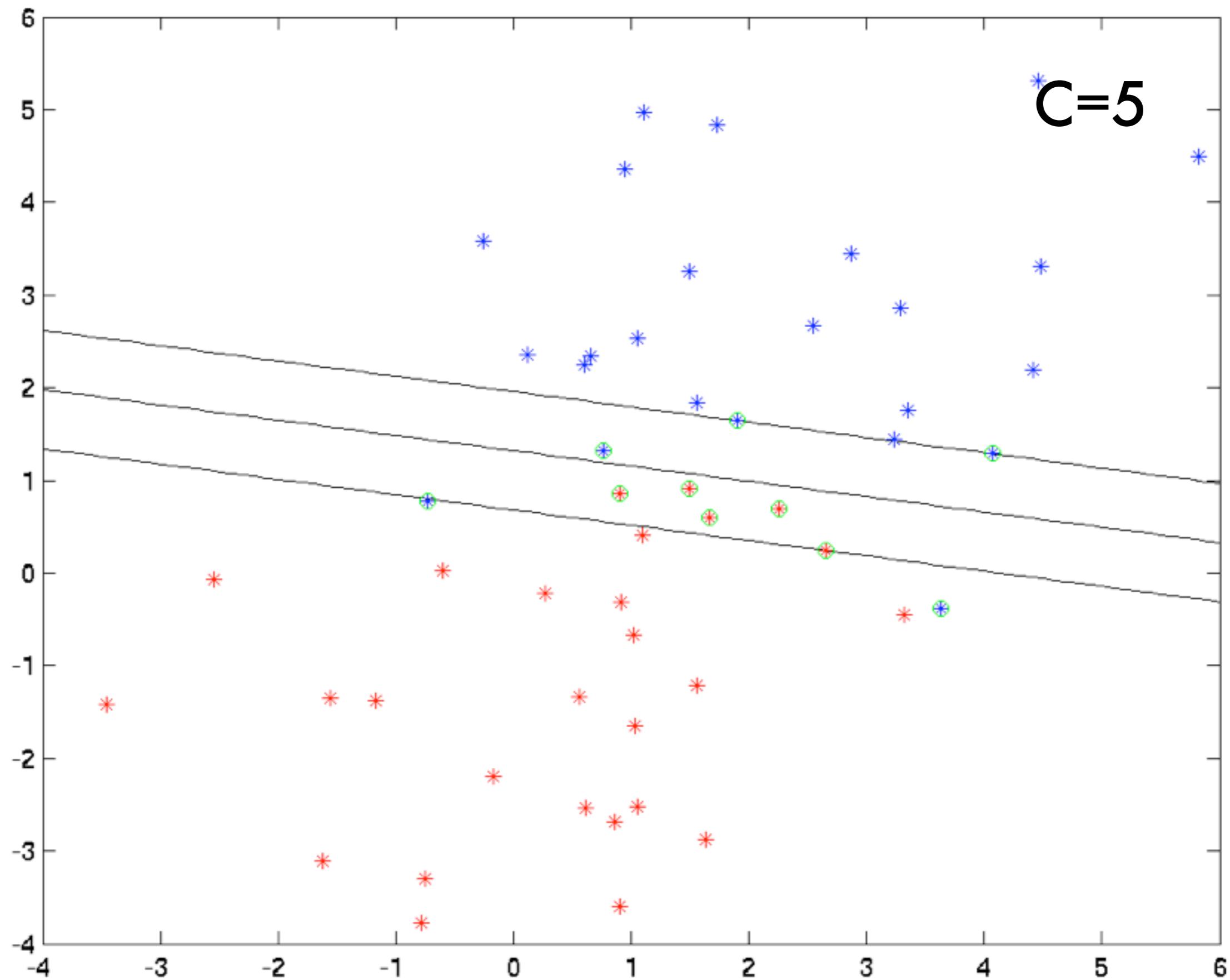




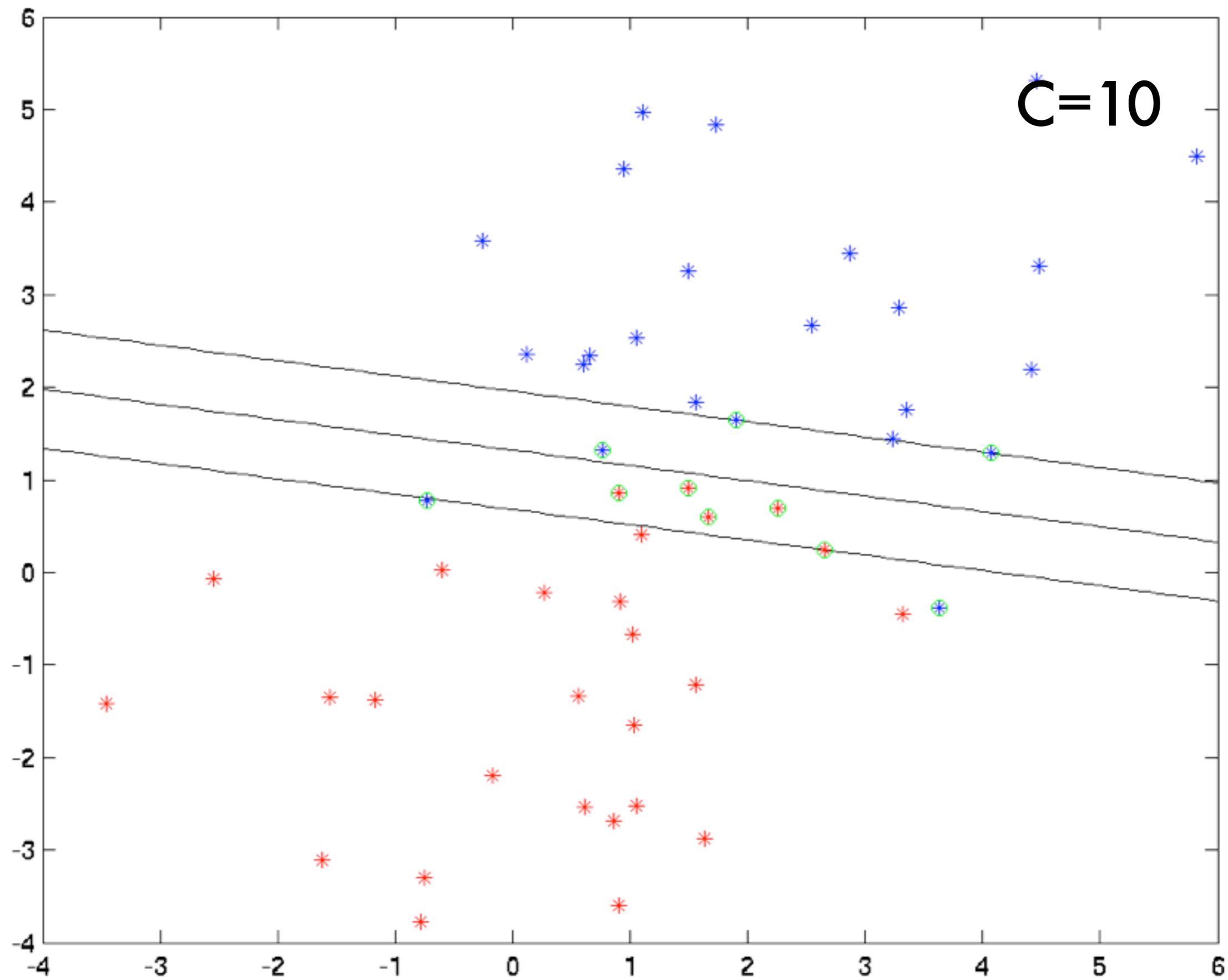
C=2



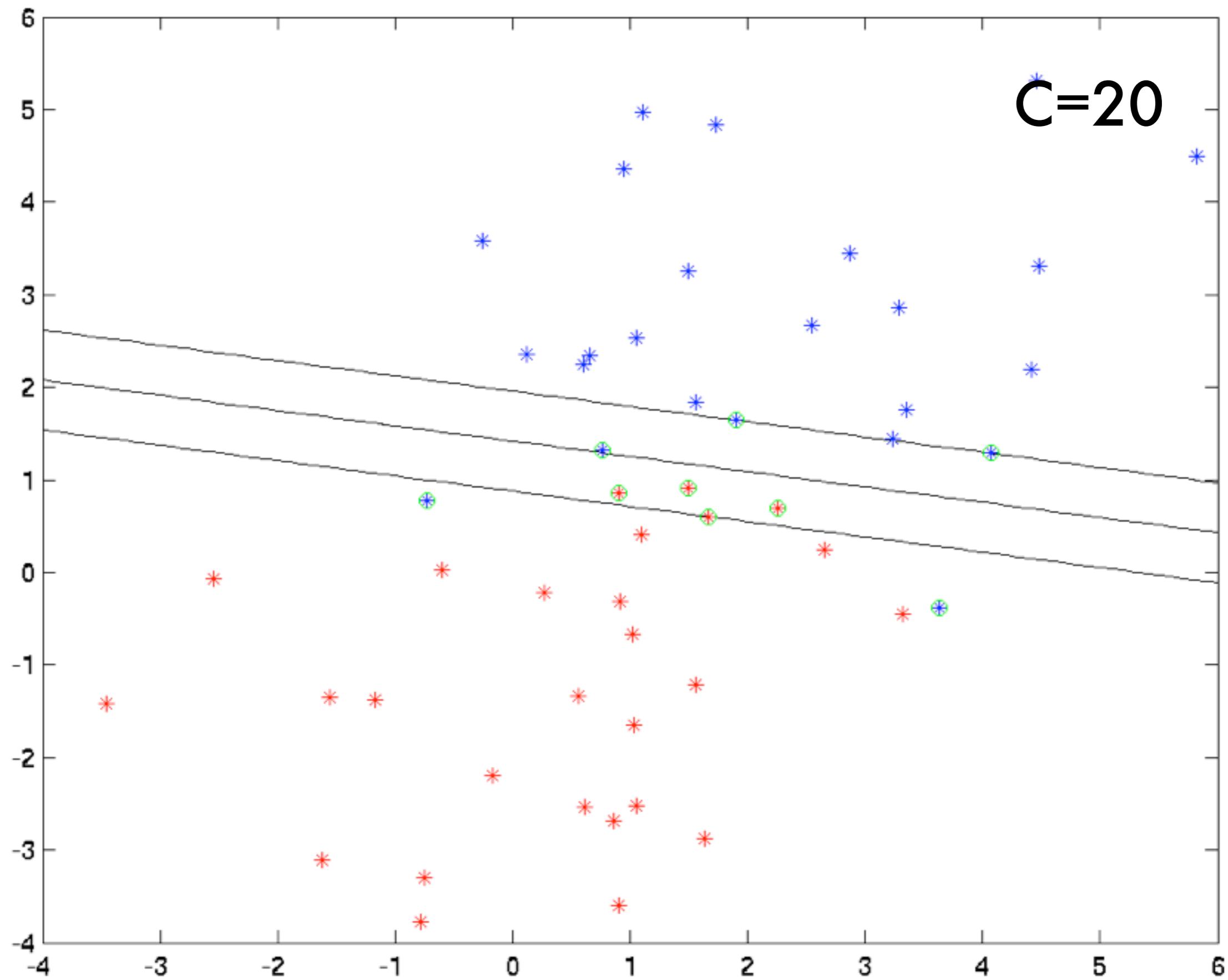
C=5



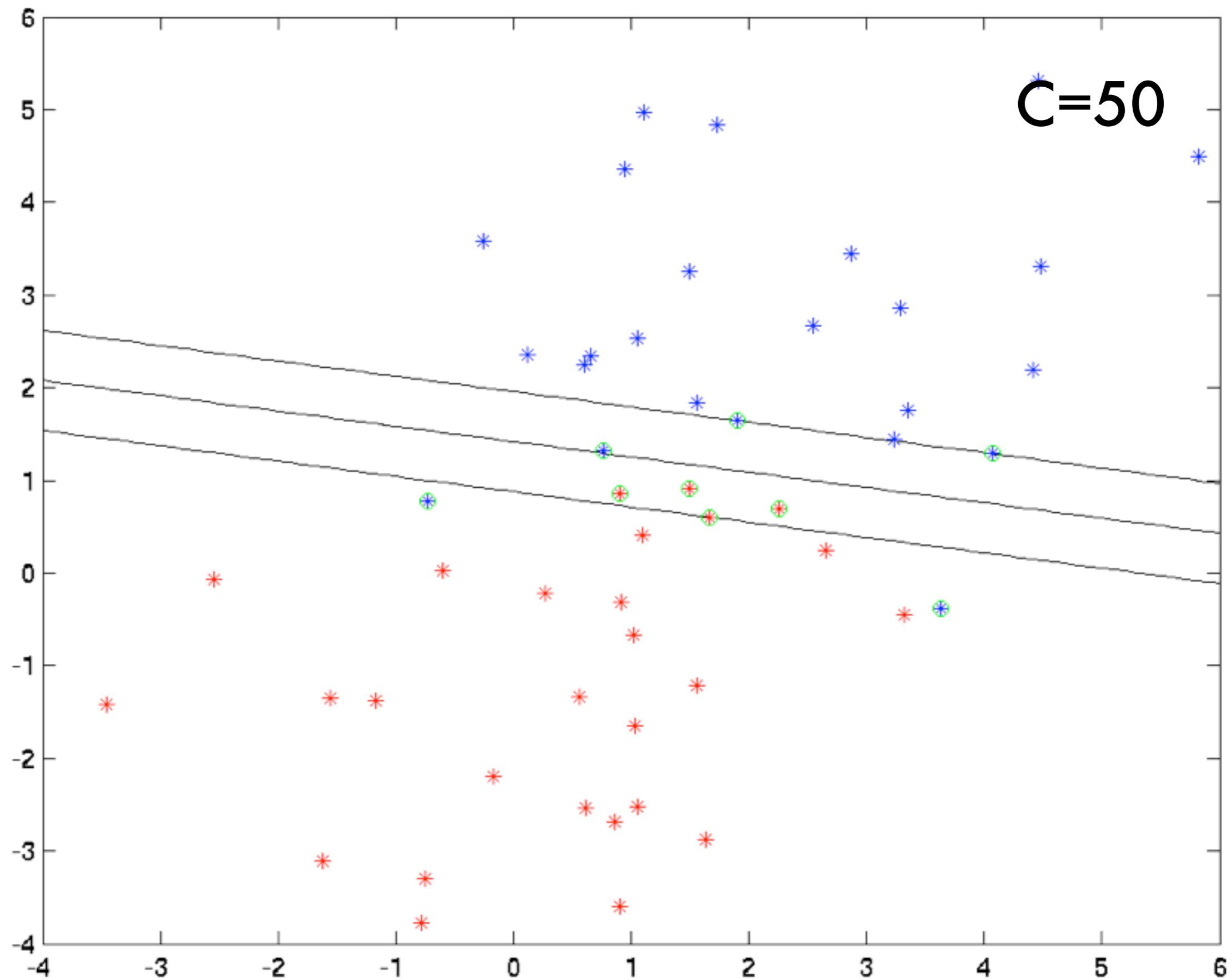
C=10



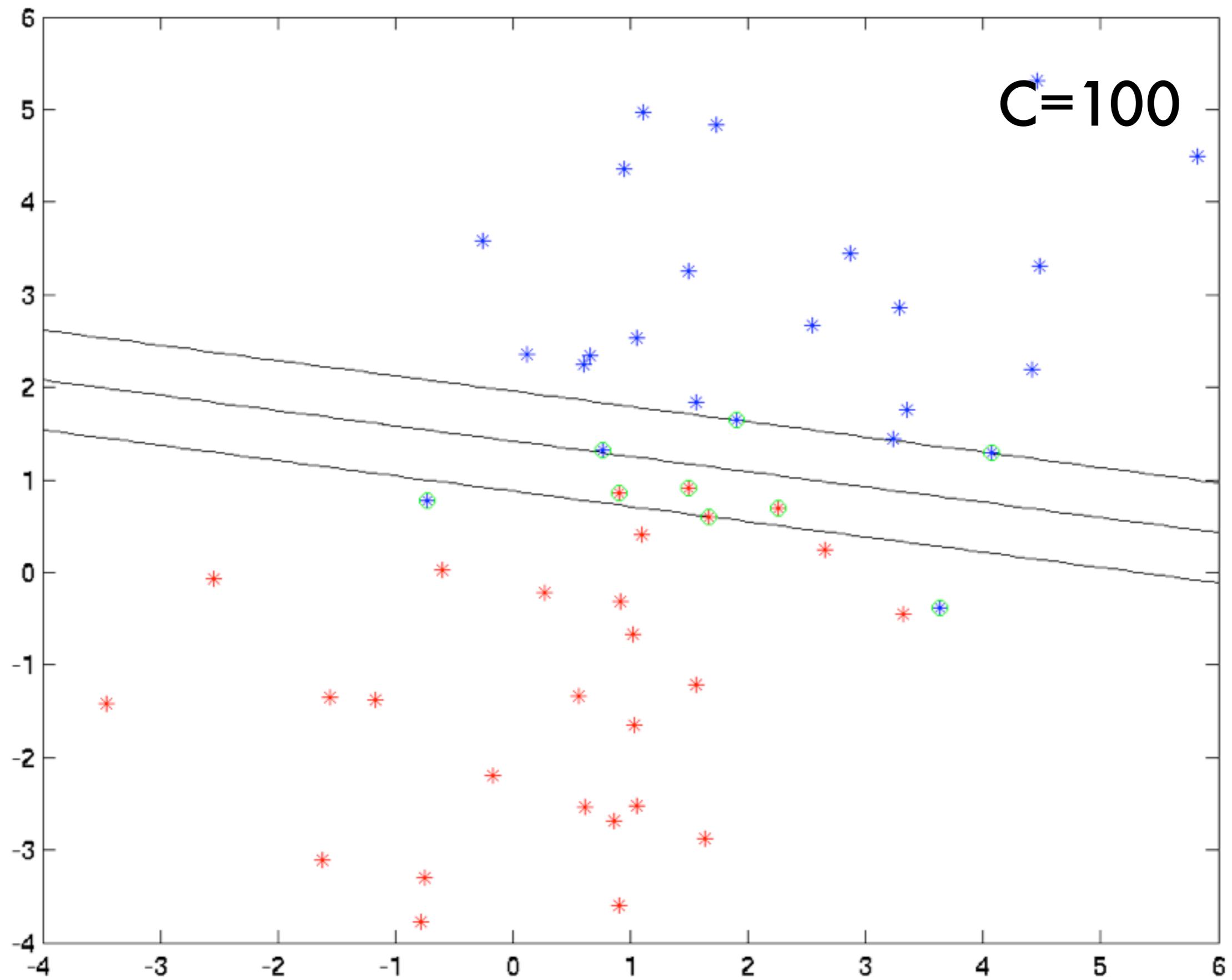
C=20



C=50



$C^* = 100$



Solving the optimization problem

- Dual problem

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

subject to $\sum_i \alpha_i y_i = 0$ 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).

Support Vector ~~Past Life~~ Regression Therapy

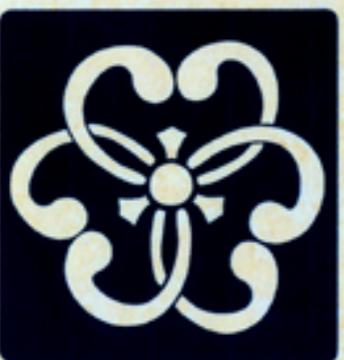
Certificate of Completion

This is to acknowledge that

CMU CLASS 10-715

has successfully completed all course requirements

7th day of October, 2015
This 25th day of March, 2012



Authorized by:

Vickie Penninger

Vickie Penninger, Reiki Master/Teacher

Risk Minimization

- Regression Loss

$$l(y, f(x)) \text{ hence } R_{\text{emp}}[f] = \frac{1}{m} \sum_{i=1}^m l(y_i, f(x_i))$$

- Penalty

$$\Omega[f] \text{ e.g. } \Omega[f] = \frac{1}{2} \|w\|^2$$

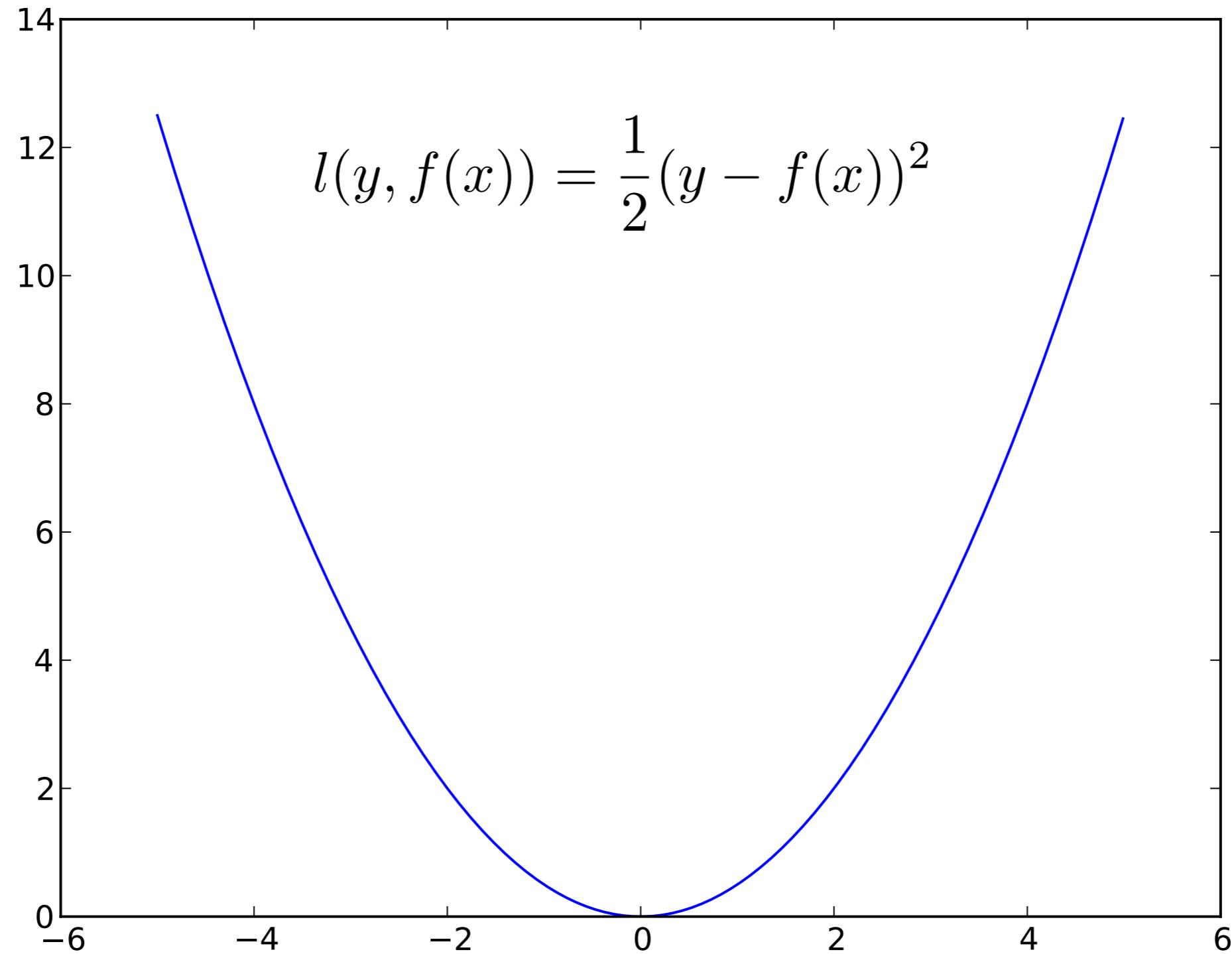
- Conversion to constrained optimization problem

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

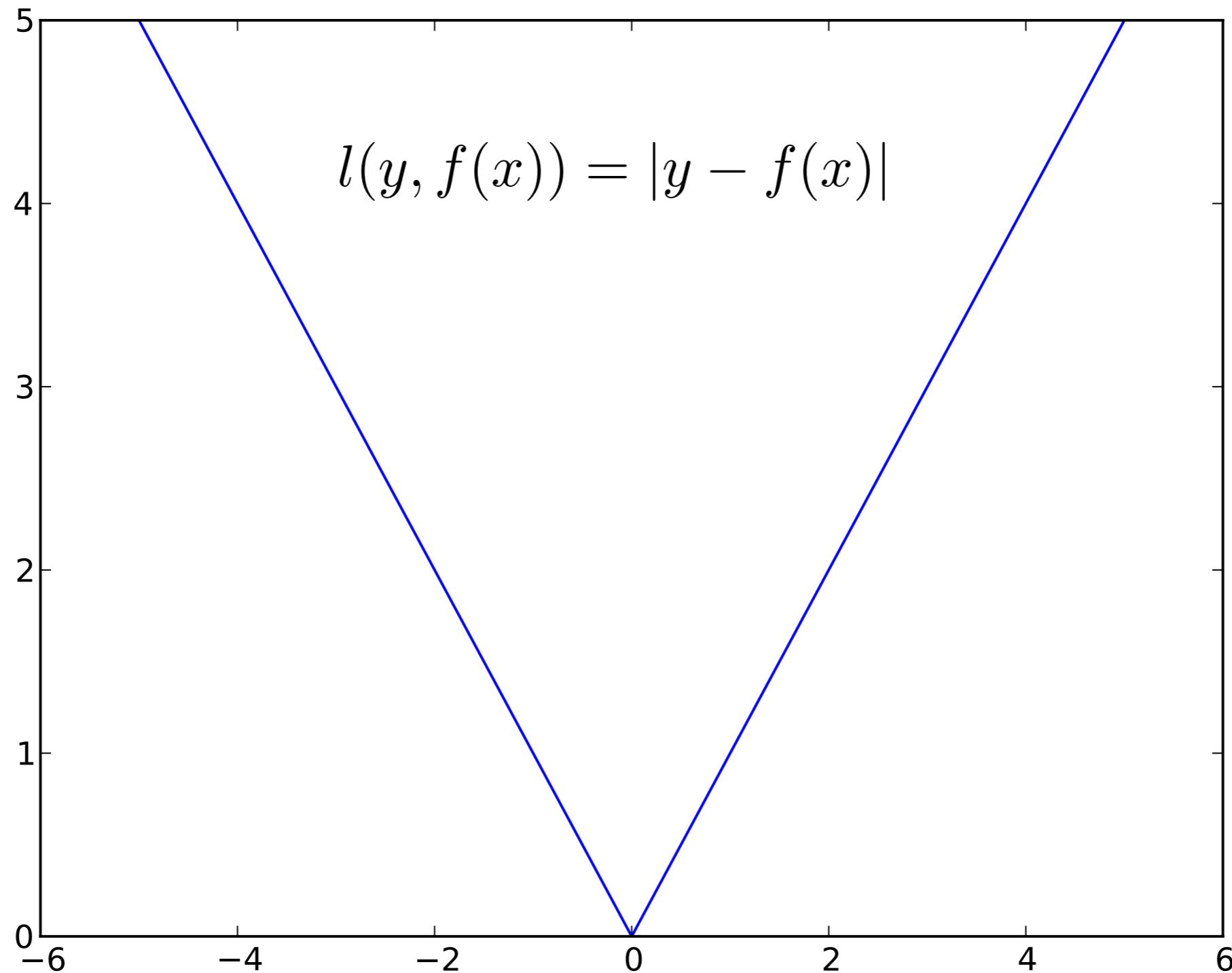
$$\text{subject to } \xi_i = y_i - [\langle w, x_i \rangle + b]$$

slack variables

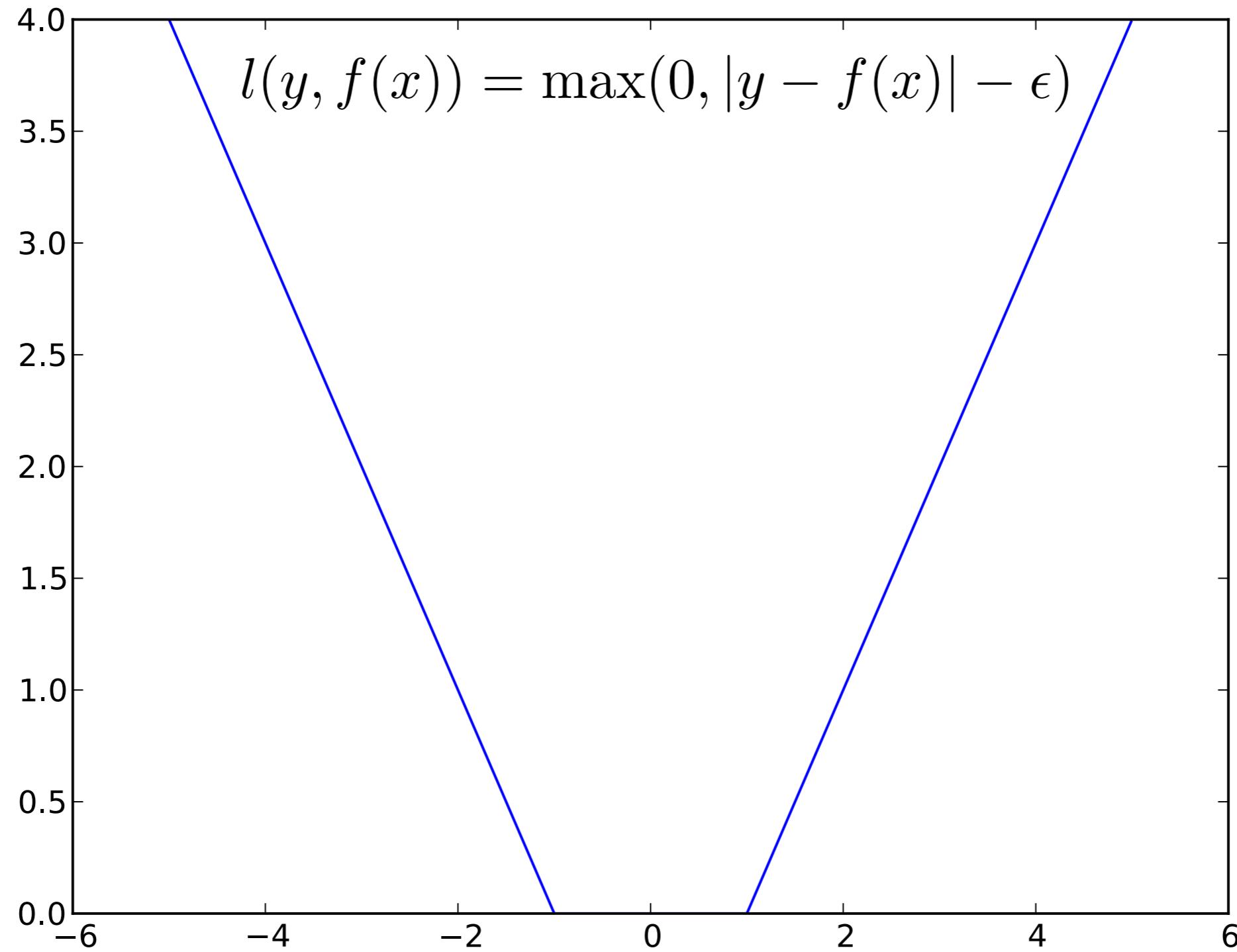
Squared loss



L1 loss



ϵ -insensitive Loss



Penalized least mean squares

- Optimization problem

$$\underset{w}{\text{minimize}} \frac{1}{2m} \sum_{i=1}^m (y_i - \langle x_i, w \rangle)^2 + \frac{\lambda}{2} \|w\|^2$$

- Solution

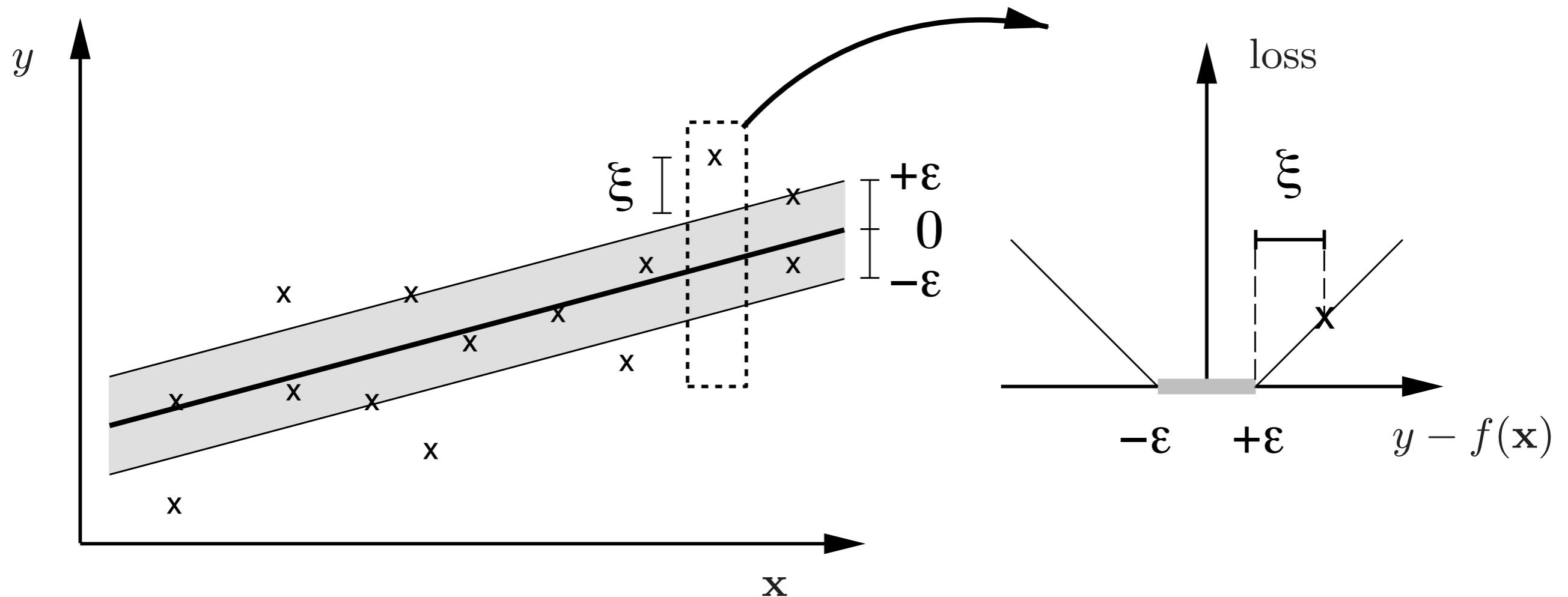
$$\begin{aligned}\partial_w [\dots] &= \frac{1}{m} \sum_{i=1}^m [x_i x_i^\top w - x_i y_i] + \lambda w \\ &= \left[\frac{1}{m} X X^\top + \lambda \mathbf{1} \right] w - \frac{1}{m} X y = 0\end{aligned}$$

$$\text{hence } w = [X X^\top + \lambda m \mathbf{1}]^{-1} X y$$

Outer product
matrix in X

Conjugate Gradient
Sherman Morrison Woodbury

SVM Regression



don't care about deviations within the tube

SVM Regression (ϵ -insensitive loss)

- Optimization Problem (as constrained QP)

$$\underset{w,b}{\text{minimize}} \quad \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m [\xi_i + \xi_i^*]$$

subject to $\langle w, x_i \rangle + b \leq y_i + \epsilon + \xi_i$ and $\xi_i \geq 0$
 $\langle w, x_i \rangle + b \geq y_i - \epsilon - \xi_i^*$ and $\xi_i^* \geq 0$

