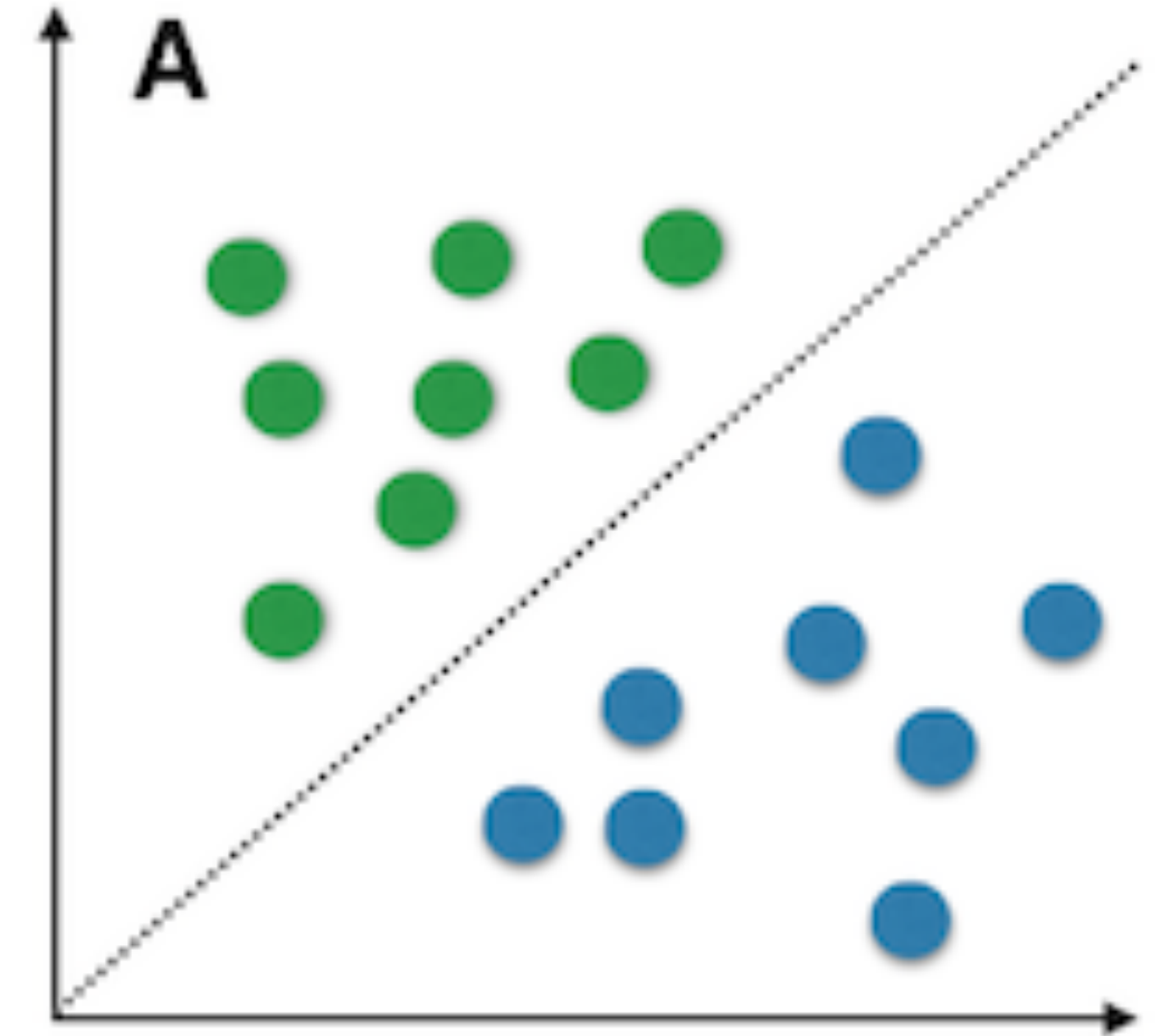# Soft & Kernel SVMs

EECE454 Intro. to Machine Learning Systems

Fall 2024
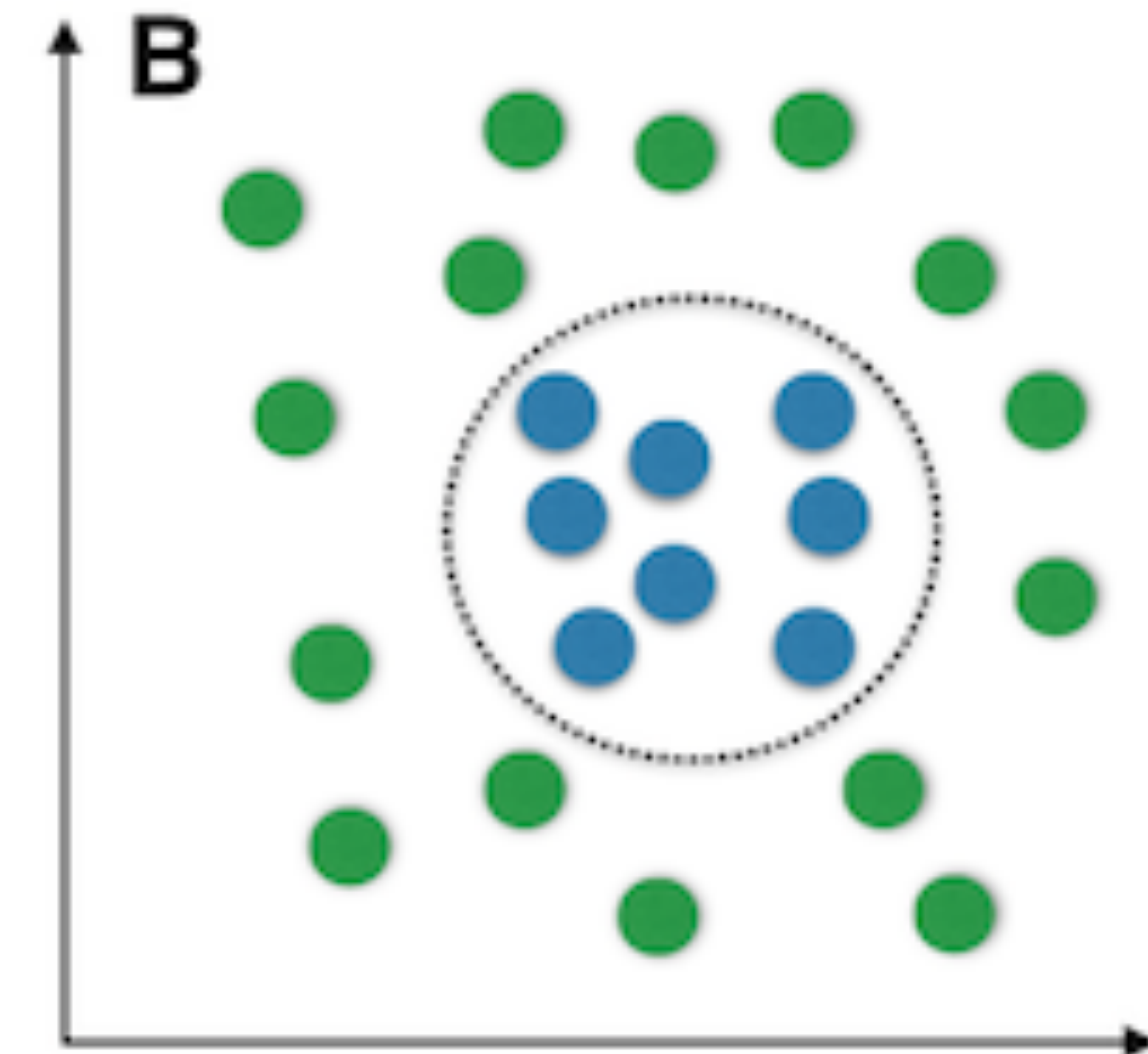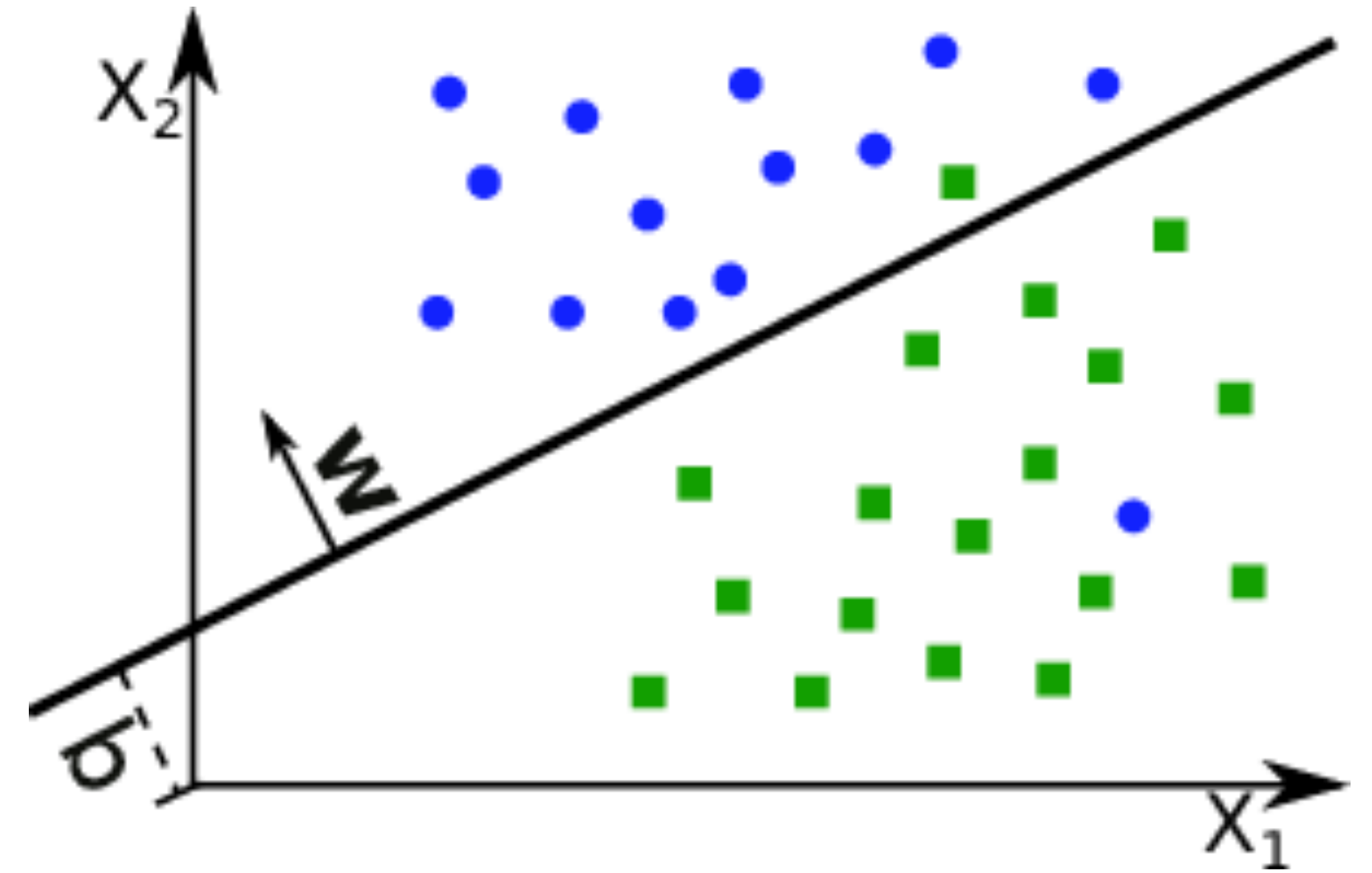
# Today

