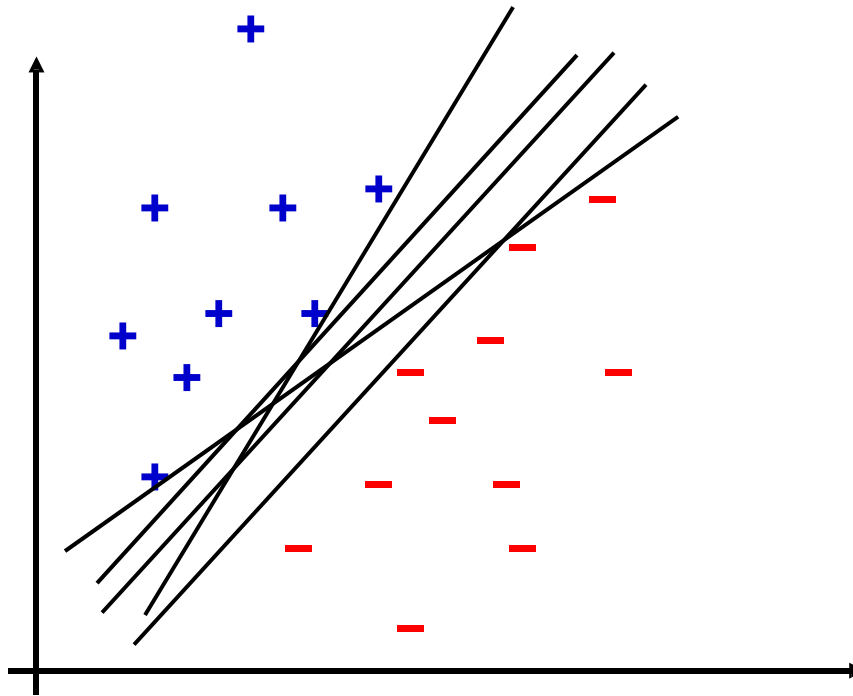# Support Vector Machines

**Key concepts**

- Functional and geometric margin of a classifier
- SVM objective: quadratic objective with linear constraints
- Constrained optimization: Lagrangian
- Primal and Dual problem, the KKT conditions
- Solution characteristics of SVM
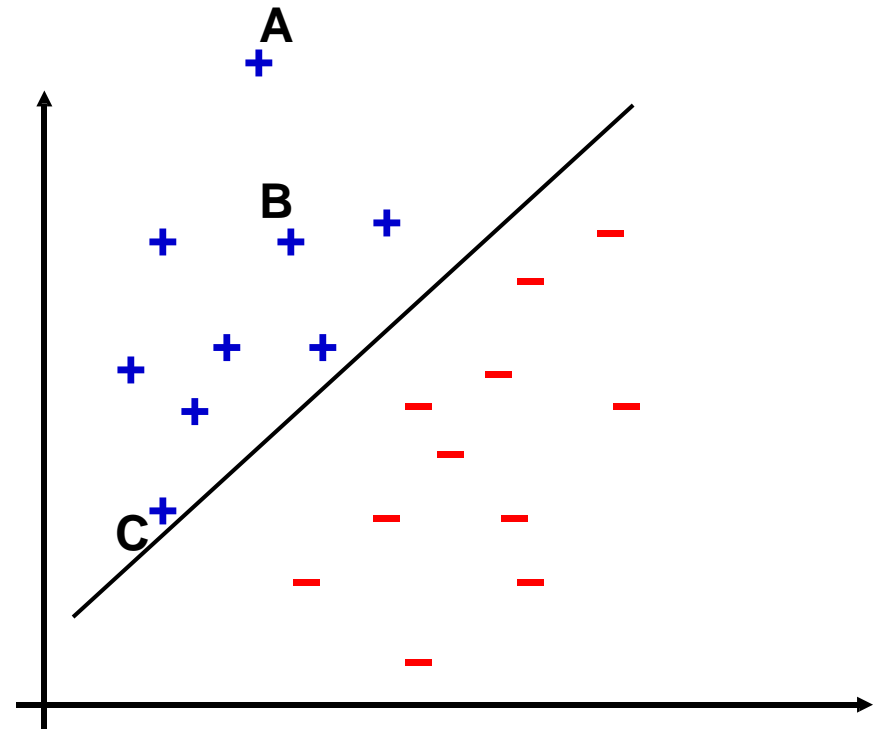- Support vectors
- Kernel SVM

# Linear Separators

- Which of the linear separators is optimal?