- Lagrange Function

$$L = \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m [\xi_i + \xi_i^*] - \sum_{i=1}^m [\eta_i \xi_i + \eta_i^* \xi_i^*] + \\ \sum_{i=1}^m \alpha_i [\langle w, x_i \rangle + b - y_i - \epsilon - \xi_i] + \sum_{i=1}^m \alpha_i^* [y_i - \epsilon - \xi_i^* - \langle w, x_i \rangle - b]$$

SVM Regression (ϵ -insensitive loss)

- First order conditions

$$\partial_w L = 0 = w + \sum_i [\alpha_i - \alpha_i^*] x_i$$

$$\partial_b L = 0 = \sum_i [\alpha_i - \alpha_i^*]$$

$$\partial_{\xi_i} L = 0 = C - \eta_i - \alpha_i$$

$$\partial_{\xi_i^*} L = 0 = C - \eta_i^* - \alpha_i^*$$

- Dual problem

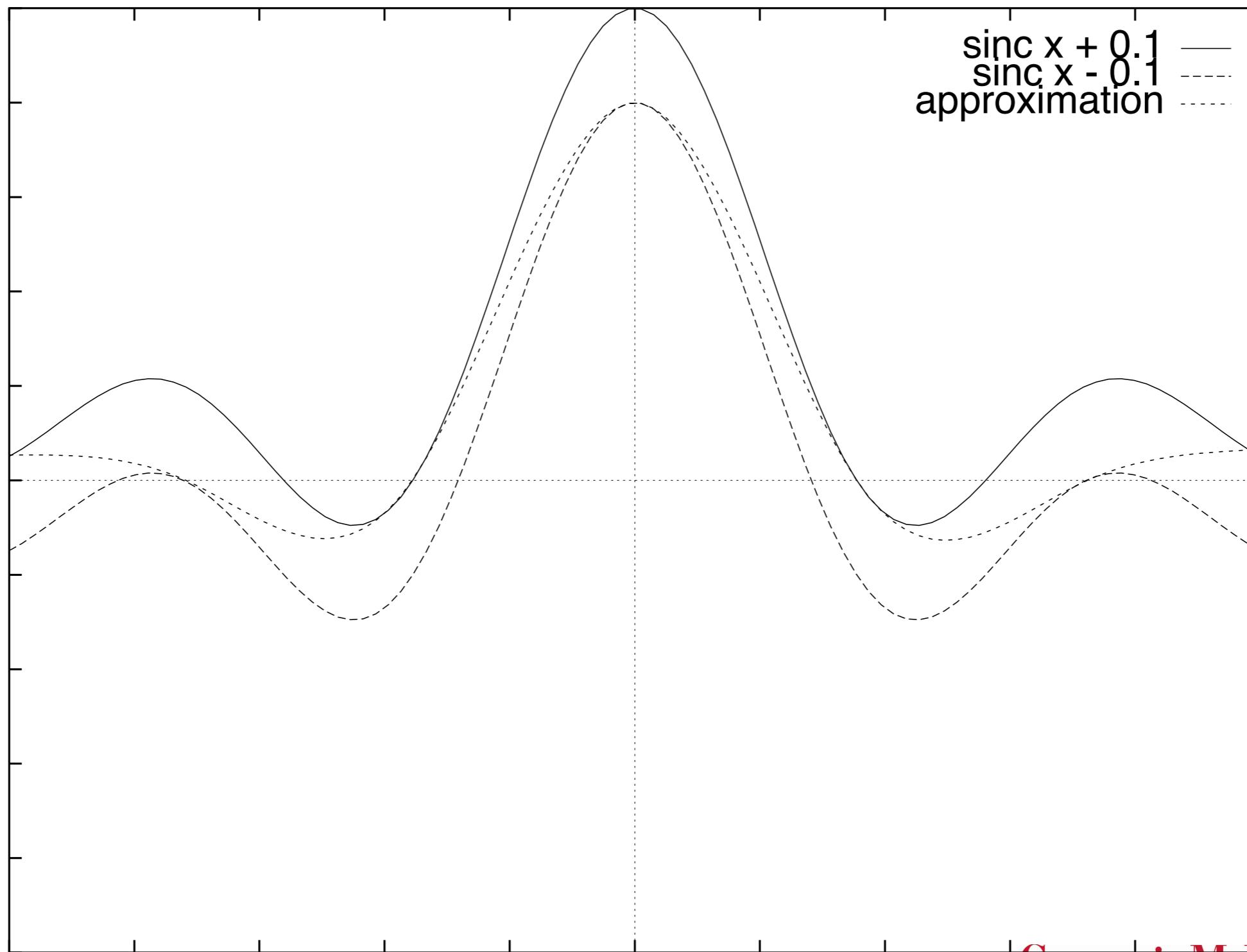
$$\underset{\alpha, \alpha^*}{\text{minimize}} \quad \frac{1}{2} (\alpha - \alpha^*)^\top K (\alpha - \alpha^*) + \epsilon \mathbf{1}^\top (\alpha + \alpha^*) + y^\top (\alpha - \alpha^*)$$

subject to $\mathbf{1}^\top (\alpha - \alpha^*) = 0$ and $\alpha_i, \alpha_i^* \in [0, C]$

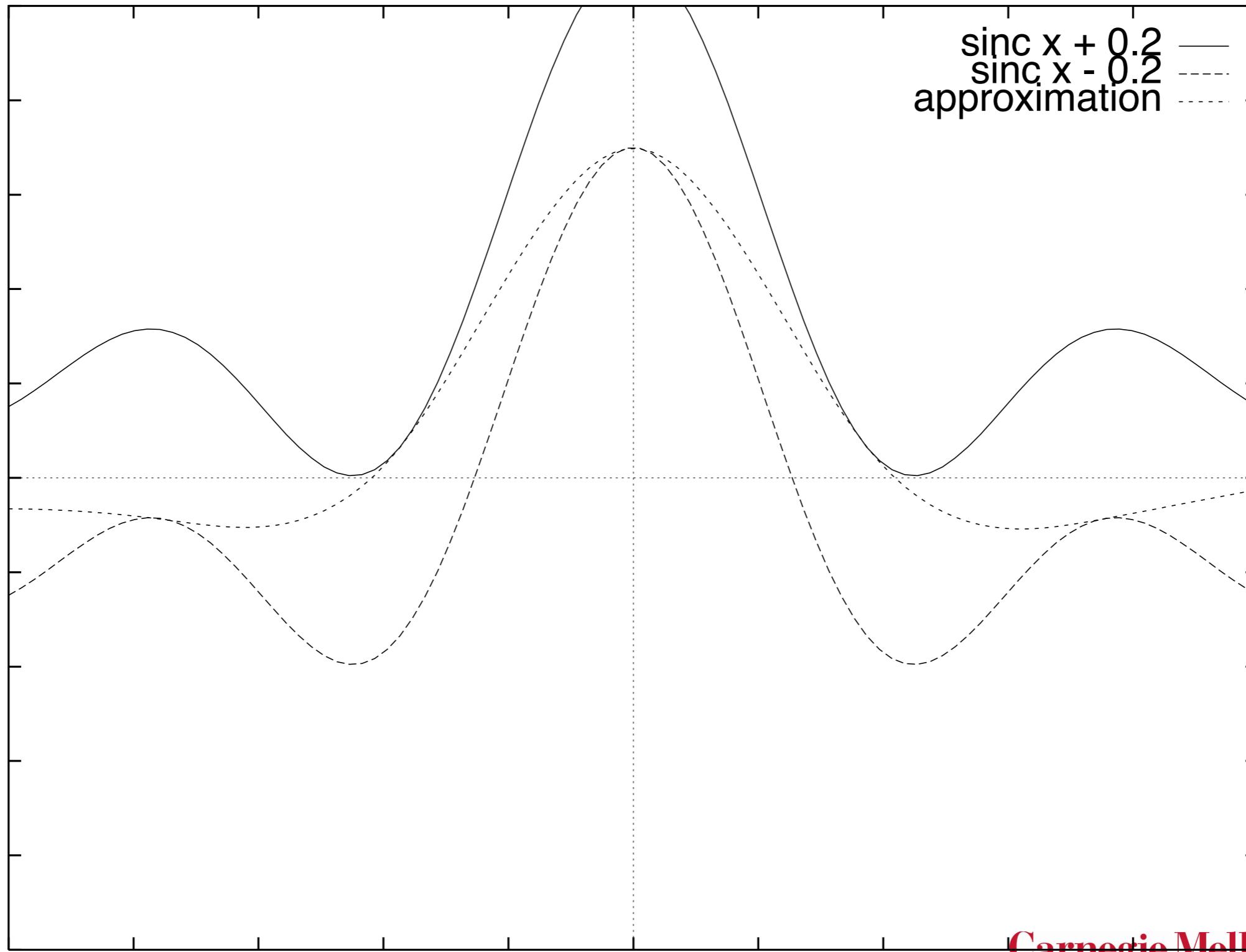
Properties

- Ignores ‘typical’ instances with small error
- Only upper or lower bound active at any time
- QP in $2n$ variables as cheap as SVM problem
- Robustness with respect to outliers
 - ℓ_1 loss yields same problem without epsilon
 - Huber’s robust loss yields similar problem but with added quadratic penalty on coefficients

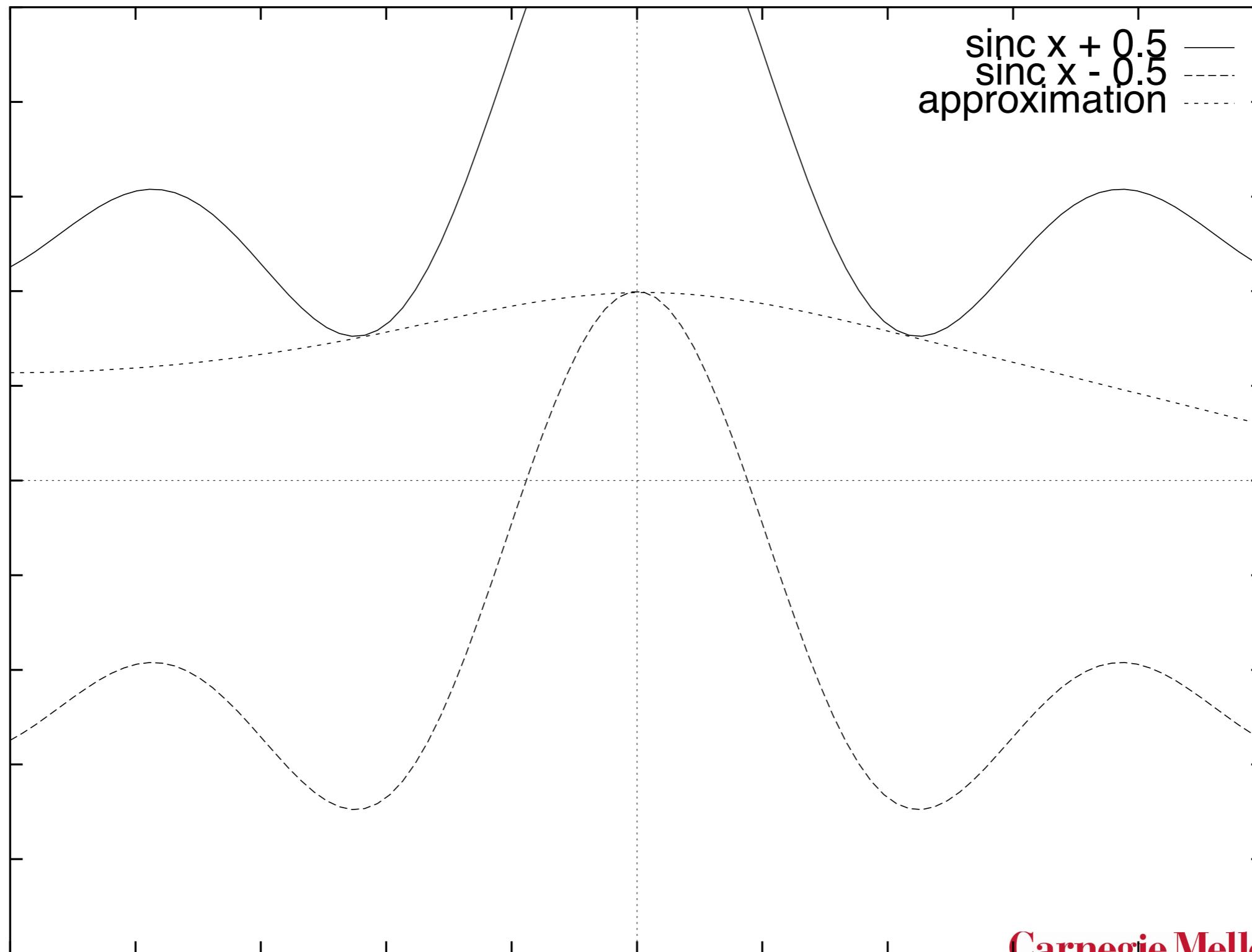
Regression example



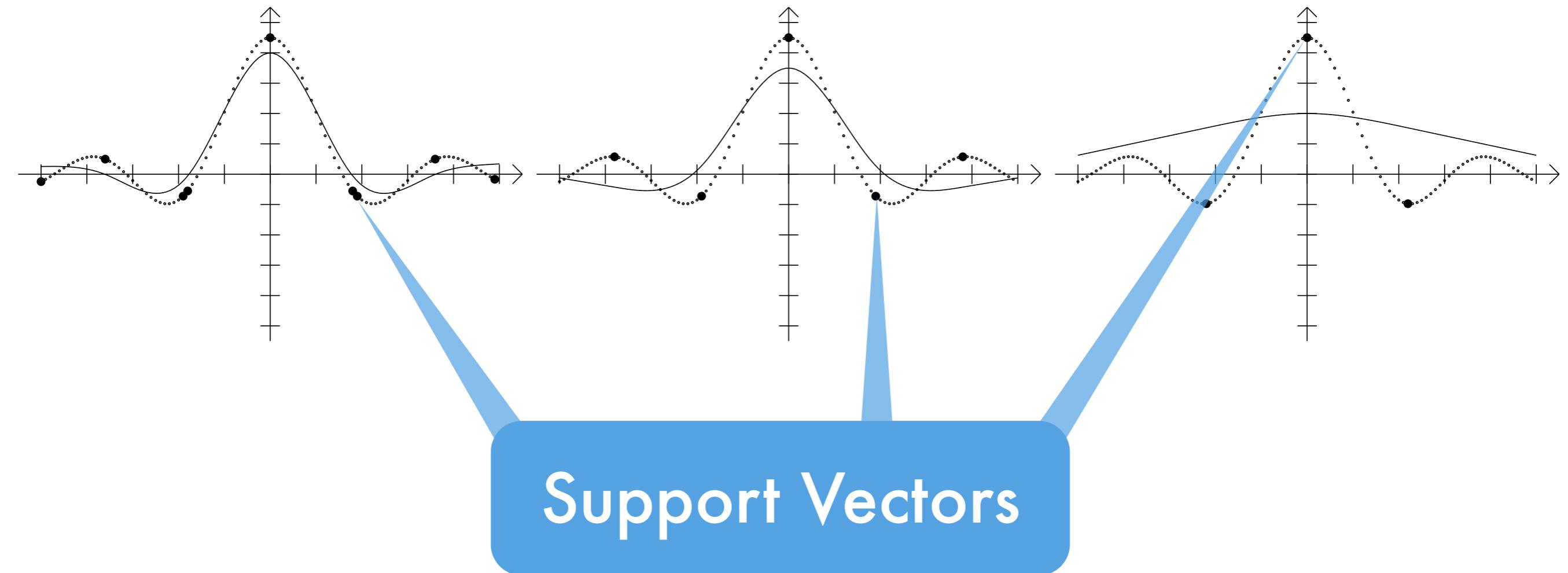
Regression example



Regression example



Regression example



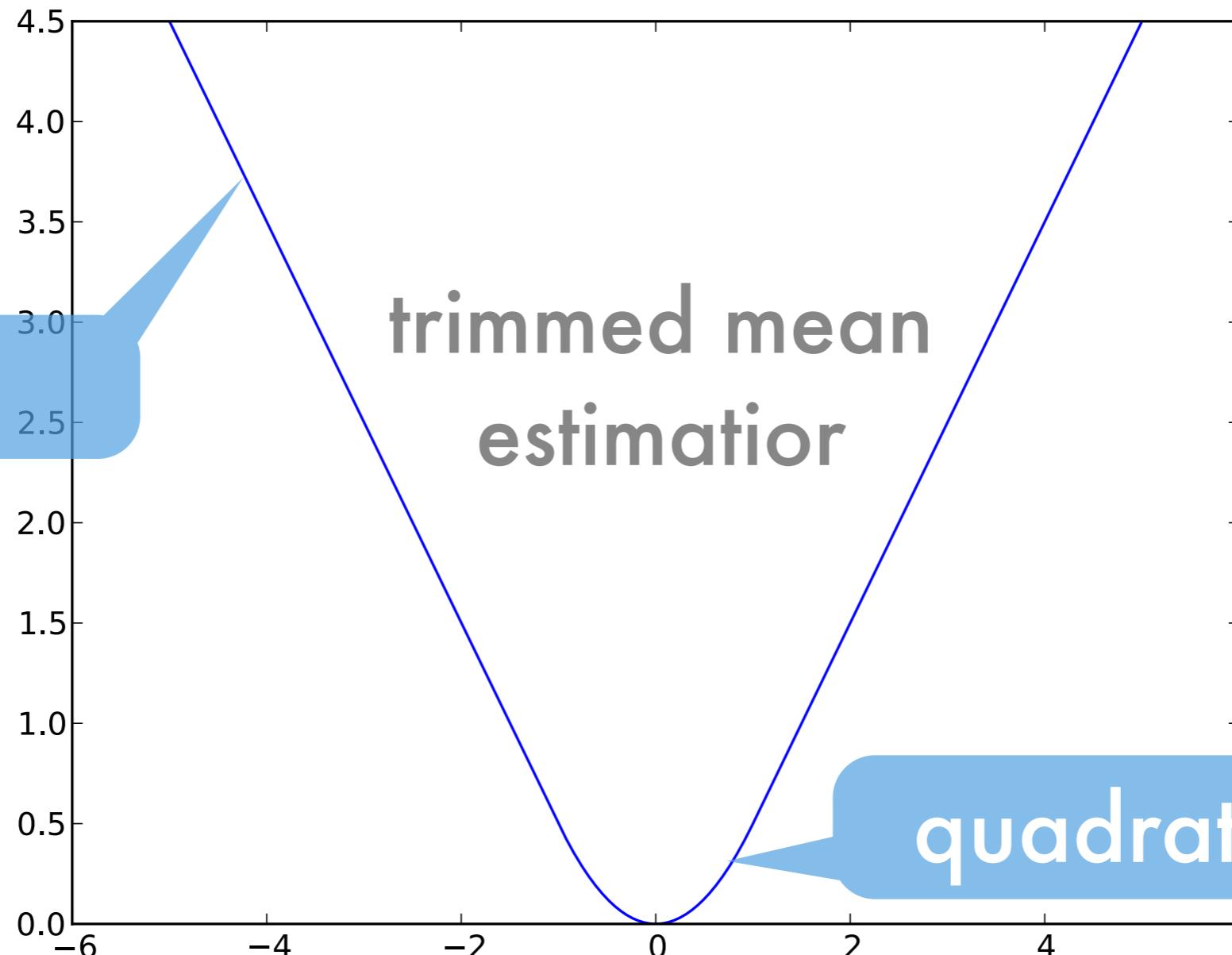
Huber's robust loss

$$l(y, f(x)) = \begin{cases} \frac{1}{2}(y - f(x))^2 & \text{if } |y - f(x)| < 1 \\ |y - f(x)| - \frac{1}{2} & \text{otherwise} \end{cases}$$

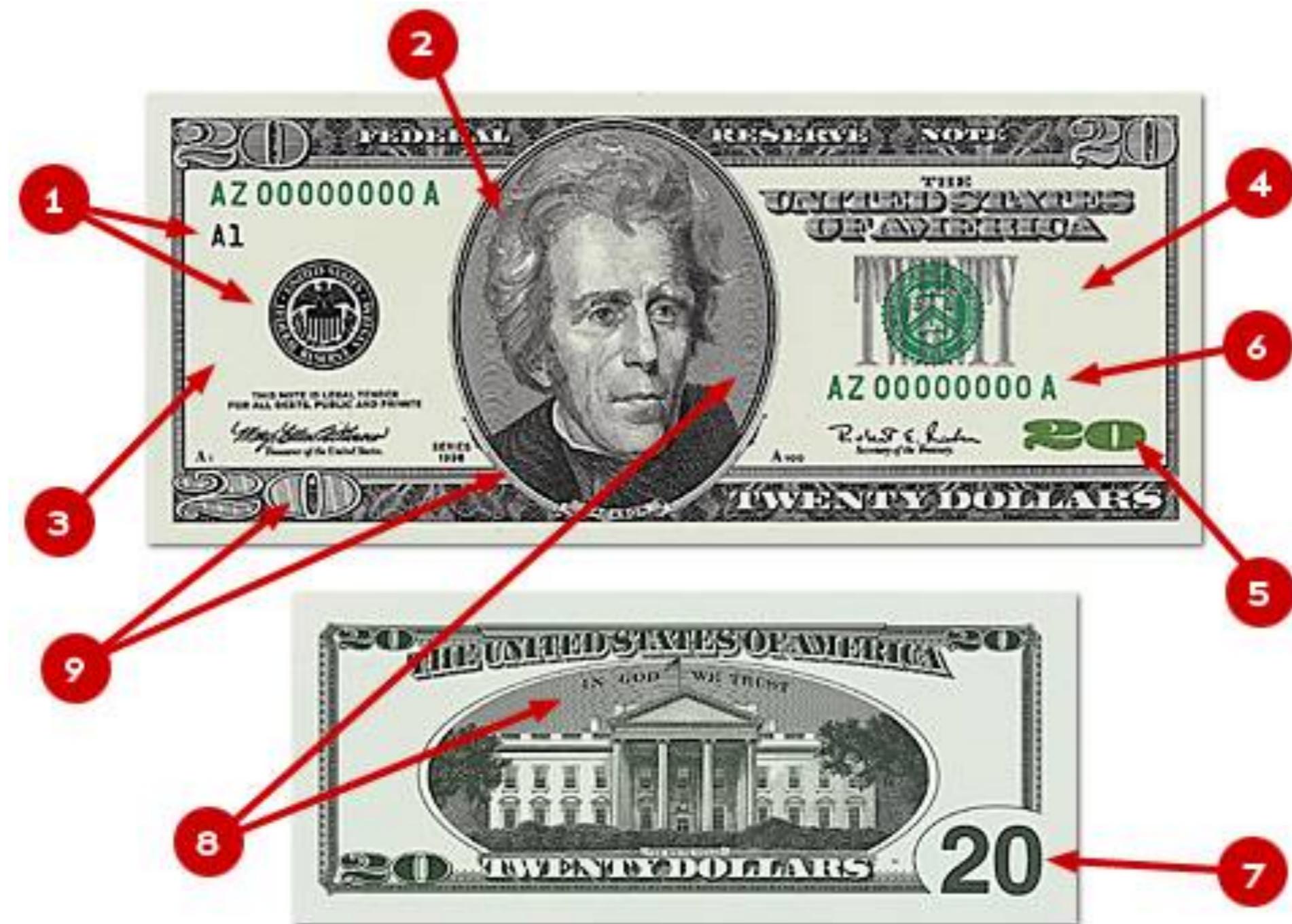
linear

trimmed mean
estimation

quadratic



Novelty Detection



Basic Idea

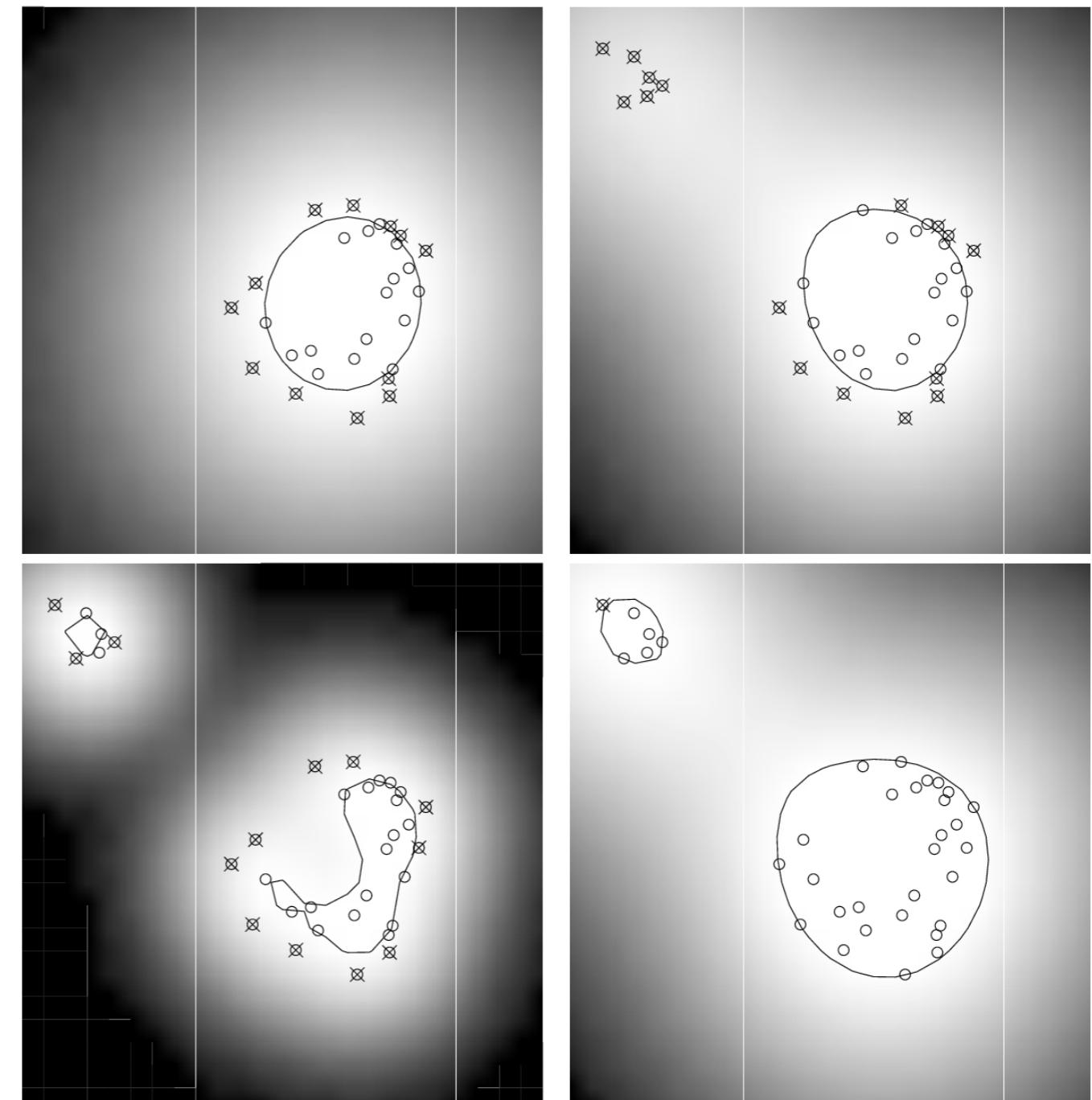
Data

Observations (x_i)
generated from
some $P(x)$, e.g.,

- network usage patterns
- handwritten digits
- alarm sensors
- factory status

Task

Find unusual events,
clean database, dis-
tinguish typical ex-
amples.



Applications

Network Intrusion Detection

Detect whether someone is trying to hack the network, downloading tons of MP3s, or doing anything else *unusual* on the network.

Jet Engine Failure Detection

You can't destroy jet engines just to see *how* they fail.

Database Cleaning

We want to find out whether someone stored bogus information in a database (typos, etc.), mislabelled digits, ugly digits, bad photographs in an electronic album.

Fraud Detection

Credit Cards, Telephone Bills, Medical Records

Self calibrating alarm devices

Car alarms (adjusts itself to where the car is parked), home alarm (furniture, temperature, windows, etc.)

Novelty Detection via

Key Idea

- Novel data is one that we don't see frequently.
- It must lie in low density regions.

Step 1: Estimate density

- Observations x_1, \dots, x_m
- Density estimate via Parzen windows

Step 2: Thresholding the density

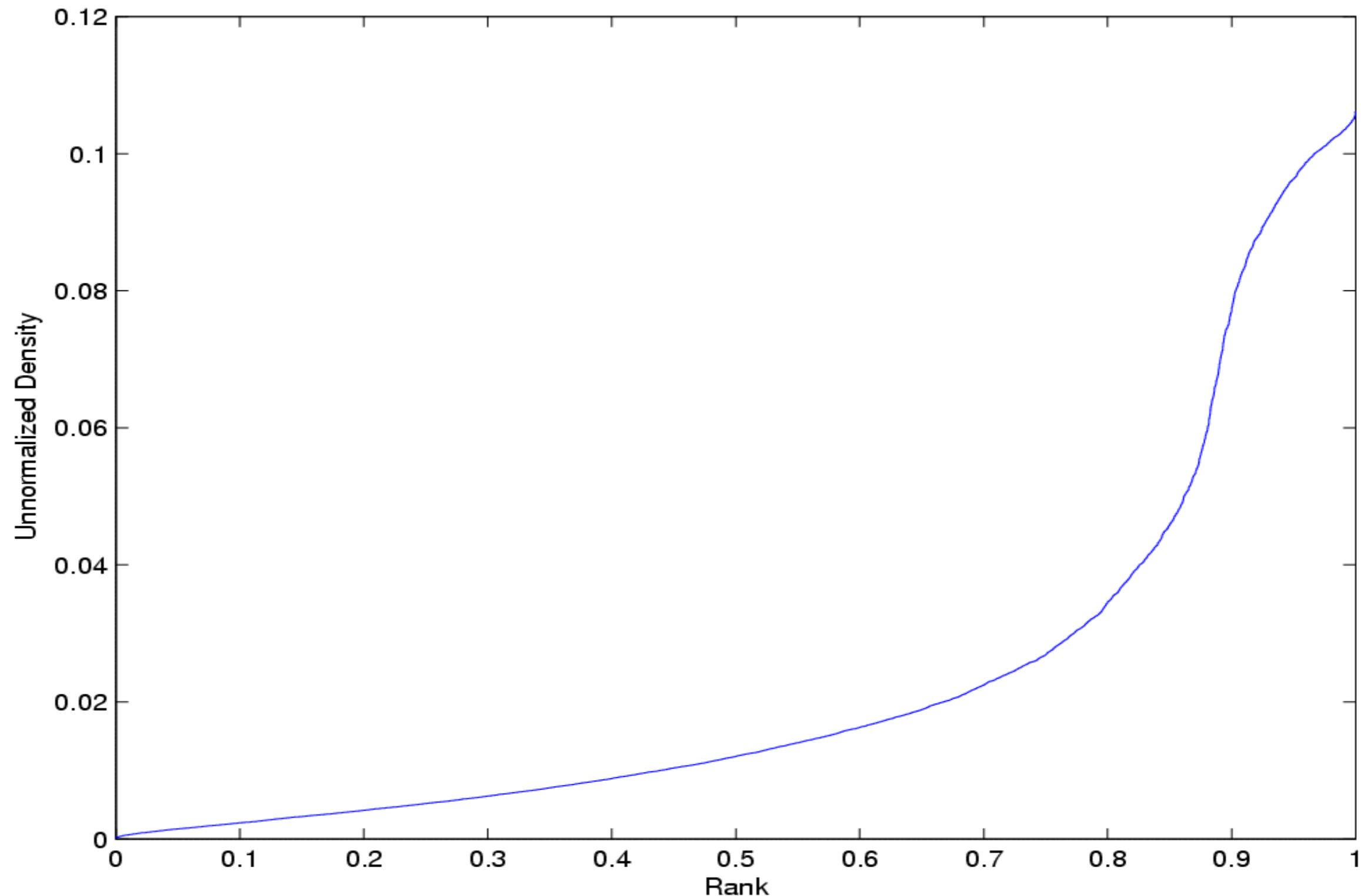
- Sort data according to density and use it for rejection
- Practical implementation: compute

$$p(x_i) = \frac{1}{m} \sum_j k(x_i, x_j) \text{ for all } i$$

and sort according to magnitude.

- Pick smallest $p(x_i)$ as novel points.

Order Statistics of Densities



Typical Data

3	4	8	6	1	1	3	6
0	0	4	7	1	4	4	2
6	0	4	3	3	7	4	1
2	6	0	0	2	1	0	0
1	7	9	3	0	6	0	0

Outliers

2 7 0 5 9 3 8 0

4 9 3 5 3 8 4 3

2 4 2 9 4 7 8 1

7 0 2 4 2 2 2 6

8 9 0 4 8 2 3

A better way

Problems

- We do not care about estimating the density properly in **regions of high density** (waste of capacity).
- We only care about the **relative density** for thresholding purposes.
- We want to eliminate a certain **fraction of observations** and tune our estimator specifically for this fraction.

Solution

- Areas of low density can be approximated as the **level set** of an auxiliary function. No need to estimate $p(x)$ directly — use proxy of $p(x)$.
- Specifically: find $f(x)$ such that x is novel if $f(x) \leq c$ where c is some constant, i.e. $f(x)$ describes the amount of novelty.