- **Last class.** Support Vector Machine

  - Linear model that maximizes the margin

  - Lagrangian dual —> Quadratic problem
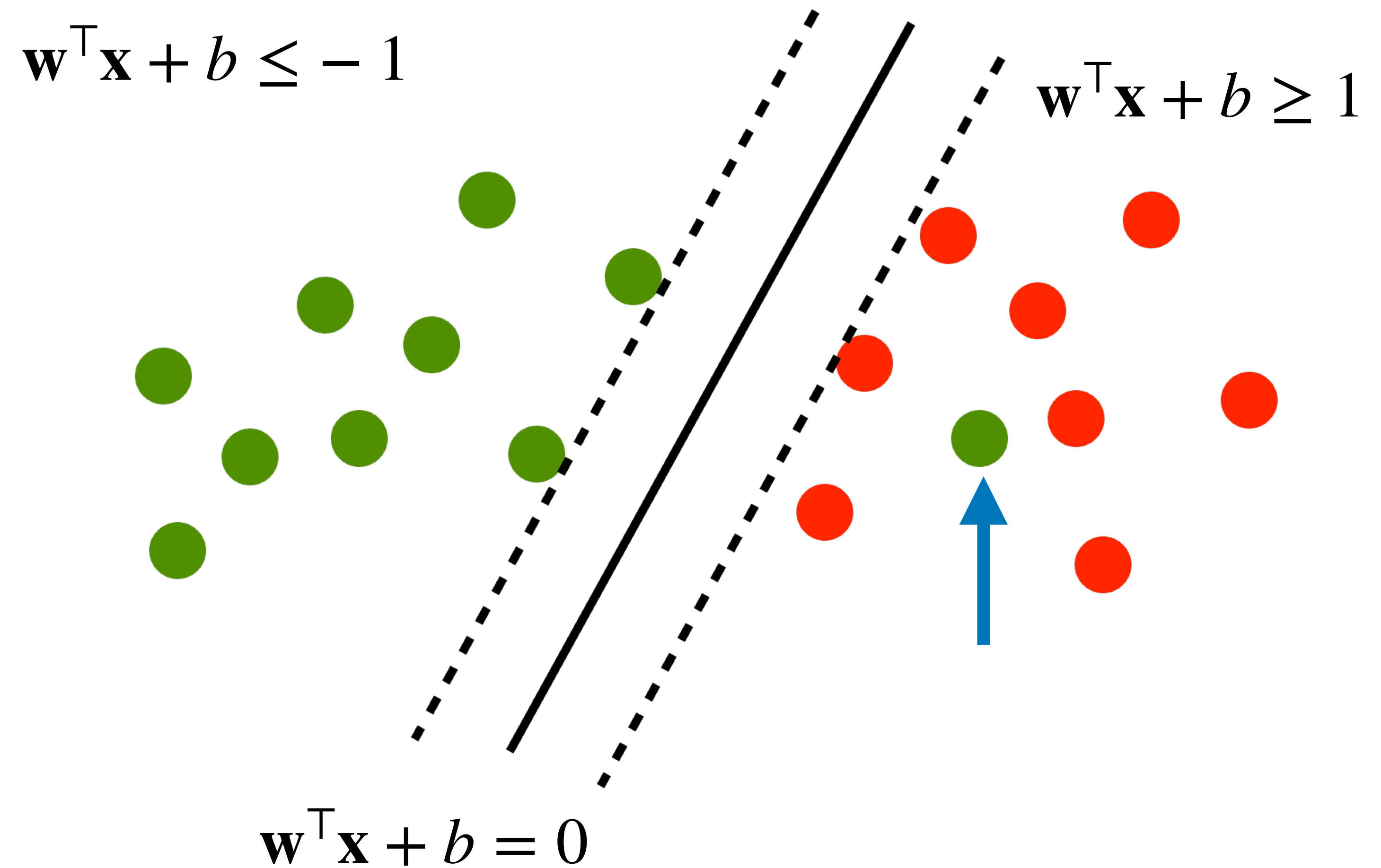
  - Required. Data is linearly separable

# Today

- **Last class.** Support Vector Machine

  - Linear model that maximizes the margin

  - Lagrangian dual —> Quadratic problem

  - Required. Data is linearly separable

- **Today.** SVMs that can handle nonseparable data

  - Soft-margin SVM

  - Kernel SVM

# Soft(-Margin) SVM

# Data with outliers

- Suppose that there exists some **outlier**

    - Then, no linear separator exists

    - Worse. finding a minimum-error separating hyperplane is NP-hard (Minsky & Papert, 1969)

- **Q.** How can we handle this situation?

$$\mathbf{w}^\top \mathbf{x} + b \leq -1$$

$$\mathbf{w}^\top \mathbf{x} + b \geq 1$$

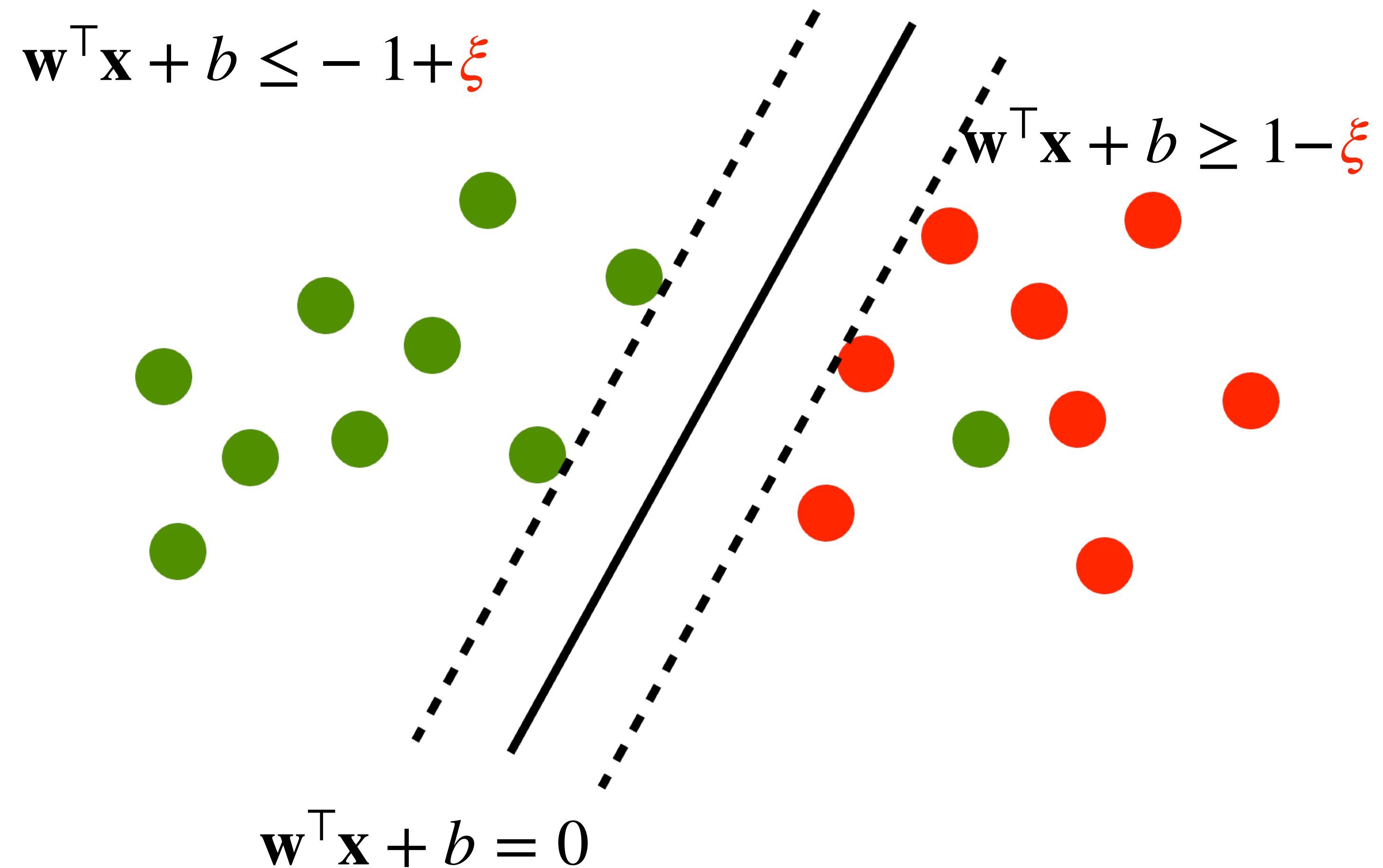$$\mathbf{w}^\top \mathbf{x} + b = 0$$

# Data with outliers

- Suppose that there exists some **outlier**

  - Then, no linear separator exists

  - Worse. finding a minimum-error separating hyperplane is NP-hard (Minsky & Papert, 1969)

- **Q.** How can we handle this situation?

  - <u>A</u>. Add some slack variable $\xi$

    - Then, aim for minimizing the slack as well

$$\mathbf{w}^\top \mathbf{x} + b \leq -1 + \xi$$

$$\mathbf{w}^\top \mathbf{x} + b \geq 1 - \xi$$

$$\mathbf{w}^\top \mathbf{x} + b = 0$$

# Formulation

- We are now solving the optimization problem

$$\ell^* = \min_{\mathbf{w}, b, \xi} \frac{\|\mathbf{w}\|^2}{2} + C \cdot \sum_i \xi_i$$

$$\text{subject to} \quad y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 - \xi_i, \qquad \xi_i \geq 0$$

# Formulation

$$\ell^* = \min_{\mathbf{w}, b, \xi} \frac{\|\mathbf{w}\|^2}{2} + C \cdot \sum_i \xi_i$$

$$\text{subject to} \quad y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 - \xi_i, \qquad \xi_i \geq 0$$

- Then, we know that the problem is always feasible

  - Constraint can be met in any case

  - For example, let $\mathbf{w} = \mathbf{0}$, $b = 0$, and $\xi_i = 1$.

# Dual Formulation

- As a dual, we get

$$\min_{\mathbf{w},b,\xi} \max_{\alpha,\eta} \left( \frac{\|\mathbf{w}\|^2}{2} + C \sum_i \xi_i - \sum_i \alpha_i \big(y_i(\mathbf{x}_i^\top \mathbf{w} + b) + \xi_i - 1\big) - \sum_i \eta_i \xi_i \right)$$

  - The optimal $(\mathbf{w}, b, \xi)$ is at the saddle point with $(\alpha, \eta)$

# Dual Formulation

- As a dual, we get

$$\min_{\mathbf{w},b,\xi} \max_{\alpha,\eta} \left( \frac{\|\mathbf{w}\|^2}{2} + C \sum_i \xi_i - \sum_i \alpha_i \big( y_i(\mathbf{x}_i^\top \mathbf{w} + b) + \xi_i - 1 \big) - \sum_i \eta_i \xi_i \right)$$

- The optimal $(\mathbf{w}, b, \xi)$ is at the saddle point with $(\alpha, \eta)$

- Derivatives for $(\mathbf{w}, b, \xi)$ needs to vanish!

