# Introduction to Machine Learning
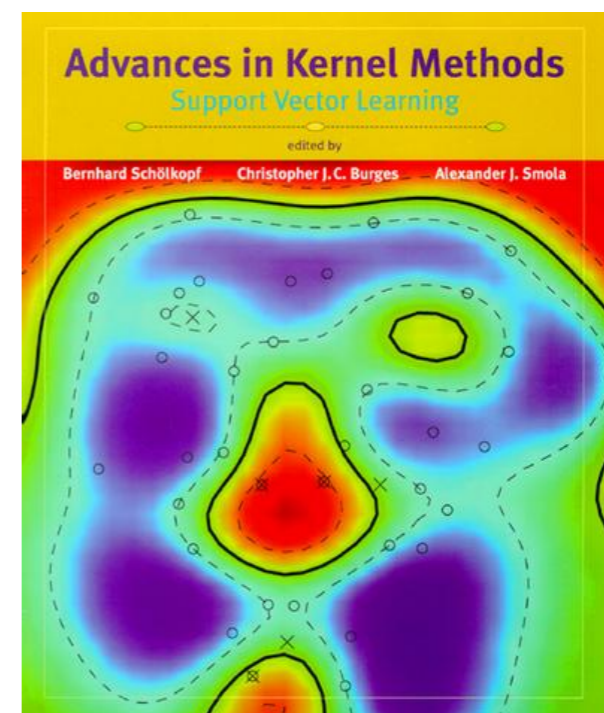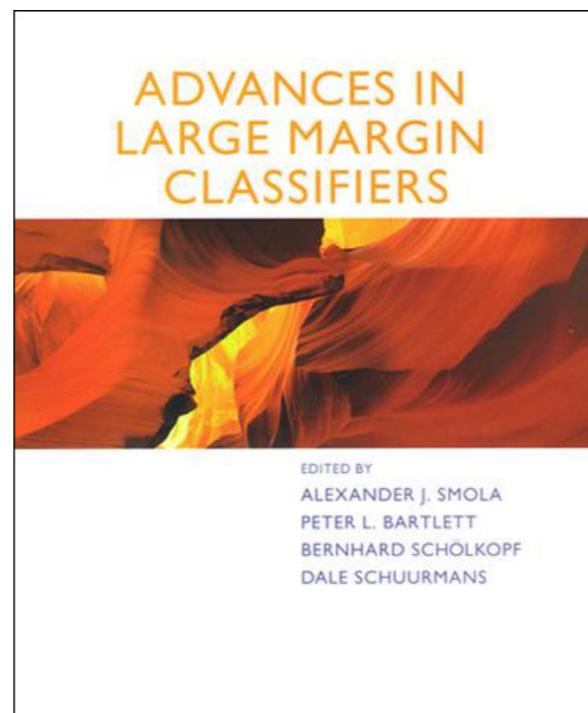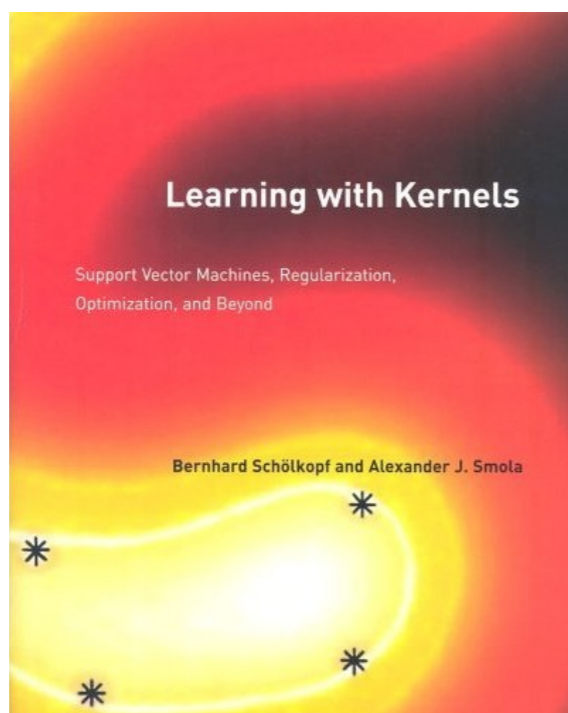## 5. Support Vector Classification

Alex Smola
Carnegie Mellon University

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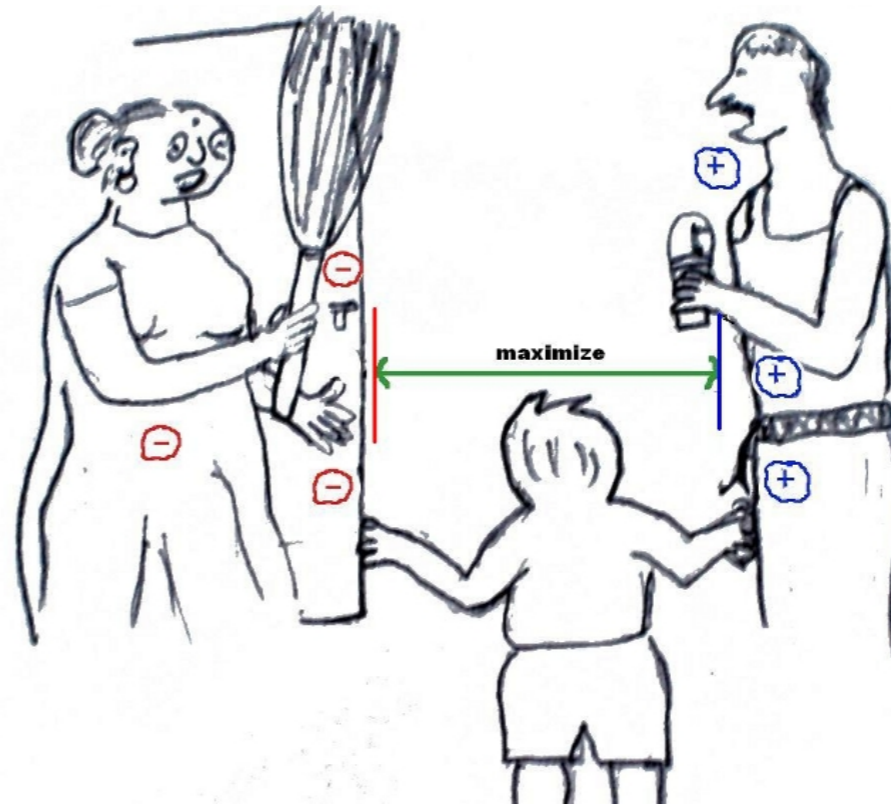10-701

Carnegie Mellon University

# Outline

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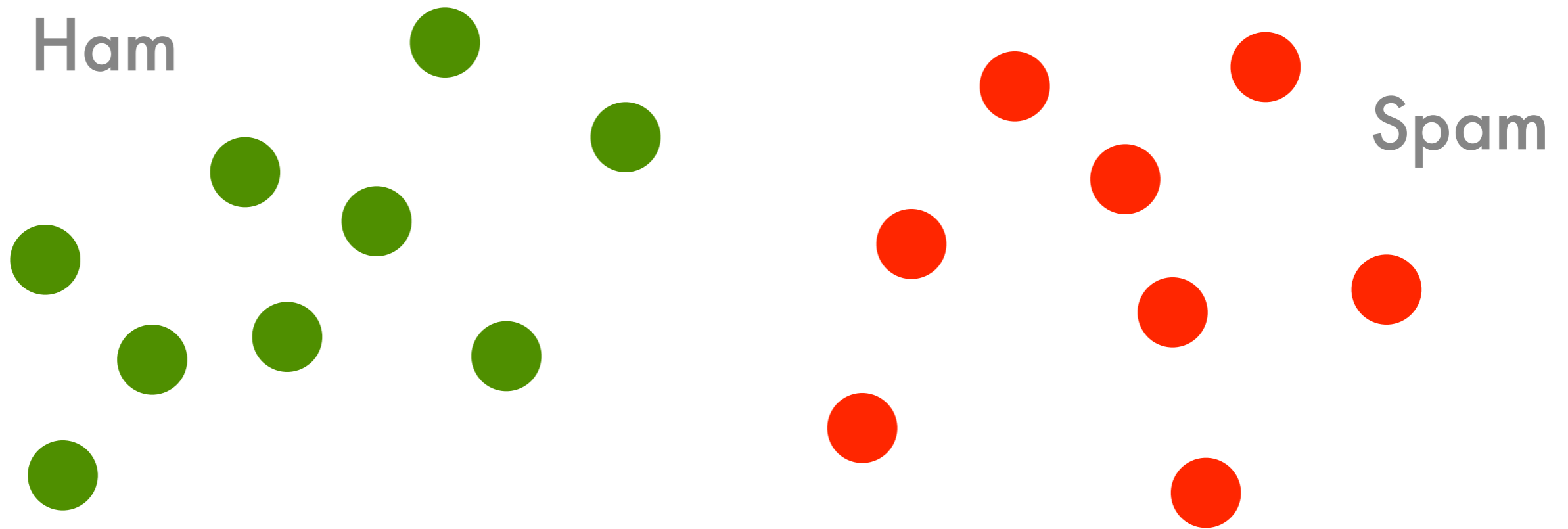