Problems with density estimation

- Exponential Family for density estimation

$$p(x|\theta) = \exp (\langle \phi(x), \theta \rangle - g(\theta))$$

- MAP estimation

$$\underset{\theta}{\text{minimize}} \sum_i g(\theta) - \langle \phi(x_i), \theta \rangle + \frac{1}{2\sigma^2} \|\theta\|^2$$

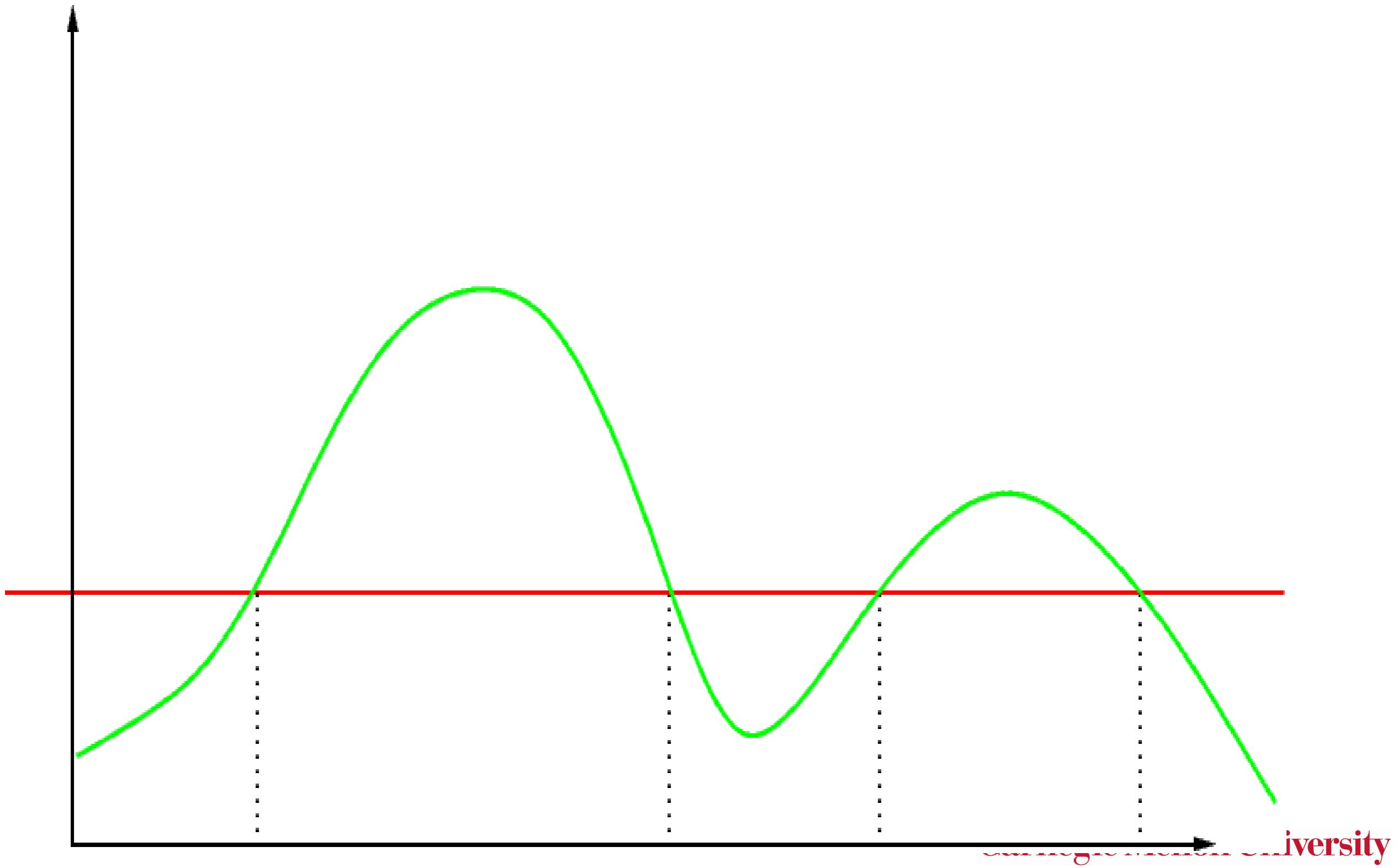
Advantages

- Convex optimization problem
- Concentration of measure

Problems

- Normalization $g(\theta)$ may be painful to compute
- For density estimation we need no normalized $p(x|\theta)$
- No need to perform particularly well in high density regions

Thresholding



Optimization Problem

Optimization Problem

$$\text{MAP} \sum_{i=1}^m -\log p(x_i|\theta) + \frac{1}{2\sigma^2} \|\theta\|^2$$

$$\text{Novelty} \sum_{i=1}^m \max \left(-\log \frac{p(x_i|\theta)}{\exp(\rho - g(\theta))}, 0 \right) + \frac{1}{2} \|\theta\|^2$$

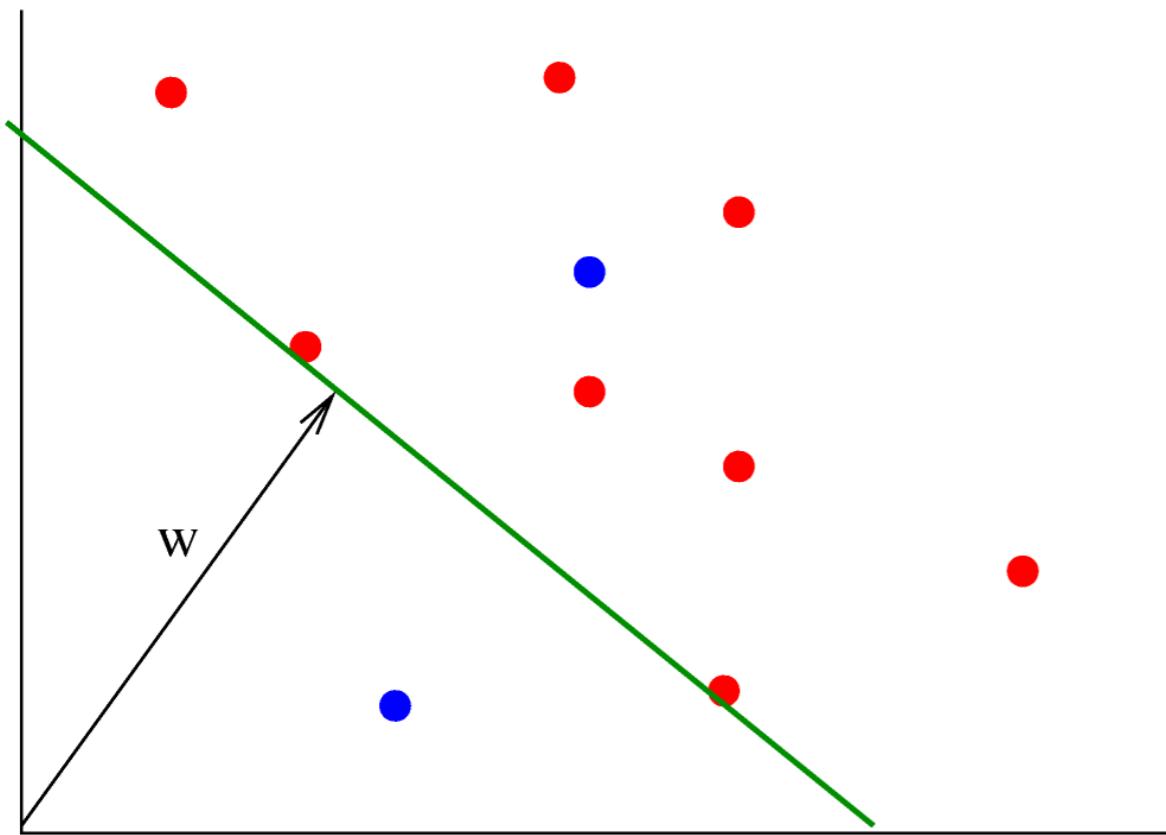
$$\sum_{i=1}^m \max(\rho - \langle \phi(x_i), \theta \rangle, 0) + \frac{1}{2} \|\theta\|^2$$

Advantages

- No normalization $g(\theta)$ needed
- No need to perform particularly well in high density regions (estimator focuses on low-density regions)
- Quadratic program

Maximum Distance

Idea Find hyperplane, given by $f(x) = \langle w, x \rangle + b = 0$ that has **maximum distance from origin** yet is still closer to the origin than the observations.



Hard Margin

minimize

$$\frac{1}{2} \|w\|^2$$

subject to

$$\langle w, x_i \rangle \geq 1$$

Soft Margin

minimize

$$\frac{1}{2} \|w\|^2 + C \sum_{i=1}^m \xi_i$$

subject to

$$\langle w, x_i \rangle \geq 1 - \xi_i$$

$$\xi_i \geq 0$$

Optimization Problem

Primal Problem

$$\begin{aligned} \text{minimize} \quad & \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m \xi_i \\ \text{subject to} \quad & \langle w, x_i \rangle - 1 + \xi_i \geq 0 \text{ and } \xi_i \geq 0 \end{aligned}$$

Lagrange Function L

- Subtract constraints, multiplied by Lagrange multipliers (α_i and η_i), from Primal Objective Function.
- Lagrange function L has **saddlepoint** at optimum.

$$L = \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m \xi_i - \sum_{i=1}^m \alpha_i (\langle w, x_i \rangle - 1 + \xi_i) - \sum_{i=1}^m \eta_i \xi_i$$

subject to $\alpha_i, \eta_i \geq 0$.

Dual Problem

Optimality Conditions

$$\begin{aligned}\partial_w L &= w - \sum_{i=1}^m \alpha_i x_i = 0 \implies w = \sum_{i=1}^m \alpha_i x_i \\ \partial_{\xi_i} L &= C - \alpha_i - \eta_i = 0 \implies \alpha_i \in [0, C]\end{aligned}$$

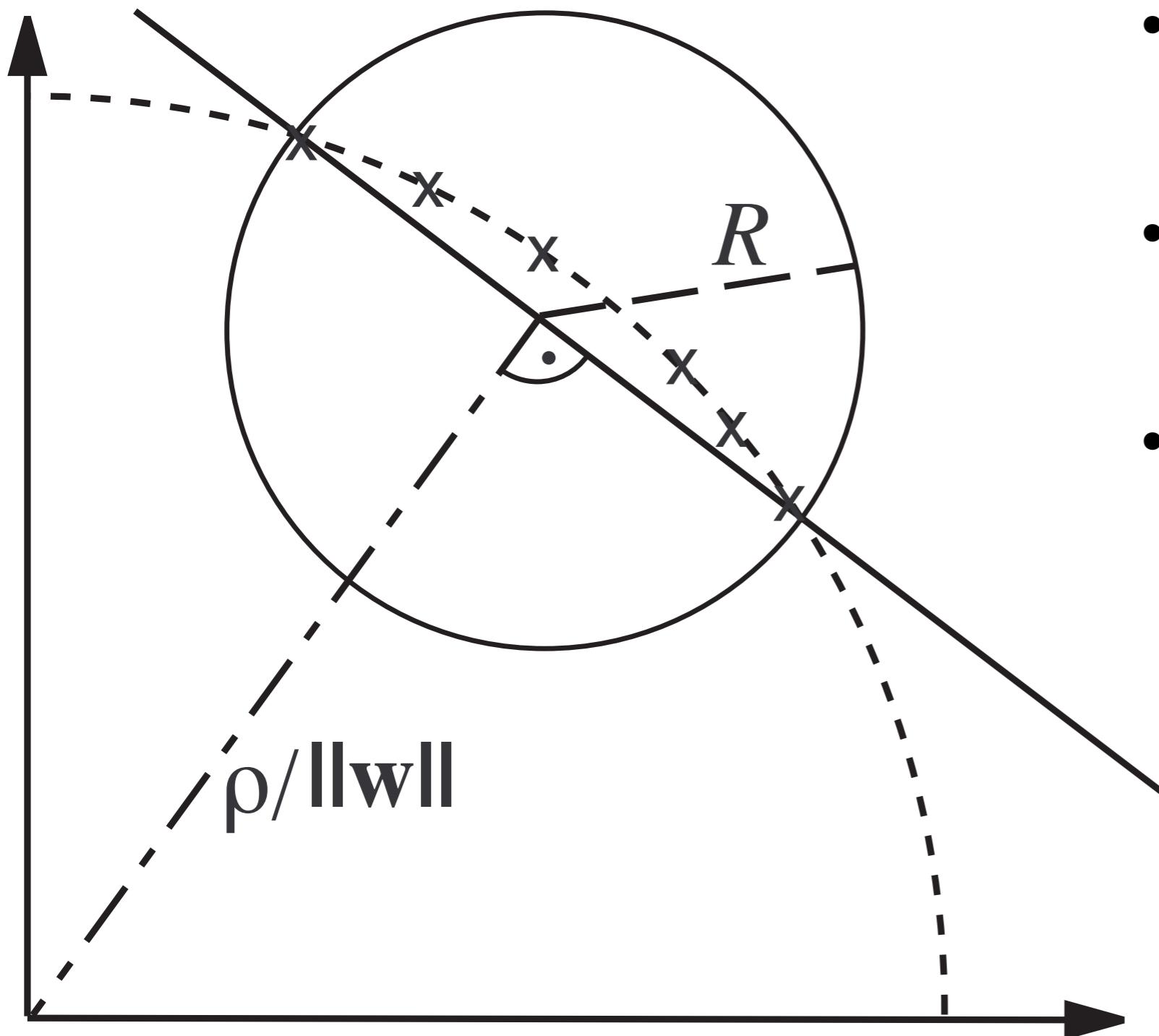
Now substitute the optimality conditions back into L .

Dual Problem

$$\begin{array}{ll}\text{minimize} & \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j \langle x_i, x_j \rangle - \sum_{i=1}^m \alpha_i \\ \text{subject to} & \alpha_i \in [0, C]\end{array}$$

All this is only possible due to the convexity of the primal problem.

Minimum enclosing ball



- Observations on surface of ball
- Find minimum enclosing ball
- Equivalent to single class SVM

Adaptive thresholds

Problem

- Depending on C , the number of novel points will vary.
- We would like to **specify the fraction** ν beforehand.

Solution

Use hyperplane separating data from the origin

$$H := \{x | \langle w, x \rangle = \rho\}$$

where the threshold ρ is **adaptive**.

Intuition

- Let the hyperplane shift by shifting ρ
- Adjust it such that the 'right' number of observations is considered novel.
- Do this automatically

Optimization Problem

Primal Problem

$$\text{minimize} \frac{1}{2} \|w\|^2 + \sum_{i=1}^m \xi_i - m\nu\rho$$

where $\langle w, x_i \rangle - \rho + \xi_i \geq 0$
 $\xi_i \geq 0$

Dual Problem

$$\text{minimize} \frac{1}{2} \sum_{i=1}^m \alpha_i \alpha_j \langle x_i, x_j \rangle$$

where $\alpha_i \in [0, 1]$ and $\sum_{i=1}^m \alpha_i = \nu m.$

The v-property theorem

- Optimization problem

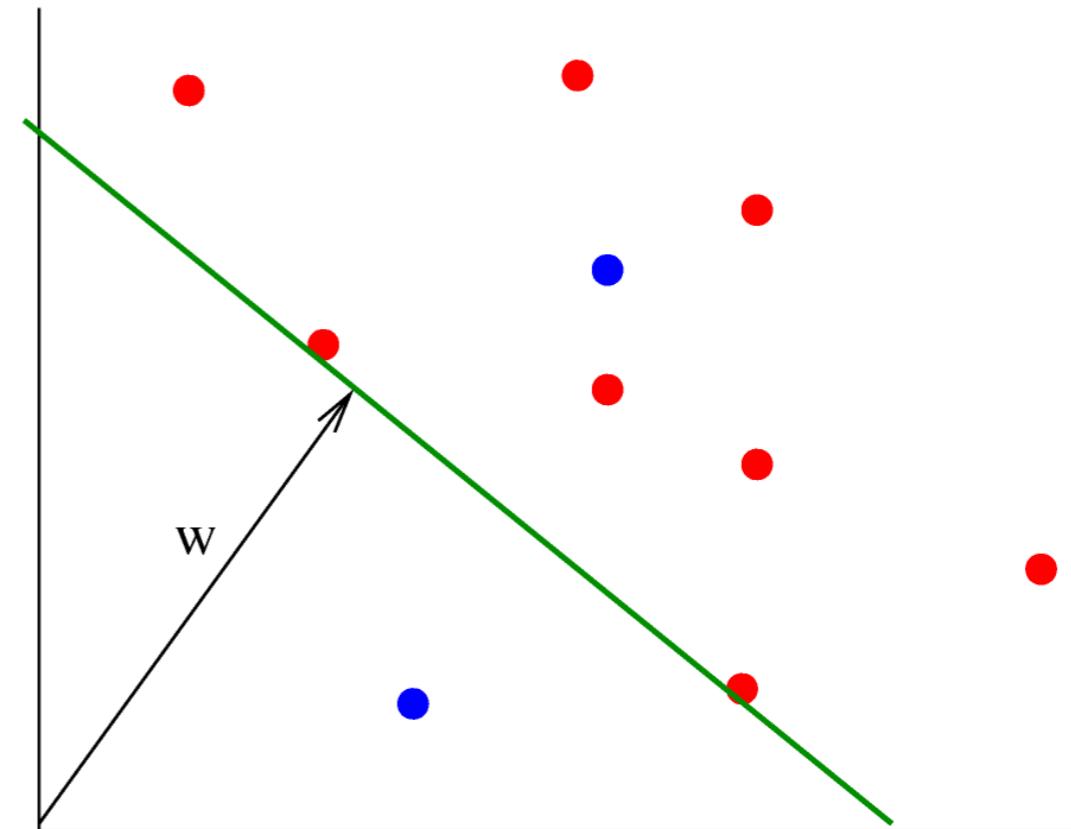
$$\underset{w}{\text{minimize}} \quad \frac{1}{2} \|w\|^2 + \sum_{i=1}^m \xi_i - m\nu\rho$$

subject to $\langle w, x_i \rangle \geq \rho - \xi_i$ and $\xi_i \geq 0$

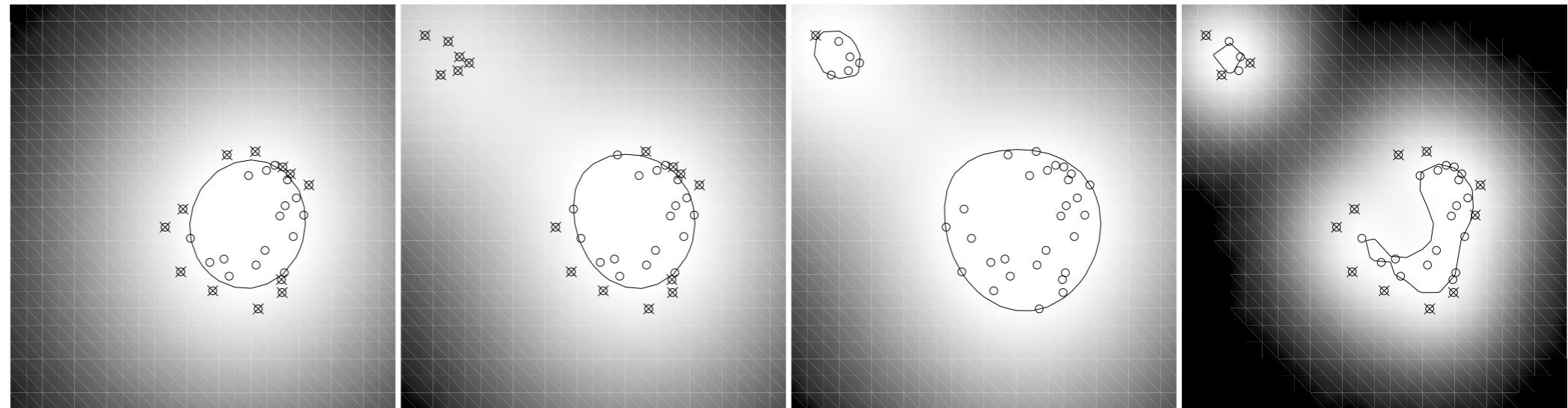
- Solution satisfies
 - At most a fraction of v points are novel
 - At most a fraction of $(1-v)$ points aren't novel
 - Fraction of points on boundary vanishes for large m (for non-pathological kernels)

Proof

- Move boundary at optimality
 - For smaller threshold m_- points on wrong side of margin contribute $\delta(m_- - \nu m) \leq 0$
 - For larger threshold m_+ points not on ‘good’ side of margin yield
$$\delta(m_+ - \nu m) \geq 0$$
 - Combining inequalities
$$\frac{m_-}{m} \leq \nu \leq \frac{m_+}{m}$$
- Margin set of measure 0



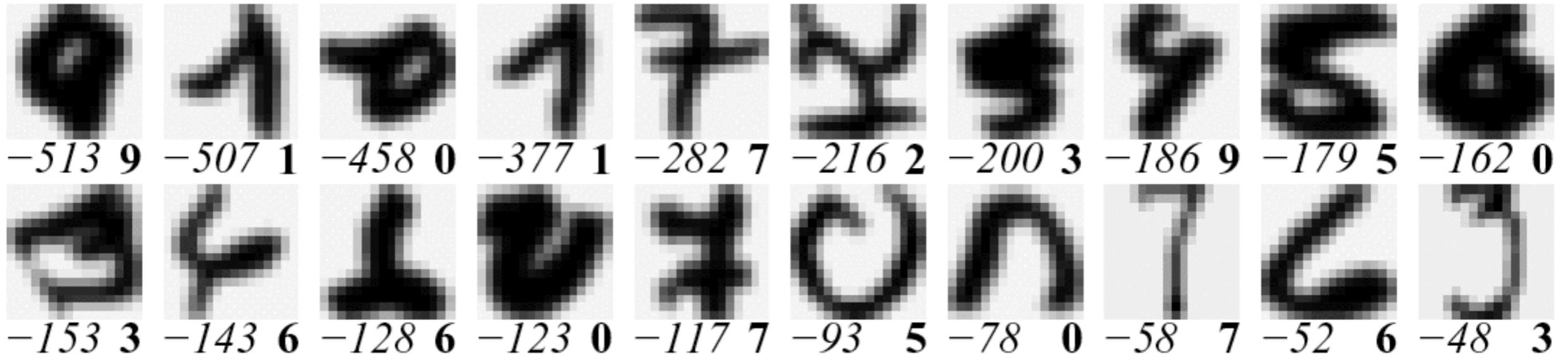
Toy example



ν , width c	0.5, 0.5	0.5, 0.5	0.1, 0.5	0.5, 0.1
frac. SVs/OLs	0.54, 0.43	0.59, 0.47	0.24, 0.03	0.65, 0.38
margin $\rho/\ \mathbf{w}\ $	0.84	0.70	0.62	0.48

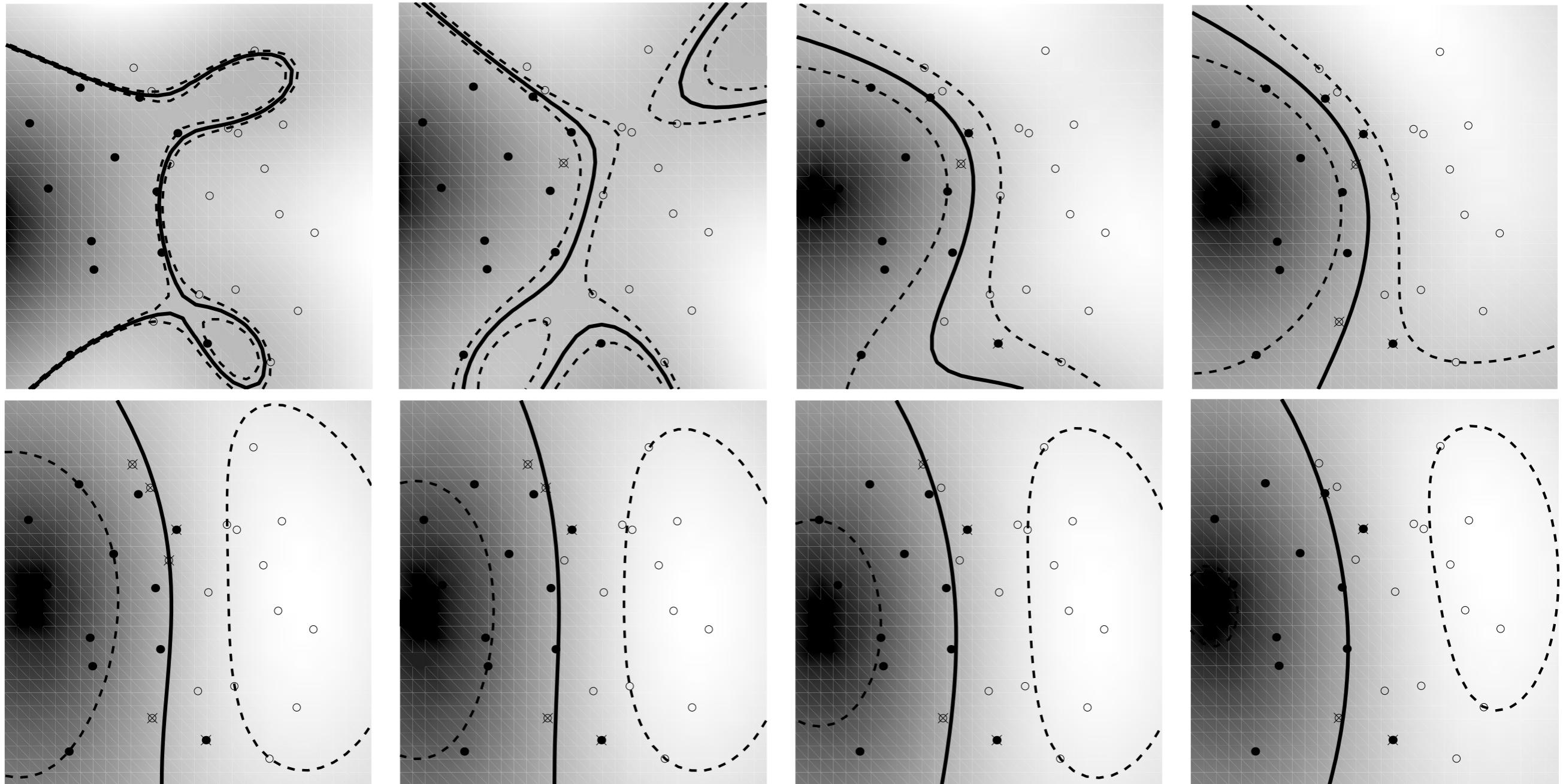
threshold and smoothness requirements

Novelty detection for OCR



- Better estimates since we only optimize in low density regions.
- Specifically tuned for small number of outliers.
- Only estimates of a level-set.
- For $\nu = 1$ we get the Parzen-windows estimator back.

Classification with the v-

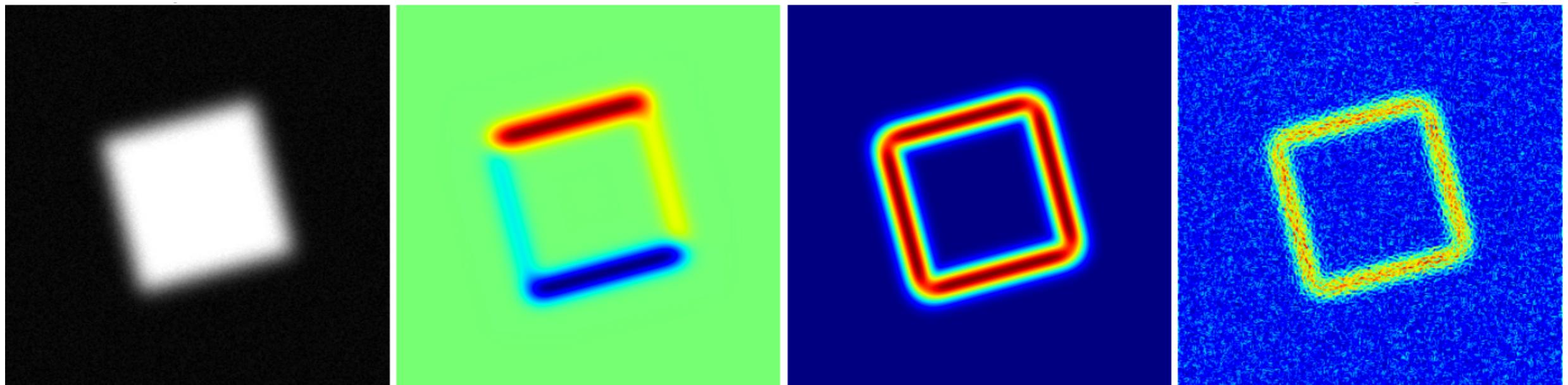


changing kernel width and threshold

Outline

- **Convex Optimization**
 - Unconstrained Optimization
 - Constrained Optimization and Duality
 - Linear and Quadratic programs
- **Support Vector Machines**
 - Classification
 - Regression
 - Novelty Detection
- **Kernels**
 - Feature Space
 - Kernel PCA
 - Kernelized SVM

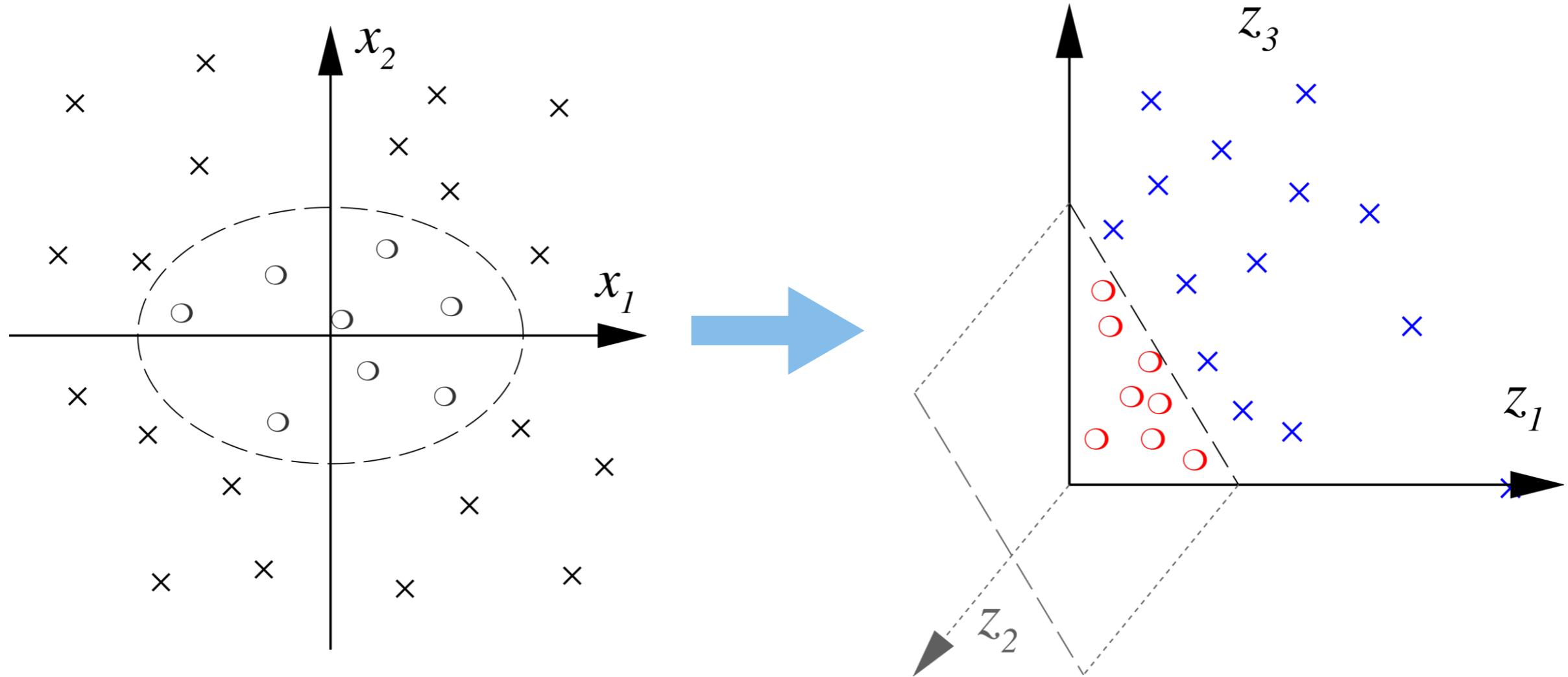
Preprocessing



Nonlinear Features

- Regression
 - We got nonlinear functions by preprocessing
- Perceptron
 - Map data into feature space $x \rightarrow \phi(x)$
 - Solve problem in this space
 - Query replace $\langle x, x' \rangle$ by $\langle \phi(x), \phi(x') \rangle$ for code
- Feature Perceptron
 - Solution in span of $\phi(x_i)$

Quadratic Features



- Separating surfaces are Circles, hyperbolae, parabolae

Constructing Features

	1	2	3	4	5	6	7	8	9	0
Loops	0	0	0	1	0	1	0	2	1	1
3 Joints	0	0	0	0	0	1	0	0	1	0
4 Joints	0	0	0	1	0	0	0	1	0	0
Angles	0	1	1	1	1	0	1	0	0	0
Ink	1	2	2	2	2	2	1	3	2	2

Delivered-To: alex.smola@gmail.com
Received: by 10.216.47.73 with SMTP id s51cs361171web;
Tue, 3 Jan 2012 14:17:53 -0800 (PST)
Received: by 10.213.17.145 with SMTP id s17mr2519891eba.147.1325629071725;
Tue, 03 Jan 2012 14:17:51 -0800 (PST)
Return-Path: <alex+caf_=alex.smola@gmail.com@smola.org>
Received: from mail-ey0-f175.google.com (mail-ey0-f175.google.com [209.85.215.175])
by mx.google.com with ESMTPS id n4si29264232eef.57.2012.01.03.14.17.51
(version=TLSv1/SSLv3 cipher=OTHER);
Tue, 03 Jan 2012 14:17:51 -0800 (PST)
Received-SPF: neutral (google.com: 209.85.215.175 is neither permitted nor denied by best
guess record for domain of alex+caf_=alex.smola@gmail.com@smola.org) client-
ip=209.85.215.175;
Authentication-Results: mx.google.com; spf=neutral (google.com: 209.85.215.175 is neither
permitted nor denied by best guess record for domain of alex
+caf_=alex.smola@gmail.com@smola.org smtp.mail=alex+caf_=alex.smola@gmail.com@smola.org;
dkim=pass (test mode) header.i=@googlemail.com
Received: by eaal1 with SMTP id l1so15092746eaa.6
for <alex.smola@gmail.com>; Tue, 03 Jan 2012 14:17:51 -0800 (PST)
Received: by 10.205.135.18 with SMTP id ie18mr5325064bkc.72.1325629071362;
Tue, 03 Jan 2012 14:17:51 -0800 (PST)
X-Forwarded-To: alex.smola@gmail.com
X-Forwarded-For: alex@smola.org alex.smola@gmail.com
Delivered-To: alex@smola.org
Received: by 10.204.65.198 with SMTP id k6cs206093bki;
Tue, 3 Jan 2012 14:17:50 -0800 (PST)
Received: by 10.52.88.179 with SMTP id bh19mr10729402vdb.38.1325629068795;
Tue, 03 Jan 2012 14:17:48 -0800 (PST)
Return-Path: <althoff.tim@googlemail.com>
Received: from mail-vx0-f179.google.com (mail-vx0-f179.google.com [209.85.220.179])
by mx.google.com with ESMTPS id dt4si11767074vdb.93.2012.01.03.14.17.48
(version=TLSv1/SSLv3 cipher=OTHER);
Tue, 03 Jan 2012 14:17:48 -0800 (PST)
Received-SPF: pass (google.com: domain of althoff.tim@googlemail.com designates
209.85.220.179 as permitted sender) client-ip=209.85.220.179;
Received: by vcbf13 with SMTP id f13so11295098vcb.10
for <alex@smola.org>; Tue, 03 Jan 2012 14:17:48 -0800 (PST)
DKIM-Signature: v=1; a=rsa-sha256; c=relaxed/relaxed;
d=googlemail.com; s=gamma;
h=mime-version:sender:date:x-google-sender-auth:message-id:subject
:from:to:content-type;
bh=WCbdZ5sXac25dpH02XcRyD0dts993hKwsAVXpGrFh0w=;
b=WK2B2+ExWnf/gvTkW6uUvKuP4XeoKn1Jq3USYTm0RARK8dSFjy0QsIHeAP9Yssxp60
7ngGoTzYqd+ZsyJfvQcLAWp1PCJhG8AMcnqWkx0NMeoFvIp2HQooZwxSOCx5ZRgY+7qX
uIbbdn4lUDXj6UFe16SpLDCkptd80Z3gr7+o=
MIME-Version: 1.0
Received: by 10.220.108.81 with SMTP id e17mr24104004vcp.67.1325629067787;
Tue, 03 Jan 2012 14:17:47 -0800 (PST)
Sender: althoff.tim@googlemail.com
Received: by 10.220.17.129 with HTTP; Tue, 3 Jan 2012 14:17:47 -0800 (PST)
Date: Tue, 3 Jan 2012 14:17:47 -0800
X-Google-Sender-Auth: 6bwi6D17HjZIkx0Eo138NZzyeHs
Message-ID: <CAFJJHDGPBW+SdZg0MdAABiAKydDk9tpeMoDiYGjoG0-WC7osg@mail.gmail.com>
Subject: CS 281B. Advanced Topics in Learning and Decision Making
From: Tim Althoff <althoff@eecs.berkeley.edu>
To: alex@smola.org
Content-Type: multipart/alternative; boundary=f46d043c7af4b07e8d04b5a7113a
--f46d043c7af4b07e8d04b5a7113a
Content-Type: text/plain; charset=ISO-8859-1

Feature Engineering for Spam Filtering

- bag of words
- pairs of words
- date & time
- recipient path
- IP number
- sender
- encoding
- links
- ... secret sauce ...