  - $\nabla_{\mathbf{w}} \mathscr{L} = \mathbf{w} - \sum \alpha_i y_i \mathbf{x}_i = 0$

  - $\nabla_b \mathscr{L} = \sum \alpha_i y_i = 0$

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

# Dual Formulation

- Doing the similar thing, we get the Lagrangian

$$-\frac{1}{2}\sum_{i,j}\alpha_i\alpha_j y_i y_j \mathbf{x}_i^\top \mathbf{x}_j + \sum_i \alpha_i - {\color{red}\sum_i \alpha_i \xi_i} + {\color{red}C\sum_i \xi_i} - {\color{red}\sum_i \eta_i \xi_i}$$

$$= -\frac{1}{2}\sum_{i,j}\alpha_i\alpha_j y_i y_j \mathbf{x}_i^\top \mathbf{x}_j + \sum_i \alpha_i$$

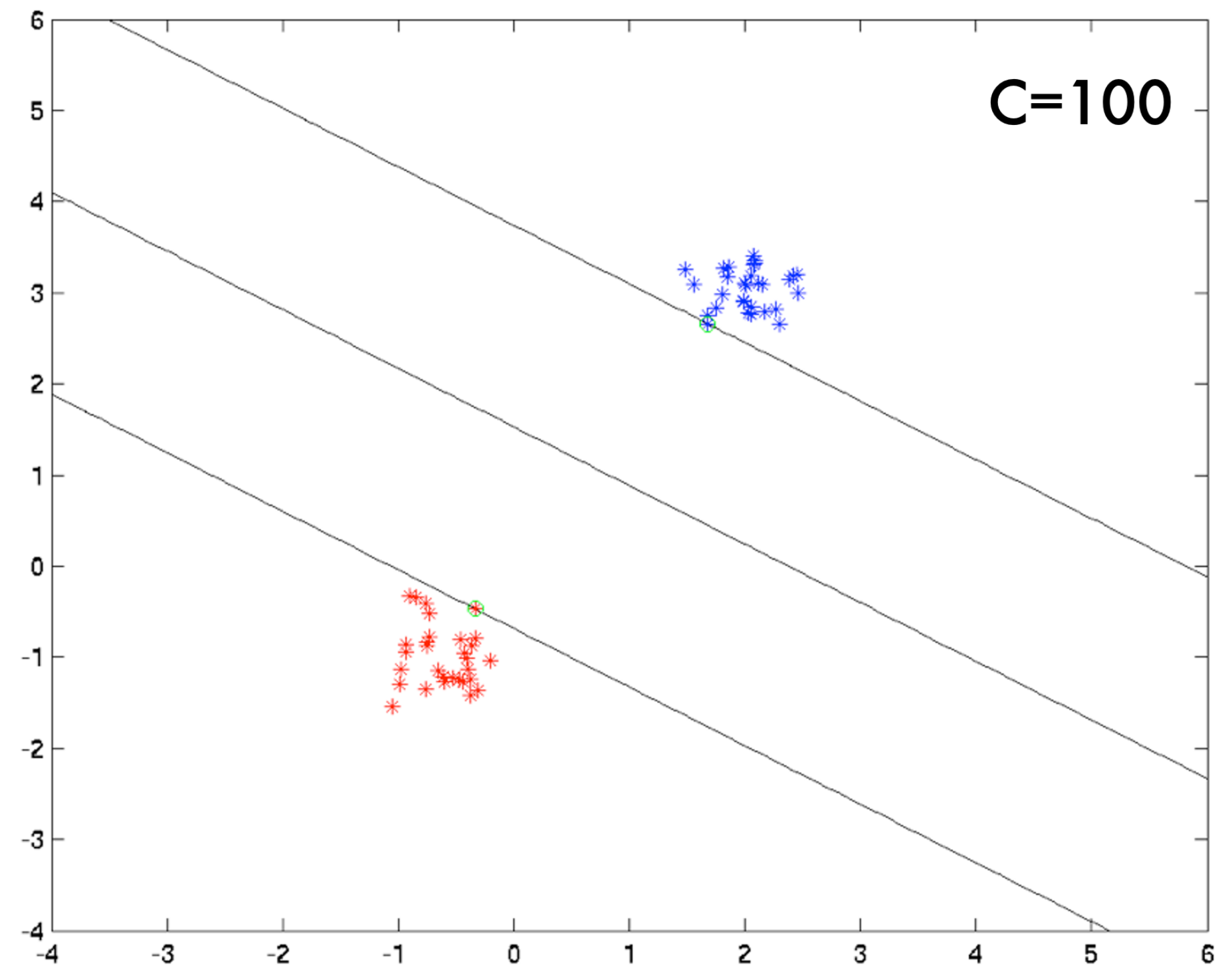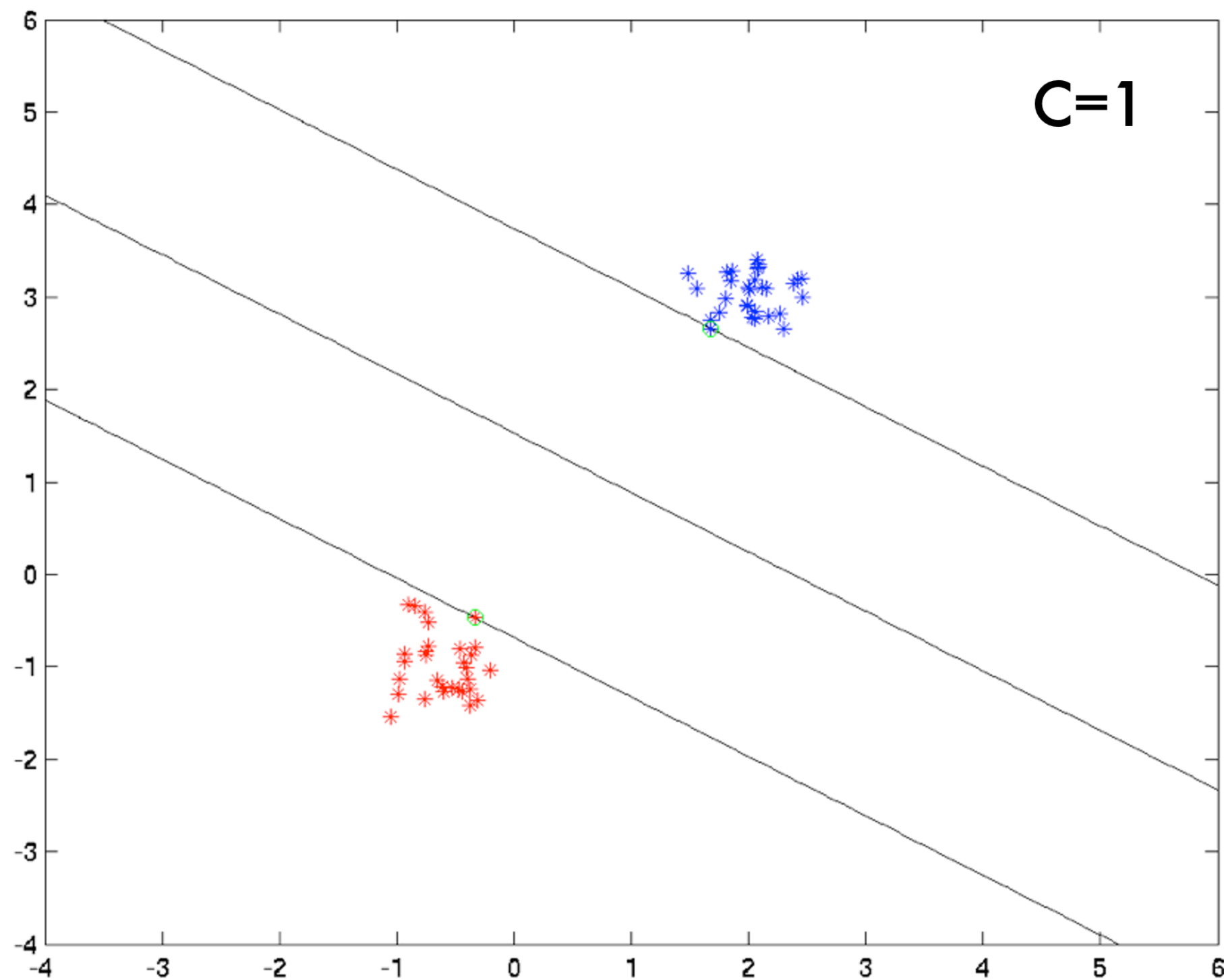# Dual Formulation

- Doing the similar thing, we get the Lagrangian

$$-\frac{1}{2}\sum_{i,j}\alpha_i\alpha_j y_i y_j \mathbf{x}_i^\top \mathbf{x}_j + \sum_i \alpha_i - \sum_i \alpha_i \xi_i + C\sum_i \xi_i - \sum_i \eta_i \xi_i$$

$$= -\frac{1}{2}\sum_{i,j}\alpha_i\alpha_j y_i y_j \mathbf{x}_i^\top \mathbf{x}_j + \sum_i \alpha_i$$