# Intuition of Margin

- Consider points A, B, and C
- We are quite confident in our prediction for A because it is far from the decision boundary.
- In contrast, we are not so confident in our prediction for C because a slight change in the decision boundary may flip the decision.



Given a training set, we would like to make all predictions correct and confident! This leads to the concept of margin.

# Functional Margin

- Given a linear classifier parameterized by $(\mathbf{w}, b)$, we define its functional margin w.r.t training example $(\mathbf{x}^i, y^i)$ as:

$$\hat{\gamma}^i = y^i(\mathbf{w}^T\mathbf{x}^i + b)$$

- If we rescale $(\mathbf{w}, b)$ by a factor $\alpha$, functional margin gets multiplied by $\alpha$

  - we can make it arbitrarily large without change anything meaningful
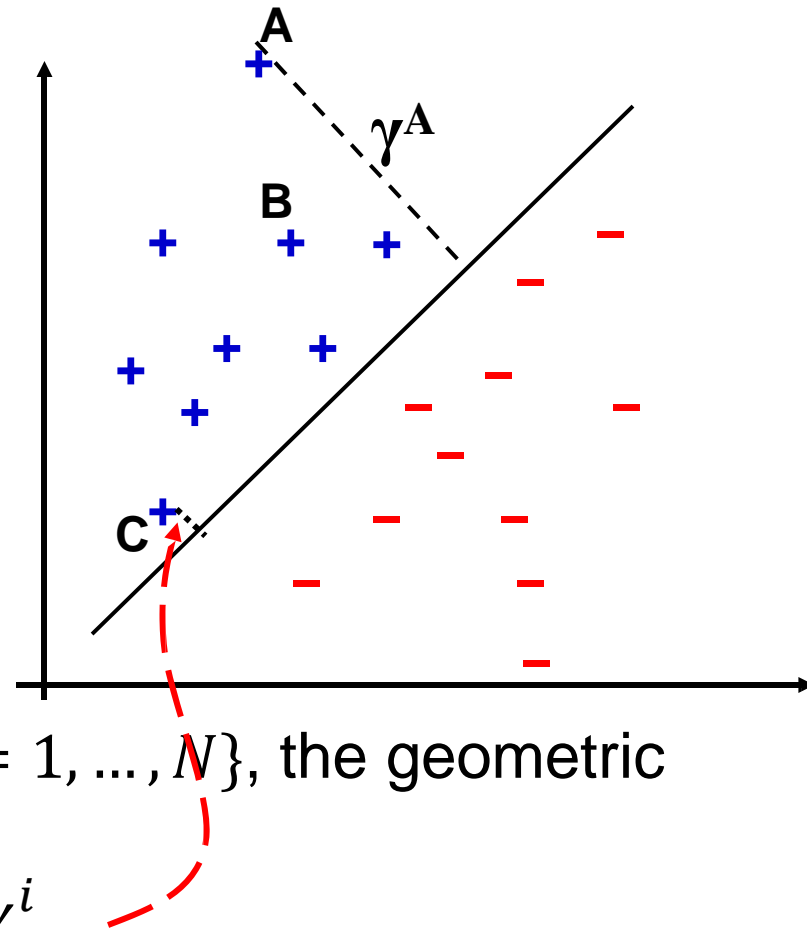  - Instead, we will look at *geometric margin*

# Geometric Margin

- The geometric margin of $(\mathbf{w}, b)$ w.r.t. $\mathbf{x}^i$ is the distance from $\mathbf{x}^i$ to the decision boundary

- This distance can be computed as

$$\gamma^i = \frac{y^i(\mathbf{w}^T\mathbf{x} + b)}{\|\mathbf{w}\|}$$



- Given training set $S = \{(\mathbf{x}^i, y^i): i = 1, ..., N\}$, the geometric margin of the classifier w.r.t. $S$ is

$$\gamma = \min_{i=1,...,N} \gamma^i$$

Points closest to the boundary are called Support vectors – we will see that these are the points that really matters

# Maximum Margin Classifier

- Given a ***linearly separable*** training set $S = \{(\mathbf{x}^i, y^i): i = 1, \ldots, N\}$, we would like to find a linear classifier with the maximum margin.