Learning with Kernels

Support Vector Machines, Regularization, Optimization, and Beyond

Bernhard Schölkopf and Alexander J. Smola

ADVANCES IN LARGE MARGIN CLASSIFIERS

EDITED BY
ALEXANDER J. SMOLA
PETER L. BARTLETT
BERNHARD SCHÖLKOPF
DALE SCHUURMANS

Advances in Kernel Methods
Support Vector Learning

edited by
Bernhard Schölkopf    Christopher J. C. Burges    Alexander J. Smola

lon University

Support Vector Machines

Carnegie Mellon University
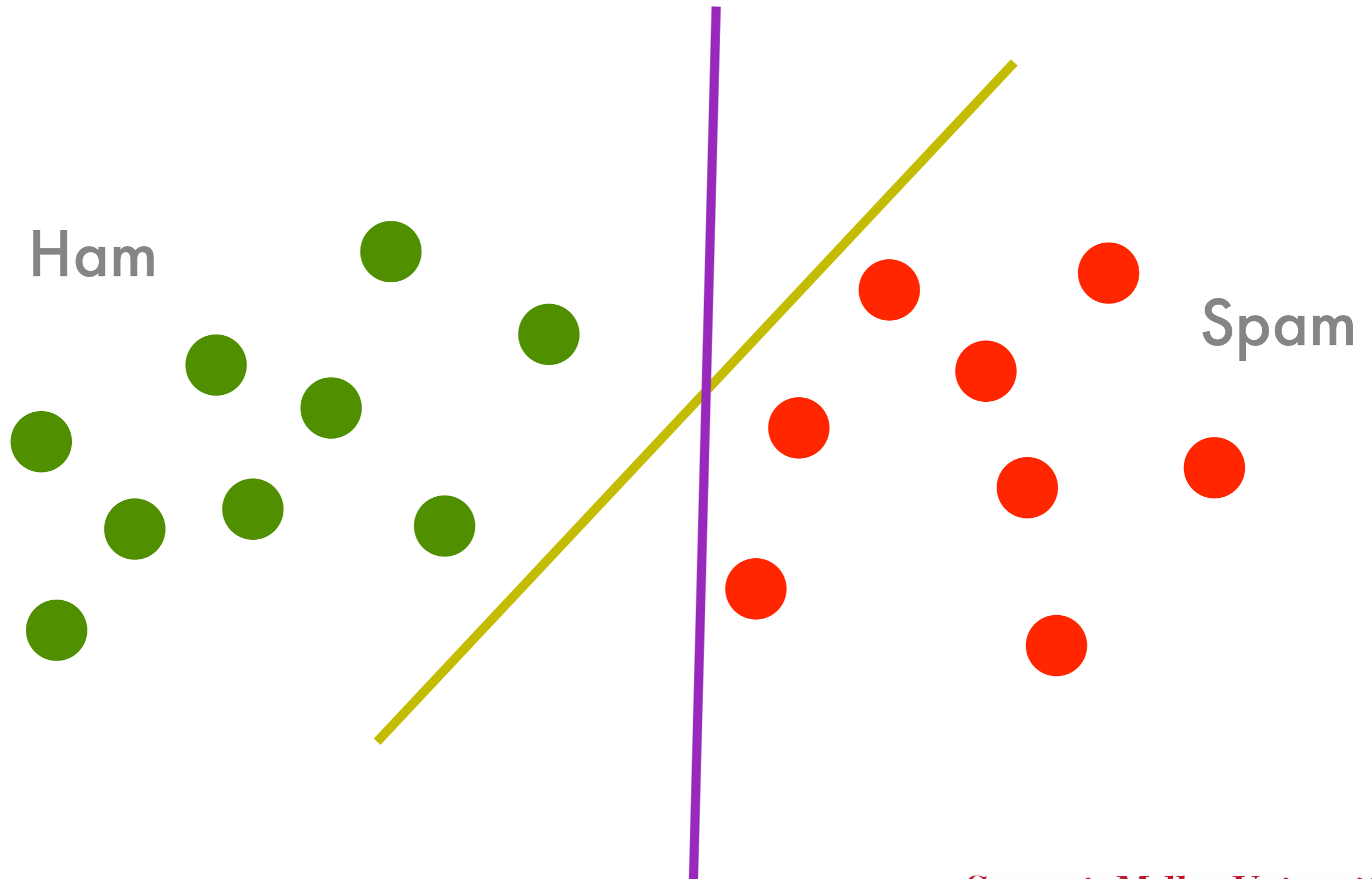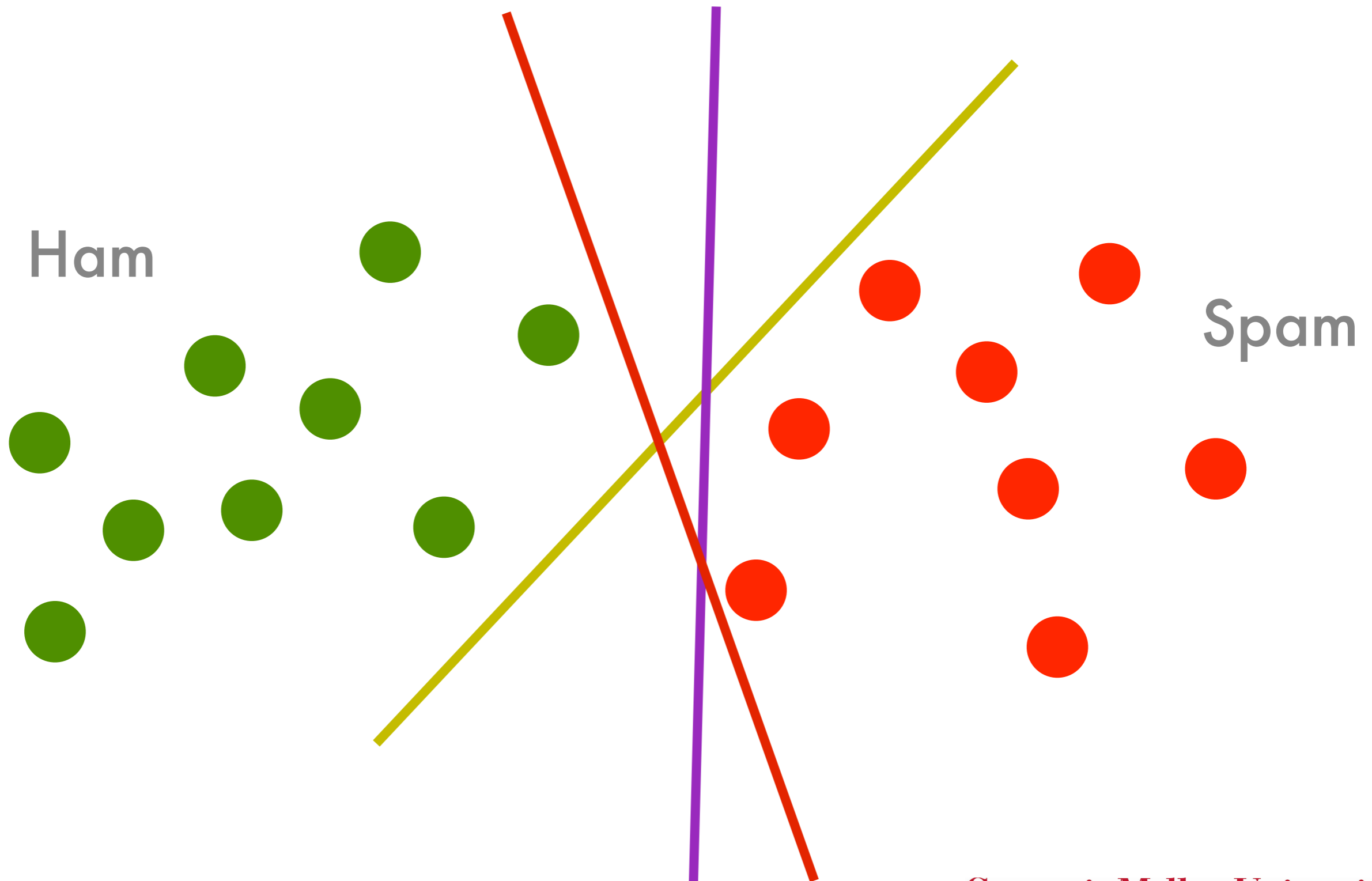
# Linear Separator