More feature engineering

- Two Interlocking Spirals
Transform the data into a radial and angular part
$$(x_1, x_2) = (r \sin \phi, r \cos \phi)$$
- Handwritten Japanese Character Recognition
 - Break down the images into strokes and recognize it
 - Lookup based on stroke order
- Medical Diagnosis
 - Physician's comments
 - Blood status / ECG / height / weight / temperature ...
 - Medical knowledge
- Preprocessing
 - Zero mean, unit variance to fix scale issue (e.g. weight vs. income)
 - Probability integral transform (inverse CDF) as alternative

The Perceptron on features

initialize $w, b = 0$

repeat

 Pick (x_i, y_i) from data

 if $y_i(w \cdot \Phi(x_i) + b) \leq 0$ then

$$w' = w + y_i \Phi(x_i)$$

$$b' = b + y_i$$

until $y_i(w \cdot \Phi(x_i) + b) > 0$ for all i

- Nothing happens if classified correctly

- Weight vector is linear combination

$$w = \sum_{i \in I} y_i \phi(x_i)$$

- Classifier is linear combination of

inner products

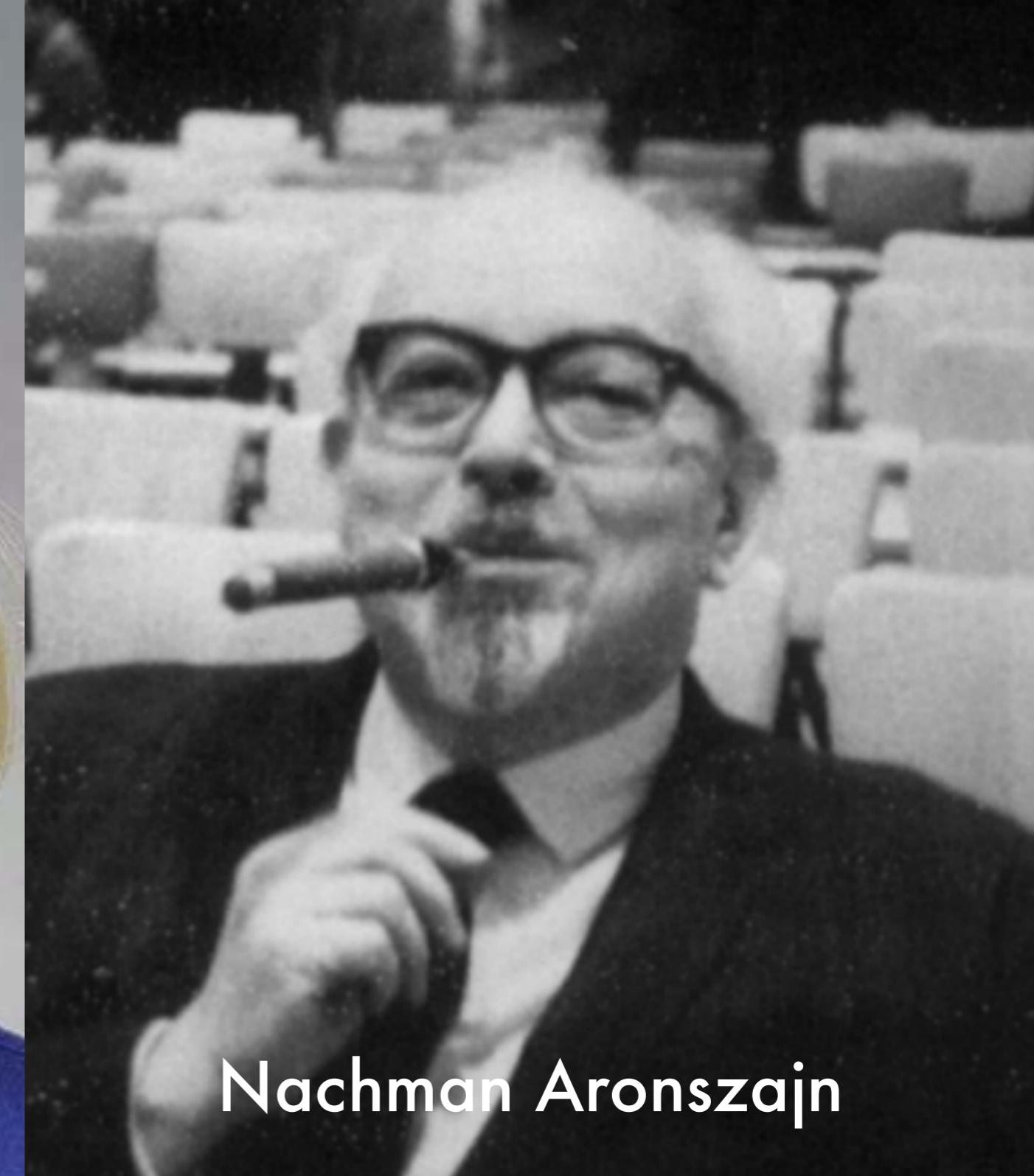
$$f(x) = \sum_{i \in I} y_i \langle \phi(x_i), \phi(x) \rangle + b$$

Problems

- Problems
 - Need domain expert (e.g. Chinese OCR)
 - Often expensive to compute
 - Difficult to transfer engineering knowledge
- Shotgun Solution
 - Compute many features
 - Hope that this contains good ones
 - Do this efficiently
- Nonlinear methods (needs lots of data & cpu)
learn the features and the classifier



Grace Wahba

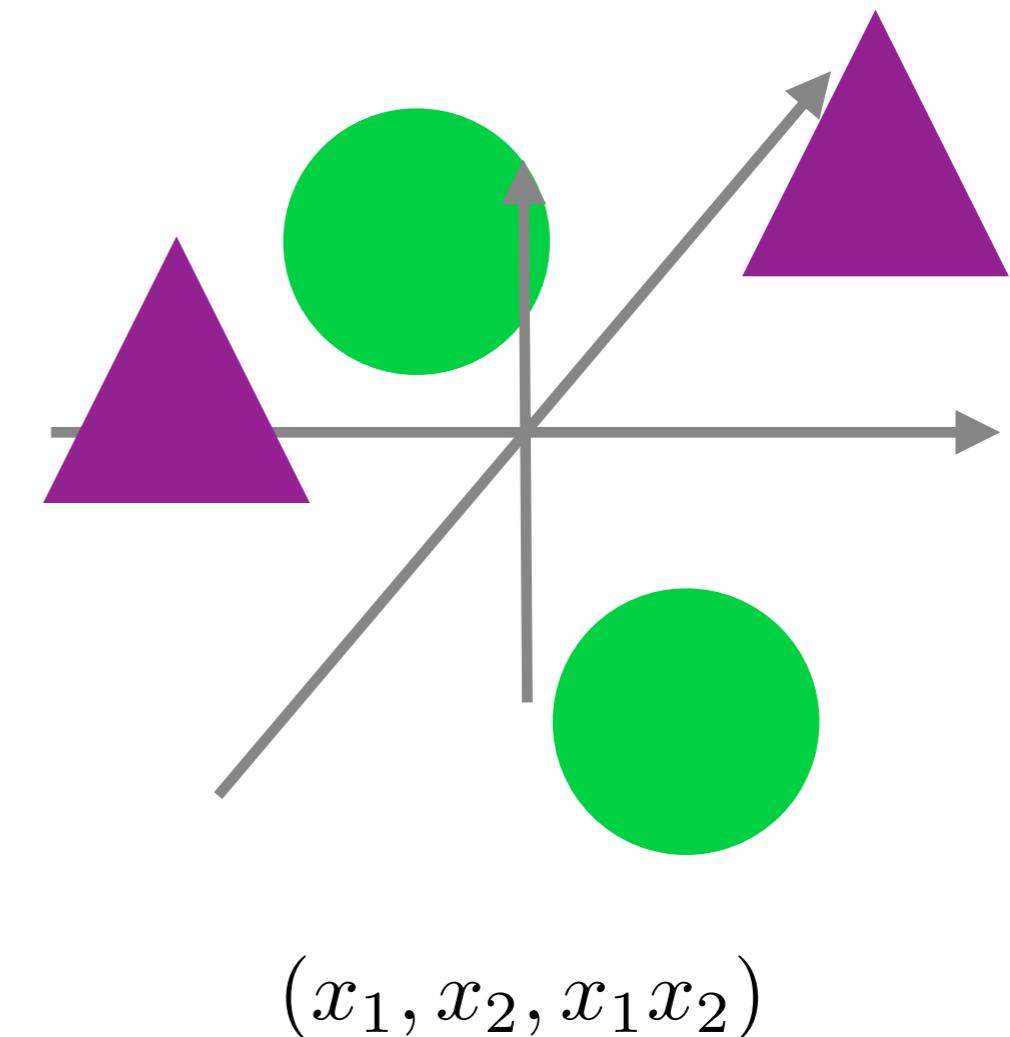
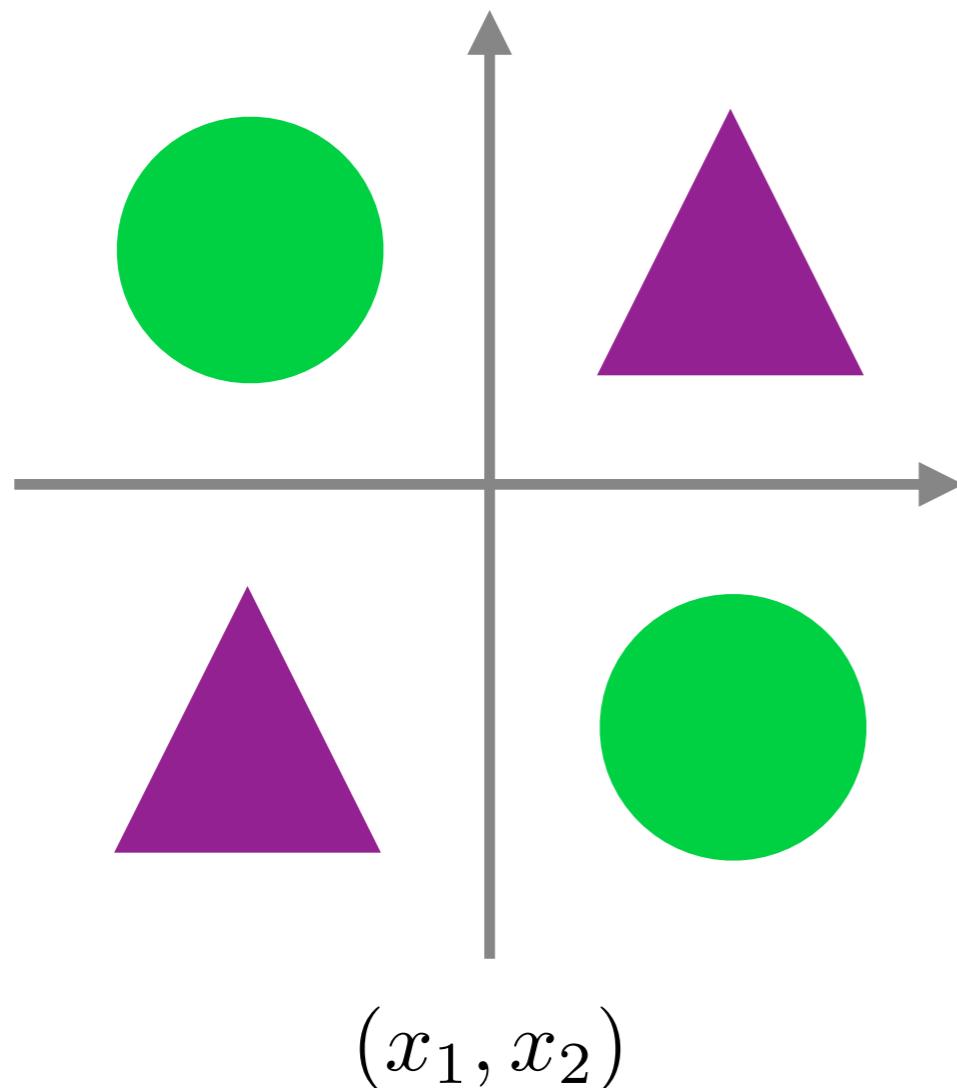


Nachman Aronszajn

Kernels

Carnegie Mellon University

Solving XOR



- XOR not linearly separable
- Mapping into 3 dimensions makes it easily solvable

Quadratic Features

Quadratic Features in \mathbb{R}^2

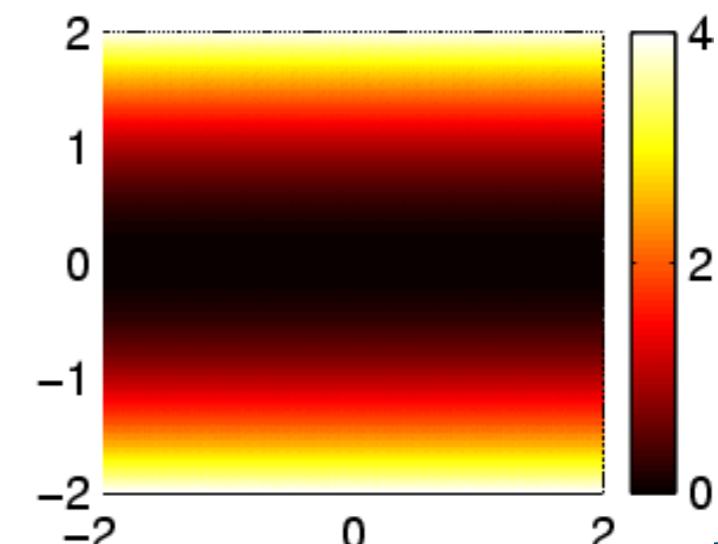
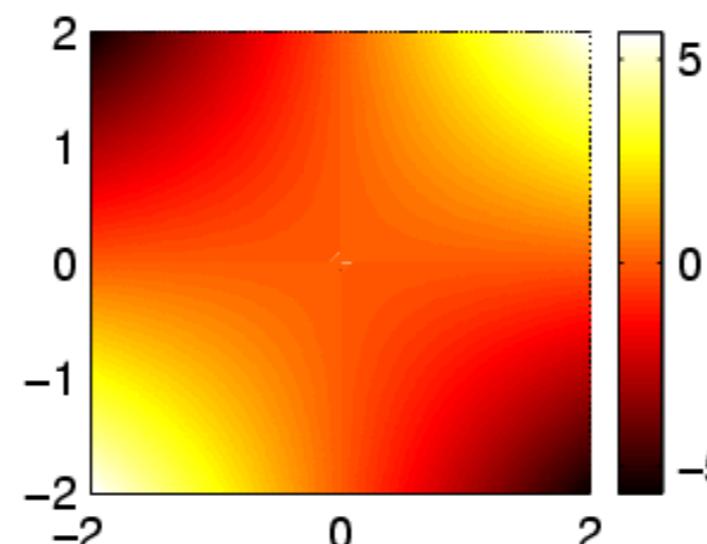
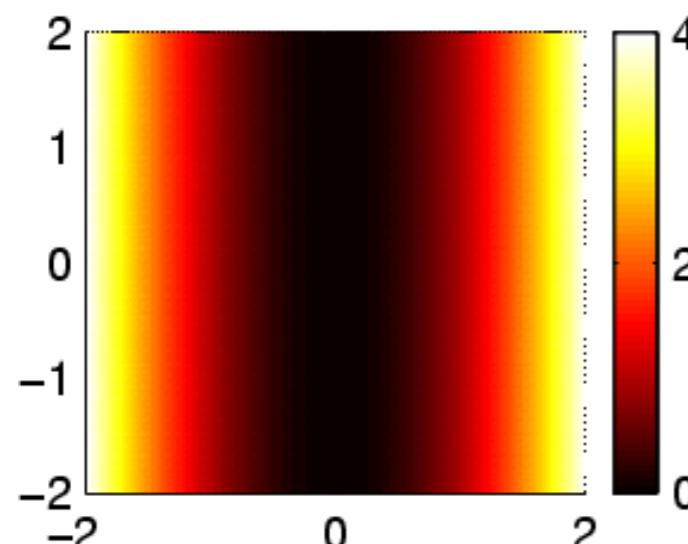
$$\Phi(x) := \left(x_1^2, \sqrt{2}x_1x_2, x_2^2 \right)$$

Dot Product

$$\begin{aligned}\langle \Phi(x), \Phi(x') \rangle &= \left\langle \left(x_1^2, \sqrt{2}x_1x_2, x_2^2 \right), \left({x'_1}^2, \sqrt{2}x'_1x'_2, {x'_2}^2 \right) \right\rangle \\ &= \langle x, x' \rangle^2.\end{aligned}$$

Insight

Trick works for any polynomials of order d via $\langle x, x' \rangle^d$.



SVM with a polynomial Kernel visualization

Created by:
Udi Aharoni

SVM with a polynomial Kernel visualization

Created by:
Udi Aharoni

Computational Efficiency

Problem

- Extracting features can sometimes be very costly.
- Example: second order features in 1000 dimensions.
This leads to $5 \cdot 10^5$ numbers. For higher order polynomial features much worse.

Solution

Don't compute the features, try to compute dot products implicitly. For some features this works . . .

Definition

A kernel function $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is a symmetric function in its arguments for which the following property holds

$$k(x, x') = \langle \Phi(x), \Phi(x') \rangle \text{ for some feature map } \Phi.$$

If $k(x, x')$ is much cheaper to compute than $\Phi(x)$. . .

The Kernel Perceptron

initialize $f = 0$

repeat

Pick (x_i, y_i) from data

if $y_i f(x_i) \leq 0$ **then**

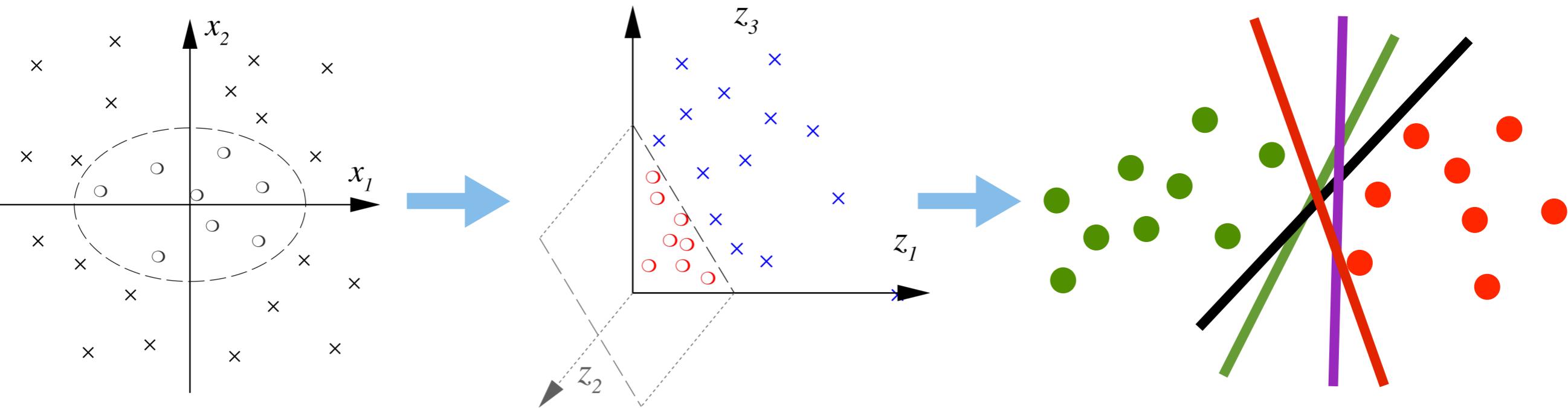
$f(\cdot) \leftarrow f(\cdot) + y_i k(x_i, \cdot) + y_i$

until $y_i f(x_i) > 0$ for all i

- Nothing happens if classified correctly
- Weight vector is linear combination $w = \sum_{i \in I} y_i \phi(x_i)$
- Classifier is linear combination of inner products

$$f(x) = \sum_{i \in I} y_i \langle \phi(x_i), \phi(x) \rangle + b = \sum_{i \in I} y_i k(x_i, x) + b$$

Processing Pipeline



- Original data
- Data in feature space (implicit)
- Solve in feature space using kernels

Polynomial Kernels

Idea

- We want to extend $k(x, x') = \langle x, x' \rangle^2$ to
$$k(x, x') = (\langle x, x' \rangle + c)^d$$
 where $c > 0$ and $d \in \mathbb{N}$.
- Prove that such a kernel corresponds to a dot product.

Proof strategy

Simple and straightforward: compute the explicit sum given by the kernel, i.e.

$$k(x, x') = (\langle x, x' \rangle + c)^d = \sum_{i=0}^m \binom{d}{i} (\langle x, x' \rangle)^i c^{d-i}$$

Individual terms $(\langle x, x' \rangle)^i$ are dot products for some $\Phi_i(x)$.

Kernel Conditions

Computability

We have to be able to compute $k(x, x')$ efficiently (much cheaper than dot products themselves).

“Nice and Useful” Functions

The features themselves have to be useful for the learning problem at hand. Quite often this means smooth functions.

Symmetry

Obviously $k(x, x') = k(x', x)$ due to the symmetry of the dot product $\langle \Phi(x), \Phi(x') \rangle = \langle \Phi(x'), \Phi(x) \rangle$.

Dot Product in Feature Space

Is there always a Φ such that k really is a dot product?

Mercer's Theorem

The Theorem

For any symmetric function $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ which is square integrable in $\mathcal{X} \times \mathcal{X}$ and which satisfies

$$\int_{\mathcal{X} \times \mathcal{X}} k(x, x') f(x) f(x') dx dx' \geq 0 \text{ for all } f \in L_2(\mathcal{X})$$

there exist $\phi_i : \mathcal{X} \rightarrow \mathbb{R}$ and numbers $\lambda_i \geq 0$ where

$$k(x, x') = \sum_i \lambda_i \phi_i(x) \phi_i(x') \text{ for all } x, x' \in \mathcal{X}.$$

Interpretation

Double integral is the continuous version of a vector-matrix-vector multiplication. For positive semidefinite matrices we have

$$\sum \sum k(x_i, x_j) \alpha_i \alpha_j \geq 0$$

Properties

Distance in Feature Space

Distance between points in feature space via

$$\begin{aligned} d(x, x')^2 &:= \|\Phi(x) - \Phi(x')\|^2 \\ &= \langle \Phi(x), \Phi(x) \rangle - 2\langle \Phi(x), \Phi(x') \rangle + \langle \Phi(x'), \Phi(x') \rangle \\ &= k(x, x) + k(x', x') - 2k(x, x') \end{aligned}$$

Kernel Matrix

To compare observations we compute dot products, so we study the matrix K given by

$$K_{ij} = \langle \Phi(x_i), \Phi(x_j) \rangle = k(x_i, x_j)$$

where x_i are the training patterns.

Similarity Measure

The entries K_{ij} tell us the overlap between $\Phi(x_i)$ and $\Phi(x_j)$, so $k(x_i, x_j)$ is a similarity measure.

Properties

K is Positive Semidefinite

Claim: $\alpha^\top K \alpha \geq 0$ for all $\alpha \in \mathbb{R}^m$ and all kernel matrices $K \in \mathbb{R}^{m \times m}$. Proof:

$$\begin{aligned}\sum_{i,j}^m \alpha_i \alpha_j K_{ij} &= \sum_{i,j}^m \alpha_i \alpha_j \langle \Phi(x_i), \Phi(x_j) \rangle \\ &= \left\langle \sum_i^m \alpha_i \Phi(x_i), \sum_j^m \alpha_j \Phi(x_j) \right\rangle = \left\| \sum_{i=1}^m \alpha_i \Phi(x_i) \right\|^2\end{aligned}$$

Kernel Expansion

If w is given by a linear combination of $\Phi(x_i)$ we get

$$\langle w, \Phi(x) \rangle = \left\langle \sum_{i=1}^m \alpha_i \Phi(x_i), \Phi(x) \right\rangle = \sum_{i=1}^m \alpha_i k(x_i, x).$$

A Counterexample

A Candidate for a Kernel

$$k(x, x') = \begin{cases} 1 & \text{if } \|x - x'\| \leq 1 \\ 0 & \text{otherwise} \end{cases}$$

This is symmetric and gives us some information about the proximity of points, yet it is not a proper kernel . . .

Kernel Matrix

We use three points, $x_1 = 1, x_2 = 2, x_3 = 3$ and compute the resulting “kernelmatrix” K . This yields

$$K = \begin{bmatrix} 1 & 1 & 0 \\ 1 & 1 & 1 \\ 0 & 1 & 1 \end{bmatrix} \text{ and eigenvalues } (\sqrt{2}-1)^{-1}, 1 \text{ and } (1-\sqrt{2}).$$

as eigensystem. Hence k is not a kernel.

Examples

Examples of kernels $k(x, x')$

Linear

$$\langle x, x' \rangle$$

Laplacian RBF

$$\exp(-\lambda \|x - x'\|)$$

Gaussian RBF

$$\exp(-\lambda \|x - x'\|^2)$$

Polynomial

$$(\langle x, x' \rangle + c)^d, c \geq 0, d \in \mathbb{N}$$

B-Spline

$$B_{2n+1}(x - x')$$

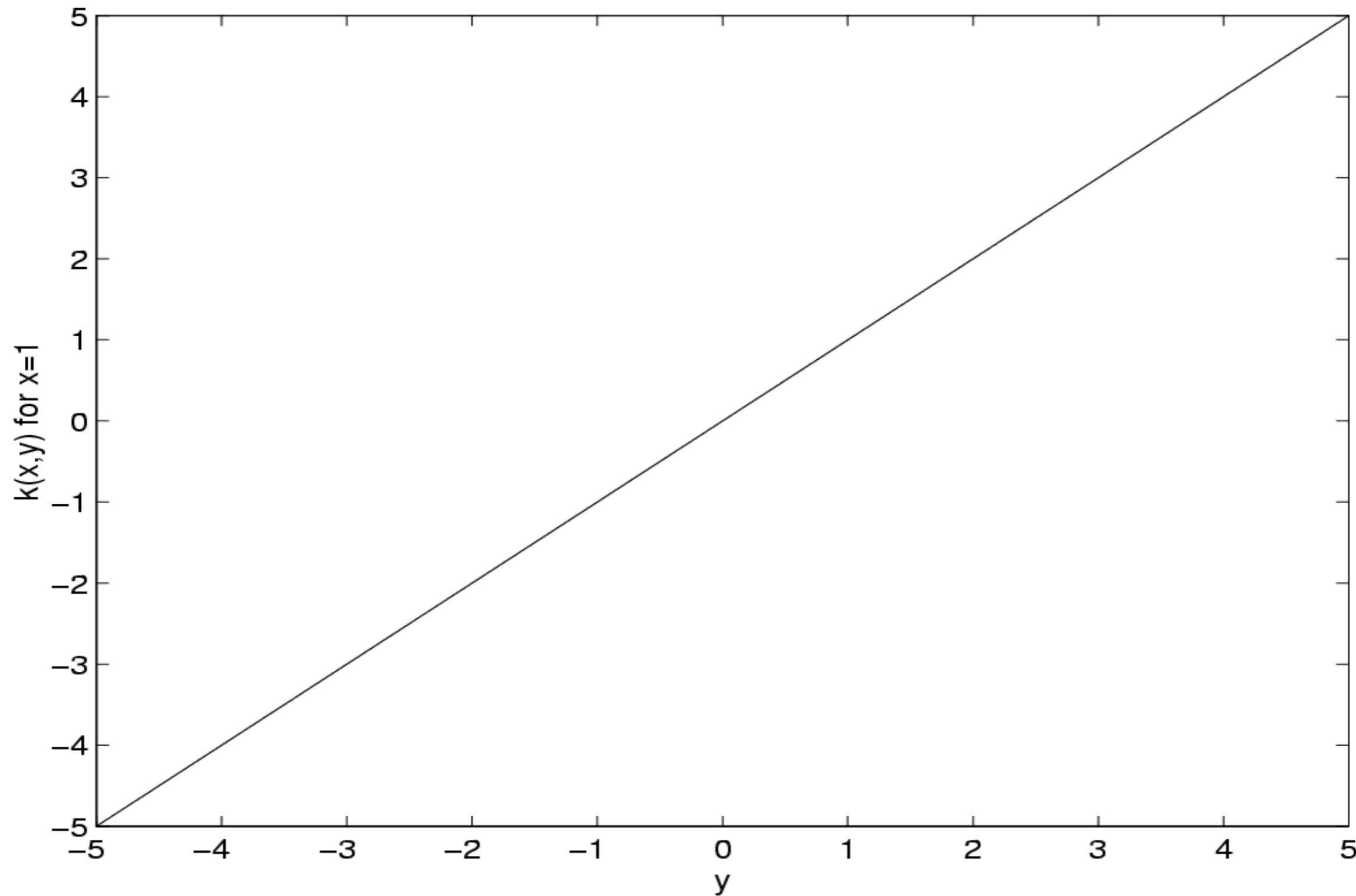
Cond. Expectation

$$\mathbf{E}_c[p(x|c)p(x'|c)]$$

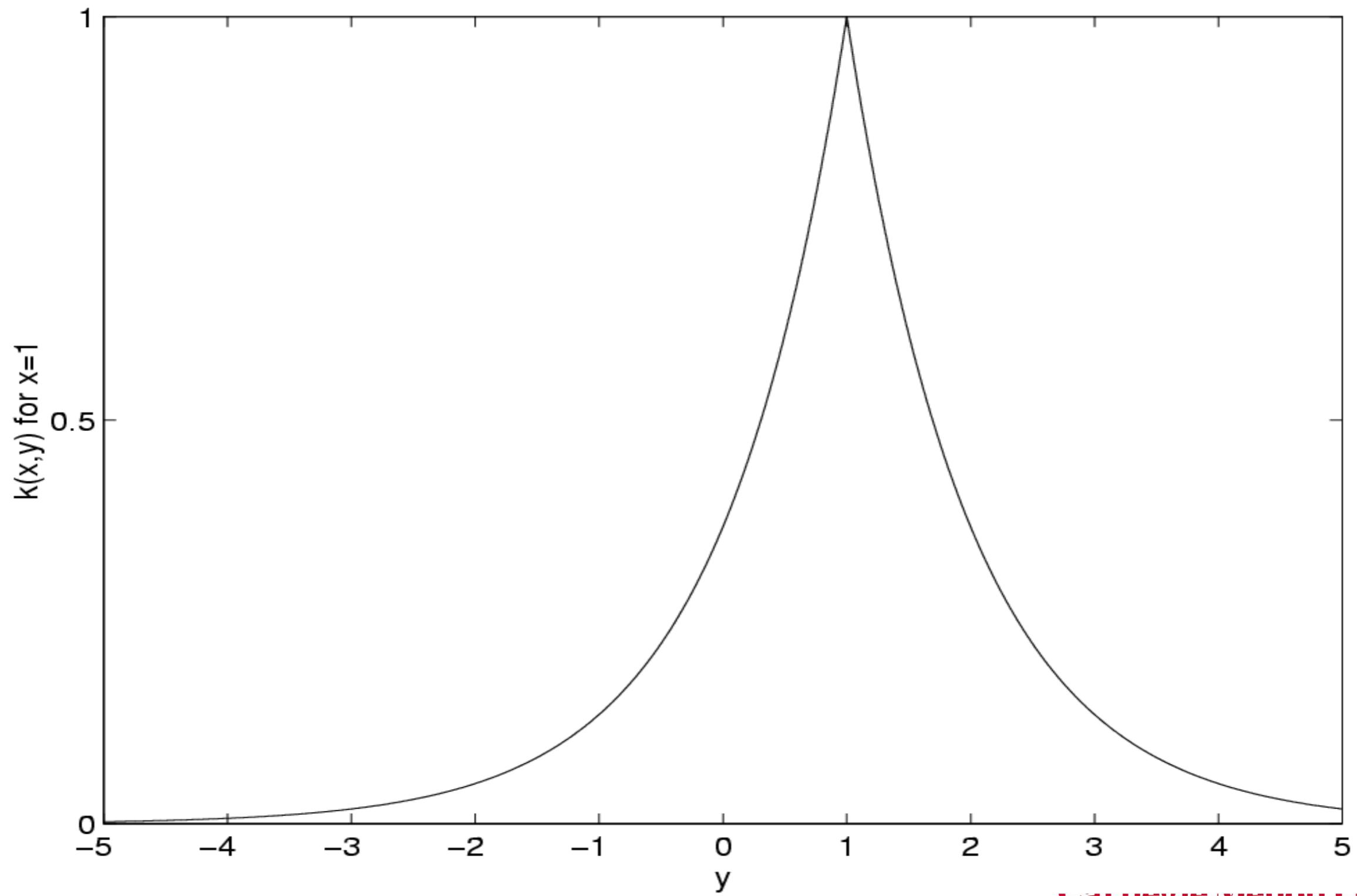
Simple trick for checking Mercer's condition

Compute the Fourier transform of the kernel and check that it is nonnegative.