- This can be represented as an optimization problem.

$$\max_{w,b,\gamma} \gamma$$

Nasty optimization problem! Let's make it look nicer!

$$\text{subject to:} \frac{y^i(\mathbf{w}^T\mathbf{x}^i + b)}{\|\mathbf{w}\|} \geq \gamma$$

- Let $\gamma' = \gamma \cdot \|w\|$, this is equivalent to

$$\max_{\mathbf{w},b,\gamma'} \frac{\gamma'}{\|\mathbf{w}\|}$$

$$\text{subject to:} \ y^i\left(\mathbf{w}^T\mathbf{x}^i + b\right) \geq \gamma' \ \ \forall i = 1, \ldots, N$$

# Maximum Margin Classifier

- Note that rescaling $\mathbf{w}$ and $b$ (by $\frac{1}{\gamma'}$) will not change the classifier, we can thus further reformulate the optimization problem

$$\max_{\mathbf{w},b,\gamma'} \frac{\gamma'}{\|\mathbf{w}\|}$$

$$\text{subject to}: \ y^i(\mathbf{w}^T\mathbf{x}^i + b) \geq \gamma', \ i = 1, ..., N$$

$$\max_{\mathbf{w},b} \frac{1}{\|\mathbf{w}\|} \ (\text{or equivalently} \min_{\mathbf{w},b}\|\mathbf{w}\|^2)$$

$$\text{subject to}: \ y^i(\mathbf{w}^T\mathbf{x}^i + b) \geq 1, \ i = 1, ..., N$$

Maximizing the geometric margin is equivalent to minimizing the magnitude of $\mathbf{w}$ subject to maintaining a functional margin of at least 1

# Solving the Optimization Problem

$$\min_{\mathbf{w},b}\|\mathbf{w}\|^2$$

$$\text{Subject to } y^i\left(\mathbf{w}^T\mathbf{x}^i + b\right) \geq 1, i = 1, \dots, N$$

- This is a ***quadratic optimization problem*** with linear constraints.
- A well-known class of mathematical programming problems, several (non-trivial) algorithms exist.
  - One can use any of them to solve for $\mathbf{w}$ and $b$
- It is useful to first formulate an equivalent dual optimization problem, which serves two purposes:
  - To show that the solution for $\mathbf{w}$ can be expressed as weighted sum of subset of training examples (aka the support vectors)
  - For applying kernel trick for nonlinear svm

# Aside: Constrained Optimization

- To solve the following optimization problem

$$\min_{x} f(x) \ s.t. \ g_i(x) \leq 0 \ \text{ for } i = 1, \ldots, m$$

- Consider the following function known as the Lagrangian

$$\mathcal{L}(x, \alpha) = f(x) + \sum_{i} \alpha_i g_i(x) \ s.t. \ \alpha_i \geq 0$$

- The original optimization problem is equivalent to solving the following:

$$\min_{x} \max_{\alpha} \mathcal{L}(x, \alpha) \quad \text{subject to } \alpha_i \geq 0$$

- By exchanging the order of min and max, we get the **dual problem**:

$$\max_{\alpha} \min_{x} \mathcal{L}(x, \alpha) \quad \text{subject to } \alpha_i \geq 0$$

# Aside: Constrained Optimization

$$\text{Primal}: f^* = \min_{x} \max_{\alpha \geq 0} L(x, \alpha)$$

$$\text{Dual}: d^* = \max_{\alpha \geq 0} \min_{x} L(x, \alpha)$$

Let $x^*$ and $\alpha^*$ be the optimal and dual solution respectively, $f^* = d^*$ if $f(x)$ is convex and $x^*$ and $\alpha^*$ satisfy the KKT conditions:

1. $\nabla L(x^*, \alpha^*) = 0$            --- zero gradient
2. $g(x^*) \leq 0$               --- primal feasibility
3. $\alpha^* \geq 0$                --- dual feasibility
4. $\alpha^* g(x^*) = 0$          --- complementary slackness

# Back to the Original Problem

$$\text{Minimize } \frac{1}{2}||\mathbf{w}||^2$$

$$\text{subject to}: 1 - y^i(\mathbf{w}^T\mathbf{x}^i + b) \leq 0, i = 1, \dots, N$$

The Lagrangian is

$$\mathcal{L}(\mathbf{w}, b, \alpha) = \frac{1}{2}\mathbf{w}^T\mathbf{w} + \sum_{i=1}^{N} \alpha_i \left(1 - y^i(\mathbf{w}^T\mathbf{x}^i + b)\right) s.t., \alpha_i \geq 0$$

- We want to solve $\max_{\alpha \geq 0} \min_{\mathbf{w}, b} \mathcal{L}(\mathbf{w}, b, \alpha)$

- Setting the gradient of $\mathcal{L}$ w.r.t. $\mathbf{w}$ and $b$ to zero:

$$\mathbf{w} - \sum_{i=1}^{N} \alpha_i y^i \mathbf{x}^i = 0 \implies \mathbf{w} = \sum_{i=1}^{N} \alpha_i y^i \mathbf{x}^i$$

$$\sum_{i=1}^{N} \alpha_i y^i = 0$$

# The Dual Problem

$$\mathcal{L}(\mathbf{w}, b, \alpha) = \frac{1}{2}\mathbf{w}^T\mathbf{w} + \sum_{i=1}^{N} \alpha_i \left(1 - y^i(\mathbf{w}^T\mathbf{x}^i + b)\right)$$

- Substitute $\mathbf{w} = \sum_{i=1}^{N} \alpha_i y^i \mathbf{x}^i$ into $\mathcal{L}$:

$$L(\alpha)$$

$$= \frac{1}{2}\sum_{i=1}^{N}\sum_{j=1}^{N} \alpha_i\alpha_j y^i y^j <\mathbf{x}^i \cdot \mathbf{x}^j> + \sum_{i=1}^{N} \alpha_i$$

$$- \sum_{i=1}^{N}\sum_{j=1}^{N} \alpha_i\alpha_j y^i y^j <\mathbf{x}^i \cdot \mathbf{x}^j> - b\underbrace{\sum_{i=1}^{N} \alpha_i y^i}_{= 0}$$

$$= \sum_{i=1}^{N} \alpha_i - \frac{1}{2}\sum_{i=1}^{N}\sum_{j=1}^{N} \alpha_i\alpha_j y^i y^j <\mathbf{x}^i \cdot \mathbf{x}^j>$$

# The Dual Problem

- The new objective function is in terms of $\alpha_i$, known as the <u>dual problem</u>

- The original problem is known as the <u>primal problem</u>

- The objective function of the dual problem needs to be maximized!
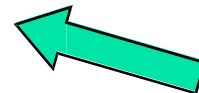
- The dual problem is therefore:

$$\max L(\mathbf{\alpha}) = \sum_{i=1}^{N} \alpha_i - \frac{1}{2} \sum_{i=1}^{N} \sum_{j=1}^{N} \alpha_i \alpha_j y^i y^j < \mathbf{x}^i \cdot \mathbf{x}^j >$$

$$\text{subject to} \quad \alpha_i \geq 0, i = 1,...,n, \qquad \sum_{i=1}^{N} \alpha_i y^i = 0$$