Ham

Spam

# Linear Separator

Ham

Spam

# Linear Separator

# Linear Separator

Ham

Spam

# Linear Separator



Ham
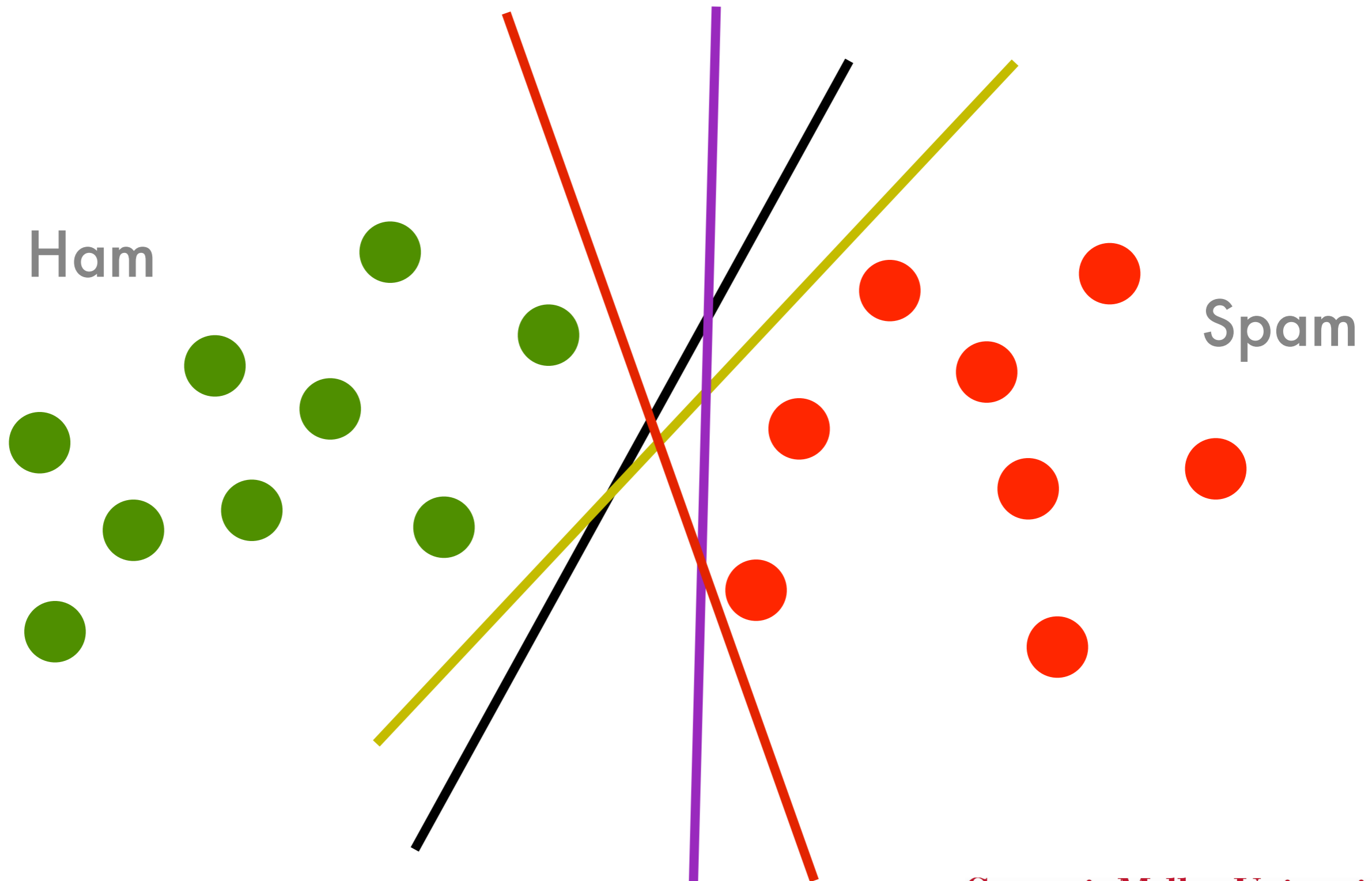
Spam
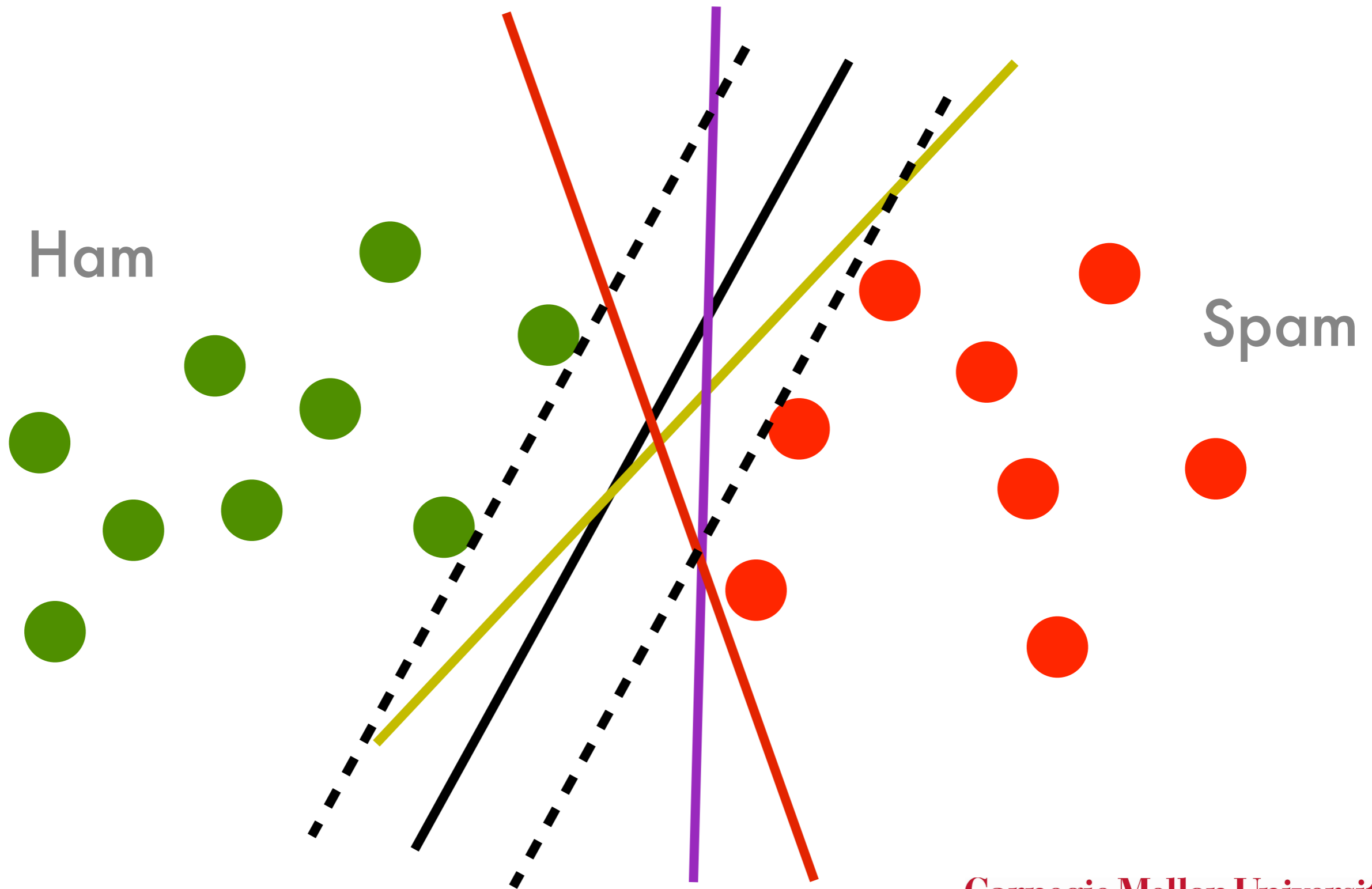
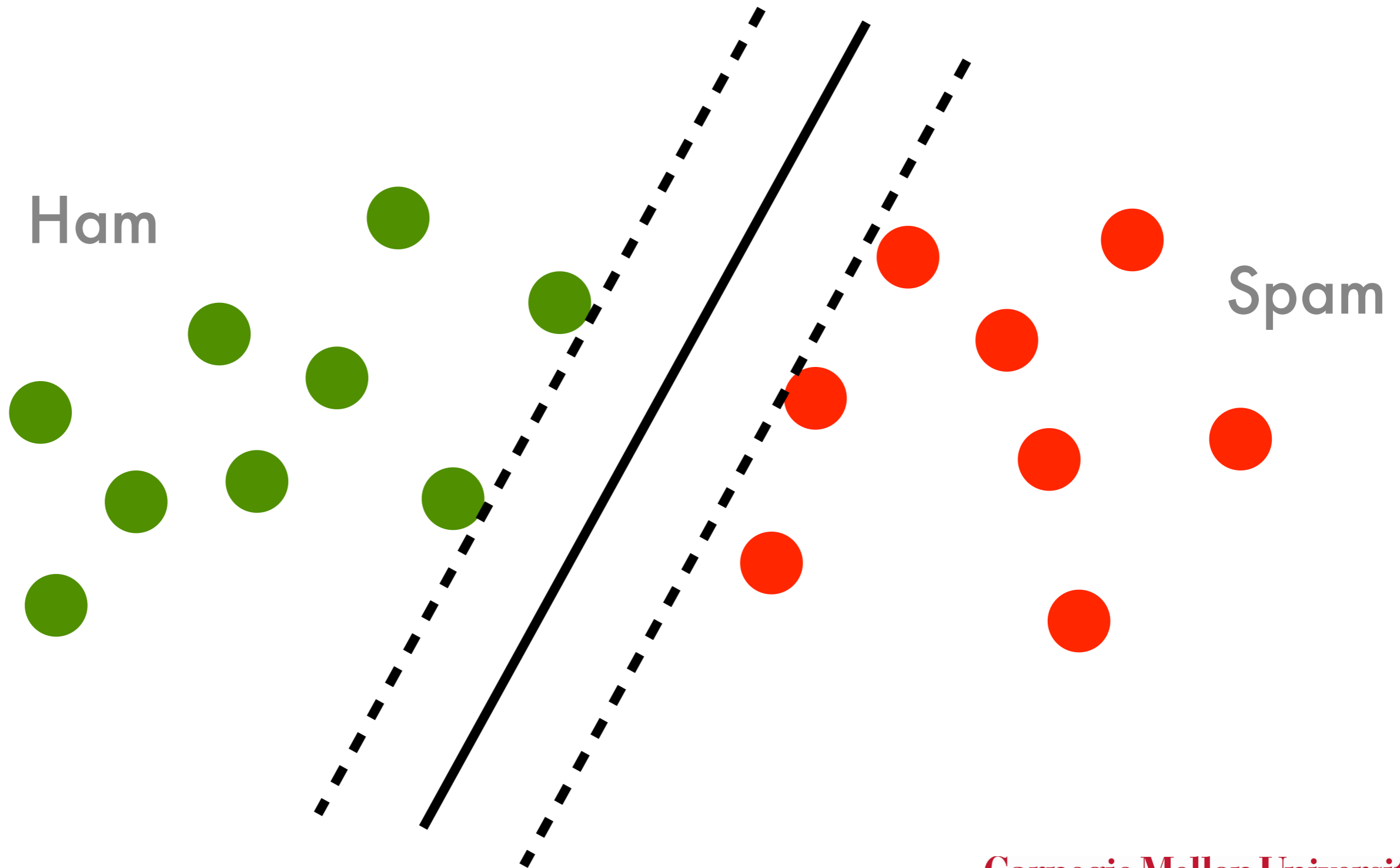# Linear Separator