- Summing up, we are solving the optimization

$$\max_{\alpha}\left(-\frac{1}{2}\sum_{i,j}\alpha_i\alpha_j y_i y_j \mathbf{x}_i^\top \mathbf{x}_j + \sum_{i=1}^{n}\alpha_i\right) \qquad \text{subject to} \qquad \sum_i \alpha_i y_i = 0 \qquad 0 \leq \alpha_i \leq C$$
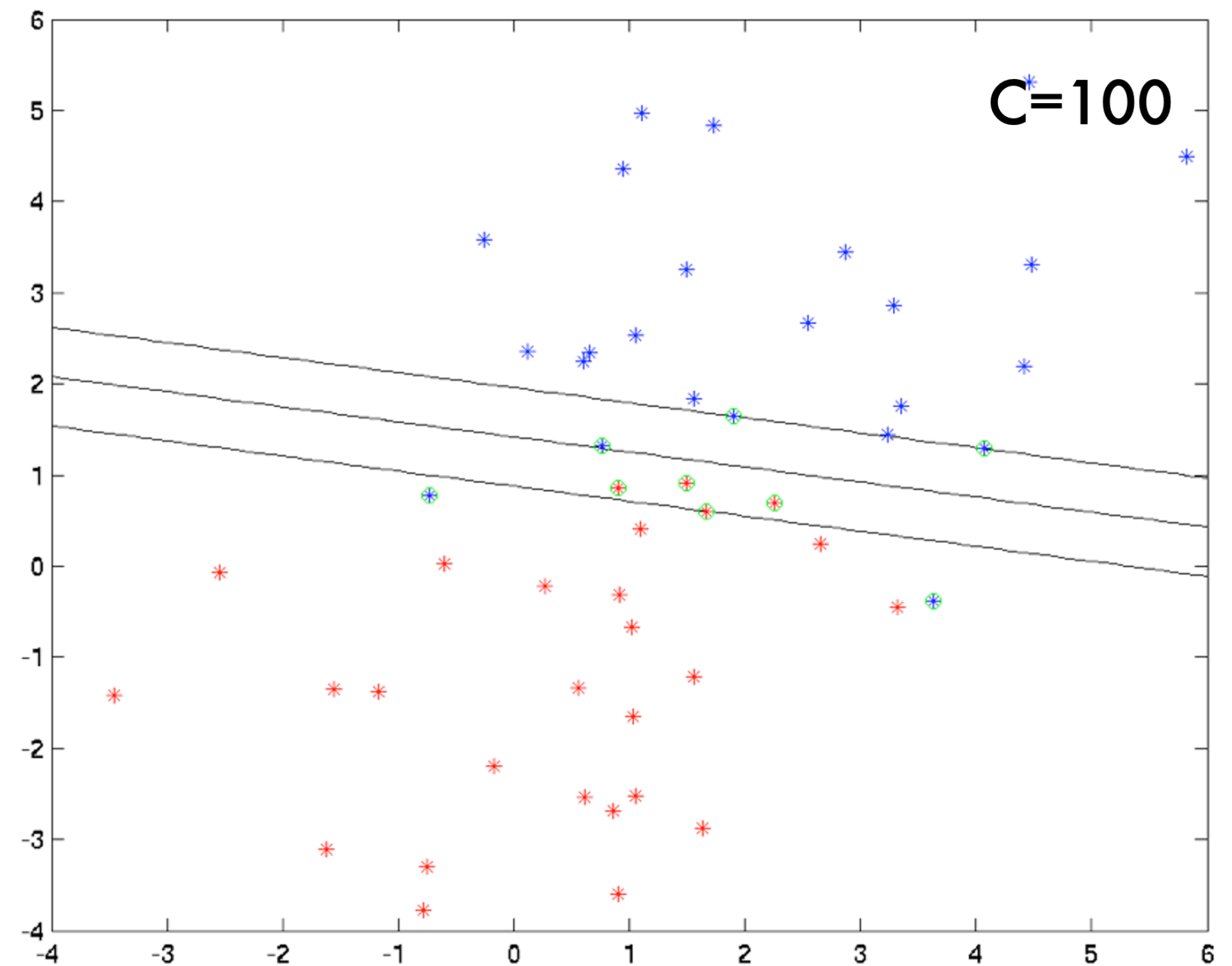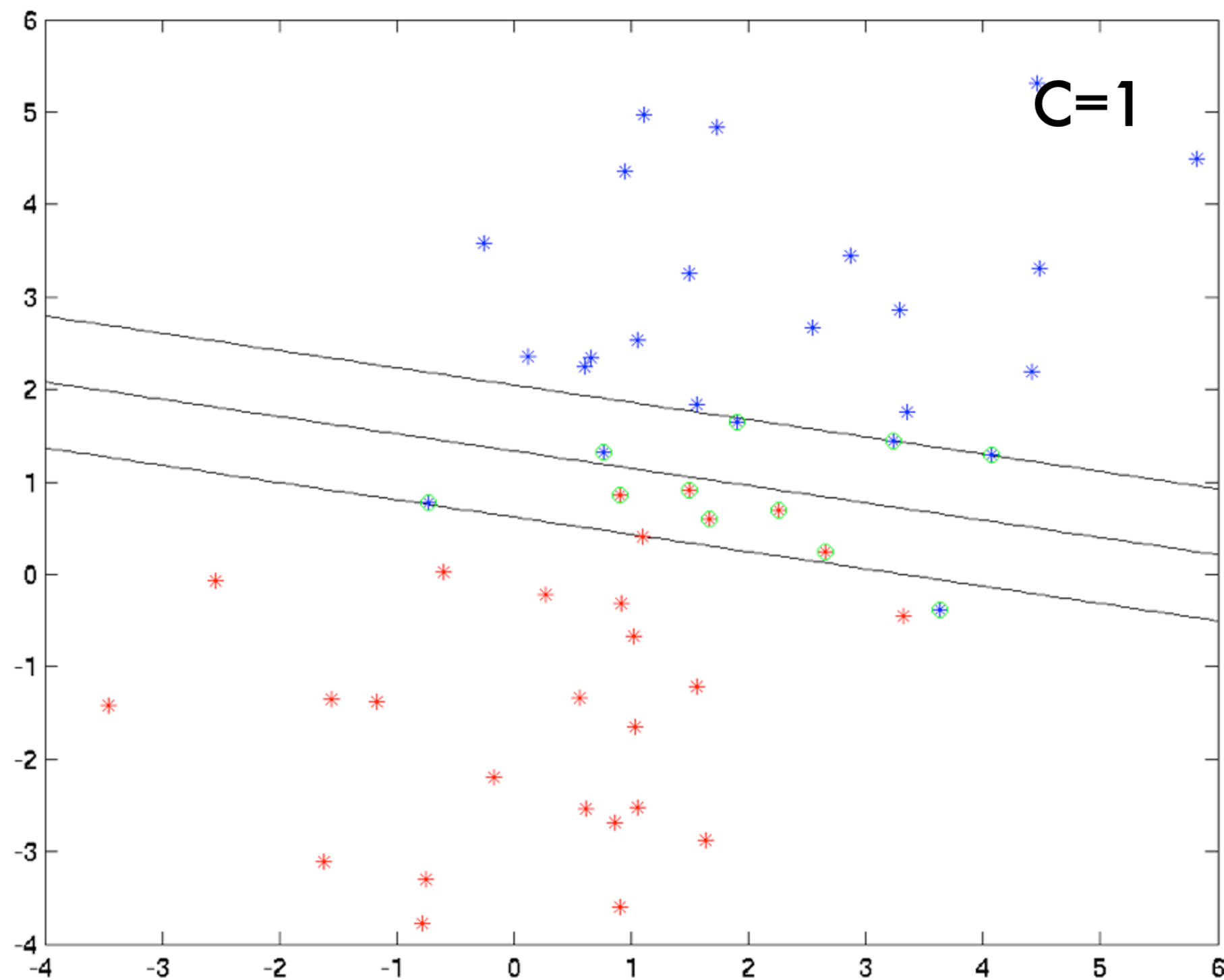
# Hyperparameter

- By increasing the hyperparameter $C$, we look for a smaller-slack solution

    - No difference when linearly separable…

# Hyperparameter

- By increasing the hyperparameter $C$, we look for a smaller-slack solution

  - No difference when linearly separable... but some difference when not

# Solving the optimization

$$\max_{\alpha} \left( -\frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \mathbf{x}_i^\top \mathbf{x}_j + \sum_{i=1}^{n} \alpha_i \right)$$
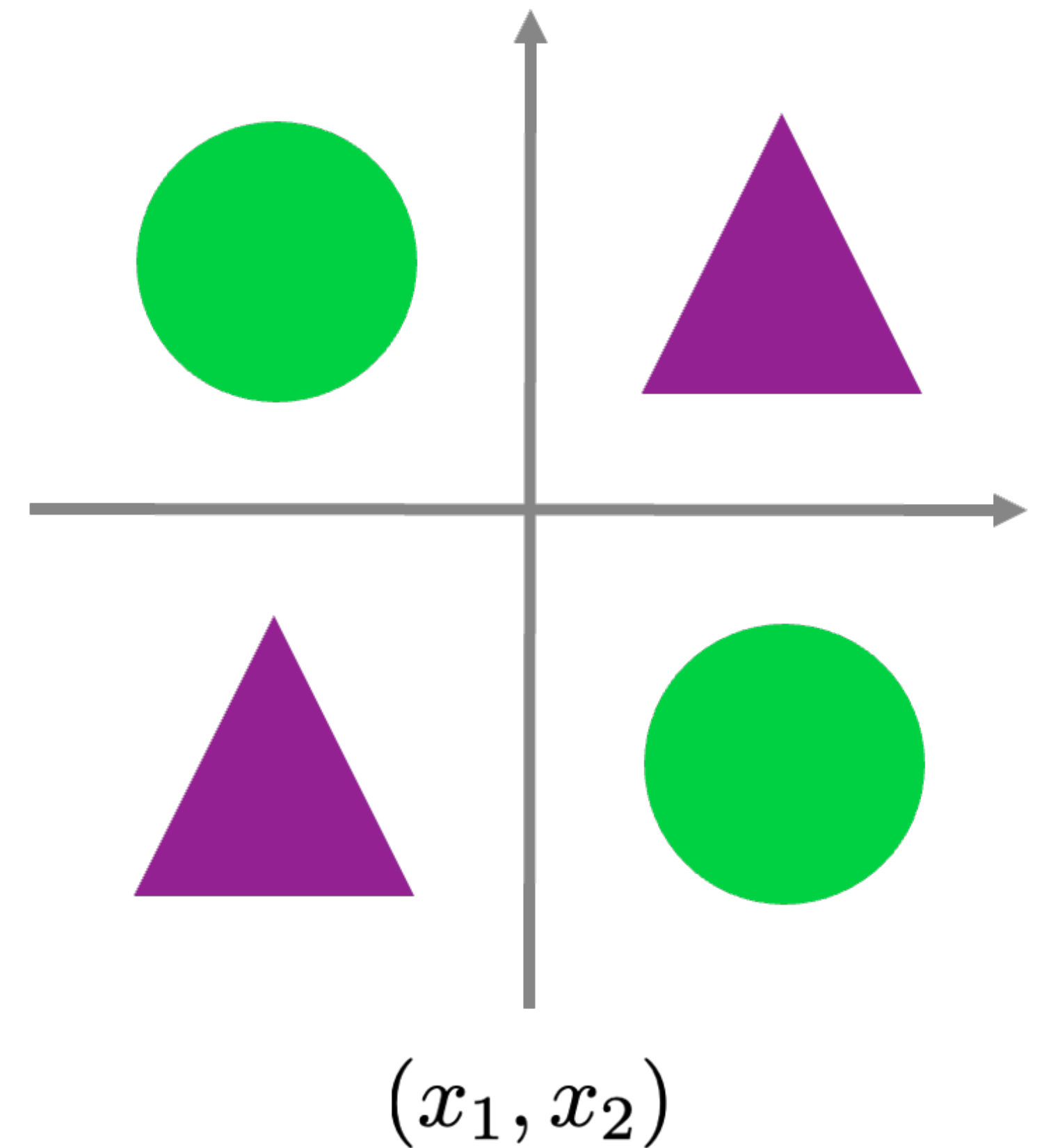
$$\text{subject to} \quad \sum_{i} \alpha_i y_i = 0 \qquad 0 \leq \alpha_i \leq C$$

- If the problem is small-scale (e.g., thousands of variables), use off-the-shelf solvers

- If the problem is large-scale, use the fact that only SVs matter, and solve in blocks

  - called "active set method"