Properties of $\alpha_i$ when we introduce the Lagrange multipliers

The result when we differentiate the original Lagrangian w.r.t. b

# The Dual Problem

$$\max L(\mathbf{\alpha}) = \sum_{i=1}^{N} \alpha_i - \frac{1}{2} \sum_{i=1}^{N} \sum_{j=1}^{N} \alpha_i \alpha_j y^i y^j < \mathbf{x}^i \cdot \mathbf{x}^j >$$

$$\text{subject to} \quad \alpha_i \geq 0, i = 1,...,n, \qquad \sum_{i=1}^{N} \alpha_i y^i = 0$$

- This is also a quadratic programming (QP) problem
  - A global maximum of $\alpha_i$ can always be found

- **w** can be recovered by $\quad \mathbf{w} = \sum_{i=1}^{N} \alpha_i y^i \mathbf{x}^i$

- b can also be recovered as well (wait for a bit)

# Characteristics of the Solution

- Many of the $\alpha_i$ are zero   --- sparse solution
- **w** is a linear combination of only <u>a small number of data points</u>
- The KKT conditions requires that:

$$\alpha_i \geq 0, i = 1, \dots, n \qquad \text{\underline{Dual feasibility}}$$
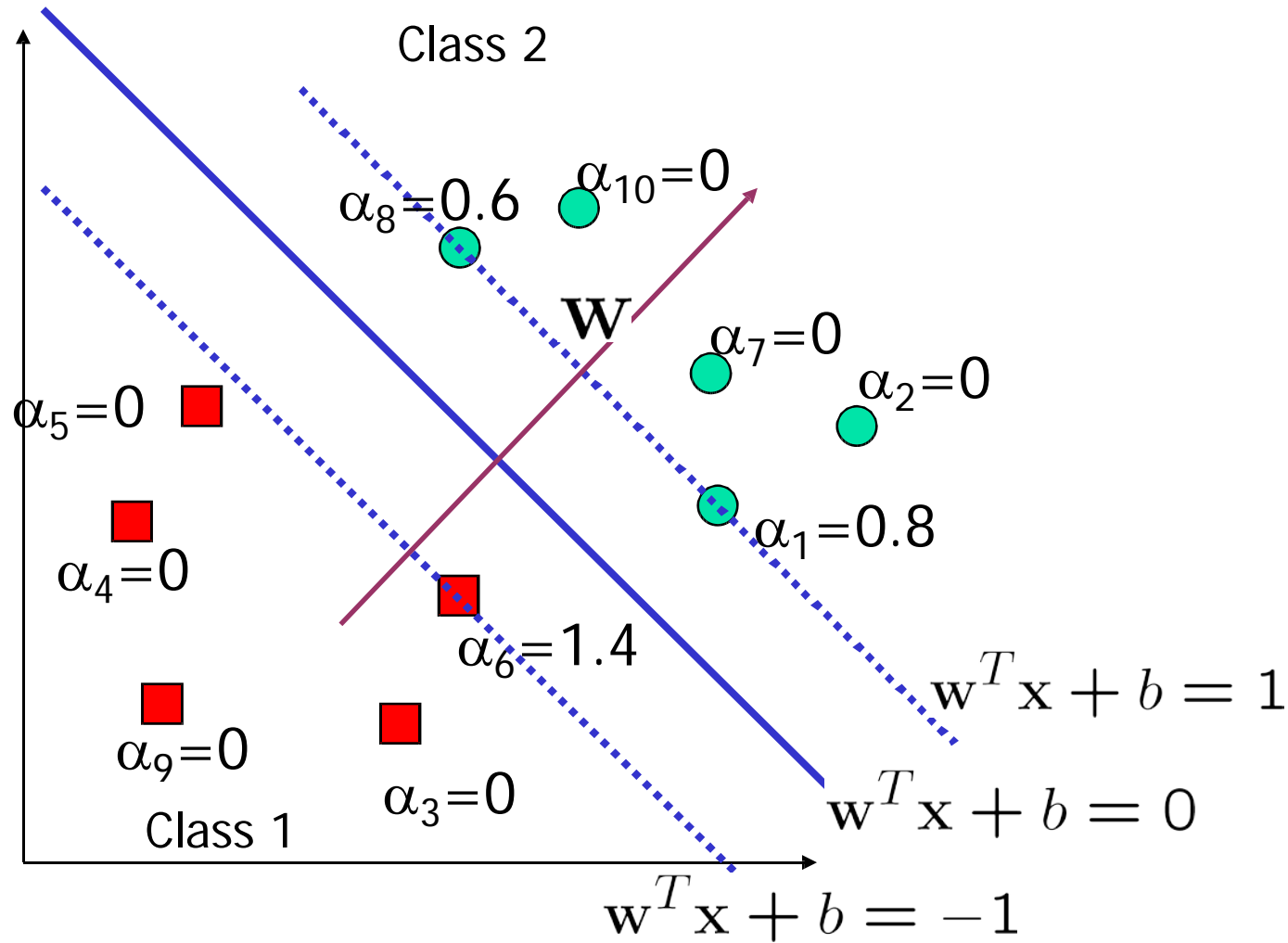
$$y^i \left( \sum_{j=1}^n \alpha_j y^j < \mathbf{x}^j \cdot \mathbf{x}^i > + b \right) \geq 1 \, , i = 1, \dots, n$$