# Linear Separator

# Large Margin Classifier
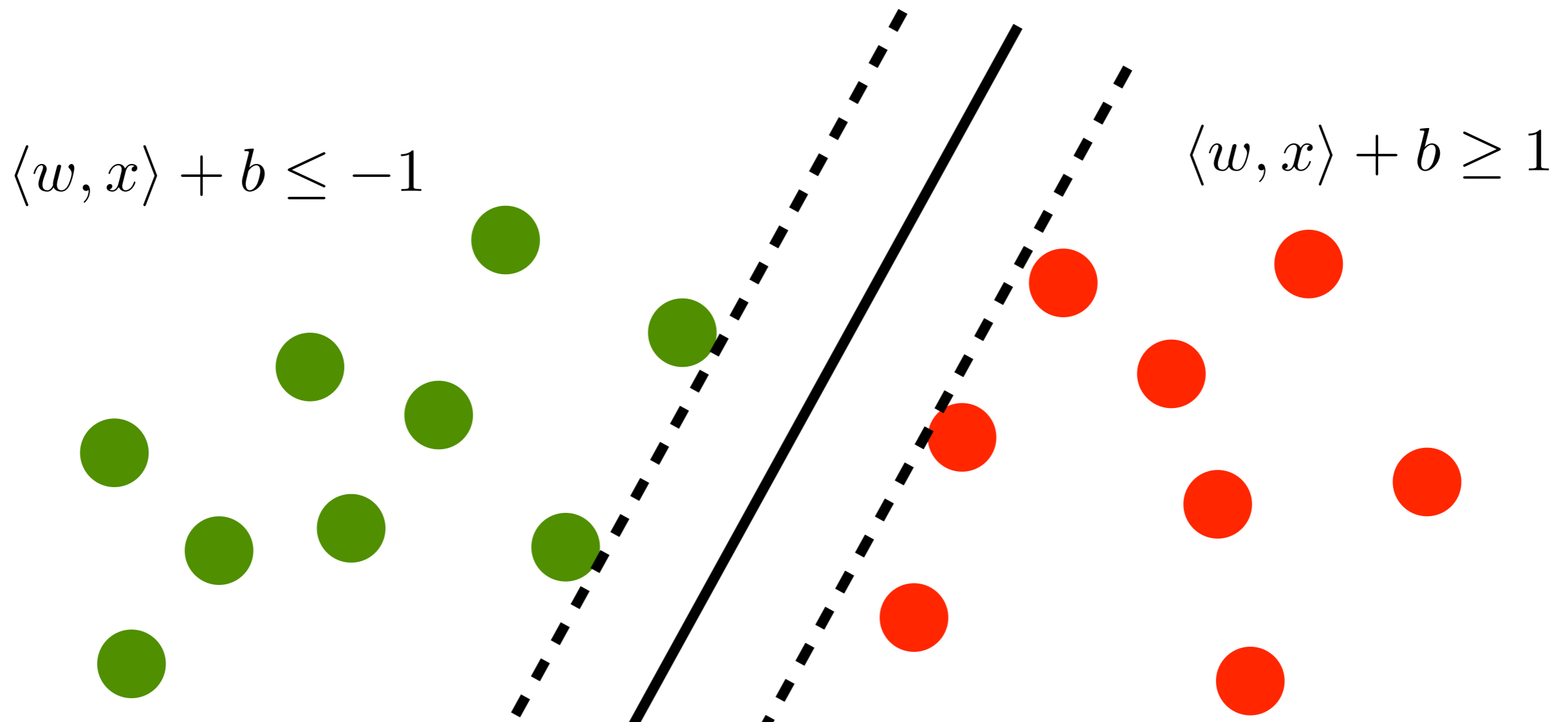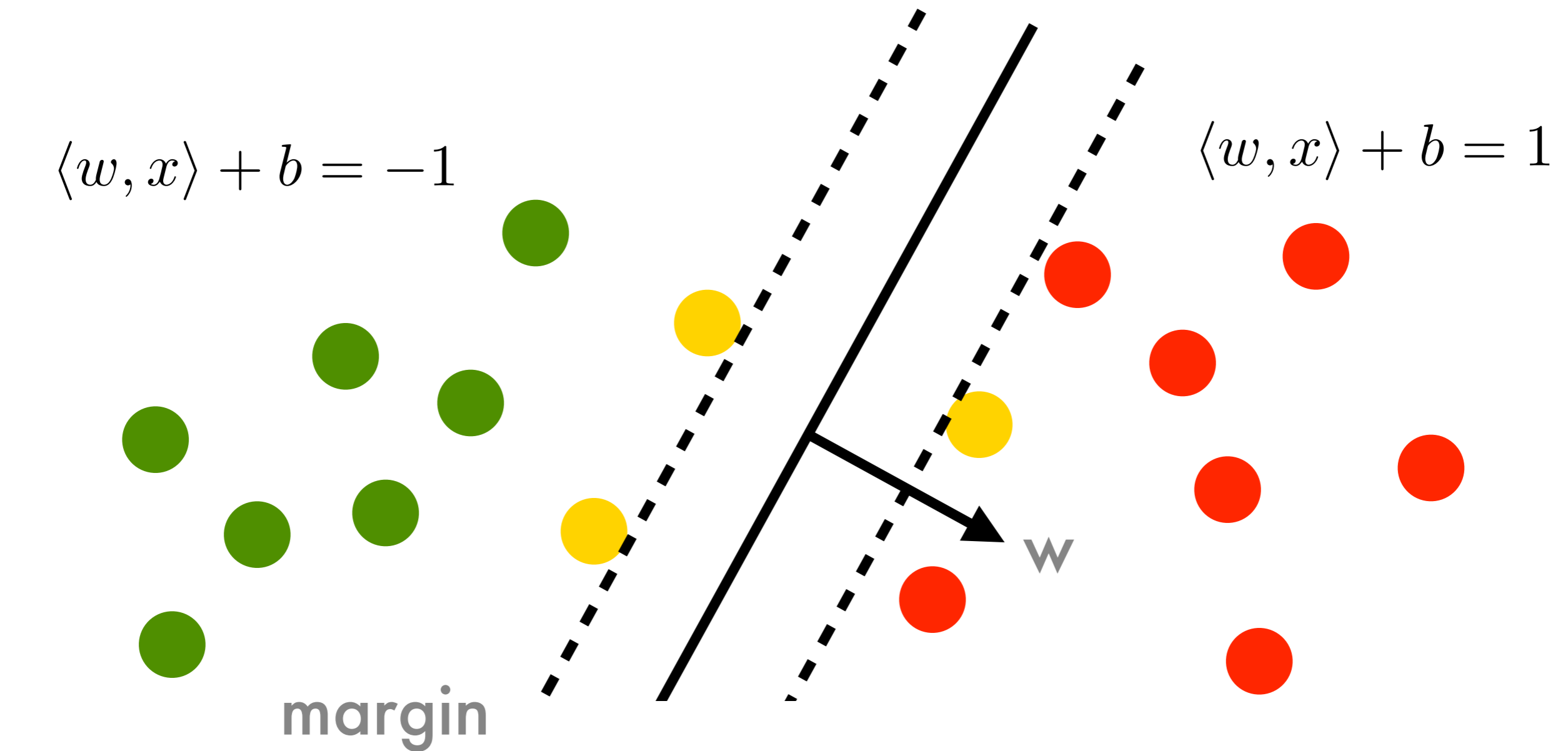
$\langle w, x \rangle + b \leq -1$