# Kernel SVM

# Nonlinear data

- Suppose that we have a data that looks like **XOR**

    - Not linearly separable

    - Thus no satisfactory linear classifier exists

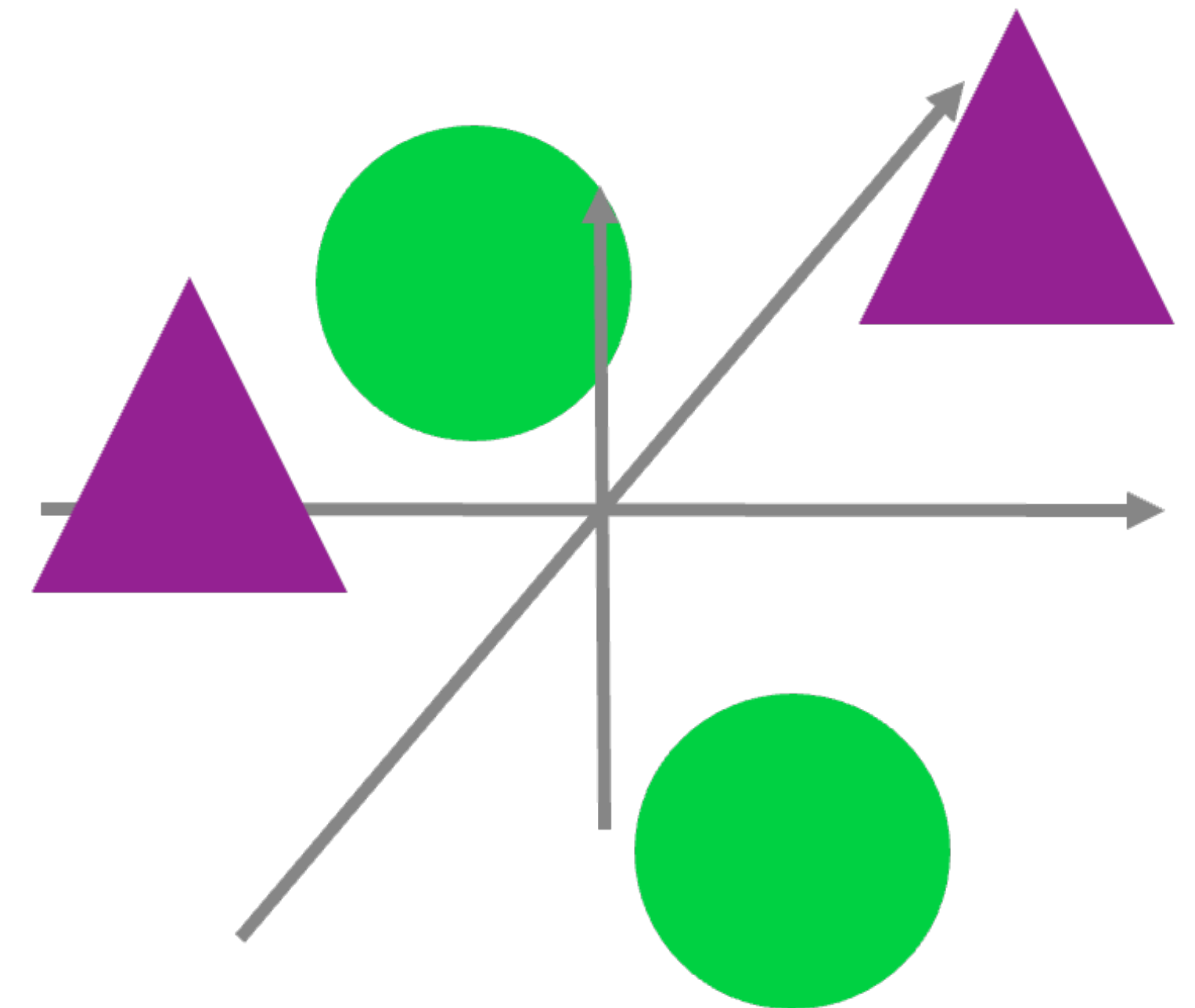- **Q.** How to handle these data?



$(x_1, x_2)$

# Nonlinear data

- Suppose that we have a data that looks like **XOR**

  - Not linearly separable

  - Thus no satisfactory linear classifier exists

- **Q.** How to handle these data?

  - <u>A</u>. Map it to a high-dimensional space

    - There exists a clean linear classifier!

$$f(\mathbf{x}) = \text{sign}\left([0 \quad 0 \quad 1] \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}\right)$$

$(x_1, x_2, \boxed{x_1 x_2})$

# More formally...

- We map the data to a high-dimensional **feature** $\Phi(\,\cdot\,) : \mathbb{R}^d \to \mathbb{R}^k$

  - Typically, $d < k$ (but not necessarily)

# More formally...

- We map the data to a high-dimensional **feature** $\Phi(\,\cdot\,) : \mathbb{R}^d \rightarrow \mathbb{R}^k$

  - Typically, $d < k$ (but not necessarily)

- Our predictor takes the form

$$f(\mathbf{x}) = \text{sign}\left( \sum_{i=1}^{n} a_i \cdot \boxed{\langle \Phi(\mathbf{x}_i), \Phi(\mathbf{x}) \rangle} + b \right)$$

  - This is quite similar to original SVMs, where

$$f(\mathbf{x}) = \text{sign}\left( \sum a_i \cdot \langle \mathbf{x}_i, \mathbf{x} \rangle + b \right)$$

# Choosing the feature

- **Question.** How should we choose $\Phi(\cdot)$?

# Choosing the feature

- **Question.** How should we choose $\Phi(\,\cdot\,)$?

  - <u>Naïve way</u>. Simply throw in many features, and let SVM choose

$$\Phi(\mathbf{x}) = (x_1, \cdots, x_d, x_1 x_2, \cdots, x_{d-1} x_d, \cdots, x_k^{100})$$

# Choosing the feature

- **Question.** How should we choose $\Phi(\,\cdot\,)$?

  - Naïve way. Simply throw in many features, and let SVM choose

$$\Phi(\mathbf{x}) = (x_1, \cdots, x_d, x_1 x_2, \cdots, x_{d-1} x_d, \cdots, x_k^{100})$$
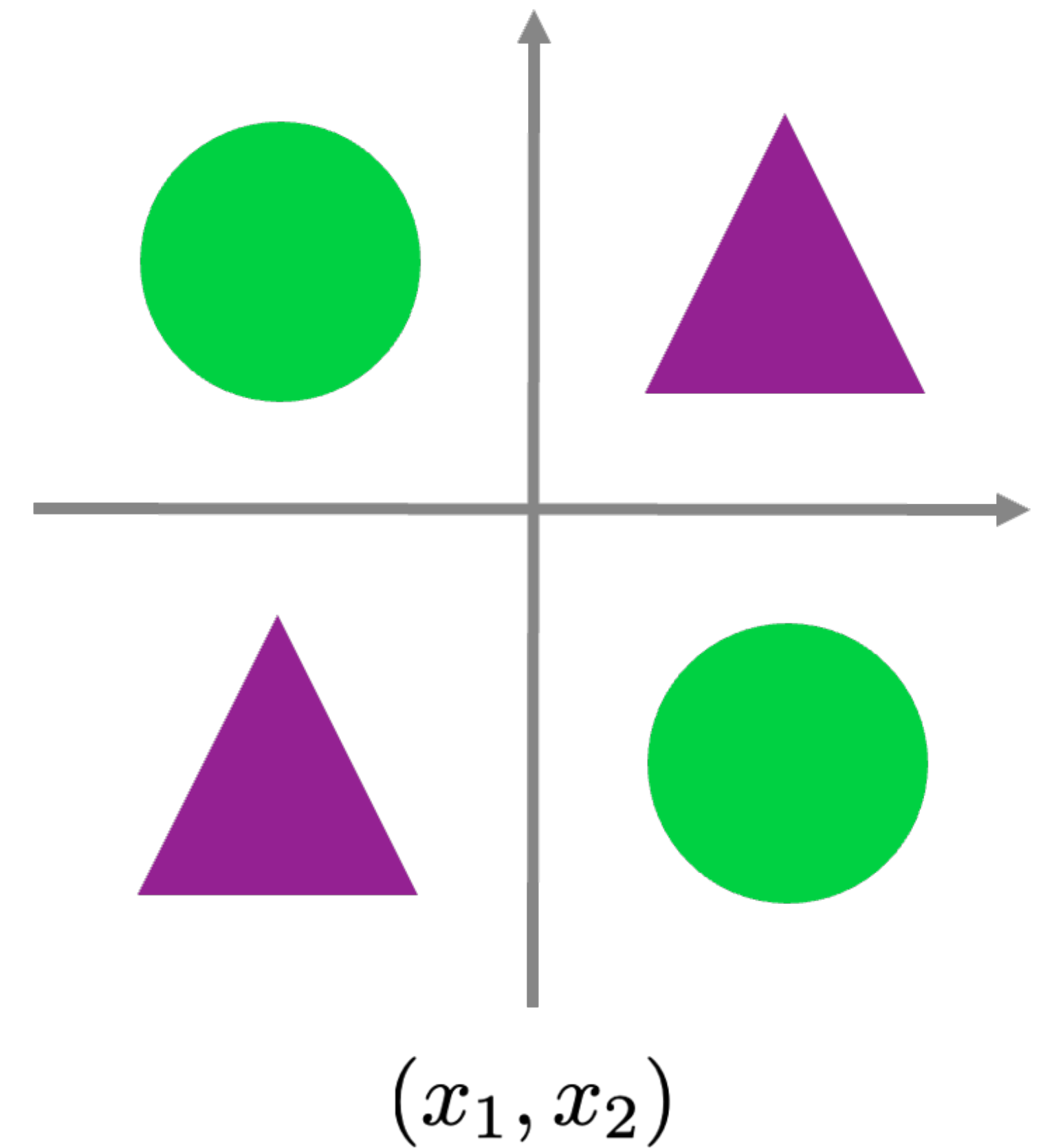
- This is **bad**!

  - overfitting

  - computation

    - computing features

    - computing inner products

# Choosing the feature

- Interestingly, some features admit **computational shortcuts**

# Choosing the feature

- Interestingly, some features admit **computational shortcuts**

  - Example. Recall the **XOR**, and think of two features.

    - $\Phi_a((x_1, x_2)) = (x_1, x_2, x_1 x_2)$

    - $\Phi_b((x_1, x_2)) = (x_1^2, x_2^2, \sqrt{2} x_1 x_2)$

  - Looks similar, but one is better than the other

  - **Question.** So which one is better?