<u>Primal feasibility: Functional margin ≥ 1</u>

$$\alpha_i \left( y^i \left( \sum_{j=1}^n \alpha_j y^j < \mathbf{x}^j \cdot \mathbf{x}^i > + b \right) - 1 \right) = 0, i = 1, \dots, n$$

<u>Complemetary slackness: $\alpha$ is nonzero only when functional margin = 1</u>

# A Geometrical Interpretation



Class 2

$\alpha_{10}=0$

$\alpha_8=0.6$

$\mathbf{W}$

$\alpha_7=0$

$\alpha_2=0$

$\alpha_5=0$

$\alpha_1=0.8$

$\alpha_4=0$

$\alpha_6=1.4$

$\mathbf{w}^T\mathbf{x}+b=1$

$\alpha_9=0$

$\alpha_3=0$

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

Class 1

$\mathbf{w}^T\mathbf{x}+b=-1$

# Support Vectors

- $\mathbf{x}^i$ with non-zero $\alpha's$ are called support vectors (SV)

- The decision boundary is determined only by the SV's

$$\mathbf{w} = \sum_{i=1}^{N} \alpha_i y^i \mathbf{x}^i$$

- Note that we know that for support vectors the functional margin = 1

- We can use this information to <u>solve for b</u>

# Classifying new examples

For classifying with a new input **x**

- Compute
$$\mathbf{w}^T\mathbf{x} + b = \sum_{i=1}^{N} \alpha_i y^i < \mathbf{x}^i \cdot \mathbf{x} > + b$$

- Note: no need to form **w** explicitly, rather, classify **x** by taking a weighted sum of **its dot products with the support vectors** (useful for generalizing from inner product to kernels)