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

linear function

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

# Large Margin Classifier

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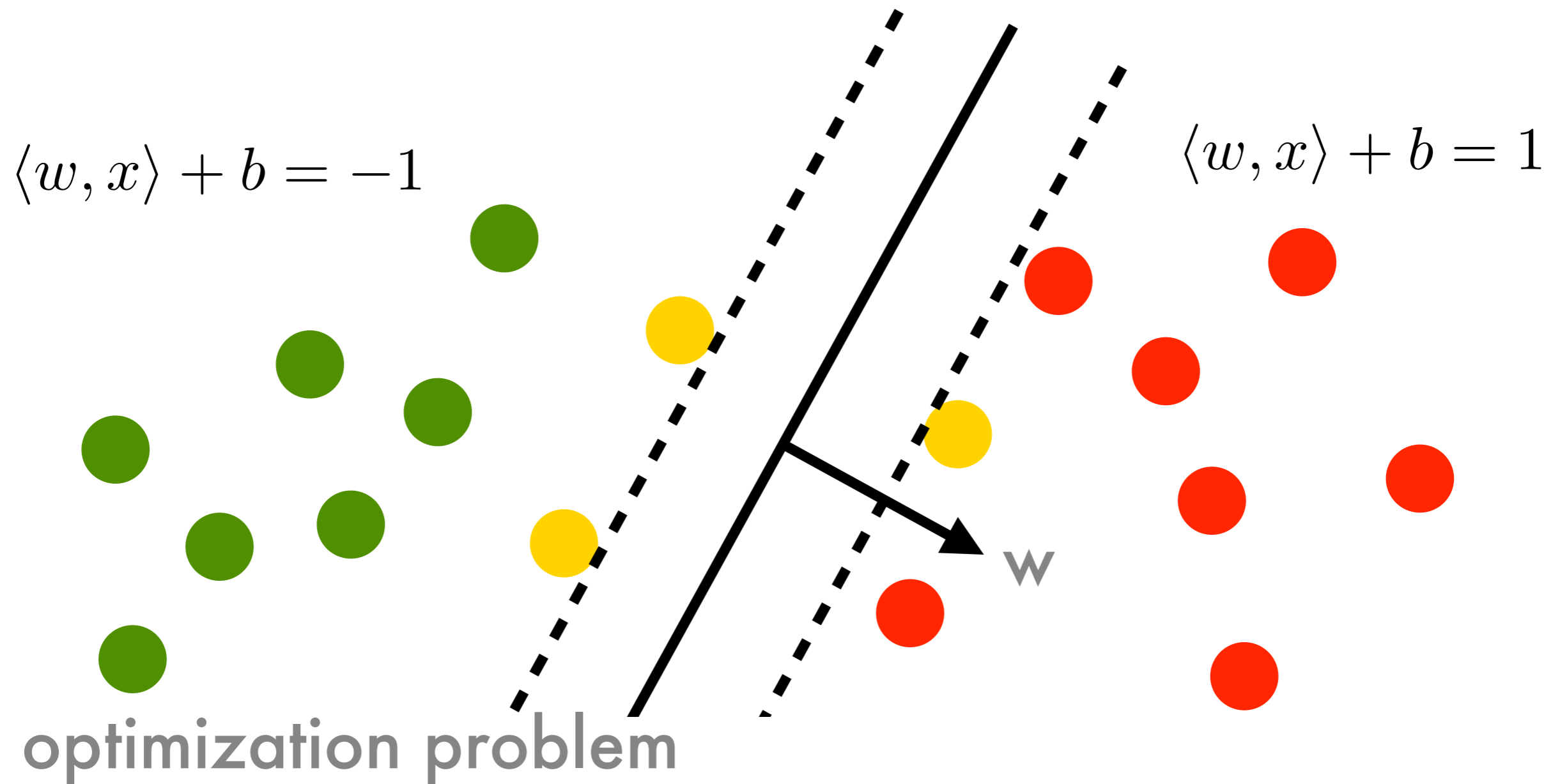
$$\langle w, x \rangle + b = 1$$

**w**

**margin**

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

# Large Margin Classifier



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

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

**w**

optimization problem

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

# Large Margin Classifier

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

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

w

**optimization problem**

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

**Carnegie Mellon University**

# Dual Problem

- **Primal optimization problem**

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

- **Lagrange function**

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

constraint

Optimality in w, b is at saddle point with α

- Derivatives in w, b need to vanish

# Dual Problem

- Lagrange function

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

- Derivatives in w, b need to vanish

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

$$\partial_b L(w, b, a) = \sum_i \alpha_i y_i = 0$$
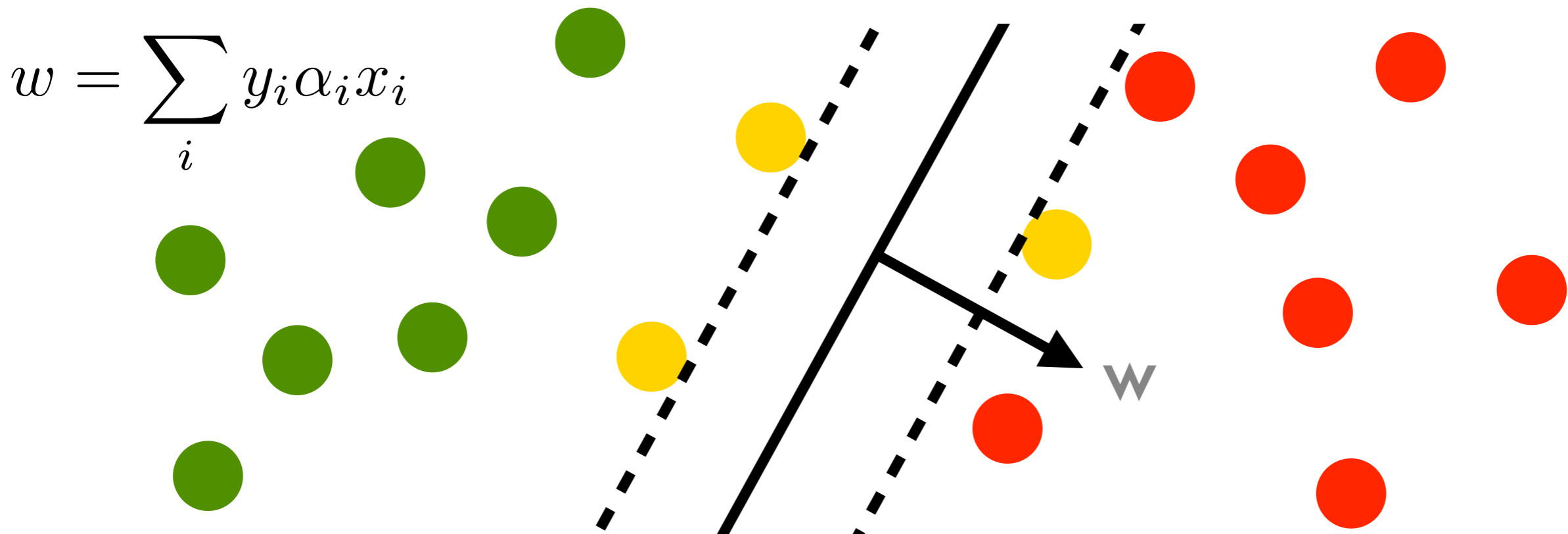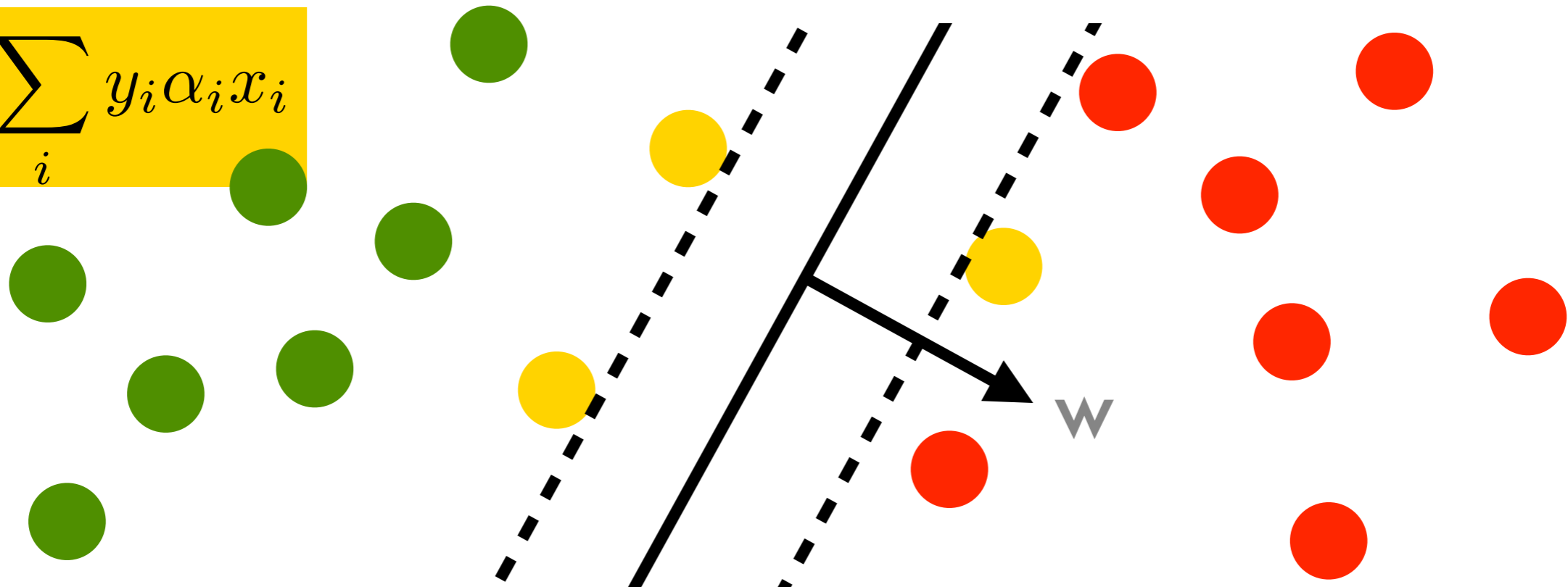
- Plugging terms back into L yields