$(x_1, x_2)$

# Choosing the feature

- Interestingly, some features admit **computational shortcuts**

  - Example. Recall the **XOR**, and think of two features.
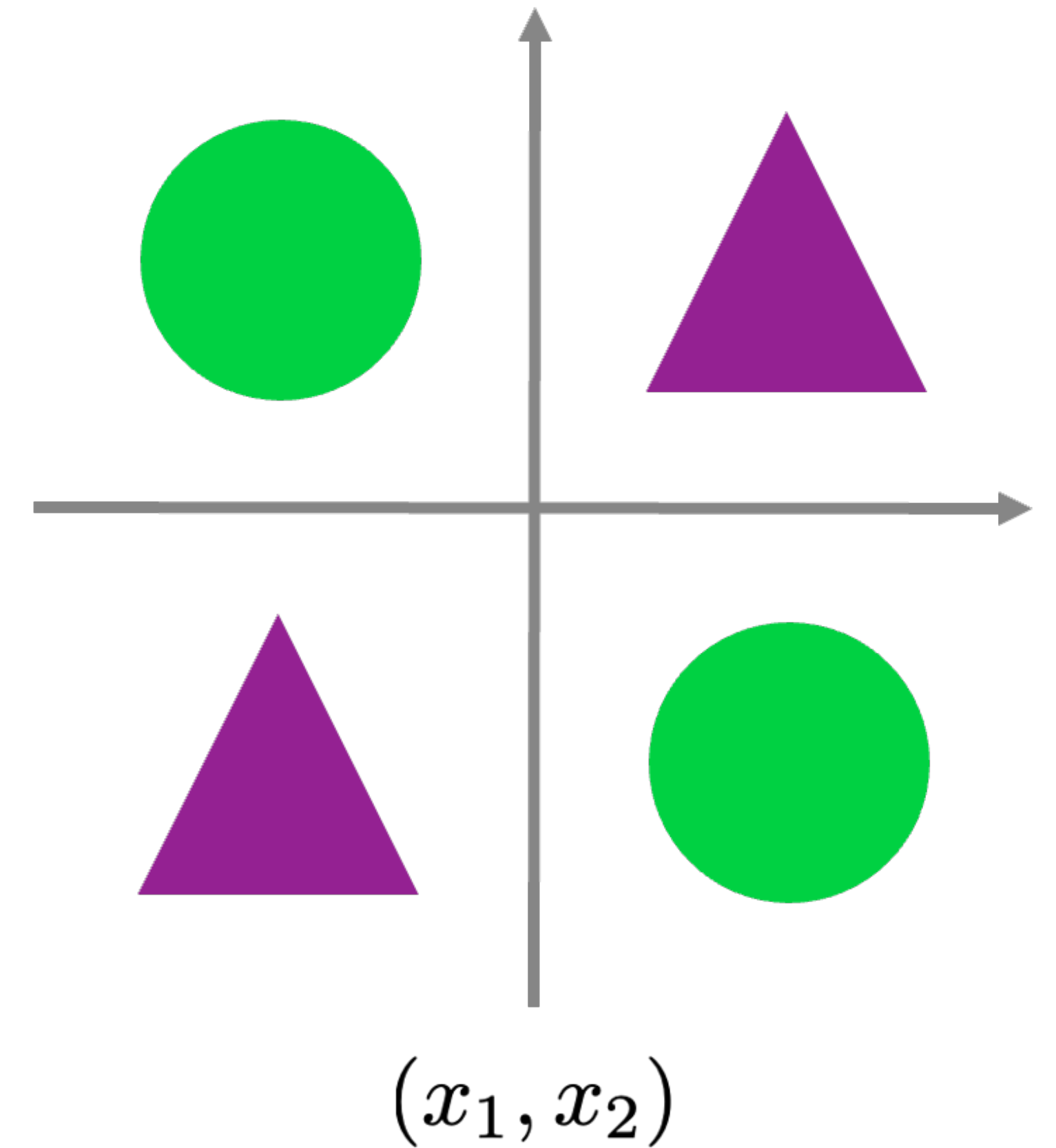
    - $\Phi_a((x_1, x_2)) = (x_1, x_2, x_1 x_2)$

    - $\Phi_b((x_1, x_2)) = (x_1^2, x_2^2, \sqrt{2} x_1 x_2)$

  - Looks similar, but one is better than the other

  - **Question.** So which one is better?

    - Answer. $\Phi_b$, for computational reasons

$(x_1, x_2)$

# Choosing the feature

- Compare the computations:

  - $\langle \Phi_a(\mathbf{x}), \Phi_a(\mathbf{y}) \rangle = x_1 y_1 + x_2 y_2 + x_1 x_2 y_1 y_2$

    - Compute 3D features $\phi_{\mathbf{x}} = \Phi_a(\mathbf{x}), \phi_{\mathbf{y}} = \Phi_a(\mathbf{y})$

    - Compute 3D inner prod $\langle \phi_{\mathbf{x}}, \phi_{\mathbf{y}} \rangle$

# Choosing the feature

- Compare the computations:

  - $\langle \Phi_a(\mathbf{x}), \Phi_a(\mathbf{y}) \rangle = x_1 y_1 + x_2 y_2 + x_1 x_2 y_1 y_2$

    - Compute 3D features $\phi_{\mathbf{x}} = \Phi_a(\mathbf{x})$, $\phi_{\mathbf{y}} = \Phi_a(\mathbf{y})$

    - Compute 3D inner prod $\langle \phi_{\mathbf{x}}, \phi_{\mathbf{y}} \rangle$

- $\langle \Phi_b(\mathbf{x}), \Phi_b(\mathbf{y}) \rangle = x_1^2 y_1^2 + x_2^2 y_2^2 + 2 x_1 x_2 y_1 y_2 \color{red}{= (\langle \mathbf{x}, \mathbf{y} \rangle)^2}$

  - Compute 2D inner prod $\langle \mathbf{x}, \mathbf{y} \rangle$

  - Take a square

    - Less memory & computation

# Kernel SVM

- **Idea.** Follow these steps.

  - Choose an easy-to-compute <span style="color:red">**similarity metric**</span> $K(\cdot\,,\cdot\,)$

  - Construct predictors of form

$$f(\mathbf{x}) = \text{sign} \left( \sum a_i \cdot K(\mathbf{x}_i, \mathbf{x}) + b \right)$$

  - Fit $a_i, b$

# Kernel SVM

- **Idea.** Follow these steps.

  - Choose an easy-to-compute **similarity metric** $K( \cdot , \cdot )$

  - Construct predictors of form

$$f(\mathbf{x}) = \text{sign} \left( \sum a_i \cdot K(\mathbf{x}_i, \mathbf{x}) + b \right)$$

  - Fit $a_i, b$

- **Question.** Is this equivalent to doing SVM with features?

    (i.e., does there always exist a $\mathbf{\Phi}$ corresponding to $K$?)

# Kernel SVM

- **Idea.** Follow these steps.

  - Choose an easy-to-compute **similarity metric** $K(\cdot, \cdot)$

  - Construct predictors of form

  $$f(\mathbf{x}) = \text{sign}\left( \sum a_i \cdot K(\mathbf{x}_i, \mathbf{x}) + b \right)$$

  - Fit $a_i, b$

- **Question.** Is this equivalent to doing SVM with features?

  (i.e., does there always exist a $\mathbf{\Phi}$ corresponding to $K$?)

  - <u>Answer</u>. Yes if $K$ is a **Mercer kernel**

# Kernel SVM

- **Definition.** A real-valued function $K(\cdot, \cdot)$ is a <span style="color:red">Mercer kernel</span> if

  - $K(\mathbf{x}, \mathbf{x}') = K(\mathbf{x}', \mathbf{x})$            (i.e., symmetric)

  - $\lim_{n \to \infty} K(\mathbf{x}^{(n)}, \mathbf{x}) \to K\left( \lim_{n \to \infty} \mathbf{x}^{(n)}, \mathbf{x} \right)$     (i.e., continuous)

  - $\sum_{i,j} \alpha_i \alpha_j K(\mathbf{x}_i, \mathbf{x}_j) \geq 0, \quad \forall \alpha_i, \alpha_j, \mathbf{x}_i, \mathbf{x}_j$     (i.e., positive-semidefinite)

# Kernel SVM

- **Definition.** A real-valued function $K(\,\cdot\,,\,\cdot\,)$ is a Mercer kernel if

  - $K(\mathbf{x}, \mathbf{x}') = K(\mathbf{x}', \mathbf{x})$                    (i.e., symmetric)

  - $\displaystyle \lim_{n \to \infty} K(\mathbf{x}^{(n)}, \mathbf{x}) \to K\left( \lim_{n \to \infty} \mathbf{x}^{(n)}, \mathbf{x} \right)$      (i.e., continuous)

  - $\displaystyle \sum_{i,j} \alpha_i \alpha_j K(\mathbf{x}_i, \mathbf{x}_j) \geq 0, \quad \forall \alpha_i, \alpha_j, \mathbf{x}_i, \mathbf{x}_j$     (i.e., positive-semidefinite)

- **Mercer's theorem.** For a Mercer kernel $K(\,\cdot\,,\,\cdot\,)$, there exists a corresponding $\Phi(\,\cdot\,)$ such that

$$K(\mathbf{x}, \mathbf{x}') = \langle \Phi(\mathbf{x}), \Phi(\mathbf{x}') \rangle$$

  - That is, we are effectively maximizing margin if we choose a nice kernel.

# Optimizing Kernel SVM

- In kernel SVM, we solve

$$\max_{\alpha} \left( -\frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j K(\mathbf{x}_i, \mathbf{x}_j) + \sum_{i=1}^{n} \alpha_i \right)$$

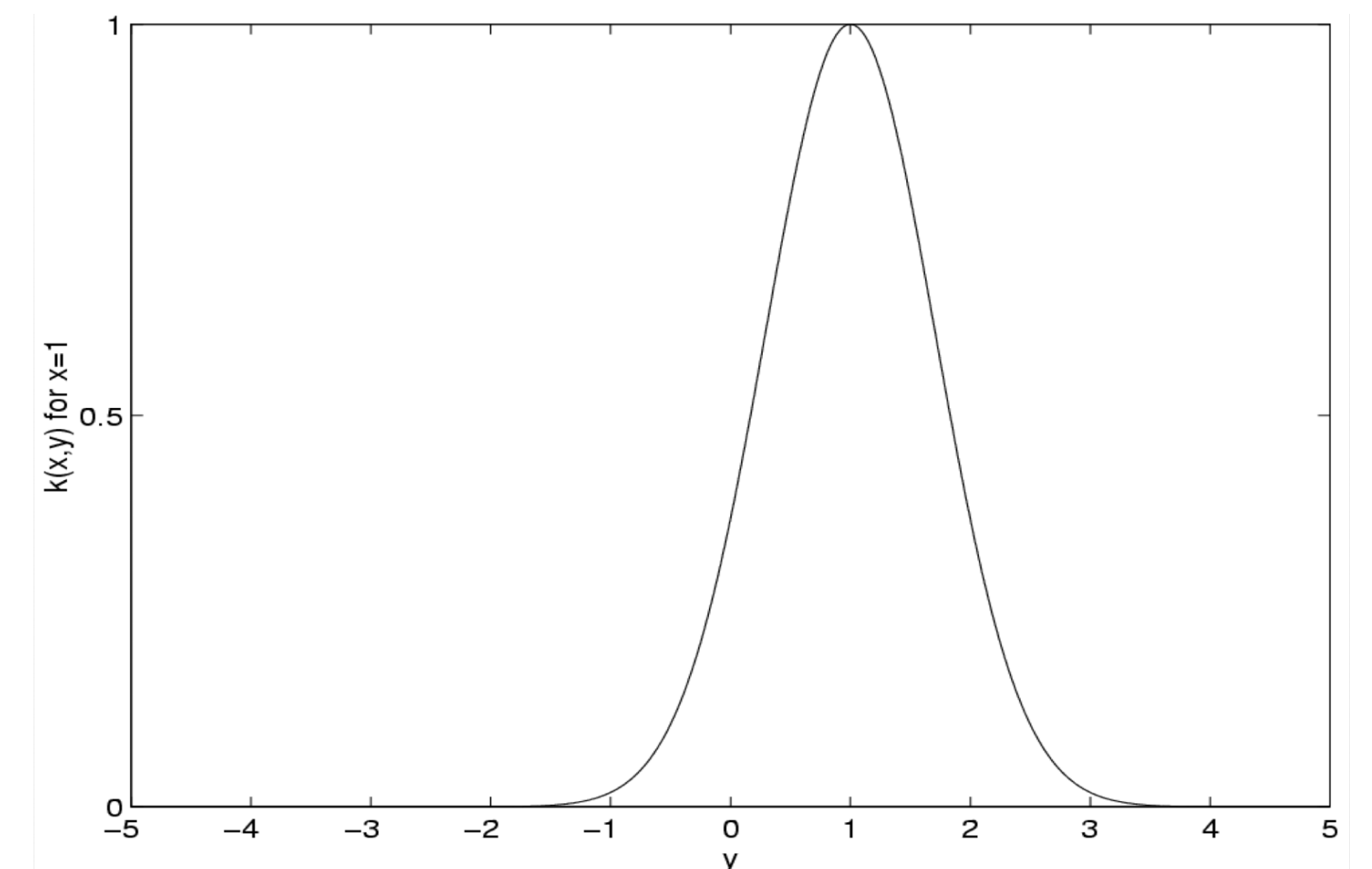- Plug in $K(\mathbf{x}, \mathbf{y}) = \mathbf{x}^\top \mathbf{y}$ to recover the original SVM