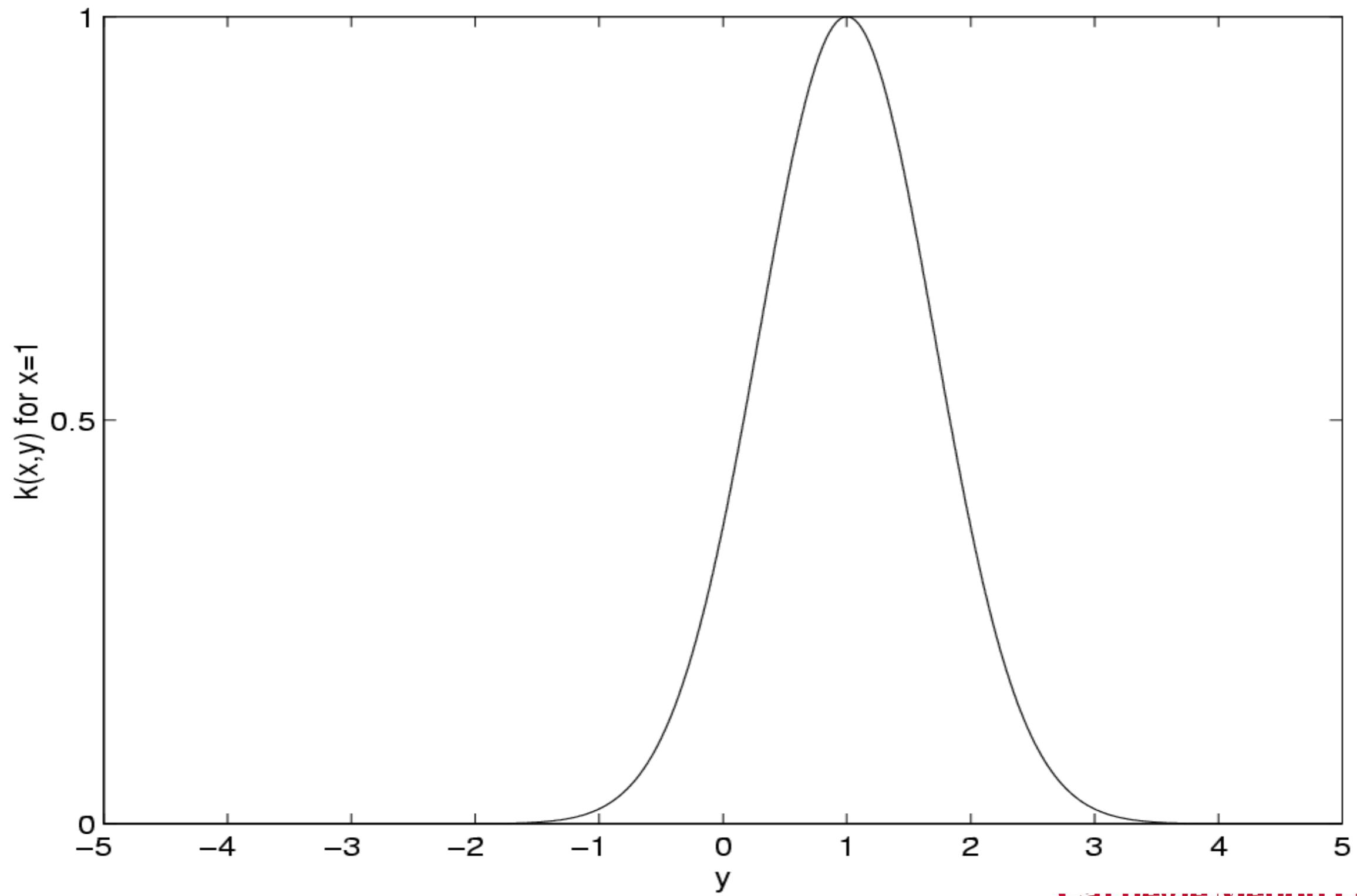
Linear Kernel



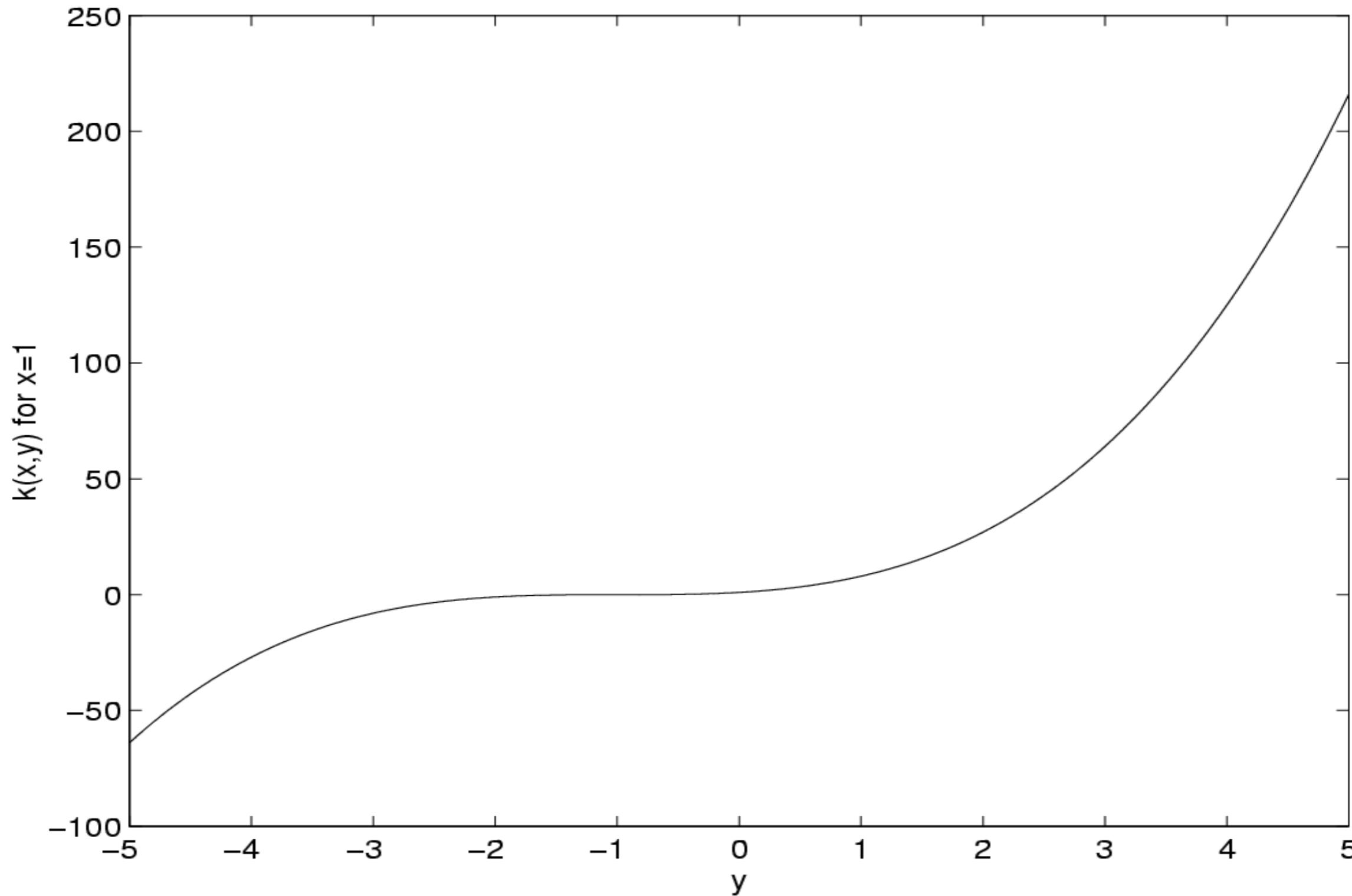
Laplacian Kernel



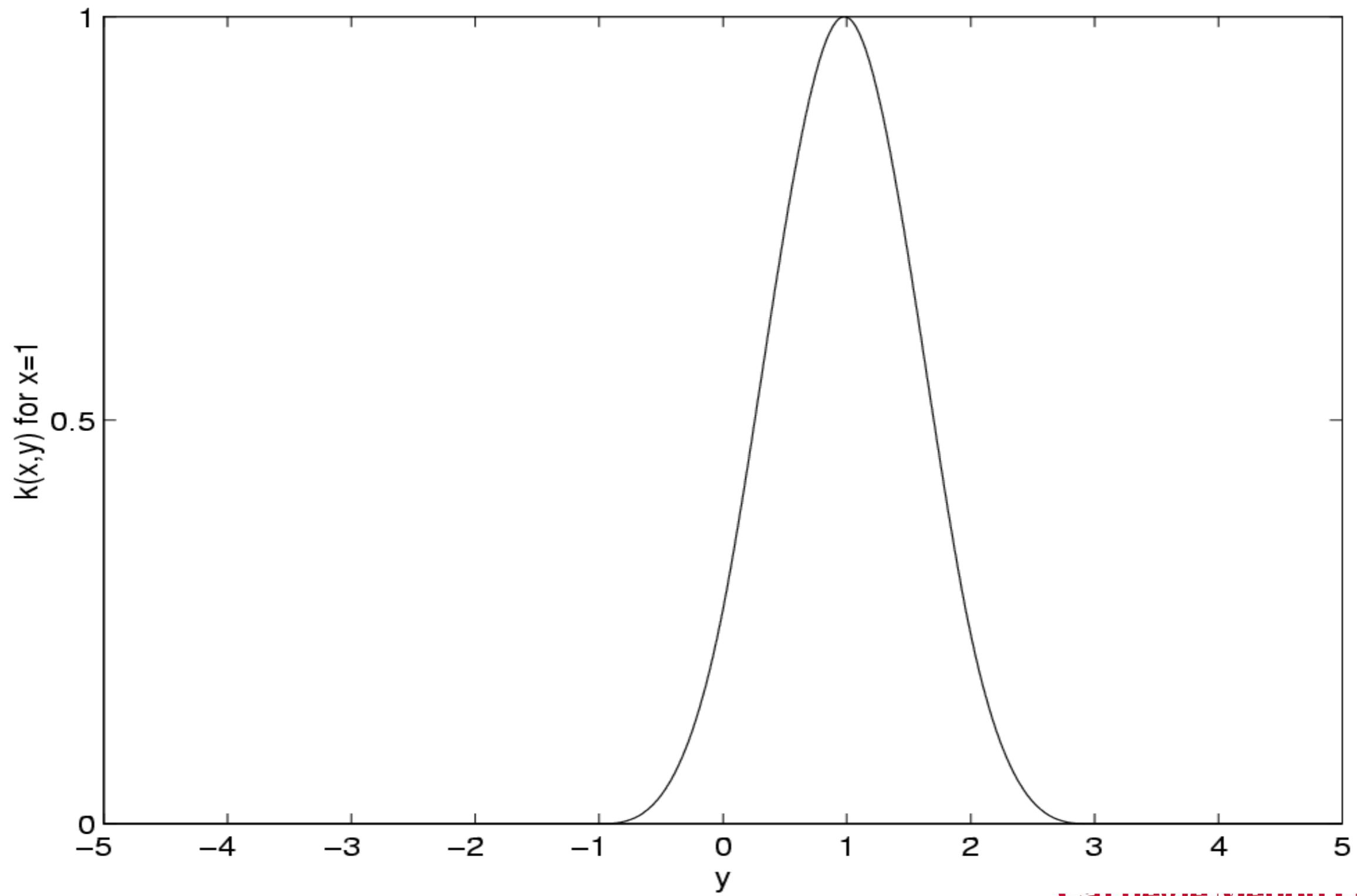
Gaussian Kernel



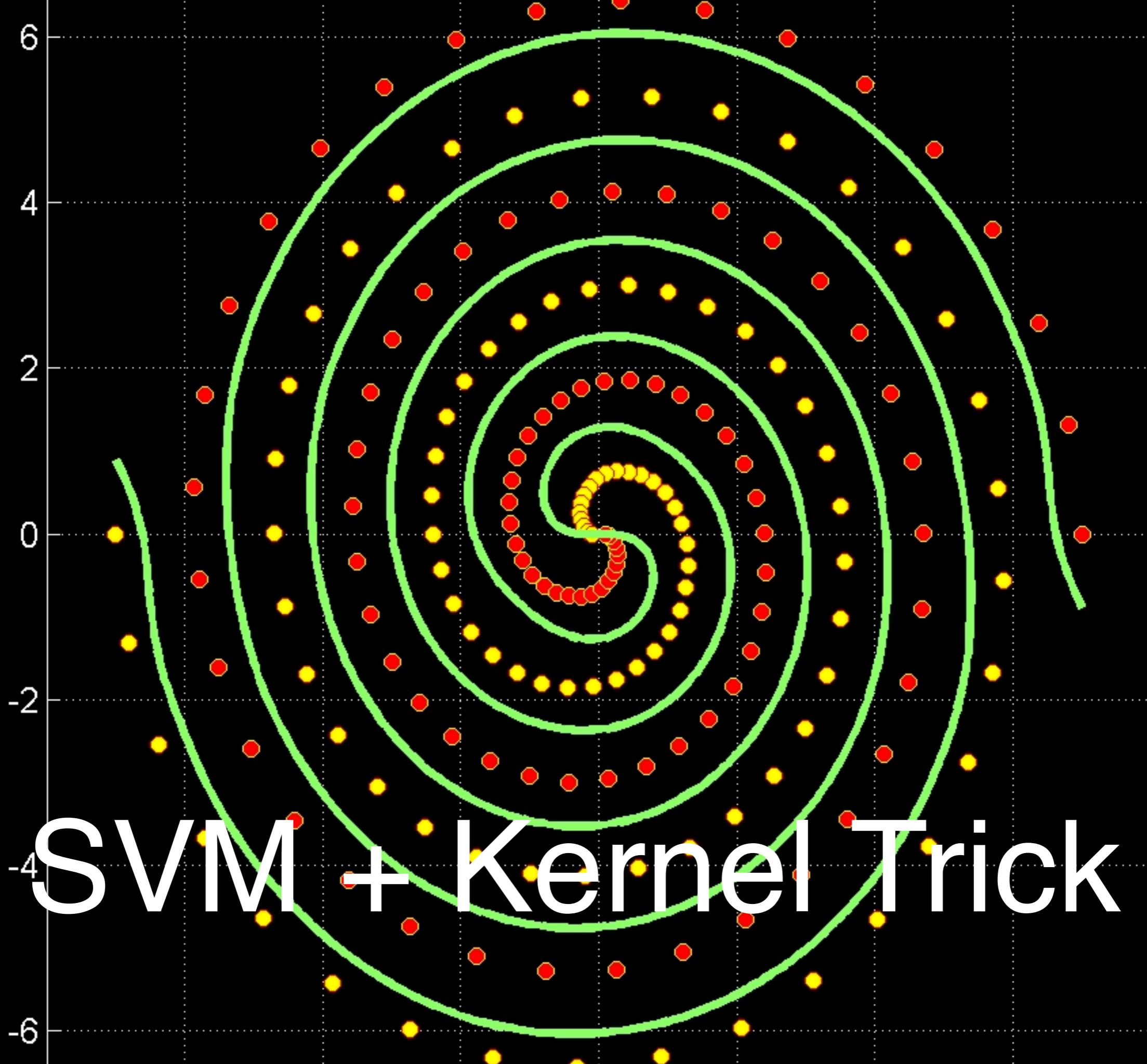
Polynomial of order 3



B₃ Spline Kernel



SVM + Kernel Trick



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

- Dual problem

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

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

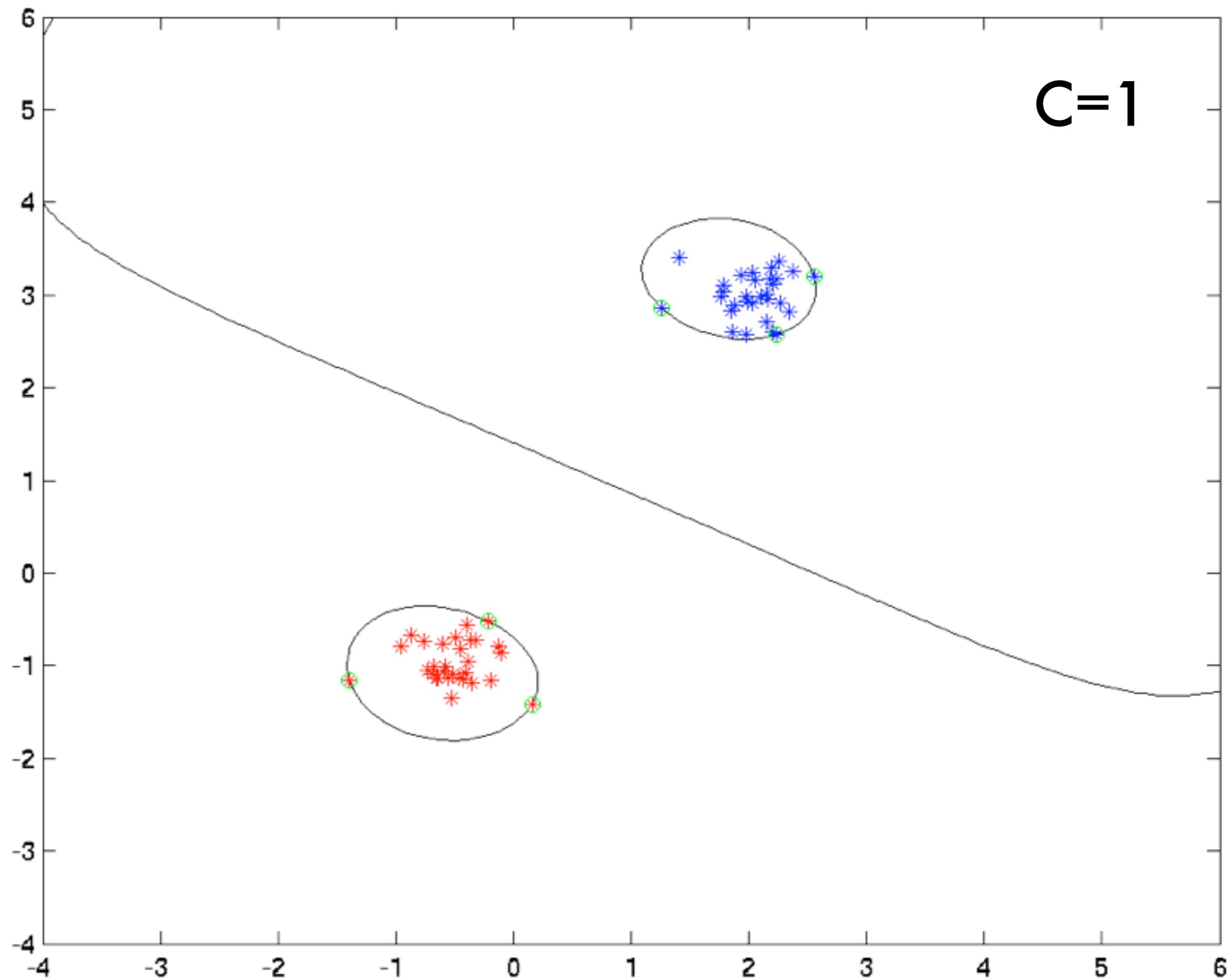
$$\underset{\alpha}{\text{maximize}} \quad -\frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j k(x_i, x_j) + \sum_i \alpha_i$$