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

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

# Support Vector Machines

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

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

w

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

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

# Support Vectors

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

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



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

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

# 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

Carnegie Mellon University

# Large Margin Classifier

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

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

linear function

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

# Large Margin Classifier

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

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

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

# 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

Carnegie Mellon University

# Large Margin Classifier

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

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



**Theorem (Minsky & Papert)**

Finding the minimum error separating hyperplane is NP hard

# Large Margin Classifier

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

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

**Theorem (Minsky & Papert)**

Finding the minimum error separating hyperplane is NP hard

# Adding slack variables

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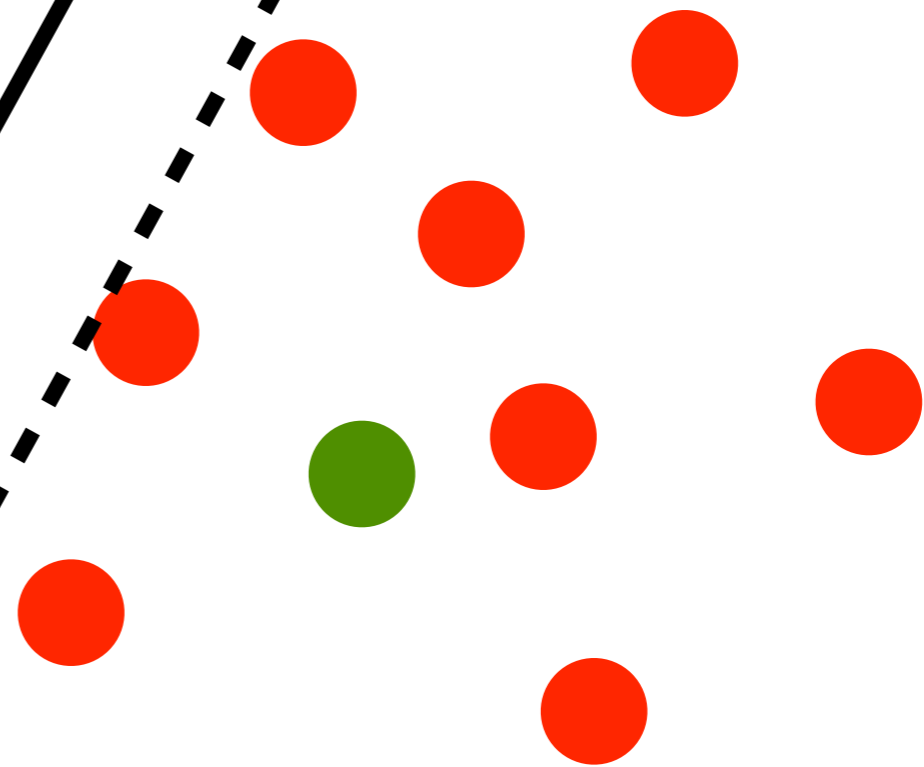
$$\langle w, x \rangle + b \geq 1 - \xi$$

**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

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

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

**minimize amount of slack**

Convex optimization problem

# 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

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

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

**minimize amount of slack**

**Convex optimization problem**

# Adding slack variables

- Hard margin problem

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

- With slack variables

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

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

Problem is always feasible. Proof:

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

# Dual Problem

- ## Primal optimization problem

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

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

- ## Lagrange function

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

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

- ## Derivatives in w,b,ξ need to vanish

# Dual Problem

- ## Lagrange function

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

- ## 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

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

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

**bound influence**

# Karush Kuhn Tucker Conditions

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

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

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



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

$$\eta_i \xi_i = 0$$

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

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

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

C=10

C=50

C=100

C=1

C=20

C=50

C=5

C=20

C=100

# Solving the optimization problem

- Dual problem

$$\underset{\alpha}{\text{maximize}} - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \langle x_i, x_j \rangle + \sum_i \alpha_i$$
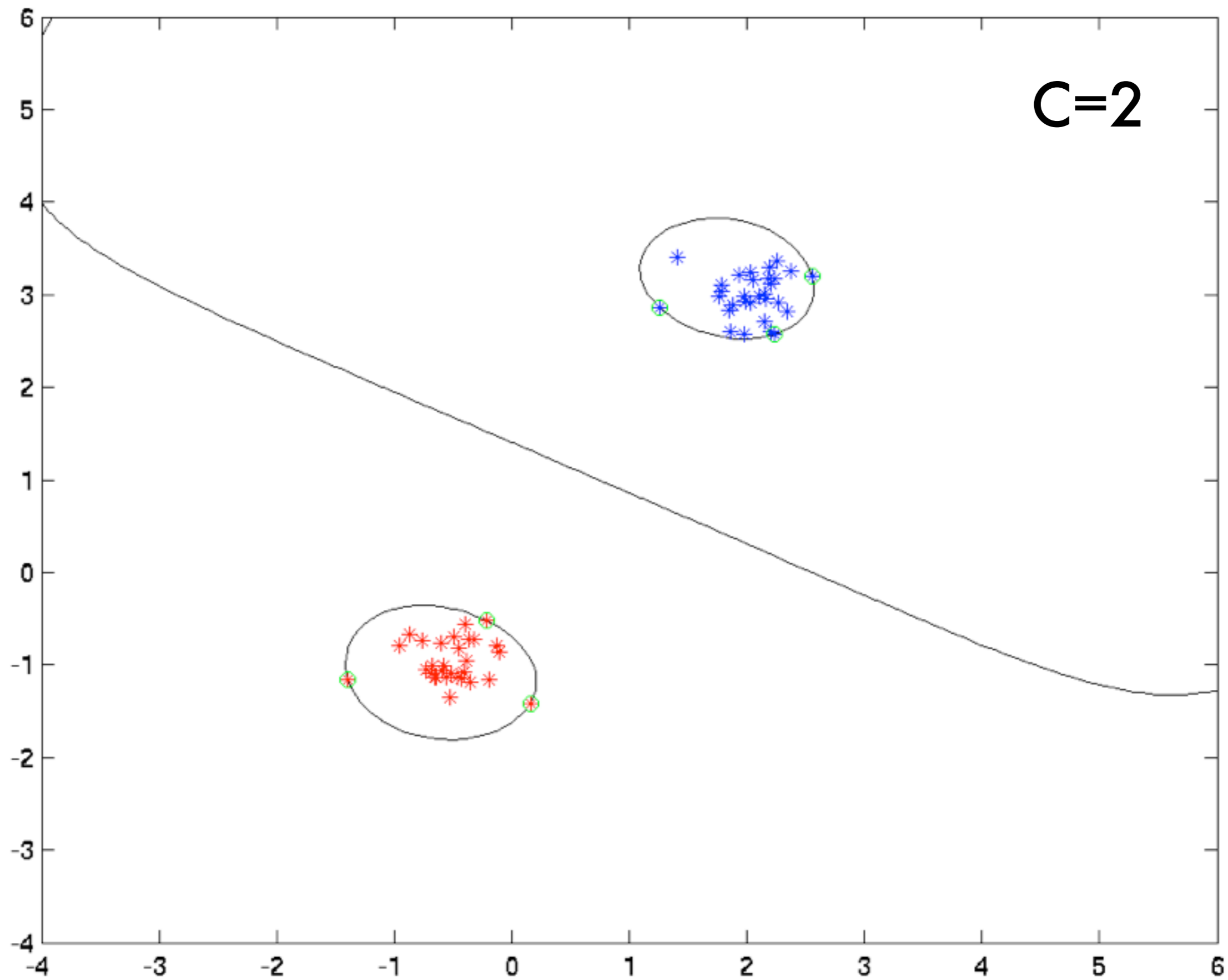
$$\text{subject to } \sum_i \alpha_i y_i = 0 \text{ and } \alpha_i \in [0, C]$$

- 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).