# Optimizing Kernel SVM

- In kernel SVM, we solve

$$\max_{\alpha} \left( -\frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j K(\mathbf{x}_i, \mathbf{x}_j) + \sum_{i=1}^{n} \alpha_i \right)$$

- Plug in $K(\mathbf{x}, \mathbf{y}) = \mathbf{x}^\top \mathbf{y}$ to recover the original SVM

- Other choices

  - Laplacian RBF $\exp(-\lambda \|\mathbf{x} - \mathbf{x}'\|_2)$

  - Gaussian RBF $\exp(-\lambda \|\mathbf{x} - \mathbf{x}'\|_2^2)$ —>

  - Polynomial $(\langle \mathbf{x}, \mathbf{x}' \rangle + c)^d$

  - B-Spline (look it up)

# Tuning Kernel SVM

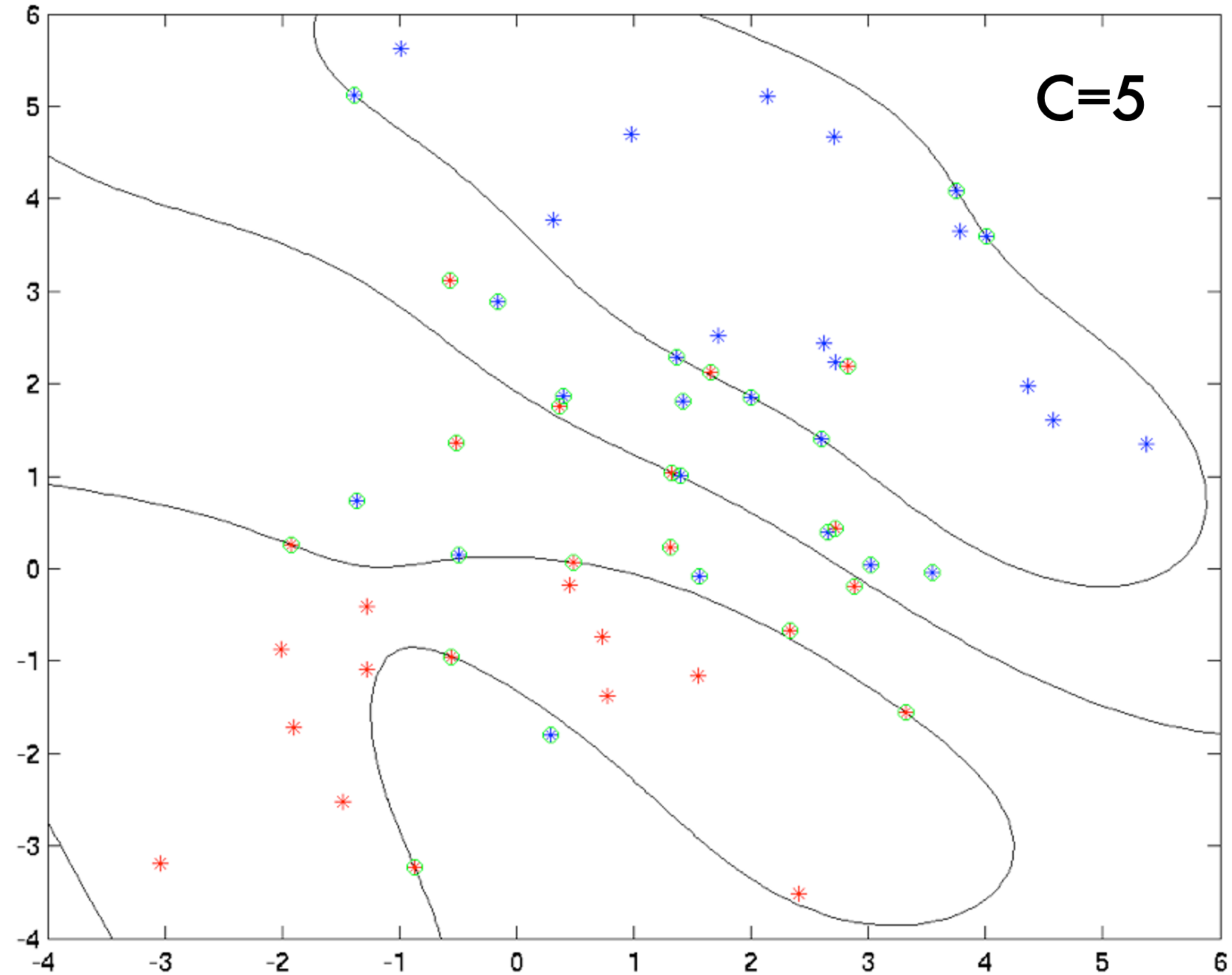- Again, we can tune hyperparameters to play with the margin
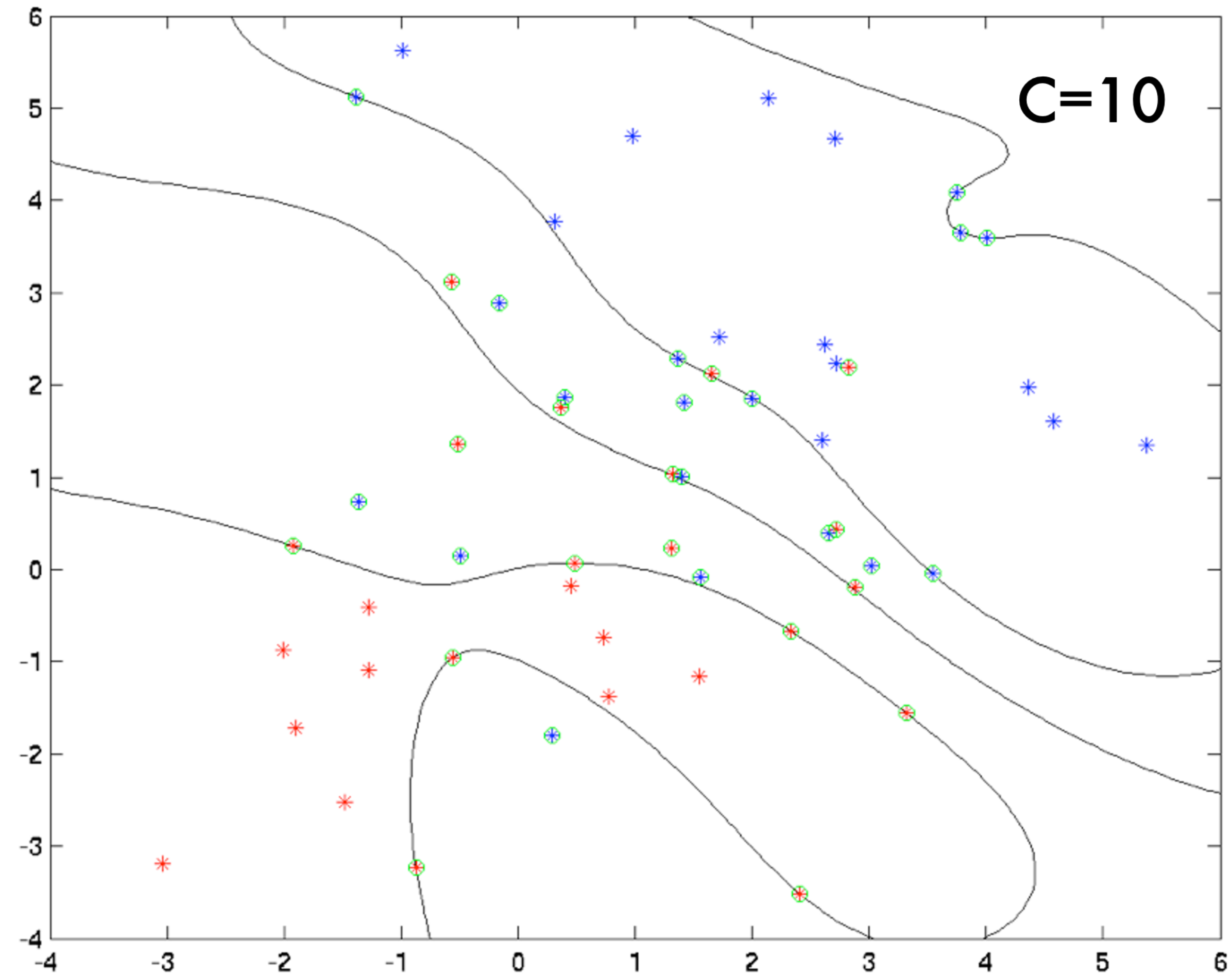
# Tuning Kernel SVM: Outliers