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

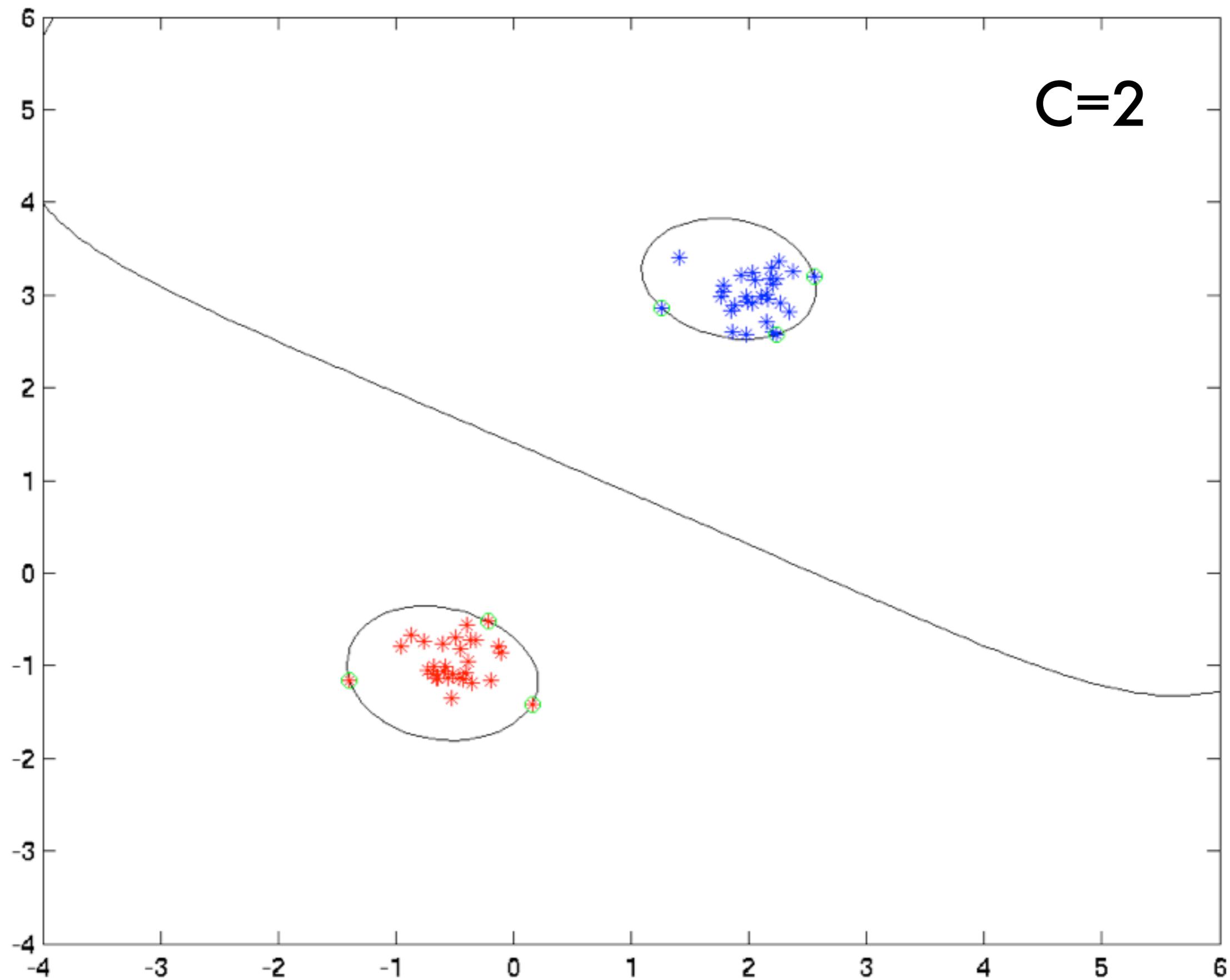
- Support vector expansion

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

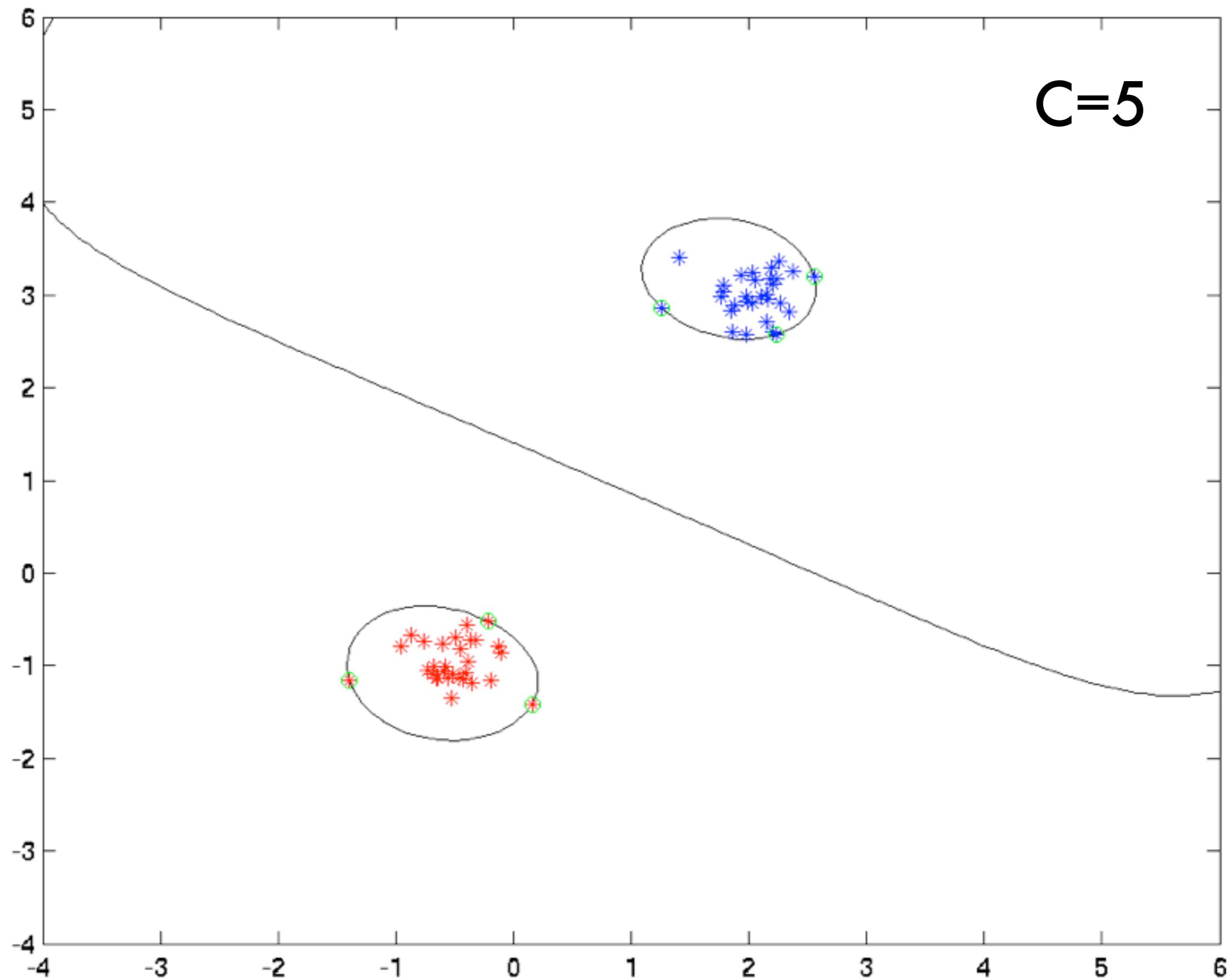
C=1



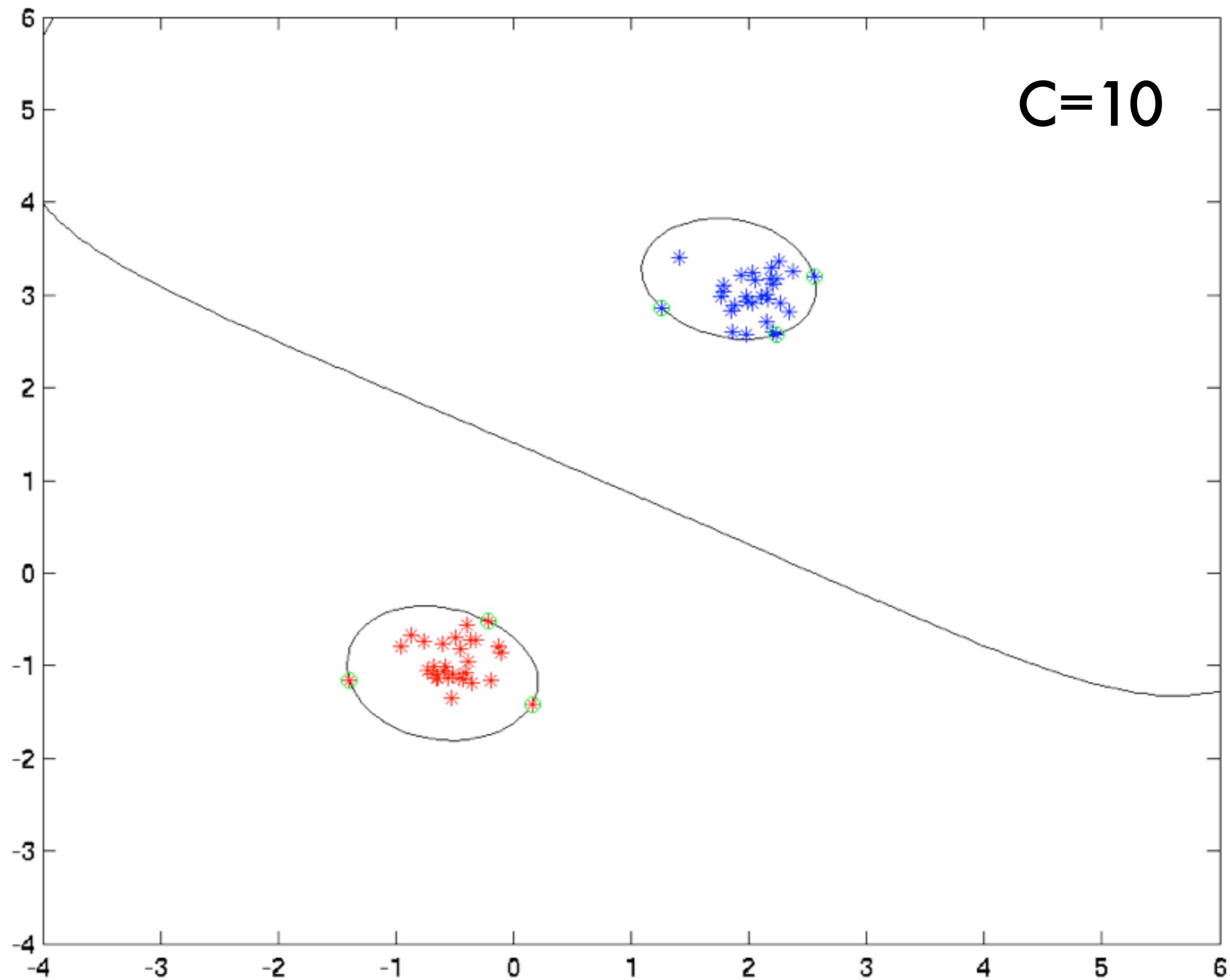
C=2



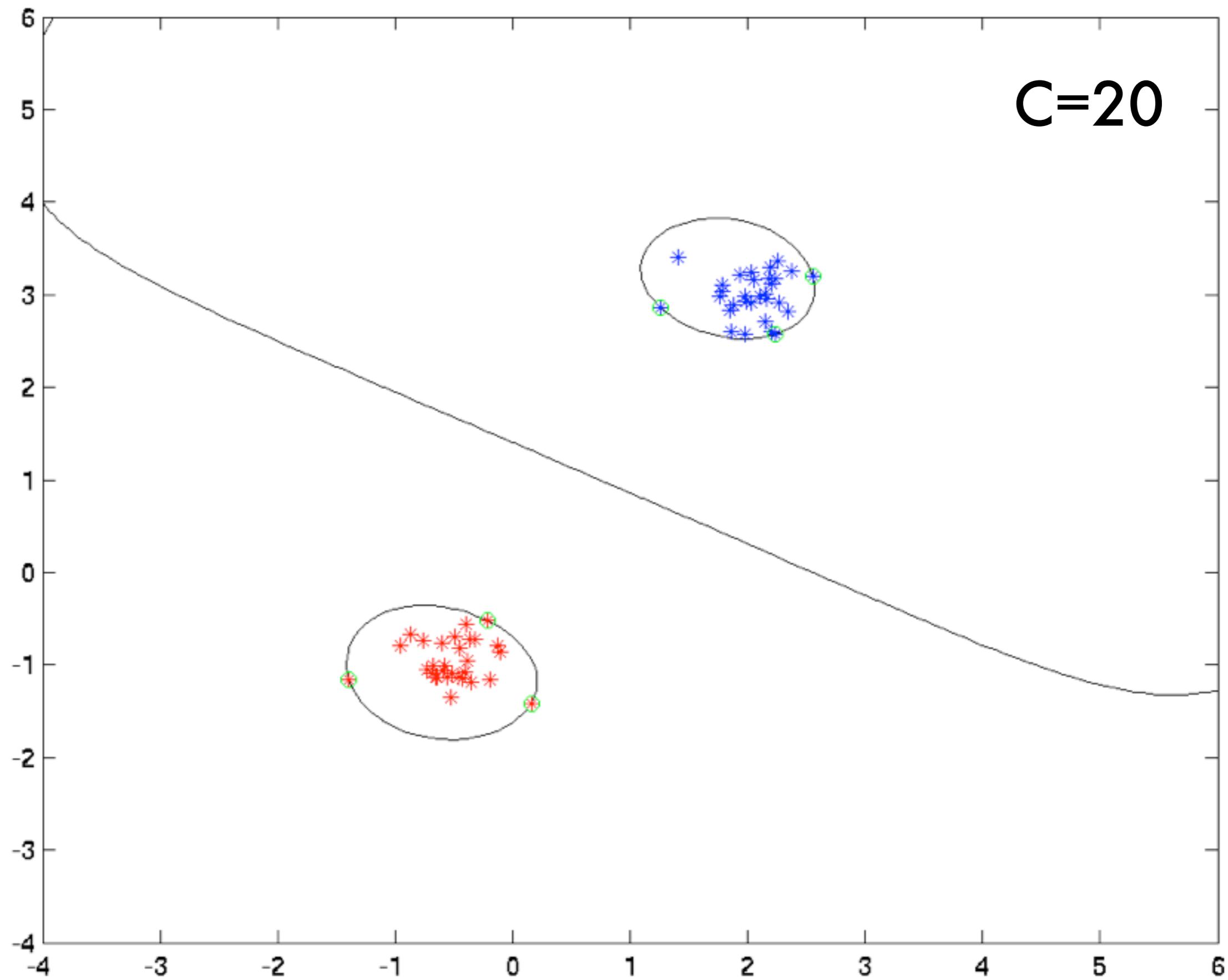
C=5



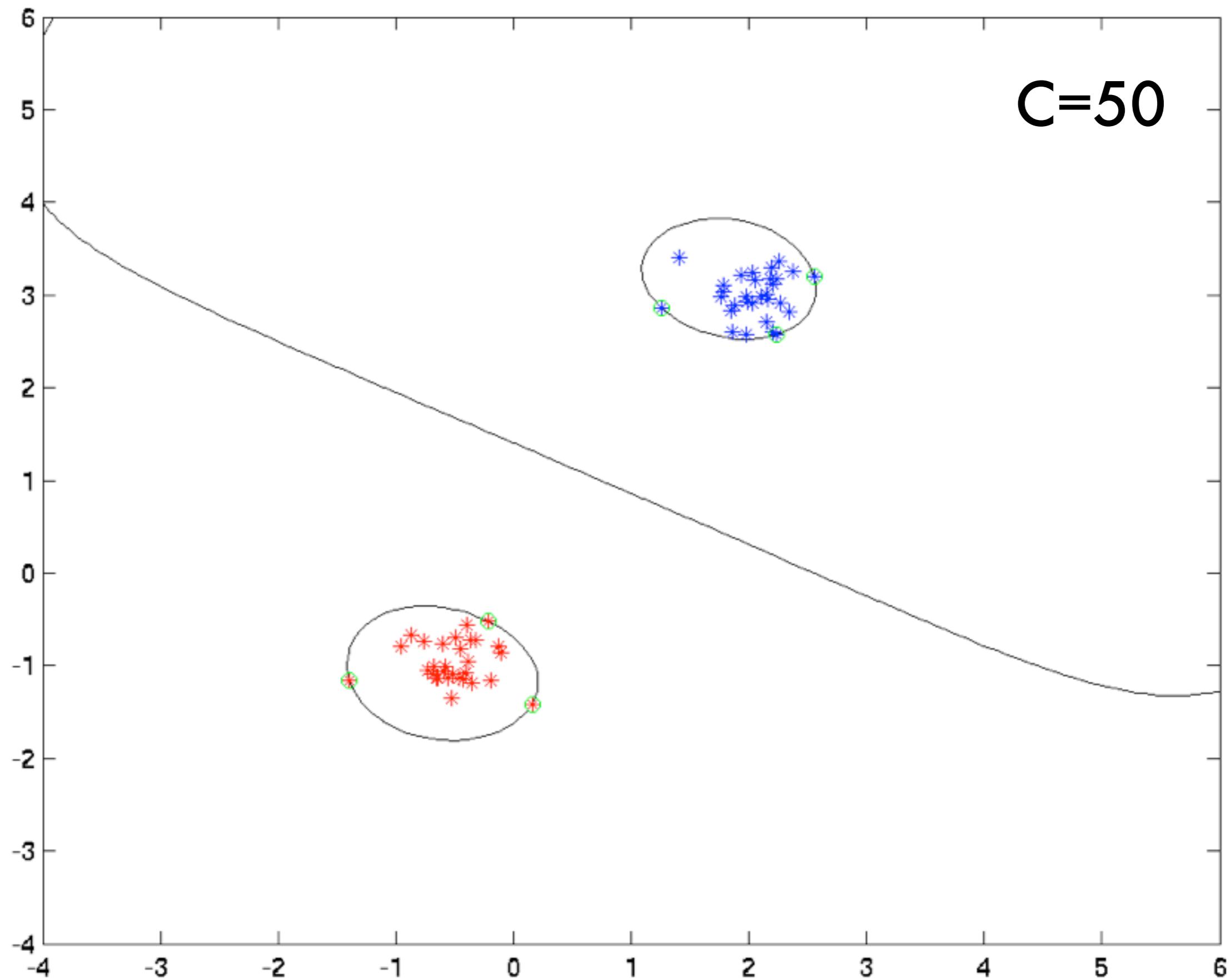
C=10



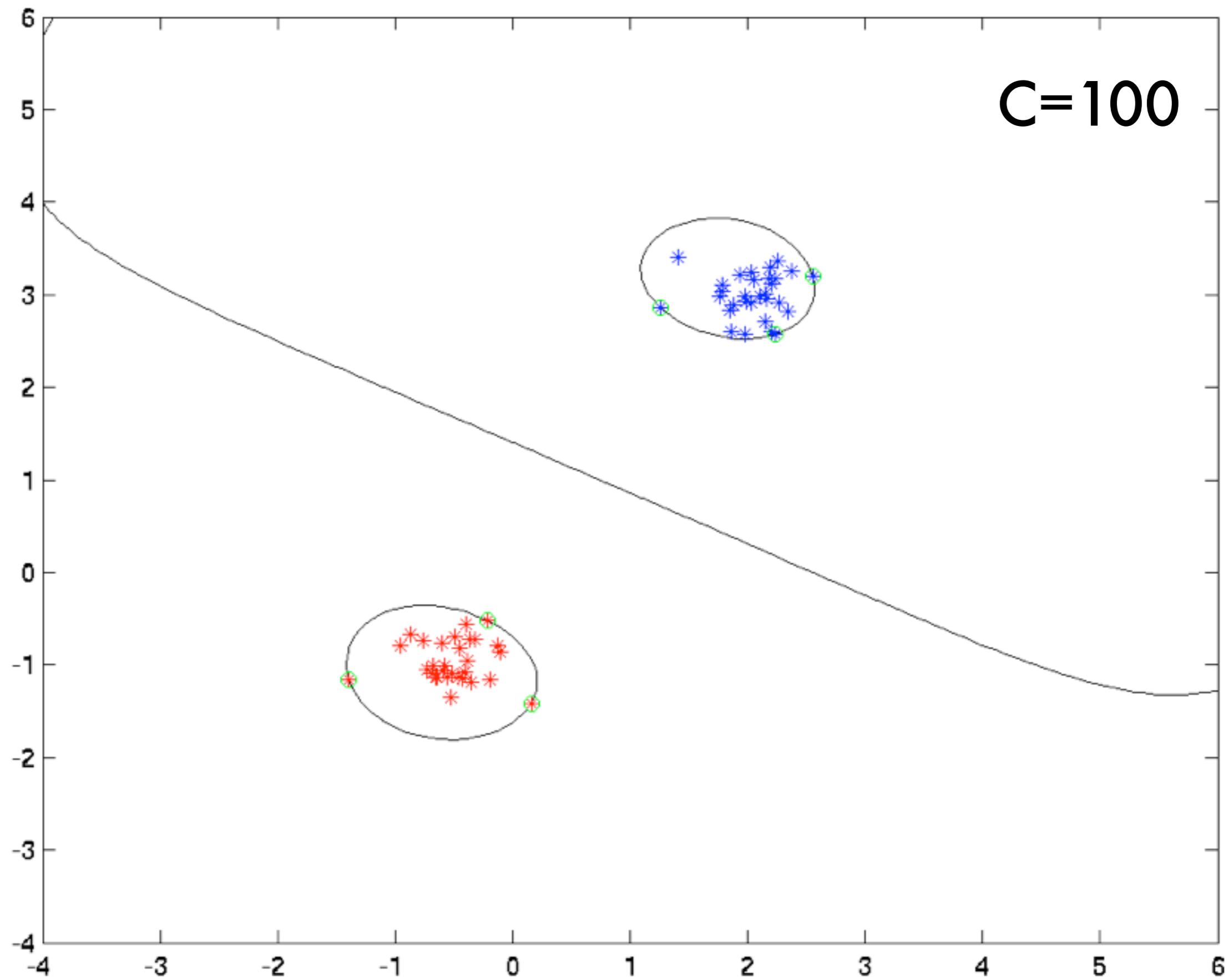
C=20



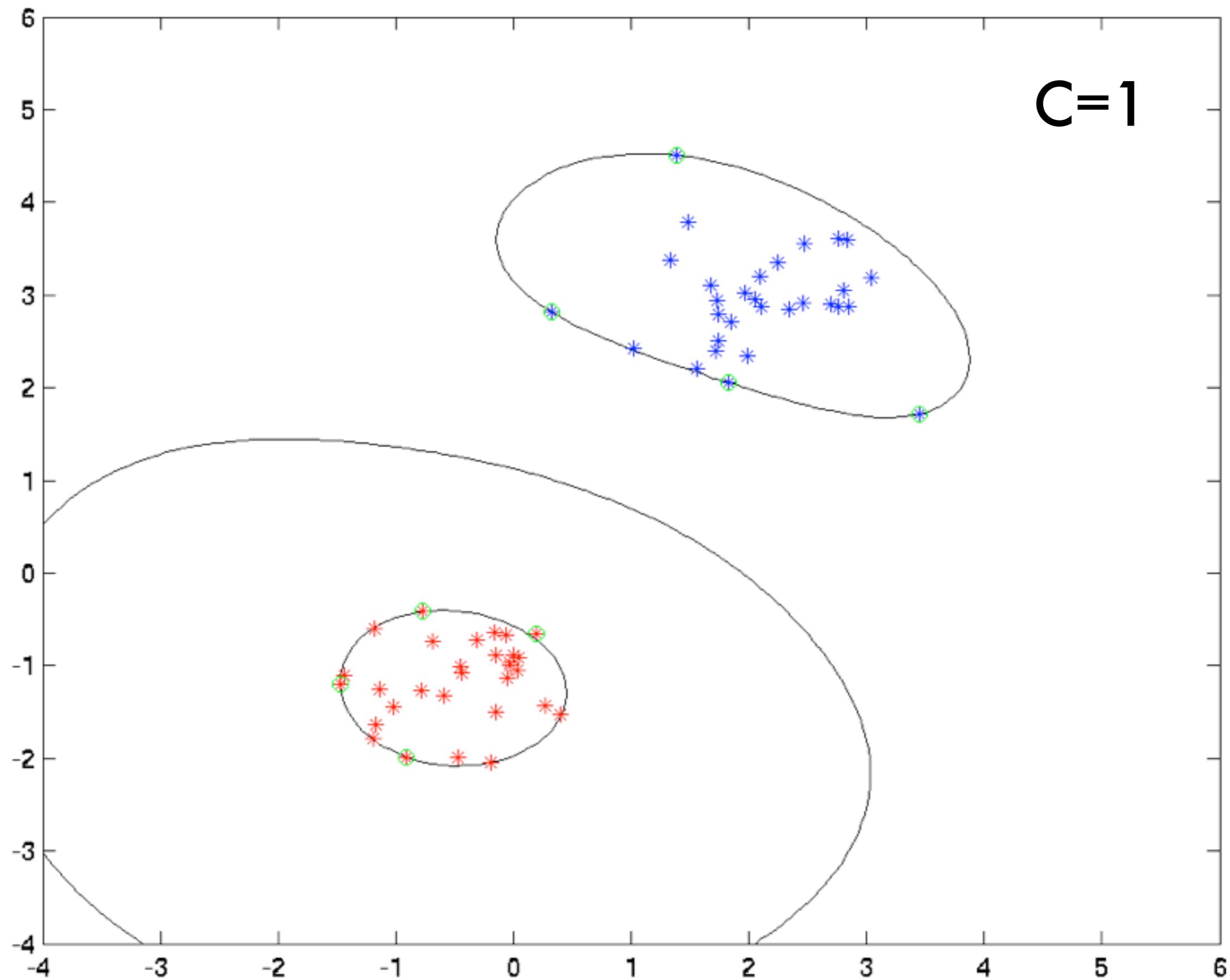
C=50



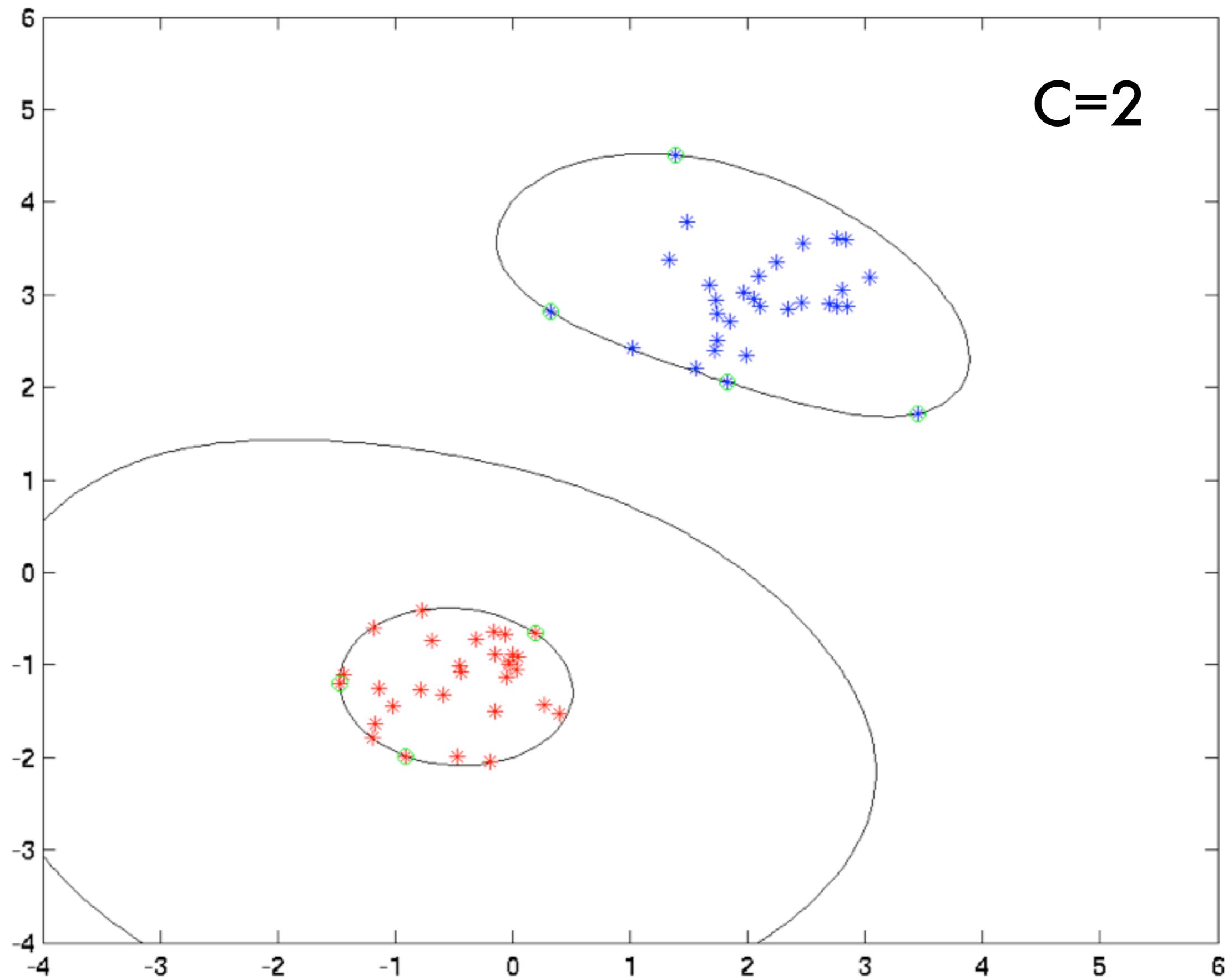
C=100



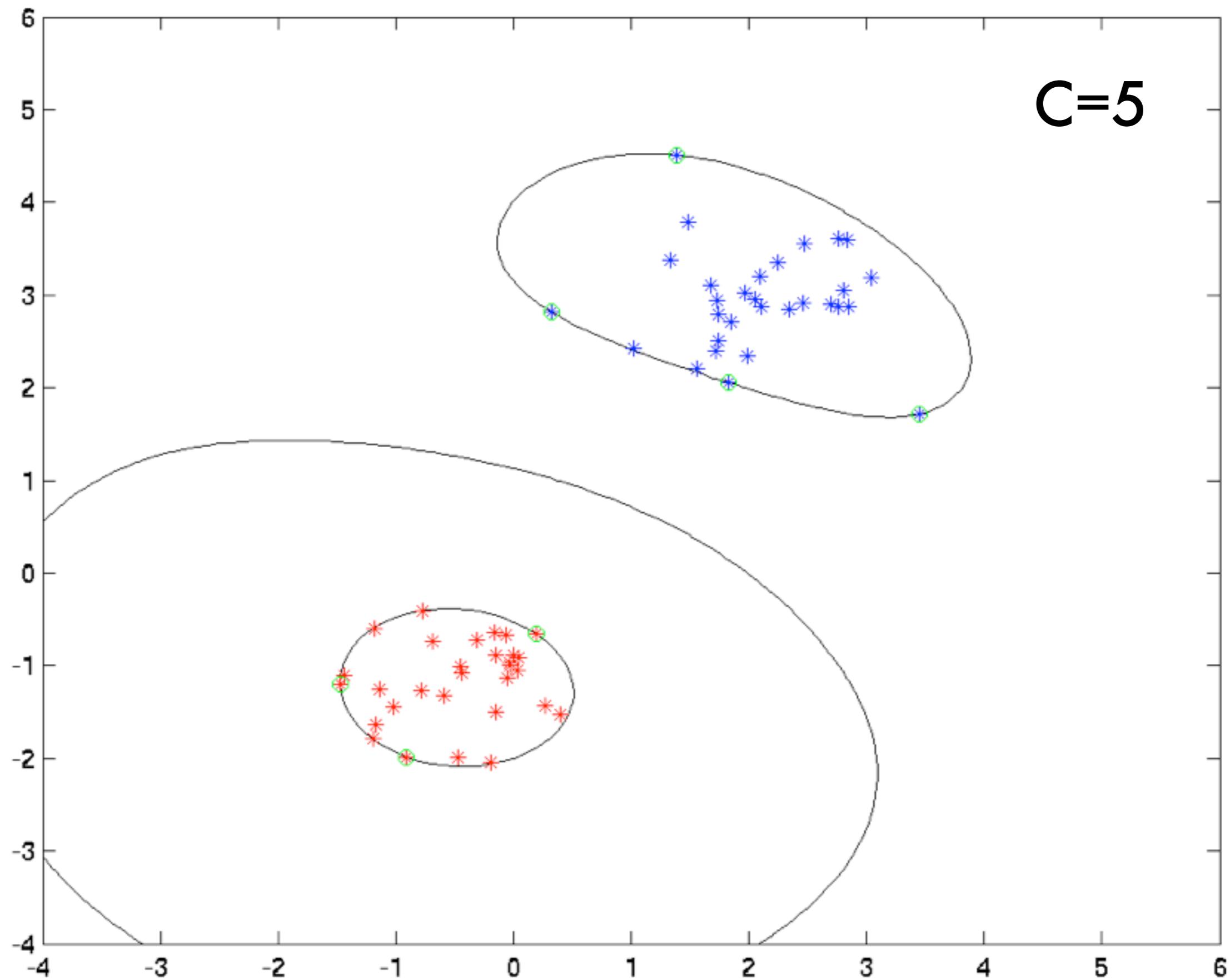
C=1



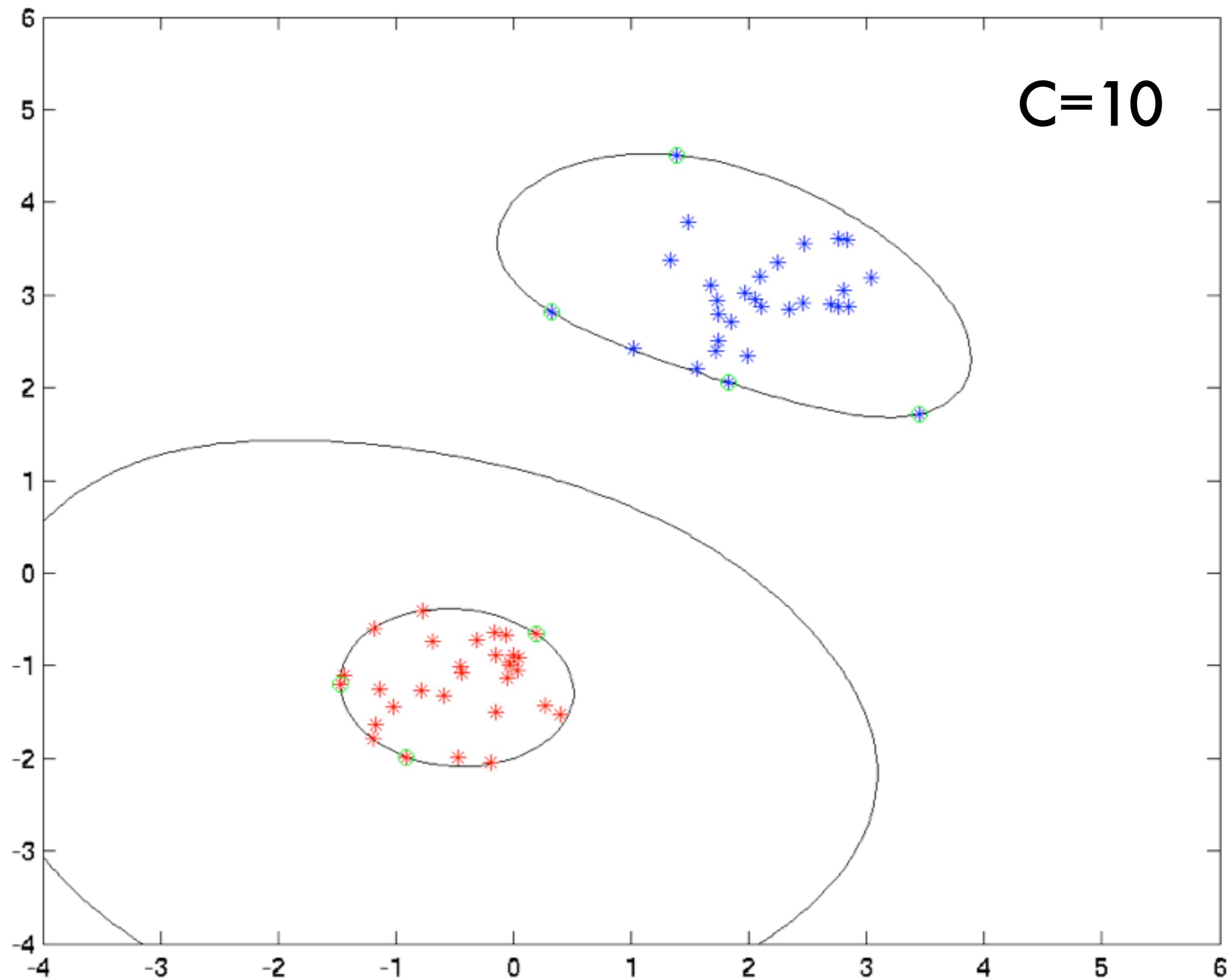
C=2



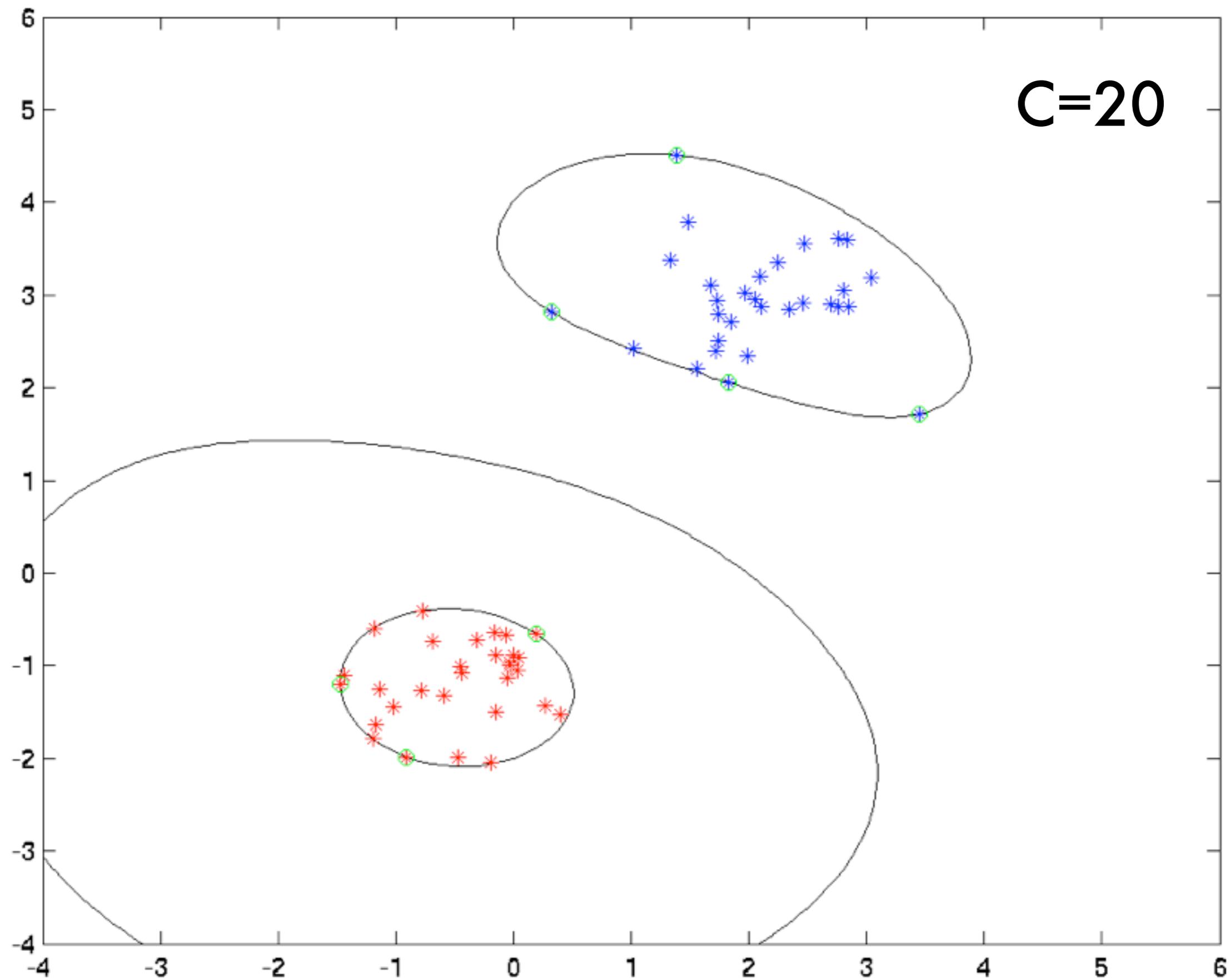
C=5



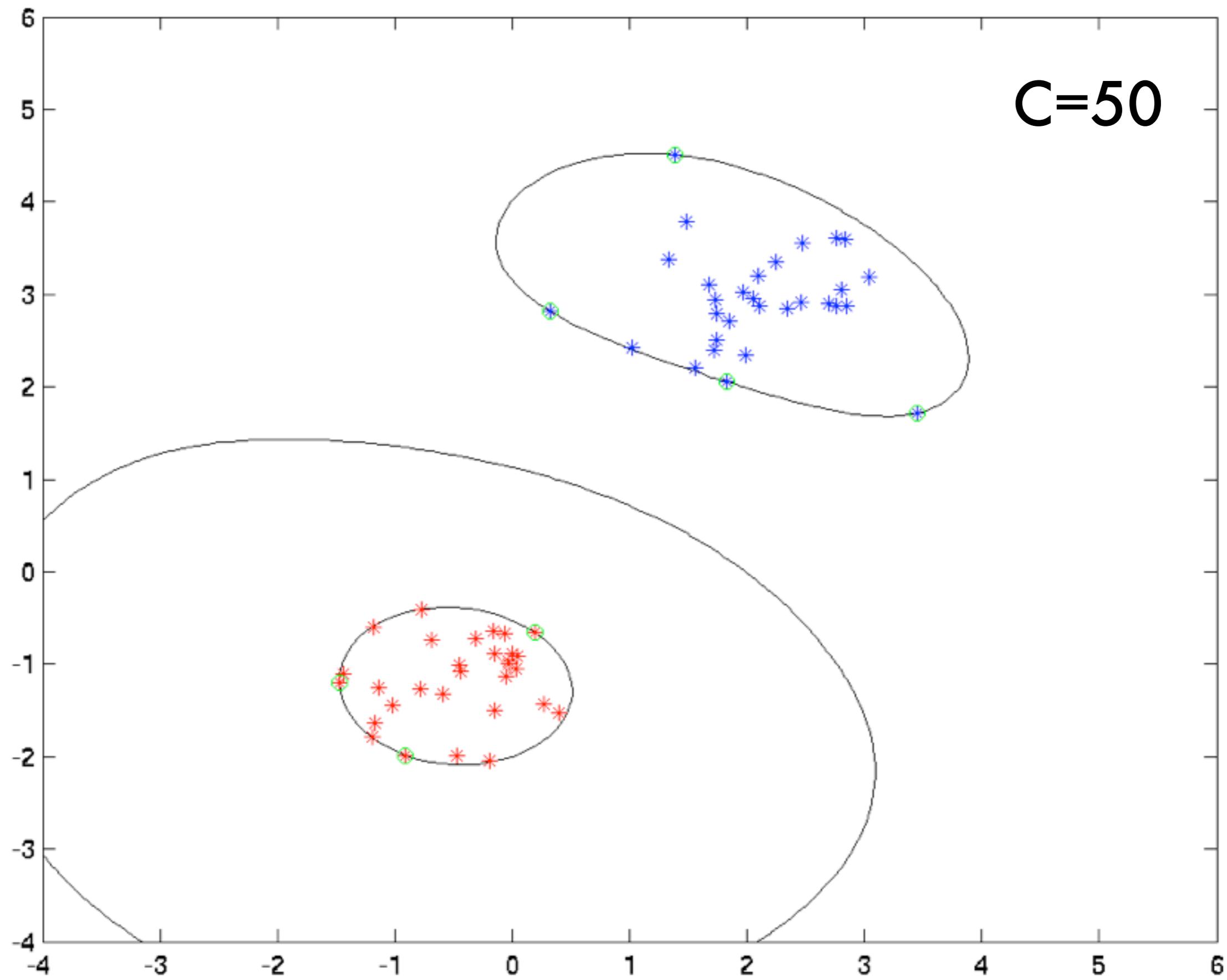
C=10



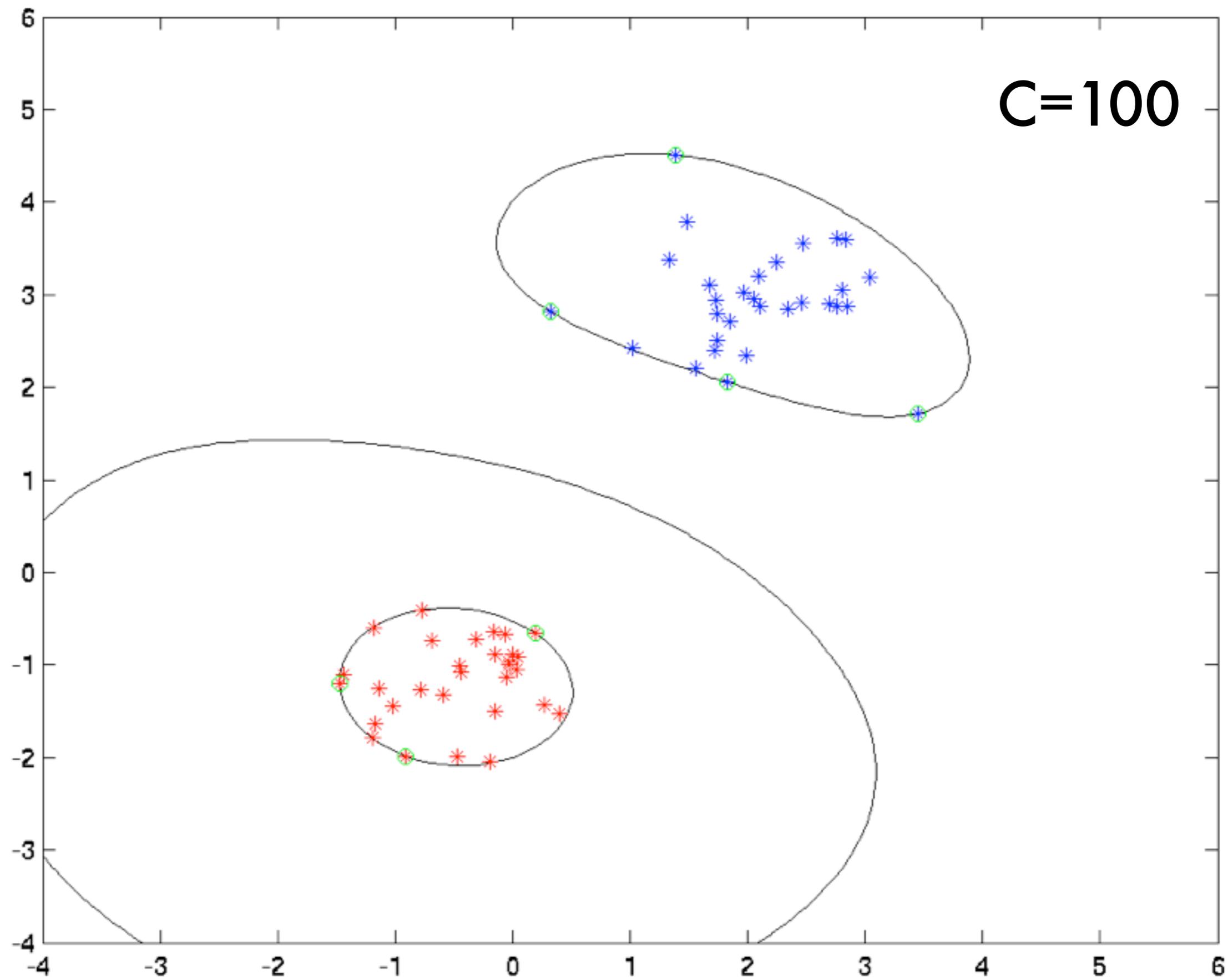
C=20



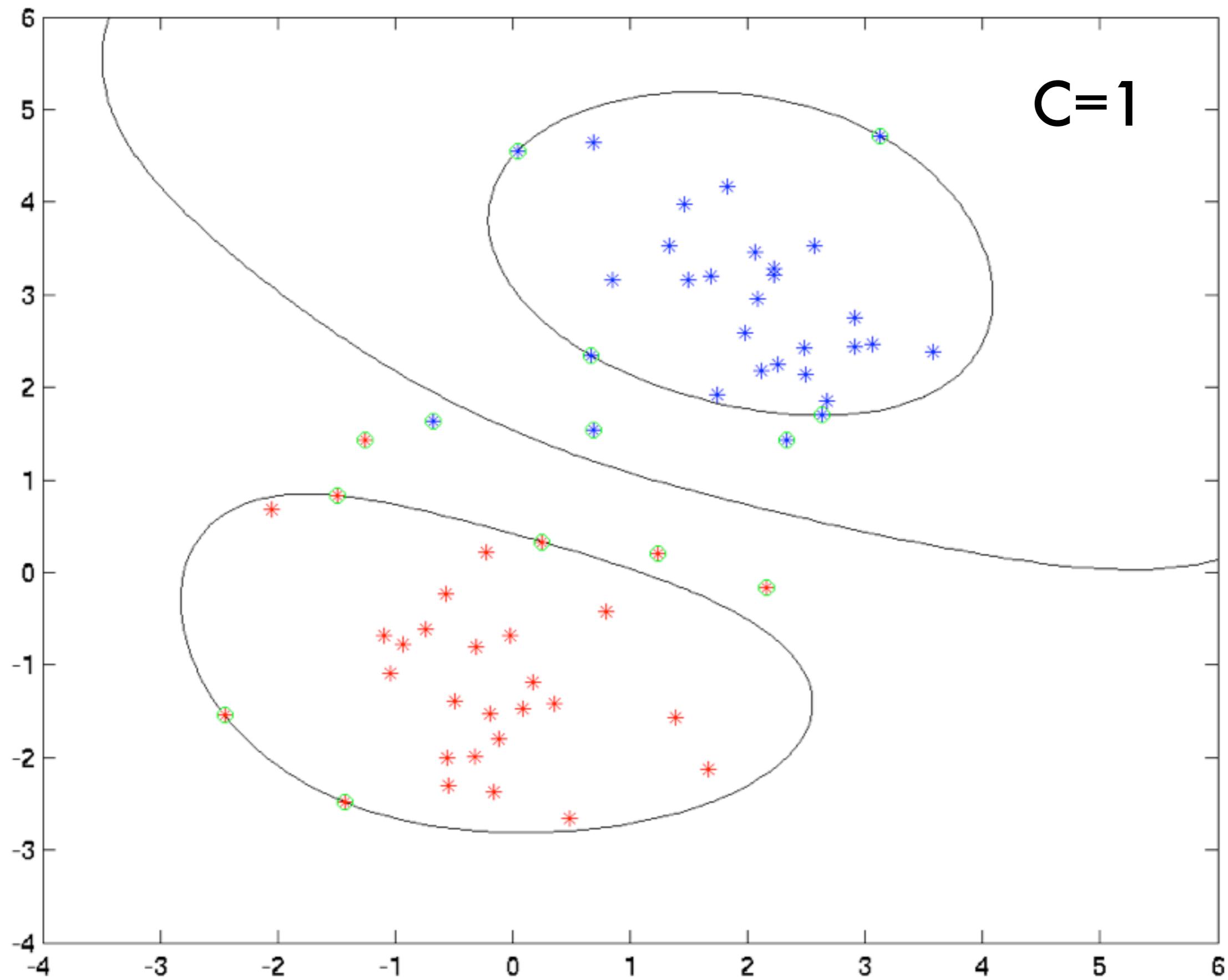
C=50



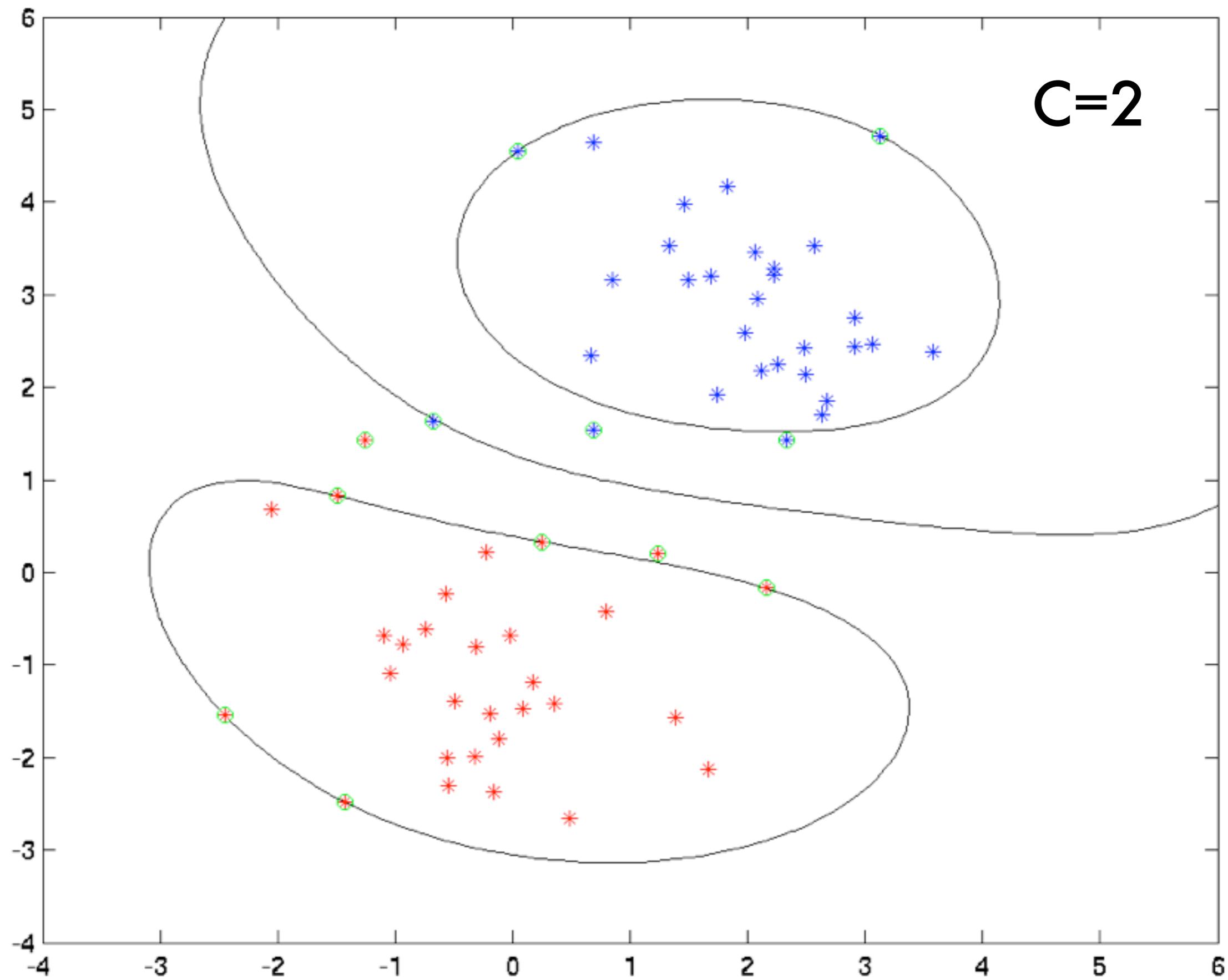
C=100



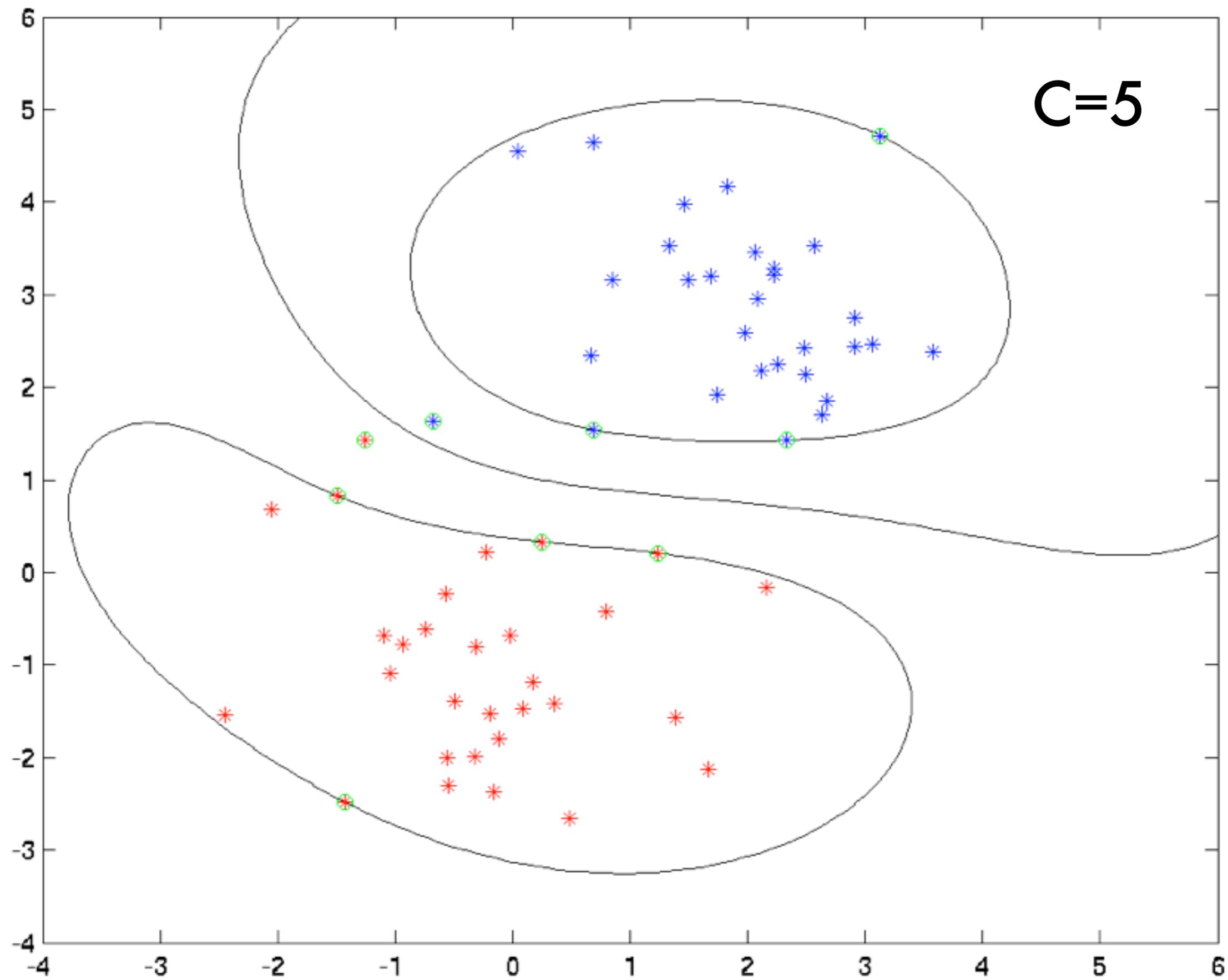
C=1



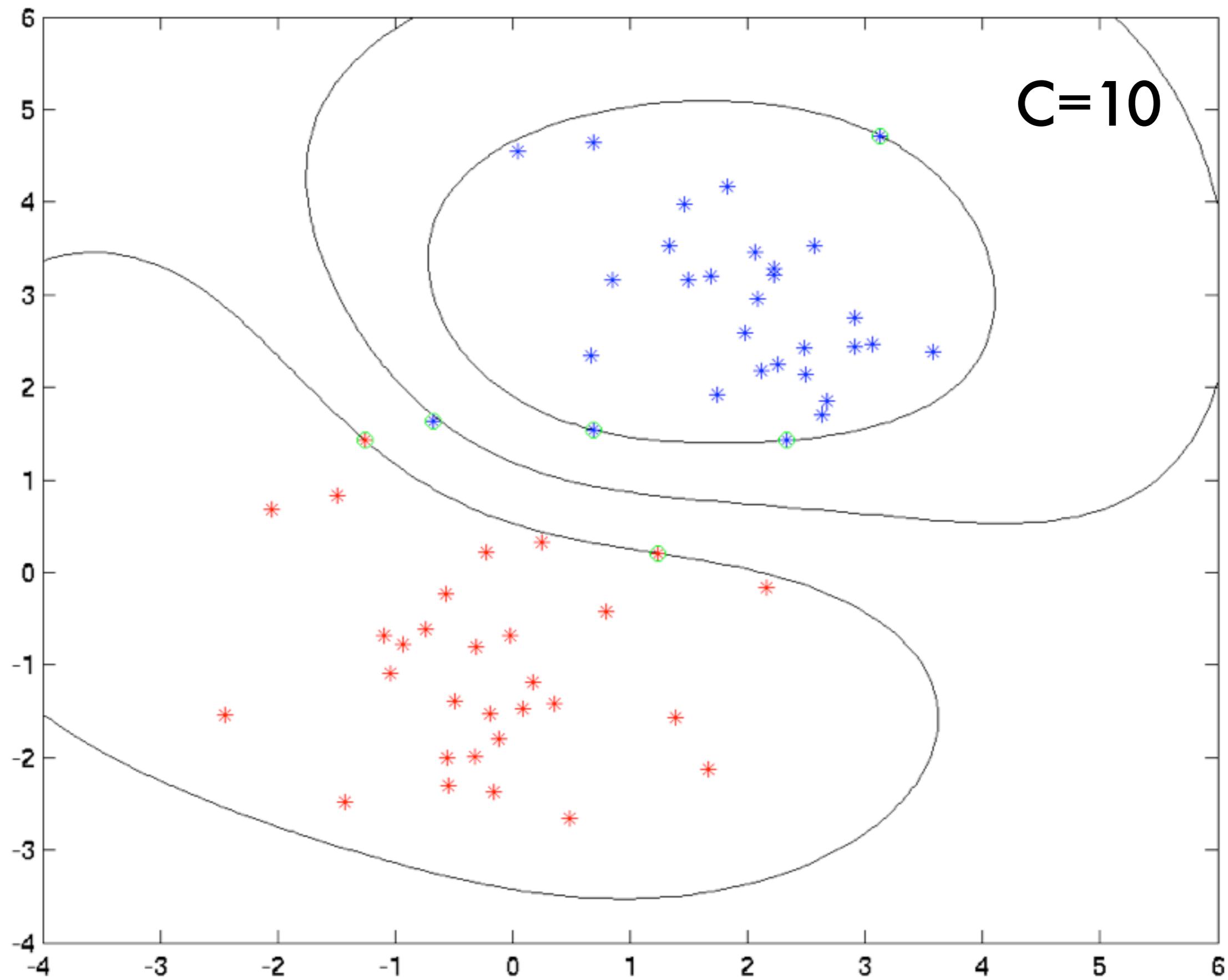
C=2

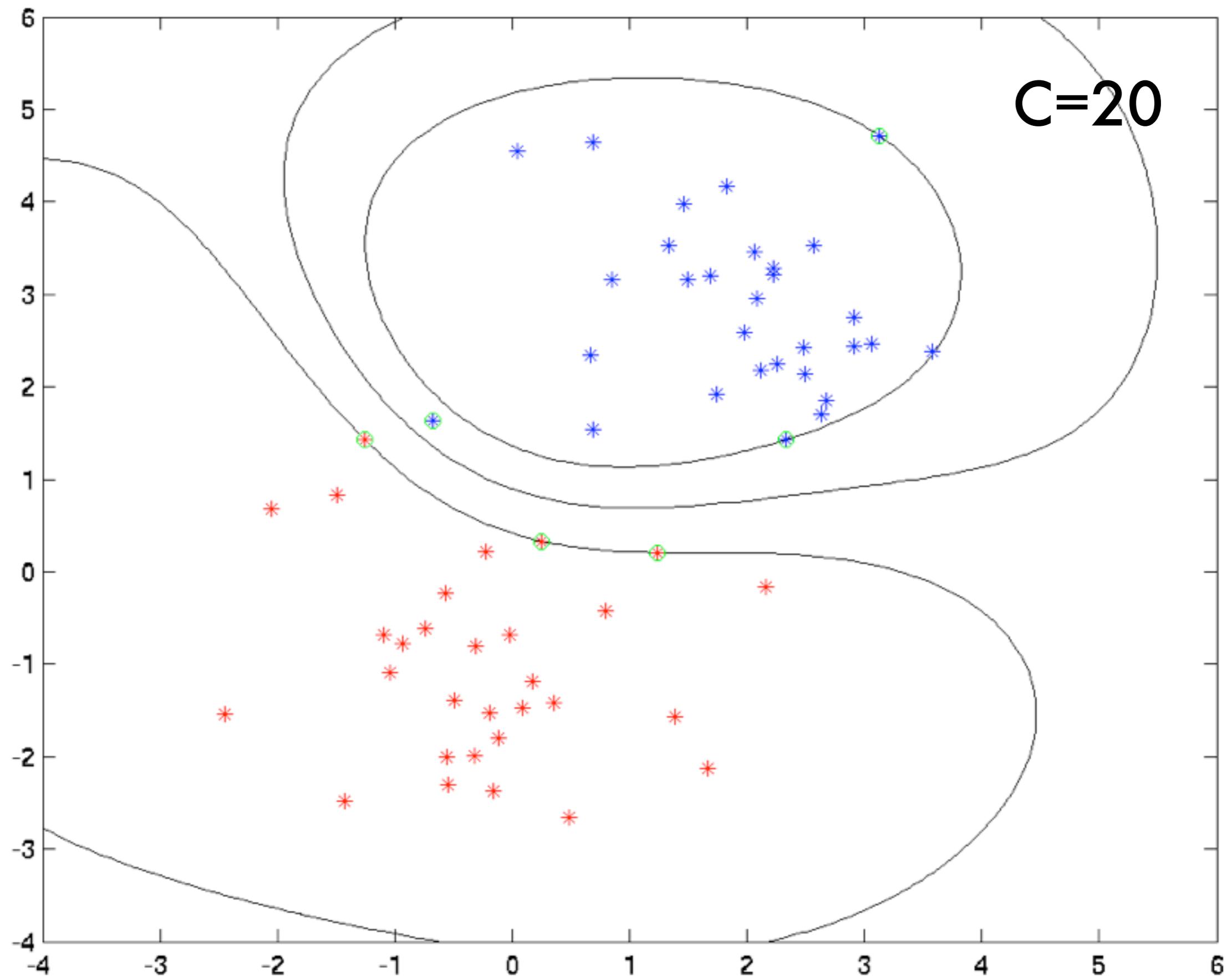


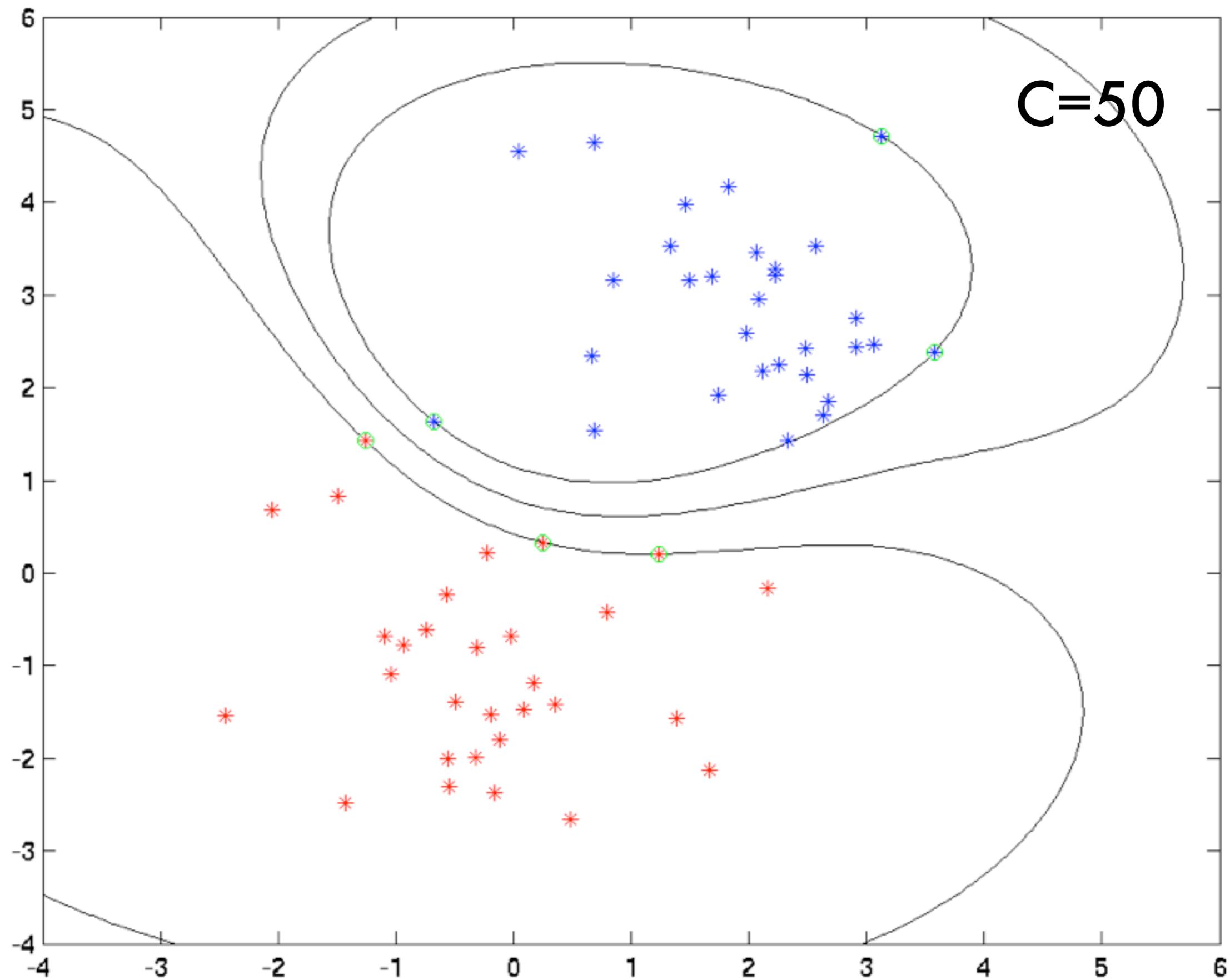
C=5

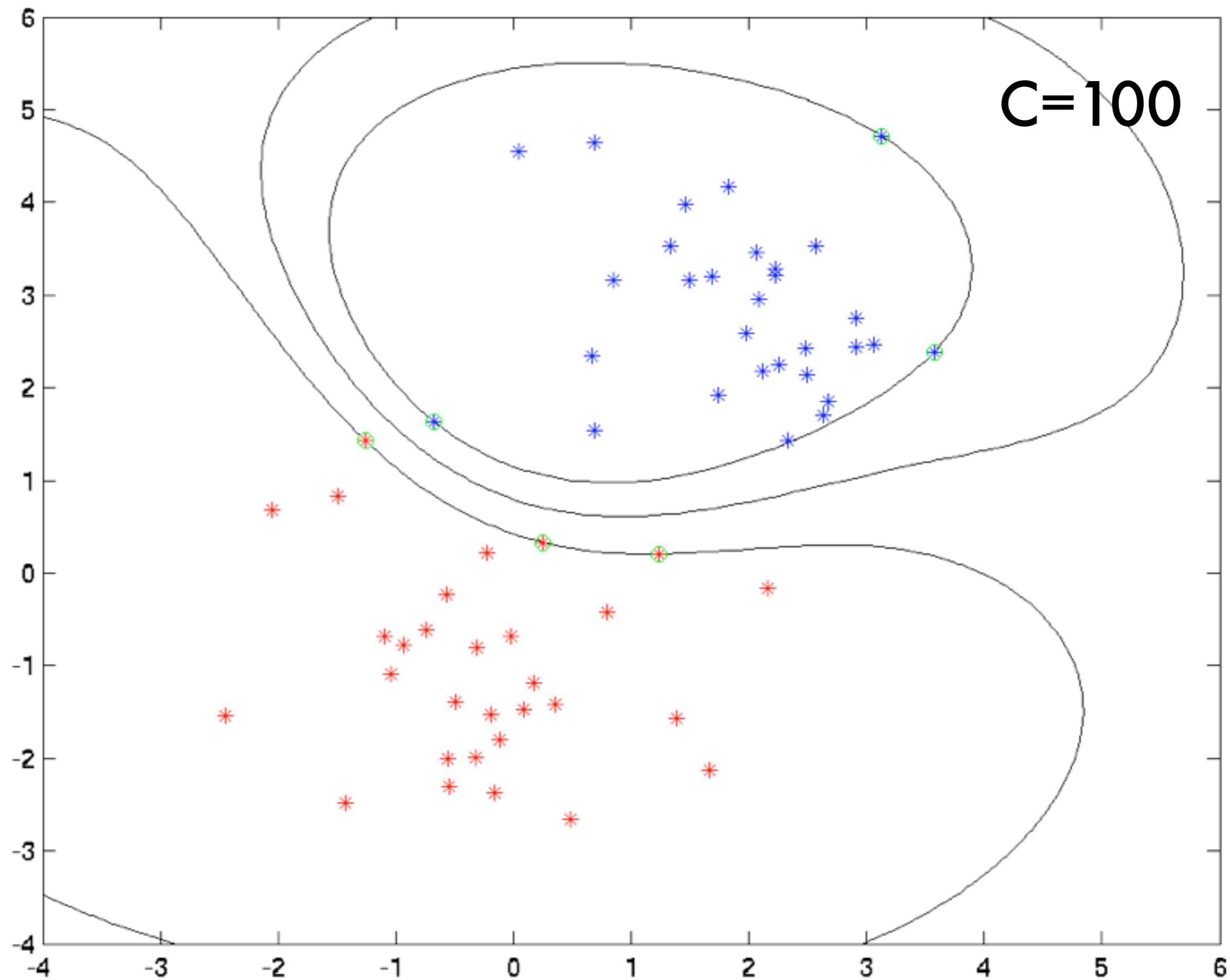


C=10

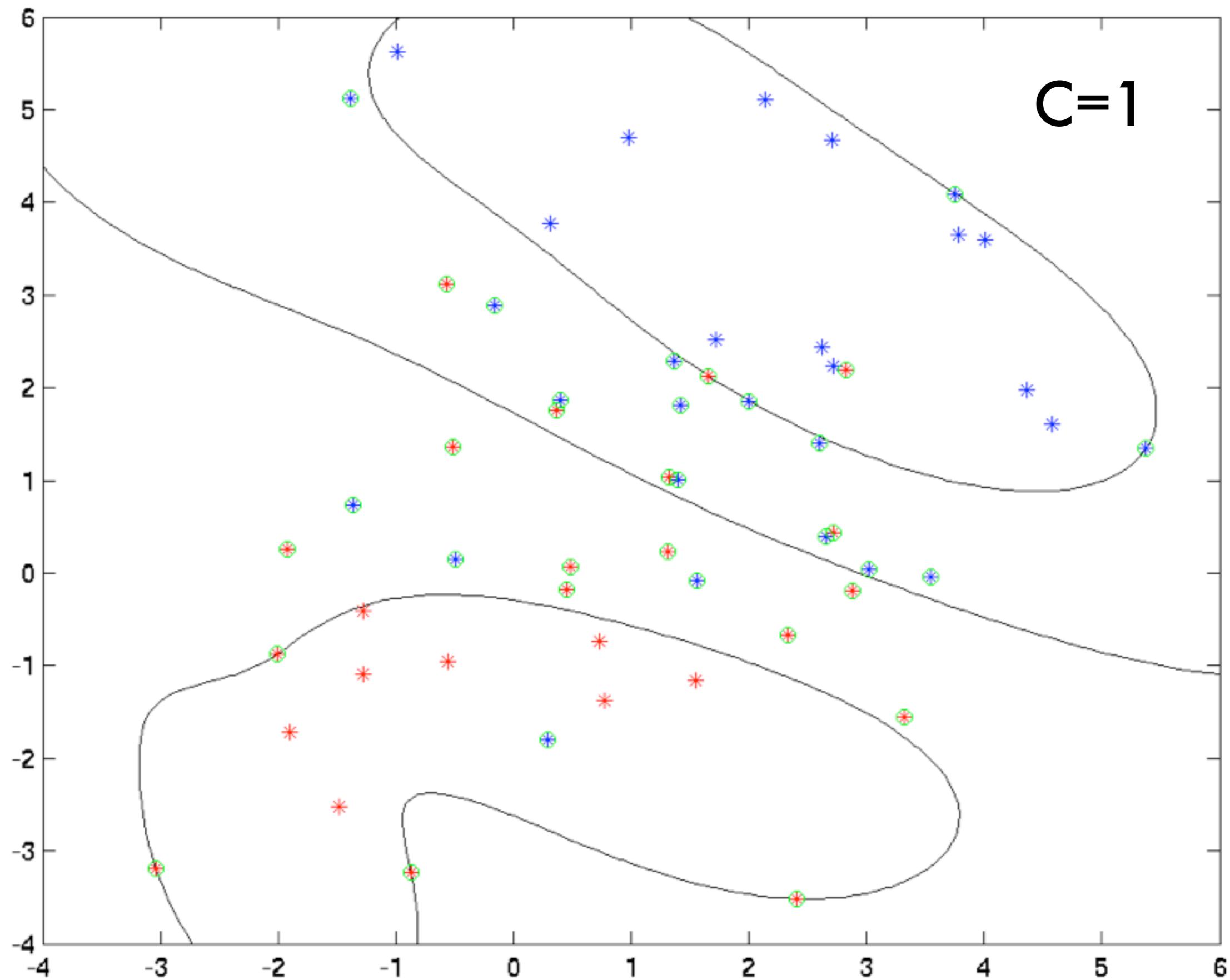




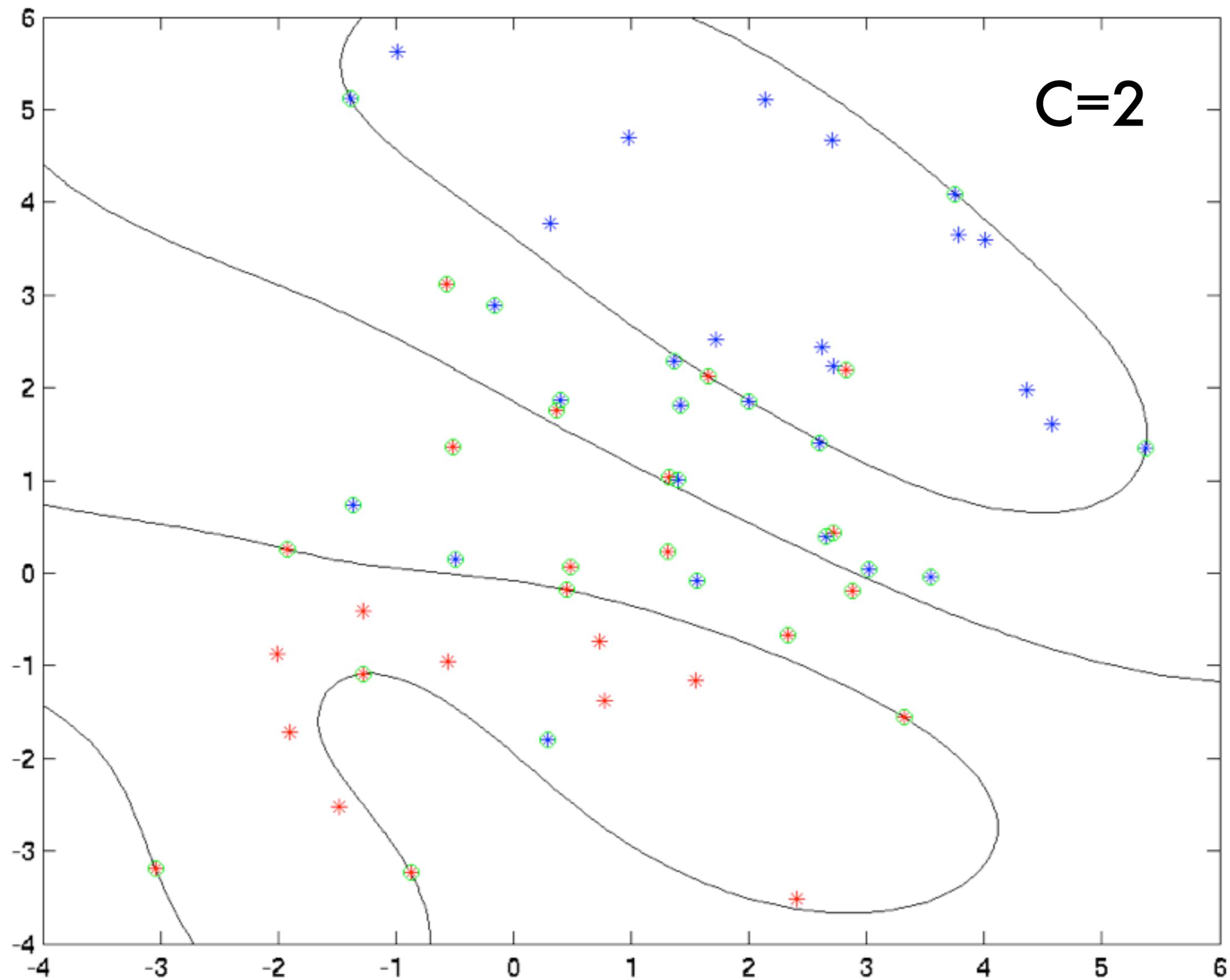