Nonlinear Separation

# The Kernel Trick

- **Linear soft margin problem**

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

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

- **Dual problem**

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

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

- **Support vector expansion**

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

# The Kernel Trick

- Linear soft margin problem

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

$$\text{subject to } y_i \left[ \langle w, \boxed{\phi(x_i)} \rangle + b \right] \geq 1 - \xi_i \text{ and } \xi_i \geq 0$$

- Dual problem

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

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

- Support vector expansion

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

C=2

C=5

C=50

C=100

C=10

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

# Risk and Loss

# Loss function point of view

- Constrained quadratic program

$$\underset{w,b}{\text{minimize}} \ \frac{1}{2} \|w\|^2 + C \sum_i \xi_i$$
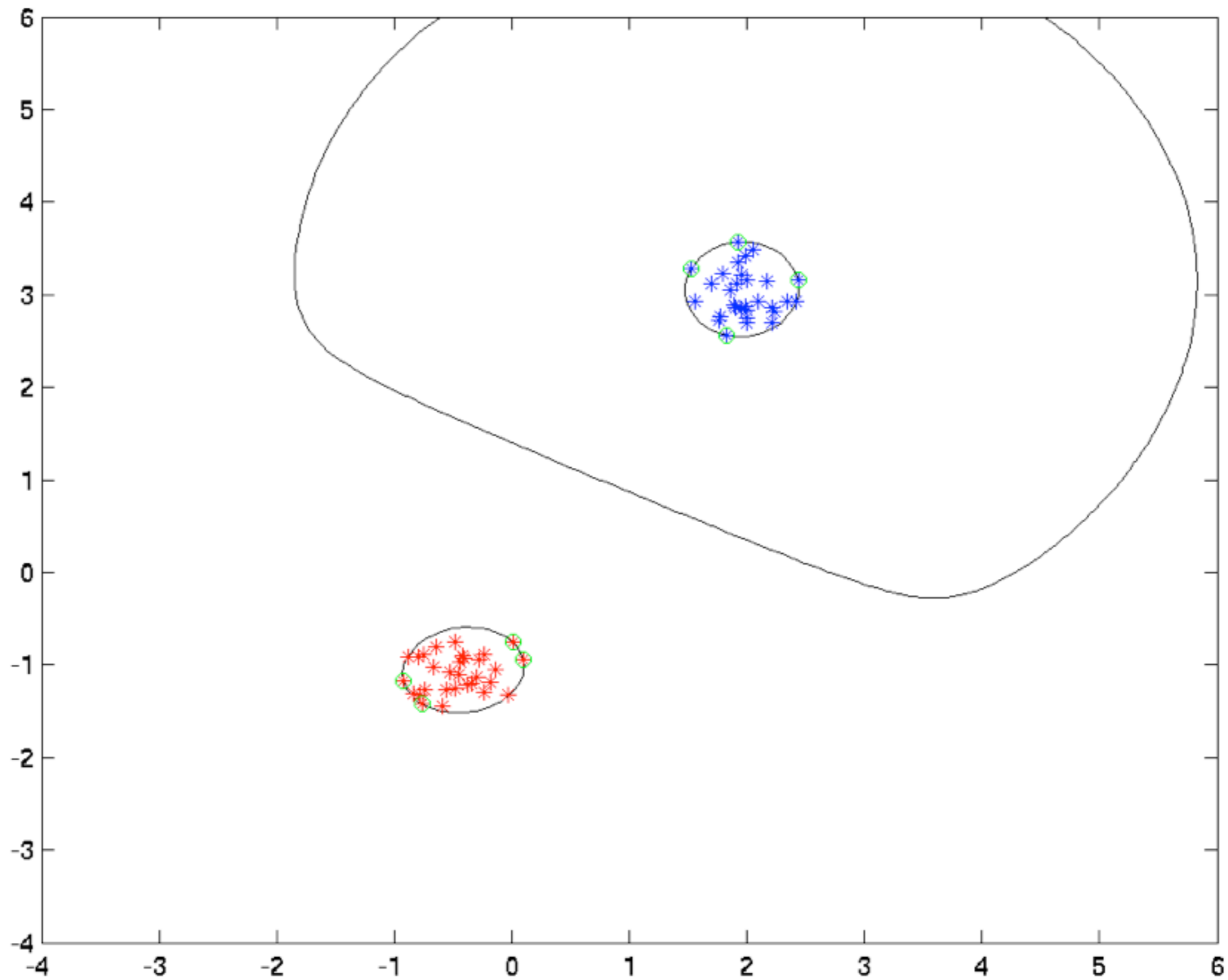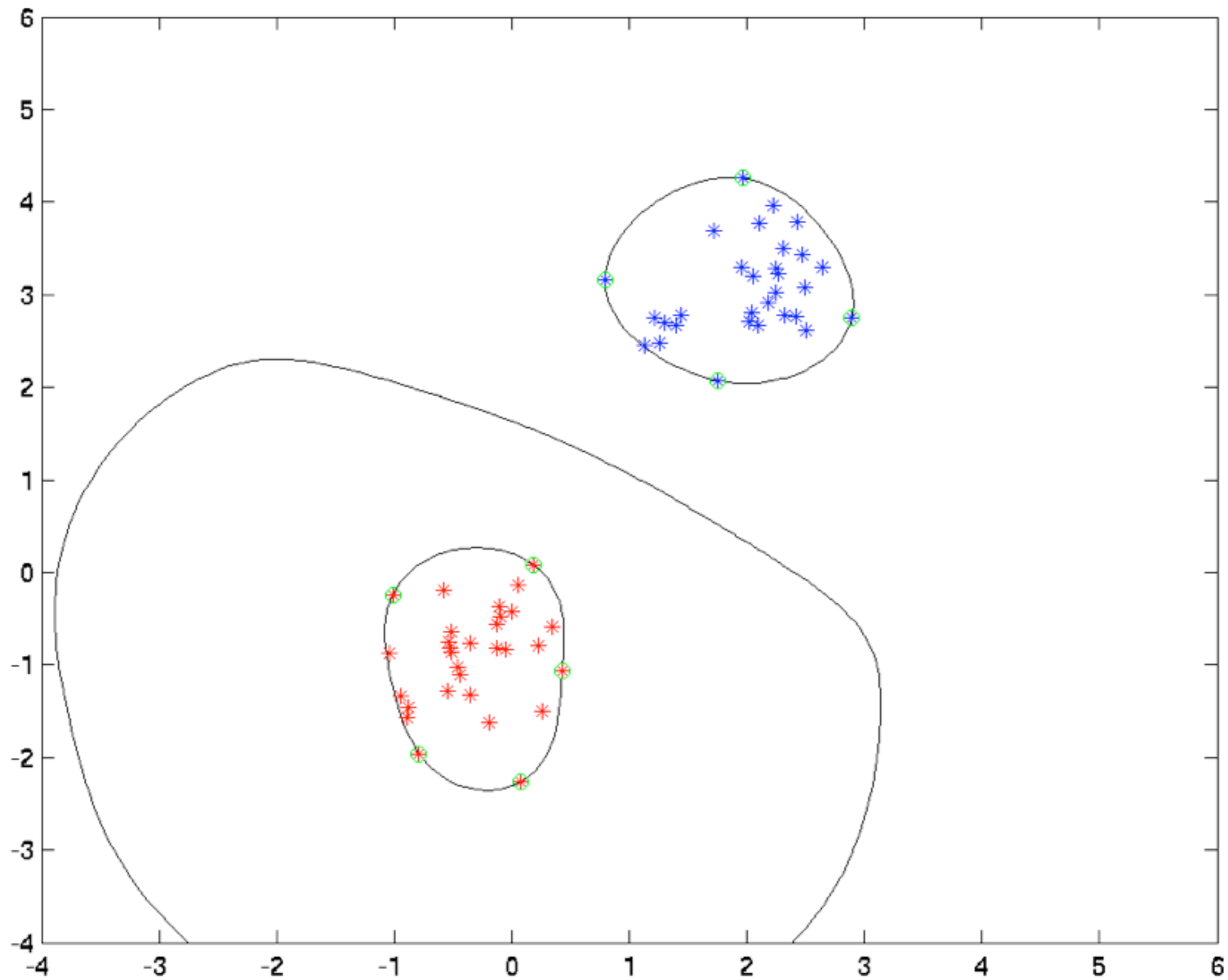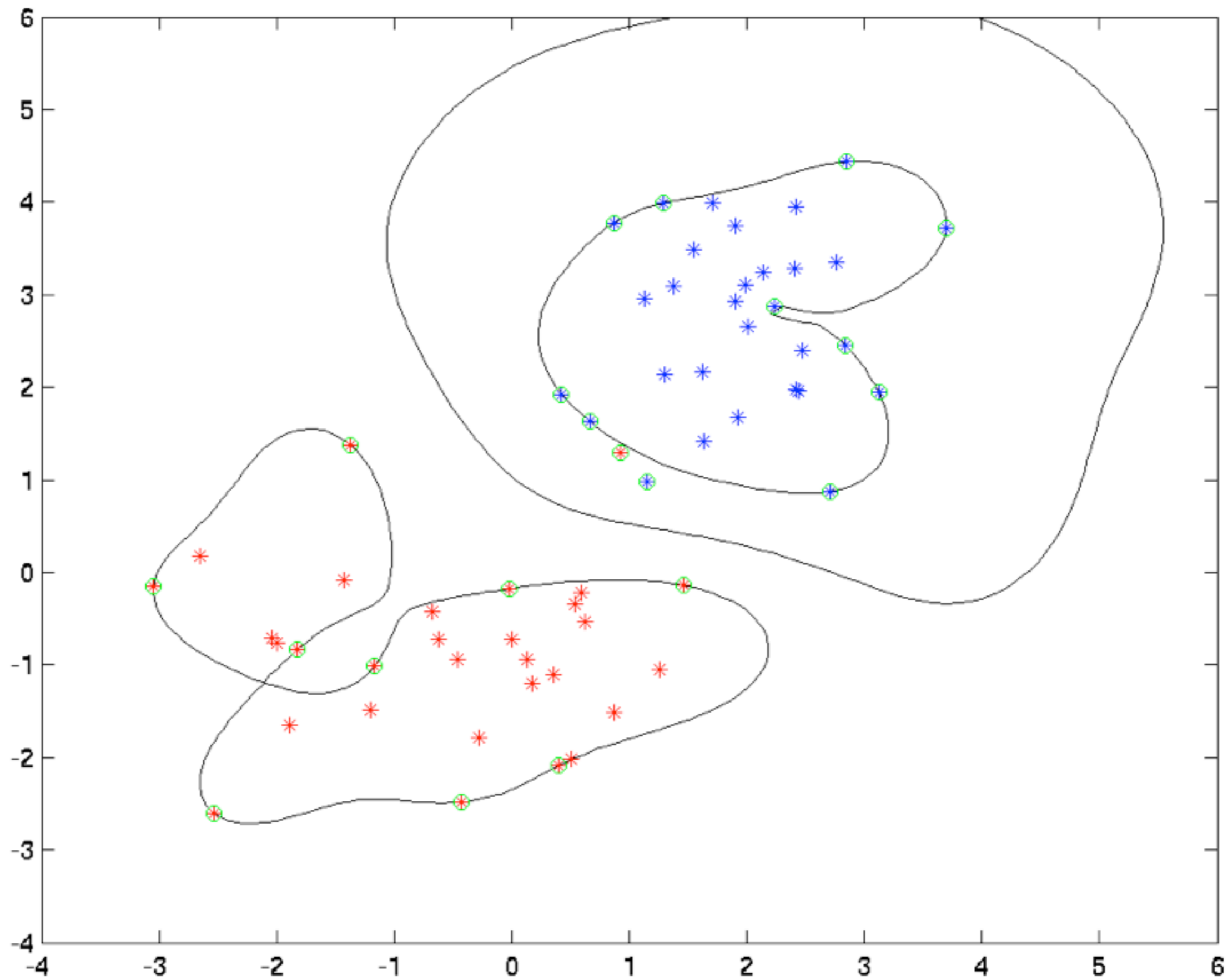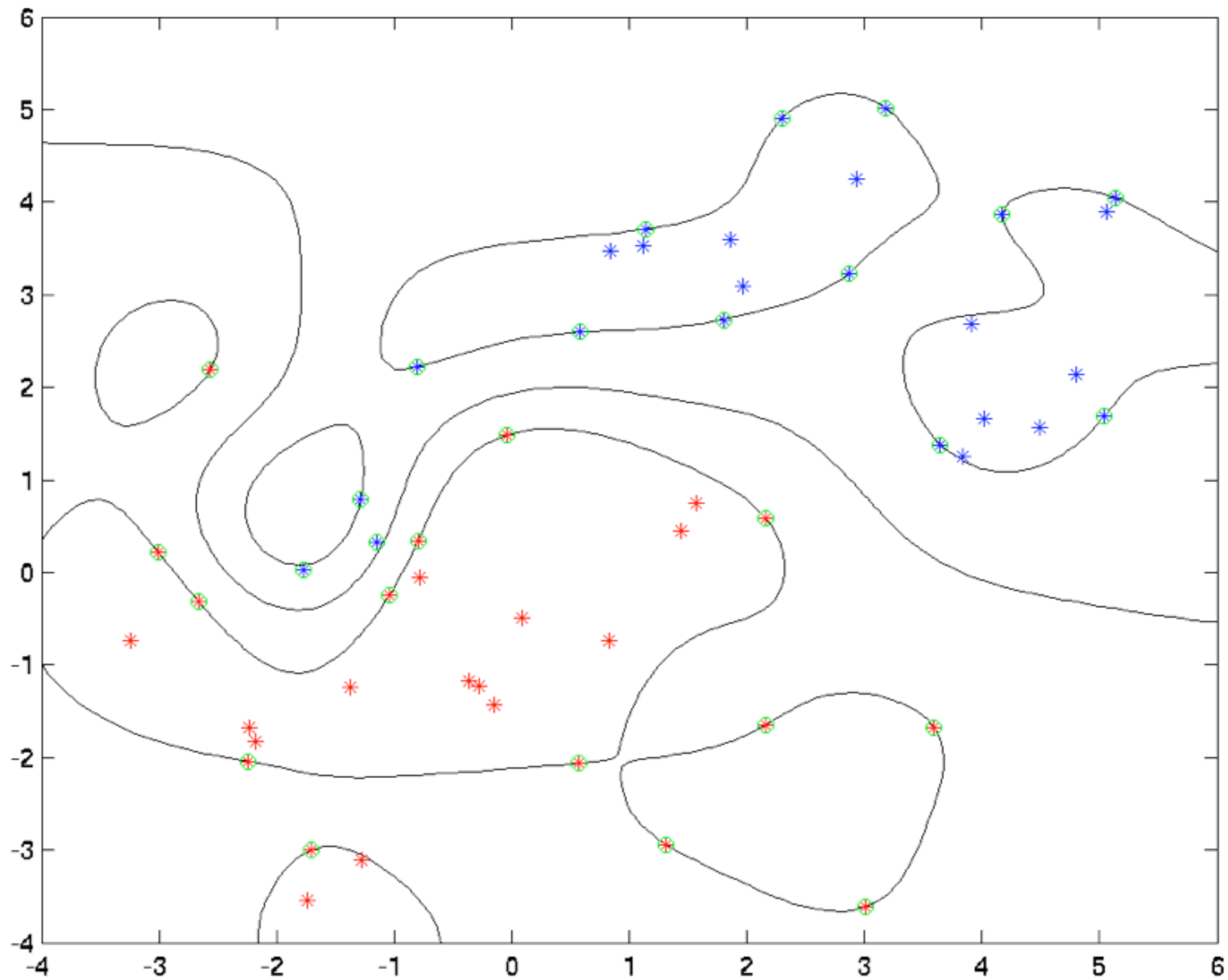
$$\text{subject to } y_i \left[ \langle w, x_i \rangle + b \right] \geq 1 - \xi_i \text{ and } \xi_i \geq 0$$

- Risk minimization setting

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

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** $\{yf(x) < 0\}$



convex upper bound

binary loss function

margin

# More loss functions

- **Logistic** $\log\left[1 + e^{-f(x)}\right]$

- **Huberized loss**

$$\begin{cases} 0 & \text{if } f(x) > 1 \\ \frac{1}{2}(1 - f(x))^2 & \text{if } f(x) \in [0, 1] \\ \frac{1}{2} - f(x) & \text{if } f(x) < 0 \end{cases}$$

- **Soft margin**

$$\max(0, 1 - f(x))$$

(asymptotically) linear

(asymptotically) 0

# Risk minimization view

- Find function f minimizing classification error

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

- Compute empirical average

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

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

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

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