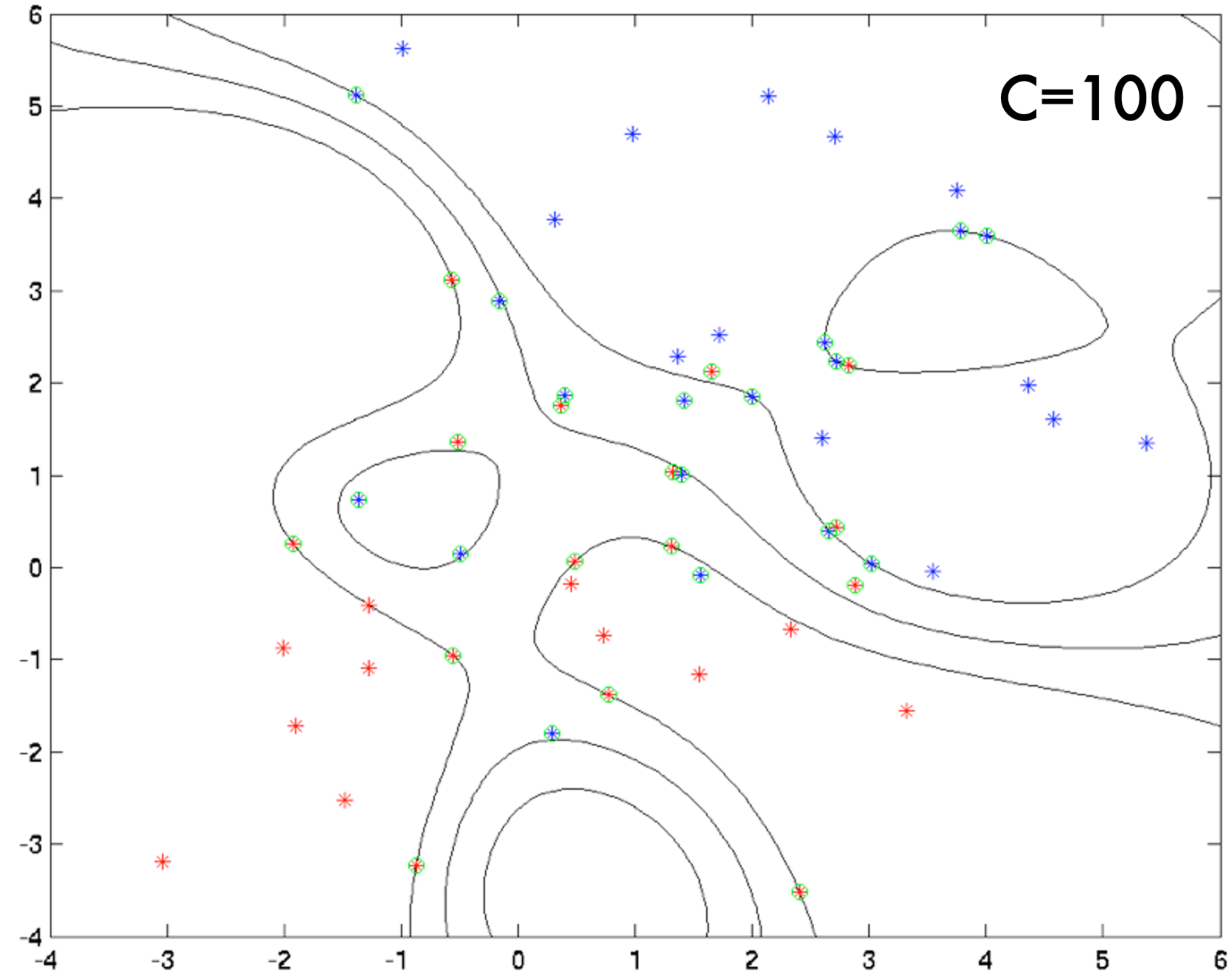# Tuning Kernel SVM: Outliers
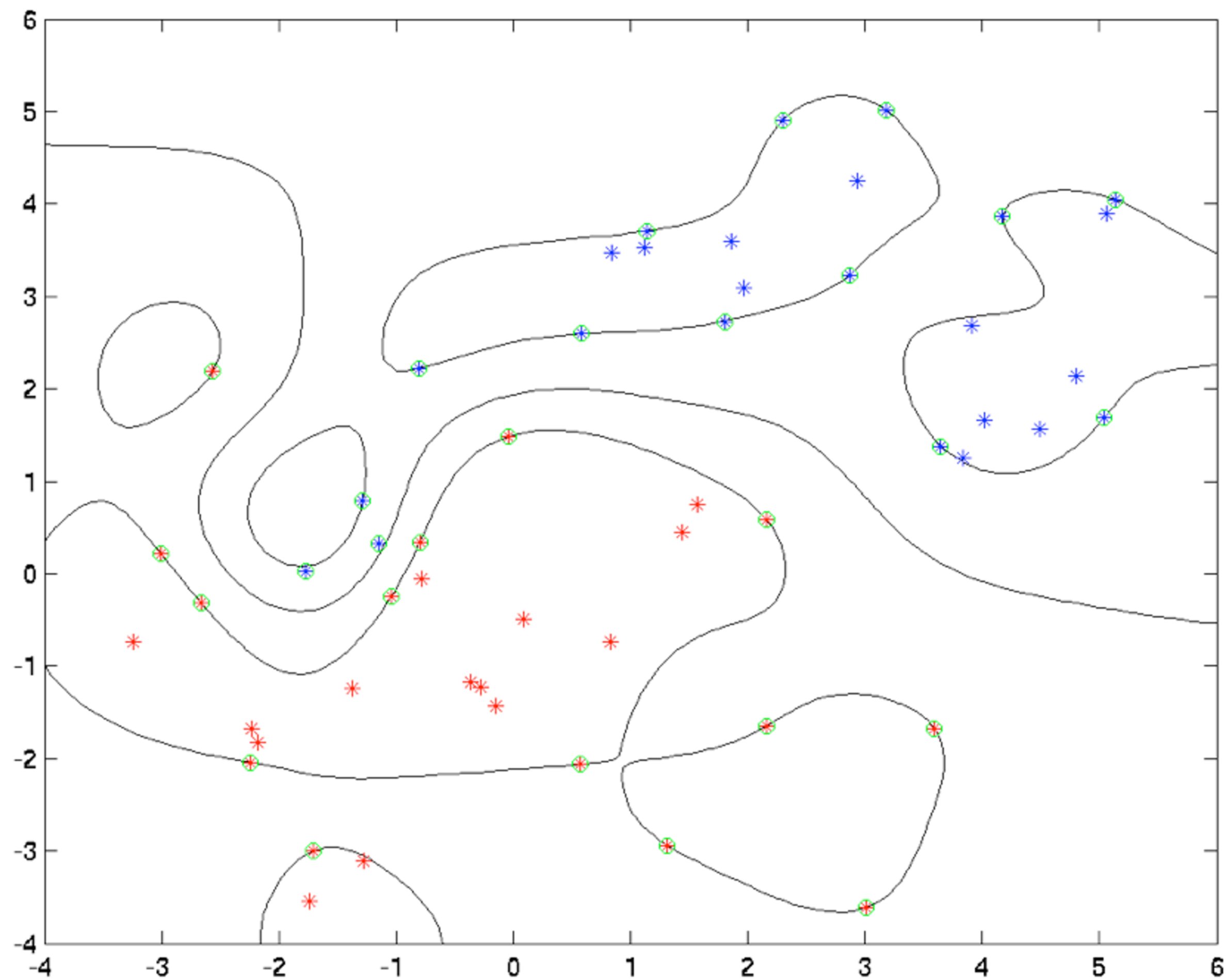
# Tuning Kernel SVM: Outliers
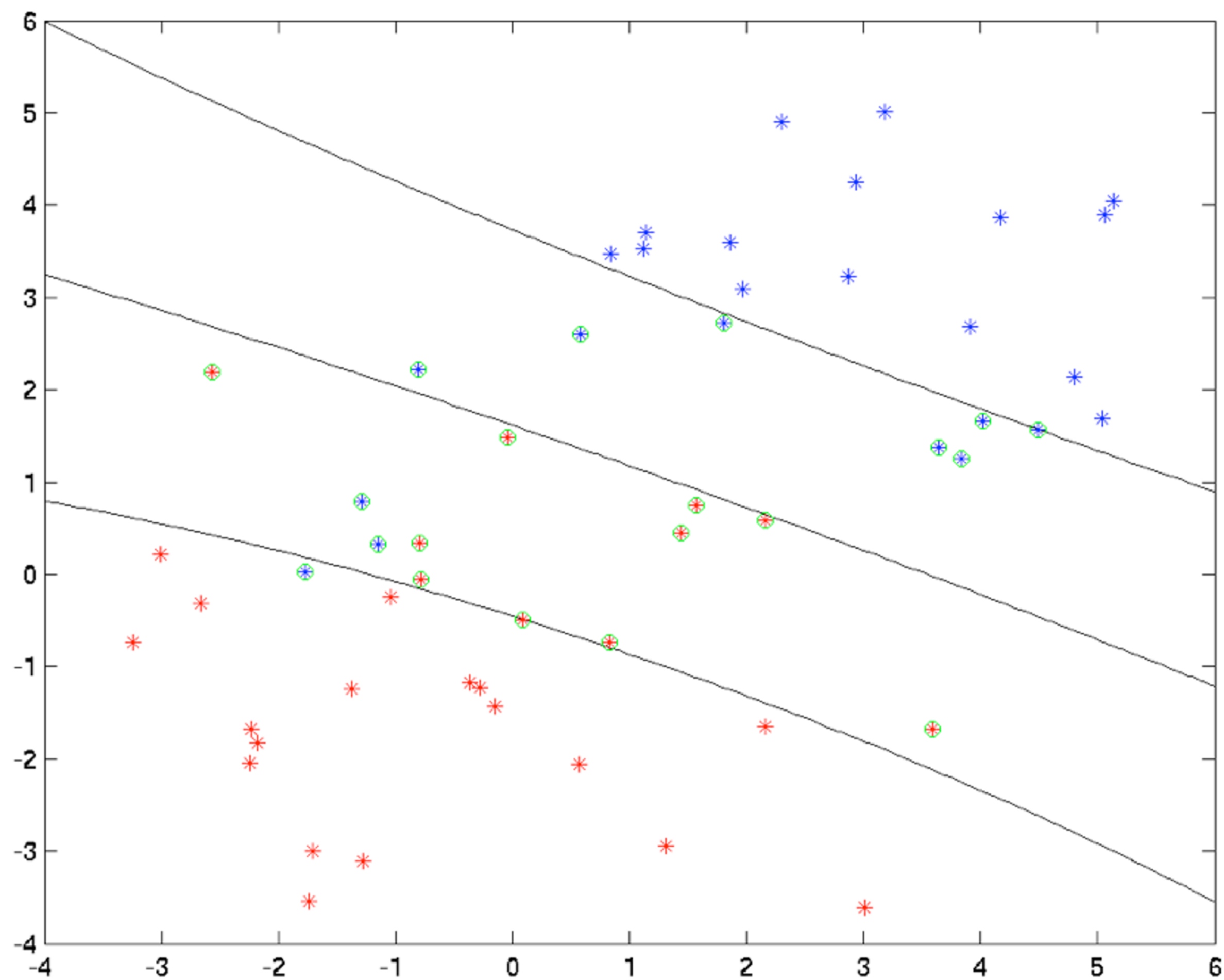
# Tuning Kernel SVM: Outliers

# Tuning Kernel SVM: Outliers

# Tuning Kernel SVM: Narrow Kernels

# Tuning Kernel SVM: Wide Kernels

# In deep learning era...

- In modern ML, we find a nice $\Phi(\,\cdot\,)$ using data + neural nets

  - Expensive, but we can afford them

  - Conduct logistic regression, instead of SVD

    - Ease of joint training

    - Also margin-maximizer (sometimes)

  - Use nice augmentations to find good similarity metric such that

    - $\Phi(\mathbf{x}) - \Phi(\mathbf{x}_{\text{aug}})$ is smaller than $\Phi(\mathbf{x}) - \Phi(\mathbf{x}')$

# Next up

- K-Means

Cheers