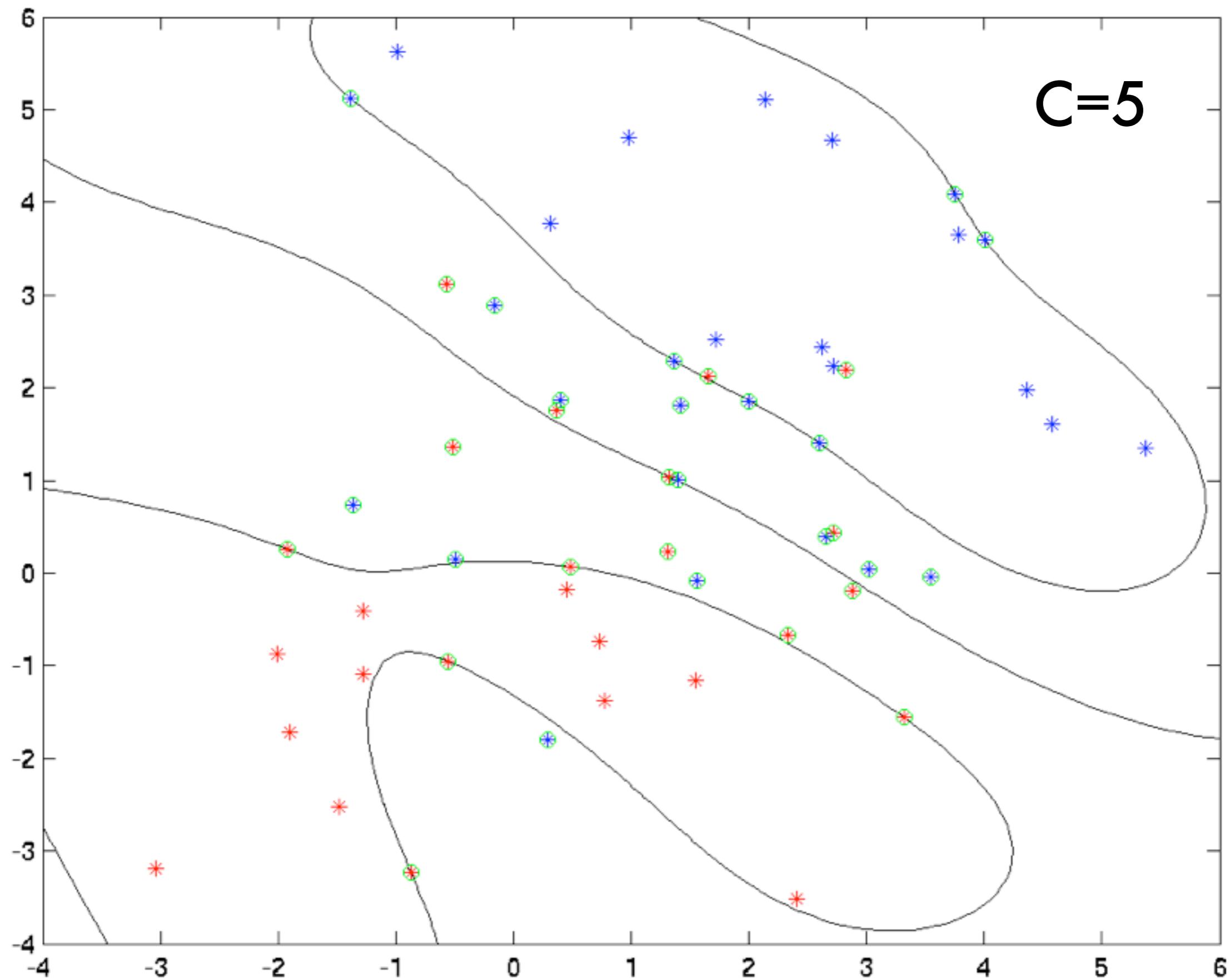
C=1



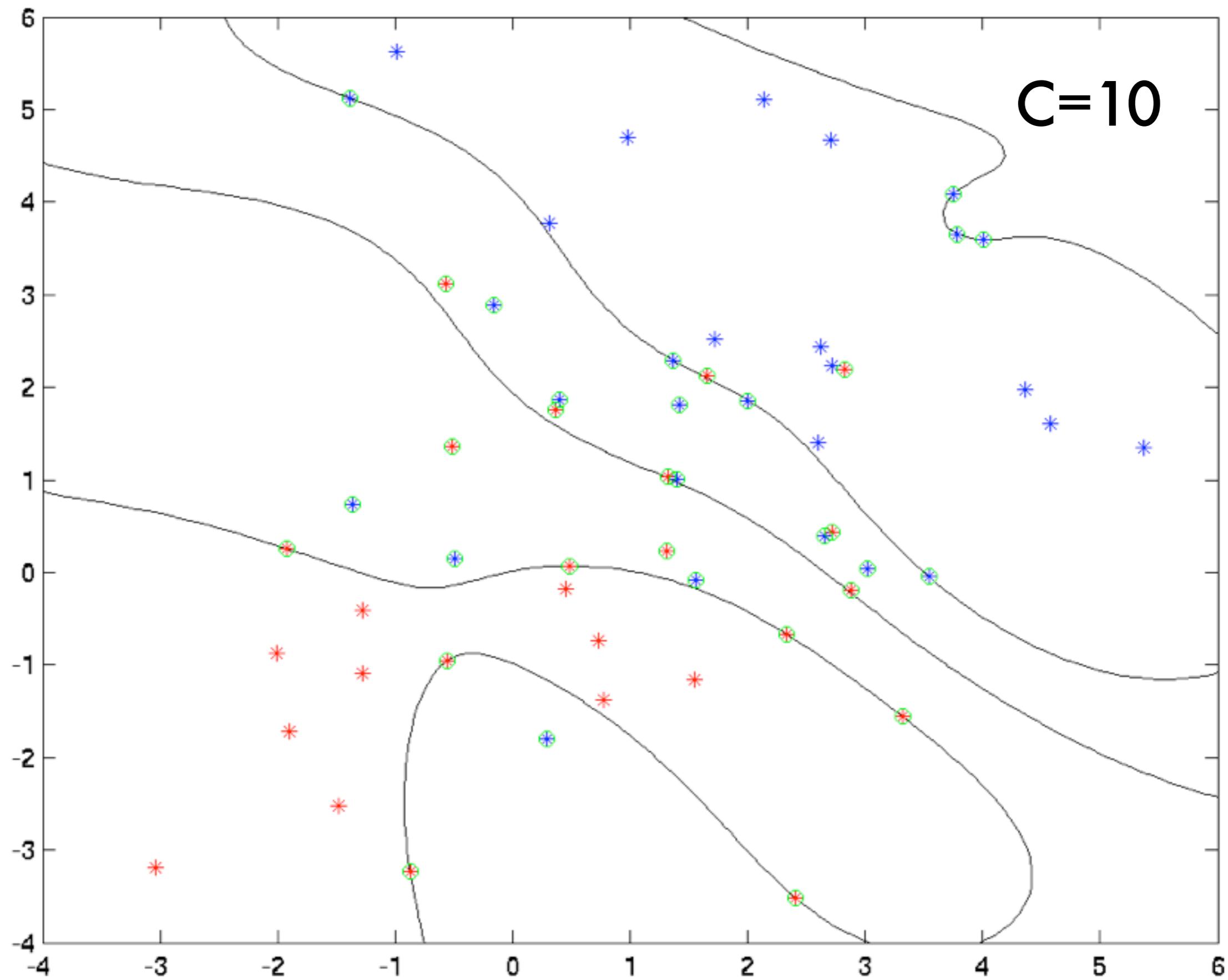
C=2



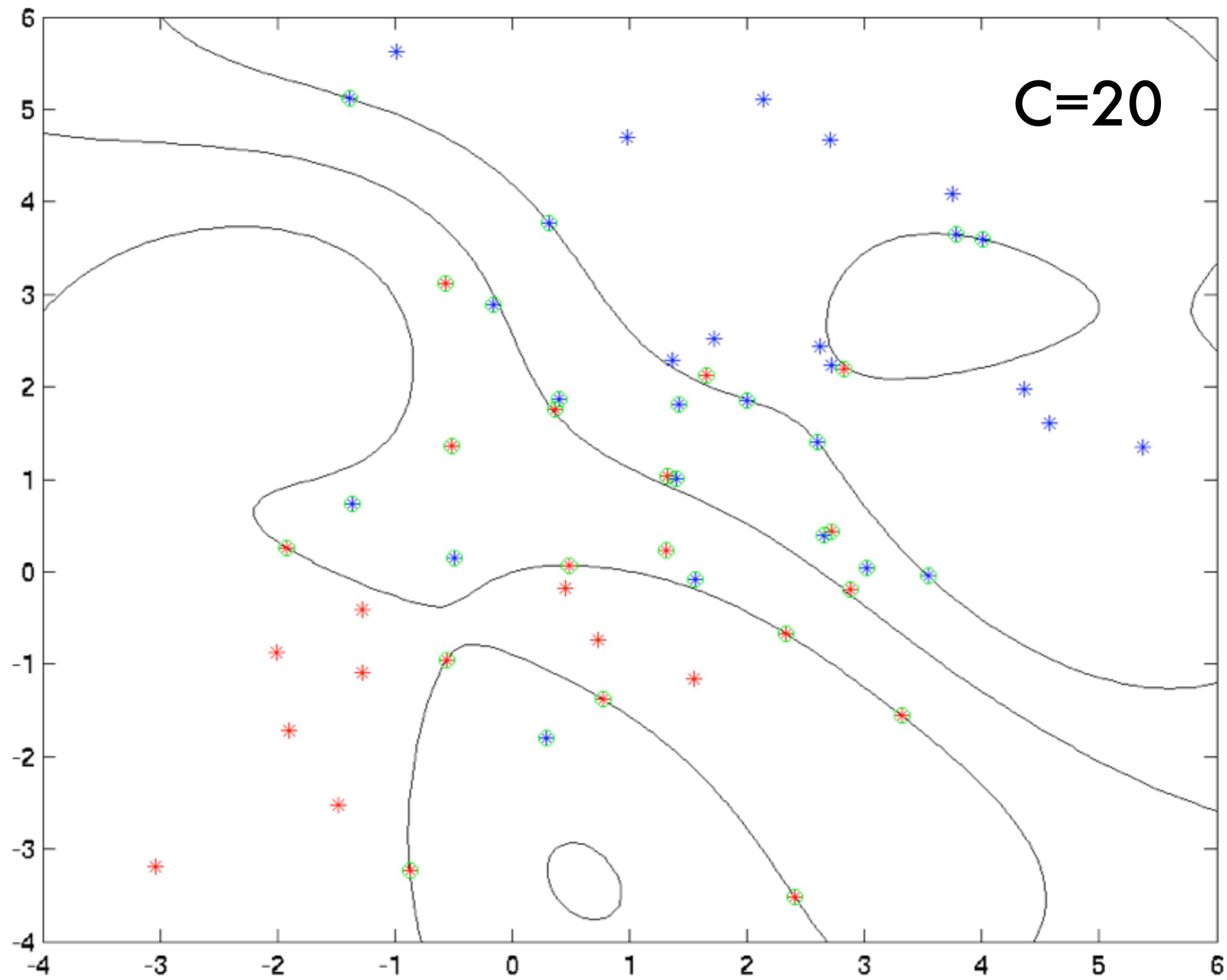
C=5



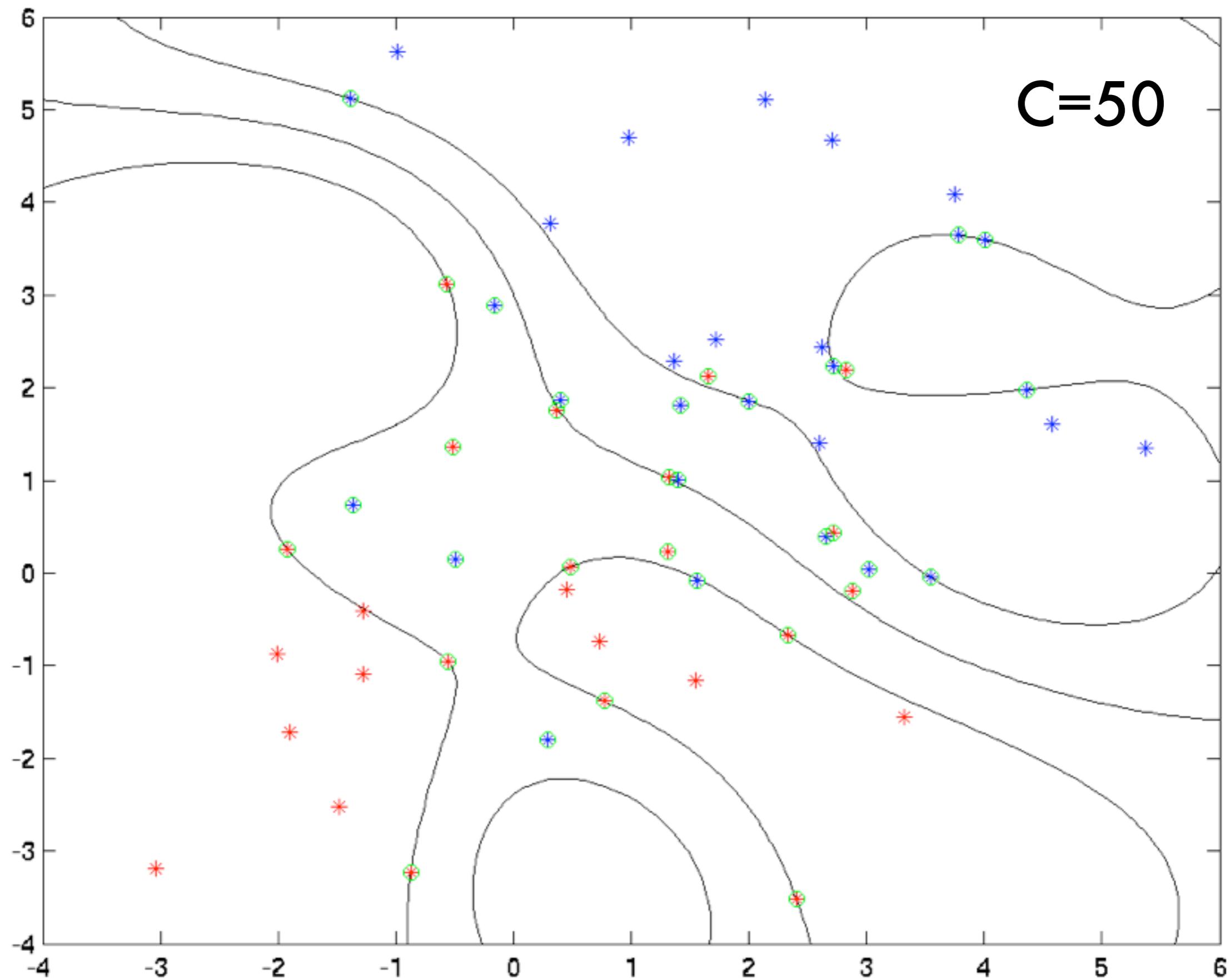
C=10



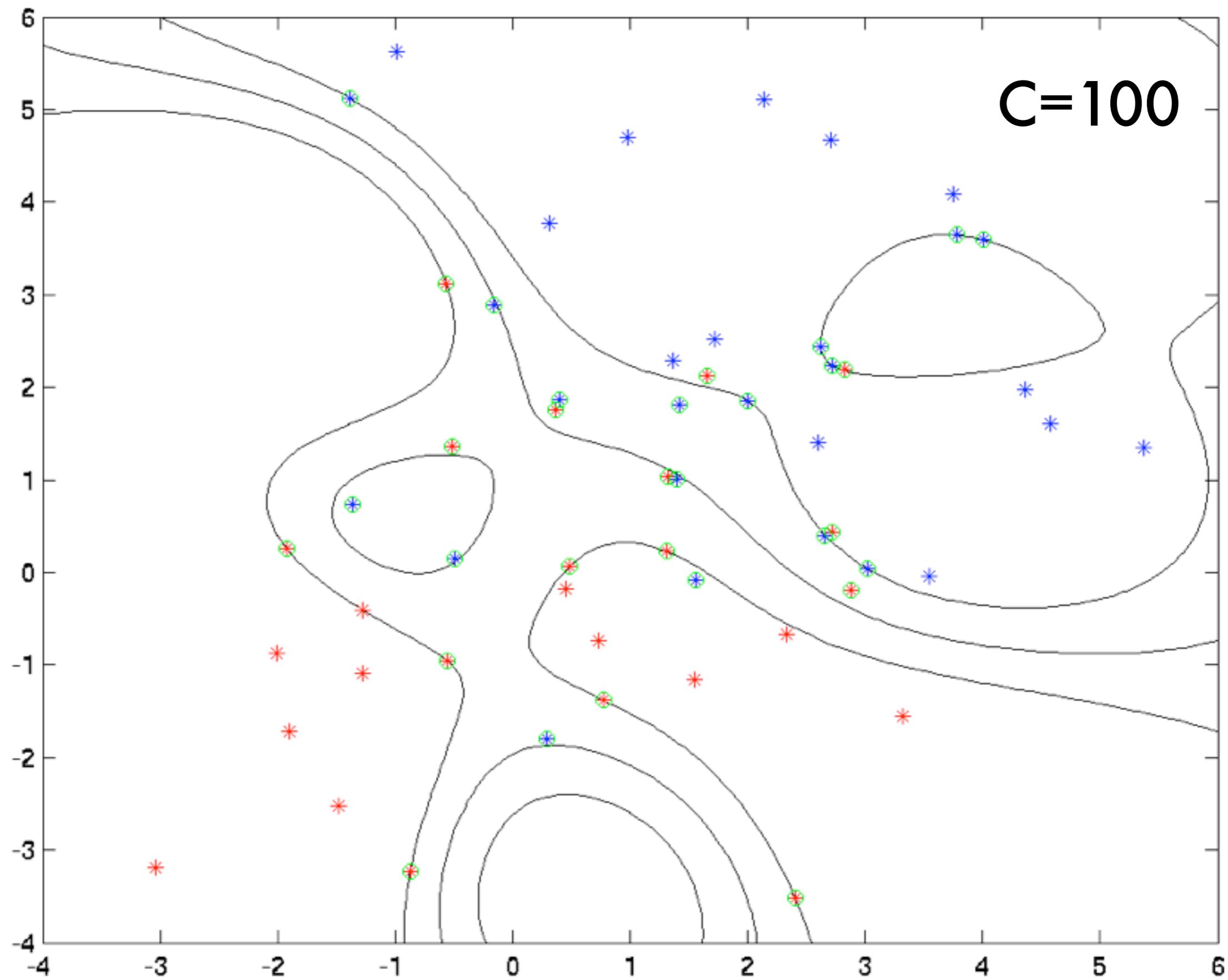
C=20



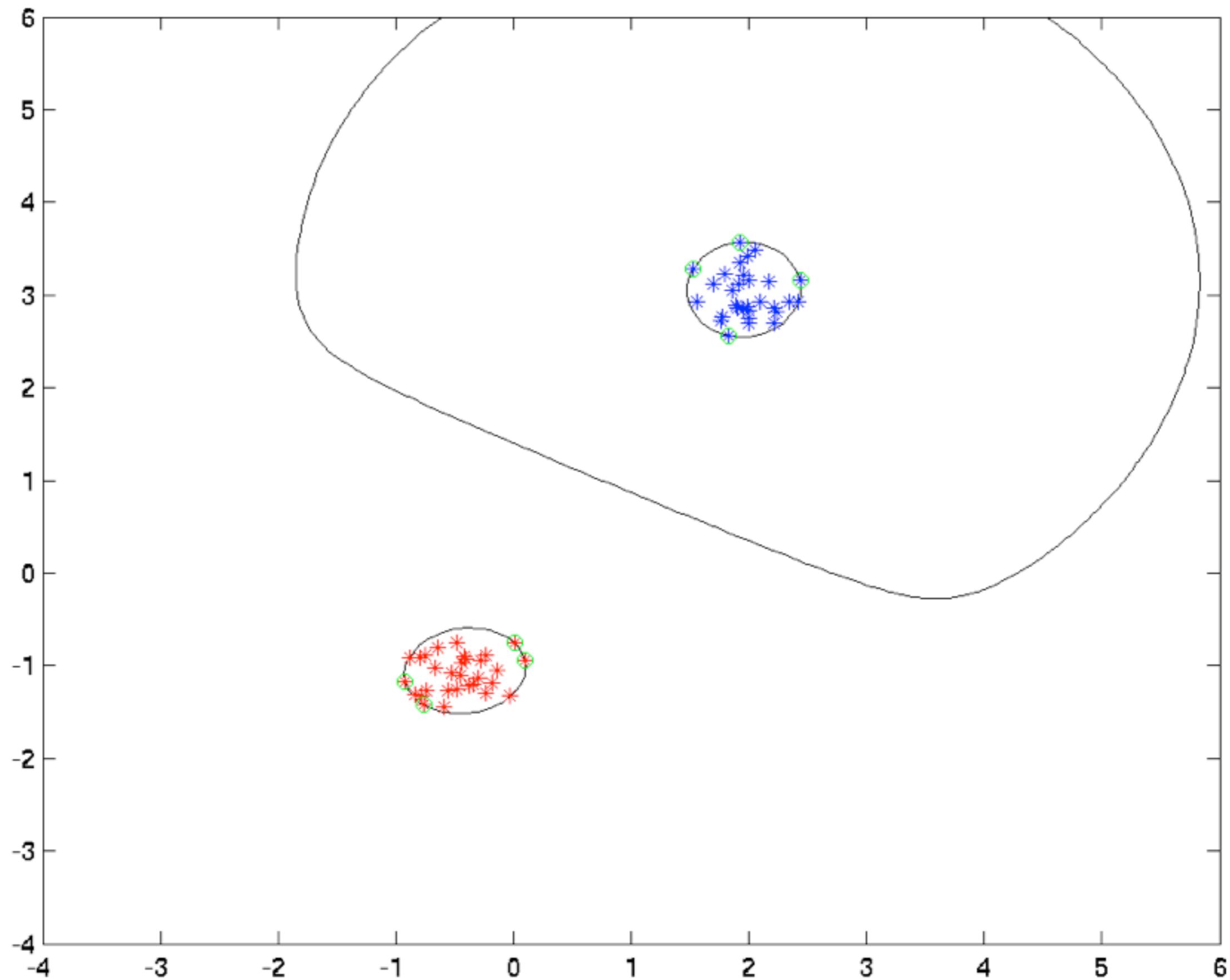
C=50

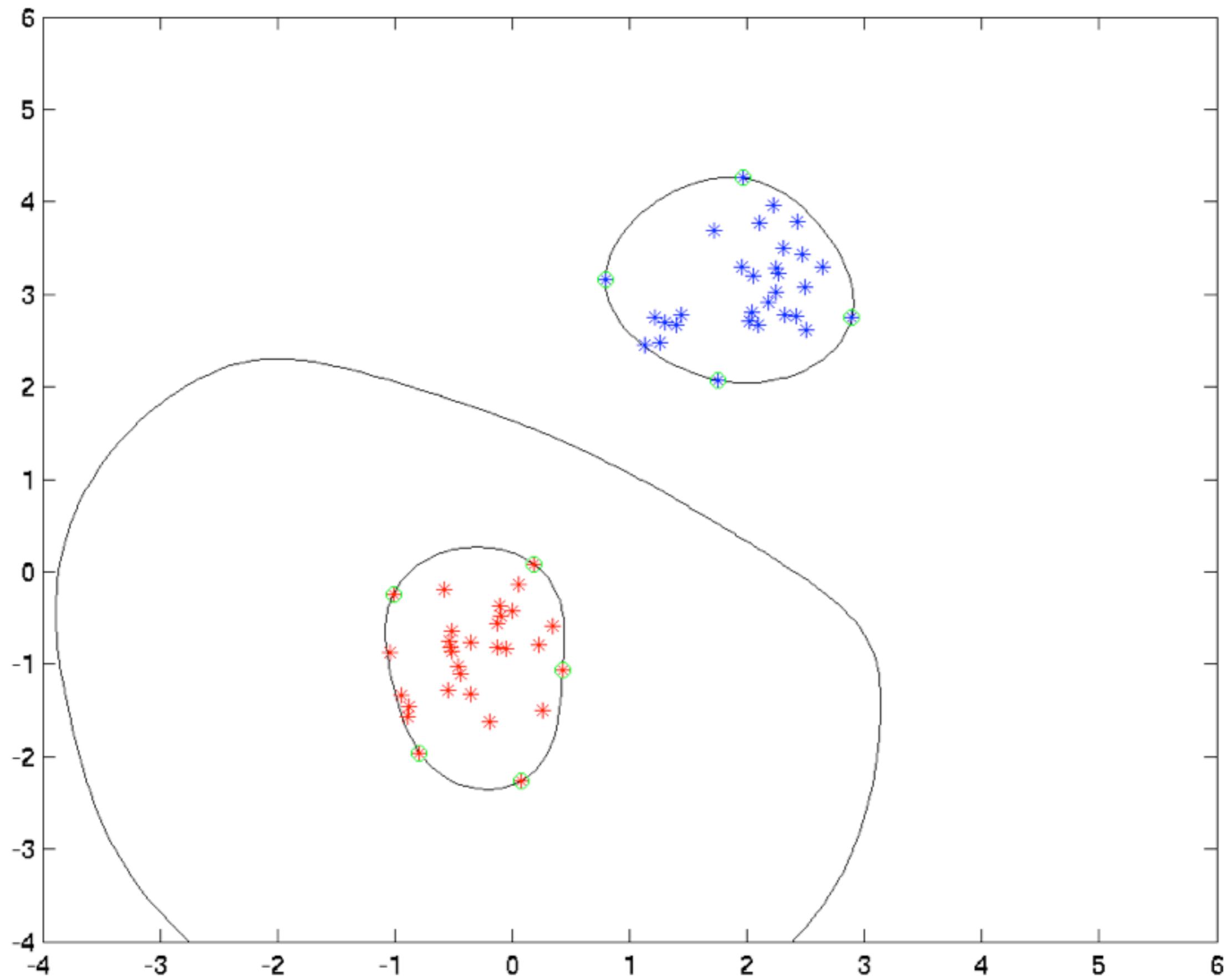


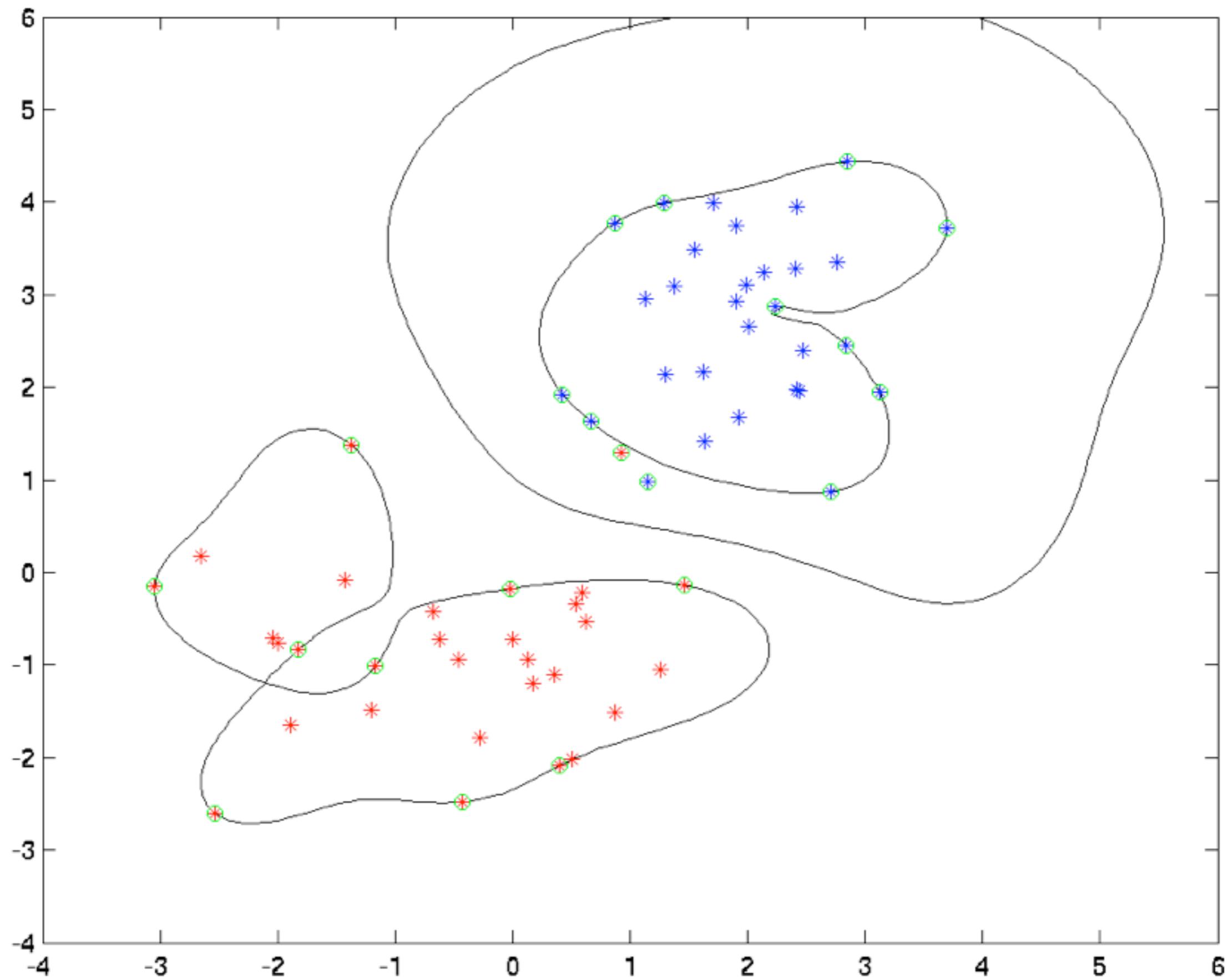
C=100

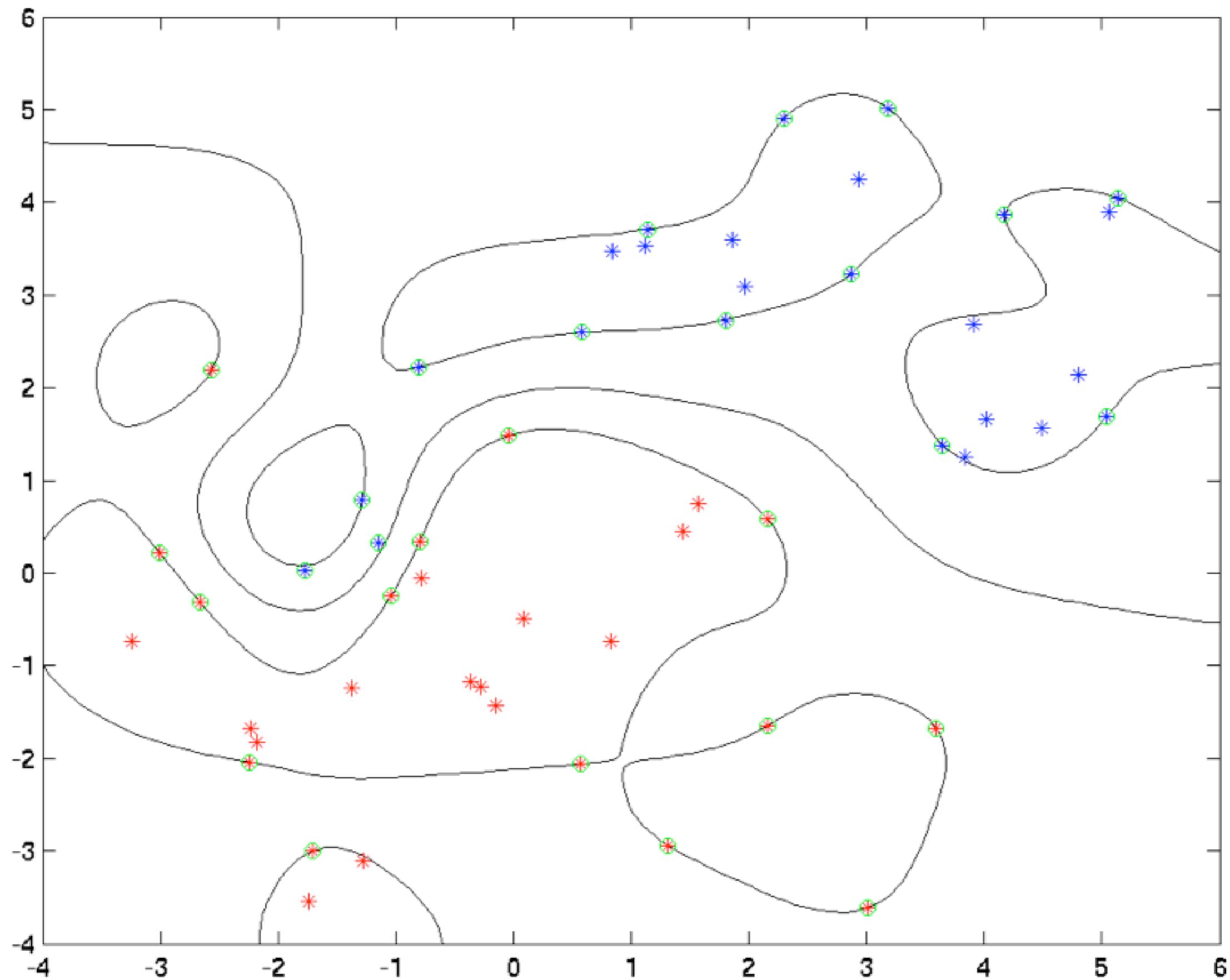


And now with a narrower kernel

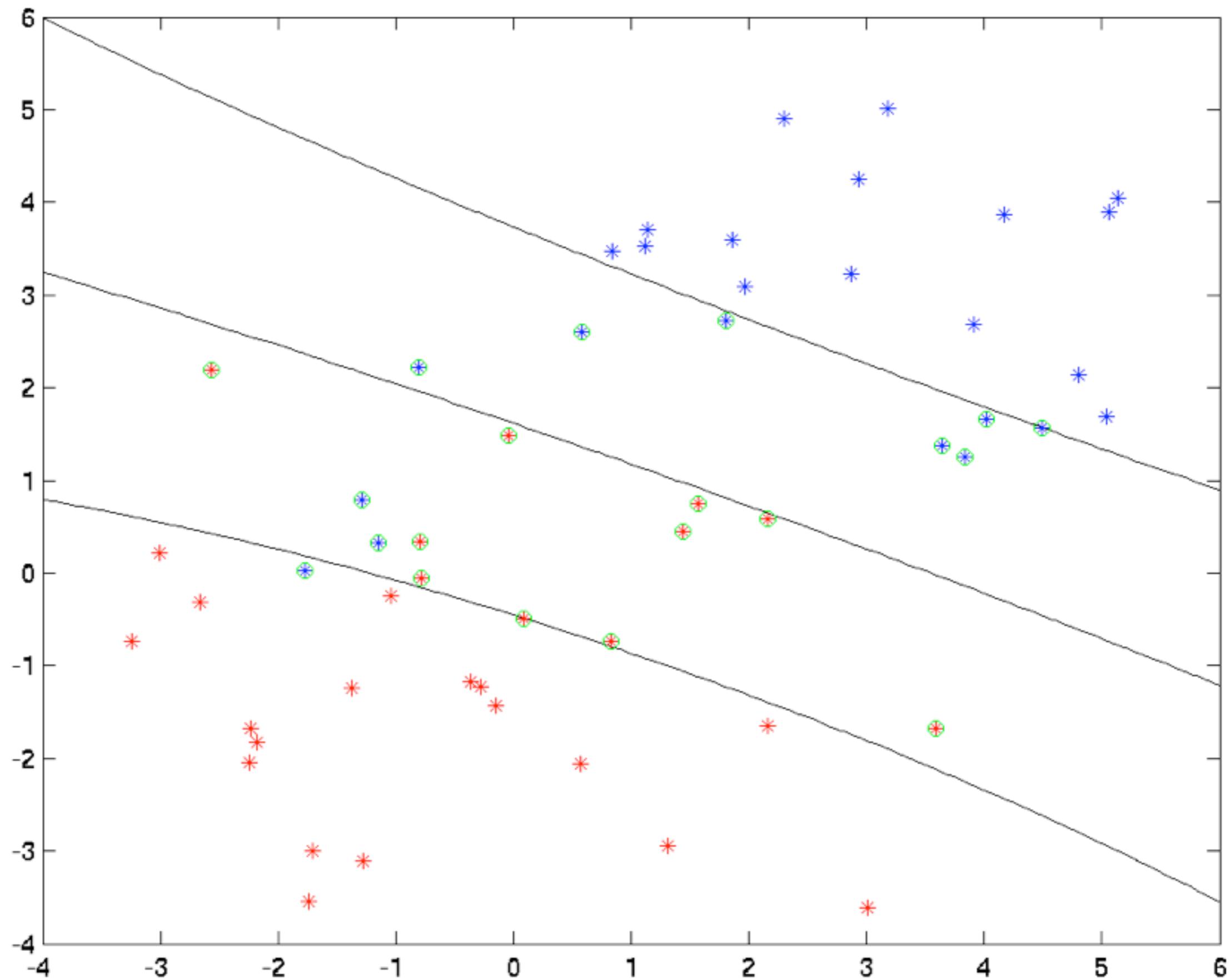




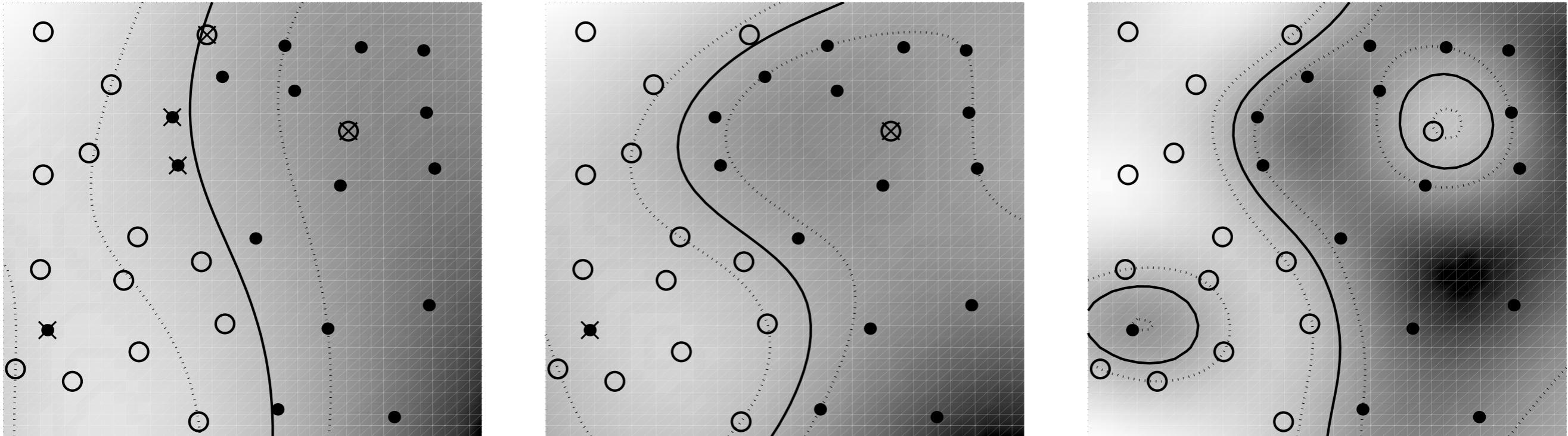




And now with a very wide kernel



Nonlinear separation



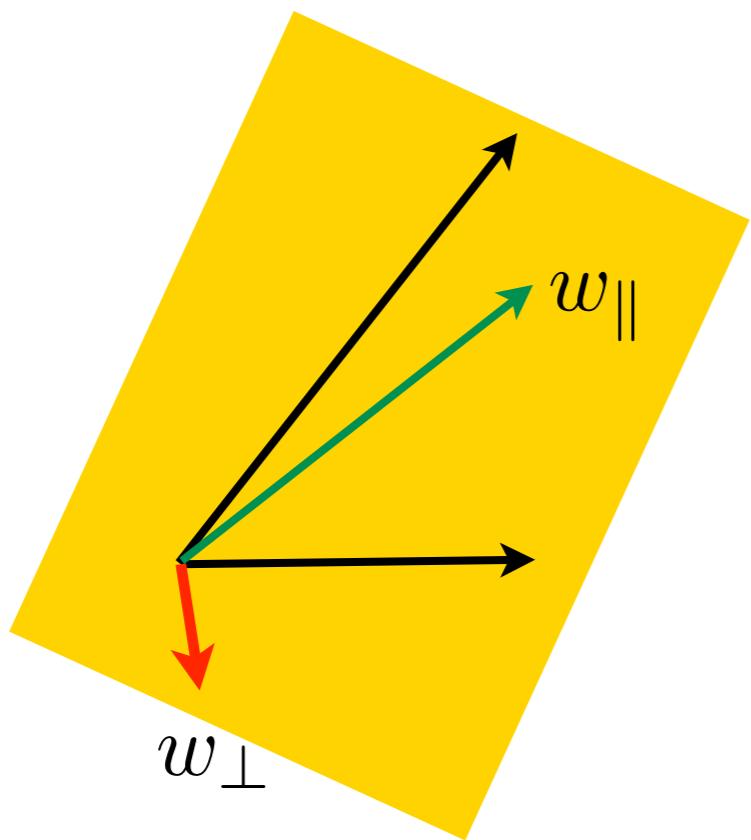
- Increasing C allows for more nonlinearities
- Decreases number of errors
- SV boundary need not be contiguous
- Kernel width adjusts function class

Penalized least mean squares ... now with kernels

- Optimization problem

$$\underset{w}{\text{minimize}} \frac{1}{2m} \sum_{i=1}^m (y_i - \langle \phi(x_i), w \rangle)^2 + \frac{\lambda}{2} \|w\|^2$$

- Representer Theorem (Kimeldorf & Wahba, 1971)



$$\|w\|^2 = \|w_{\parallel}\|^2 + \|w_{\perp}\|^2$$

empirical
risk dependent

Penalized least mean squares ... now with kernels

- Optimization problem

$$\underset{w}{\text{minimize}} \frac{1}{2m} \sum_{i=1}^m (y_i - \langle \phi(x_i), w \rangle)^2 + \frac{\lambda}{2} \|w\|^2$$

- Representer Theorem (Kimeldorf & Wahba, 1971)
 - Optimal solution is in span of data $w = \sum \alpha_i \phi(x_i)$
 - Proof - risk term only depends on data via $\phi(x_i)$
 - Regularization ensures that orthogonal part is 0
- Optimization problem in terms of w

$$\underset{\alpha}{\text{minimize}} \frac{1}{2m} \sum_{i=1}^m \left(y_i - \sum_j K_{ij} \alpha_j \right)^2 + \frac{\lambda}{2} \sum_{i,j} \alpha_i \alpha_j K_{ij}$$

solve for $\alpha = (K + m\lambda I)^{-1} y$ as linear system

Many more applications

- (Kernel) Principal Component Analysis
- (Kernel) Independent Component Analysis
- (Kernel) Linear Discriminant Analysis
- (Kernel) Canonical Correlation Analysis
- (Kernel) Two-Sample Test
- (Kernel) Graphical Models
- (Kernel) Covariate Shift Correction
- (Kernel) ...

Outline

- **Convex Optimization**
 - Unconstrained Optimization
 - Constrained Optimization and Duality
 - Linear and Quadratic programs
- **Support Vector Machines**
 - Classification
 - Regression
 - Novelty Detection
- **Kernels**
 - Feature Space
 - Kernel PCA
 - Kernelized SVM