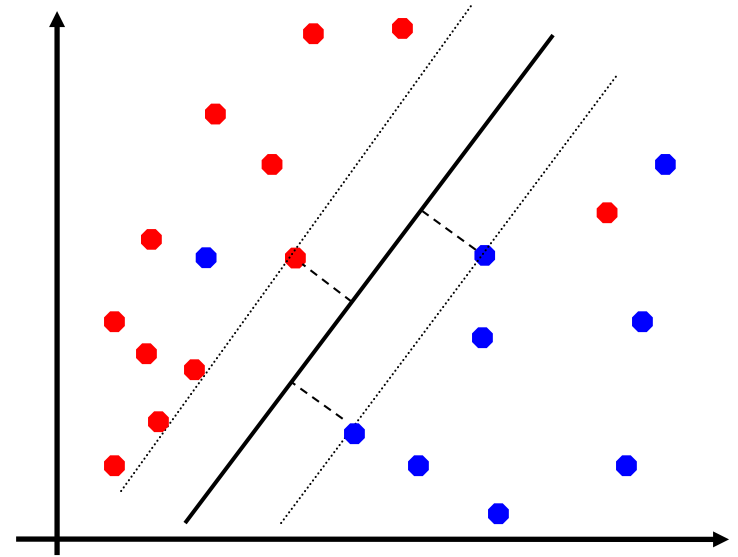# Solving the QP optimization problem

- Many approaches have been proposed for QP
  - Loqo, cplex, etc. (see http://www.numerical.rl.ac.uk/qp/qp.html)
- Early work focuses on "interior-point" methods
  - Start with an initial solution that can violate the constraints
  - Improve this solution by optimizing the objective function and/or reducing the amount of constraint violation
- Stochastic sub-gradient descent has been shown to lead to extremely efficient primal solver for large scale problems
- In practice, one can just regard the QP solver as a "black-box" without bothering how it works, but depending on the scale of the problem some solvers might be more appropriate than others

# Non-separable Data

*What if the data is not linearly separable?*

– The solution does not exist

– i.e., the set of linear constraints are not satisfiable

– But we should still be able to find a good decision boundary

**Solution:**

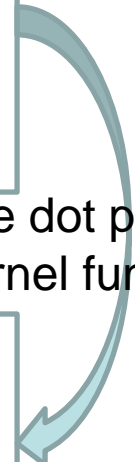- Project the data onto higher dimensional space
- Via kernel function

# Kernel SVM

Linear SVM:

$$\max L(\mathbf{\alpha}) = \sum_{i=1}^{N} \alpha_i - \frac{1}{2} \sum_{i=1}^{N} \sum_{j=1}^{N} \alpha_i \alpha_j y^i y^j < \mathbf{x}^i \cdot \mathbf{x}^j >$$

$$\text{subject to} \quad \alpha_i \geq 0, i = 1,...,n, \qquad \sum_{i=1}^{N} \alpha_i y^i = 0$$
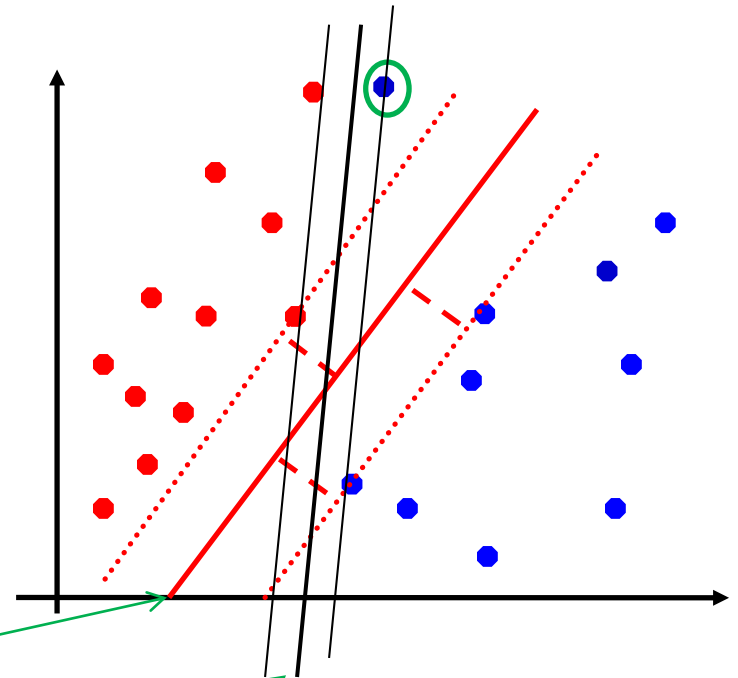
Replace dot product with kernel function

Kernel SVM:

$$\max L(\mathbf{\alpha}) = \sum_{i=1}^{N} \alpha_i - \frac{1}{2} \sum_{i=1}^{N} \sum_{j=1}^{N} \alpha_i \alpha_j y^i y^j \; K(\mathbf{x}^i, \mathbf{x}^j)$$

$$\text{subject to} \quad \alpha_i \geq 0, i = 1,...,n, \qquad \sum_{i=1}^{N} \alpha_i y^i = 0$$

# Maximum margin overfits to outliers

*Consider the blue point circled out. It is an outlier that is labeled as blue but really should belong to red*
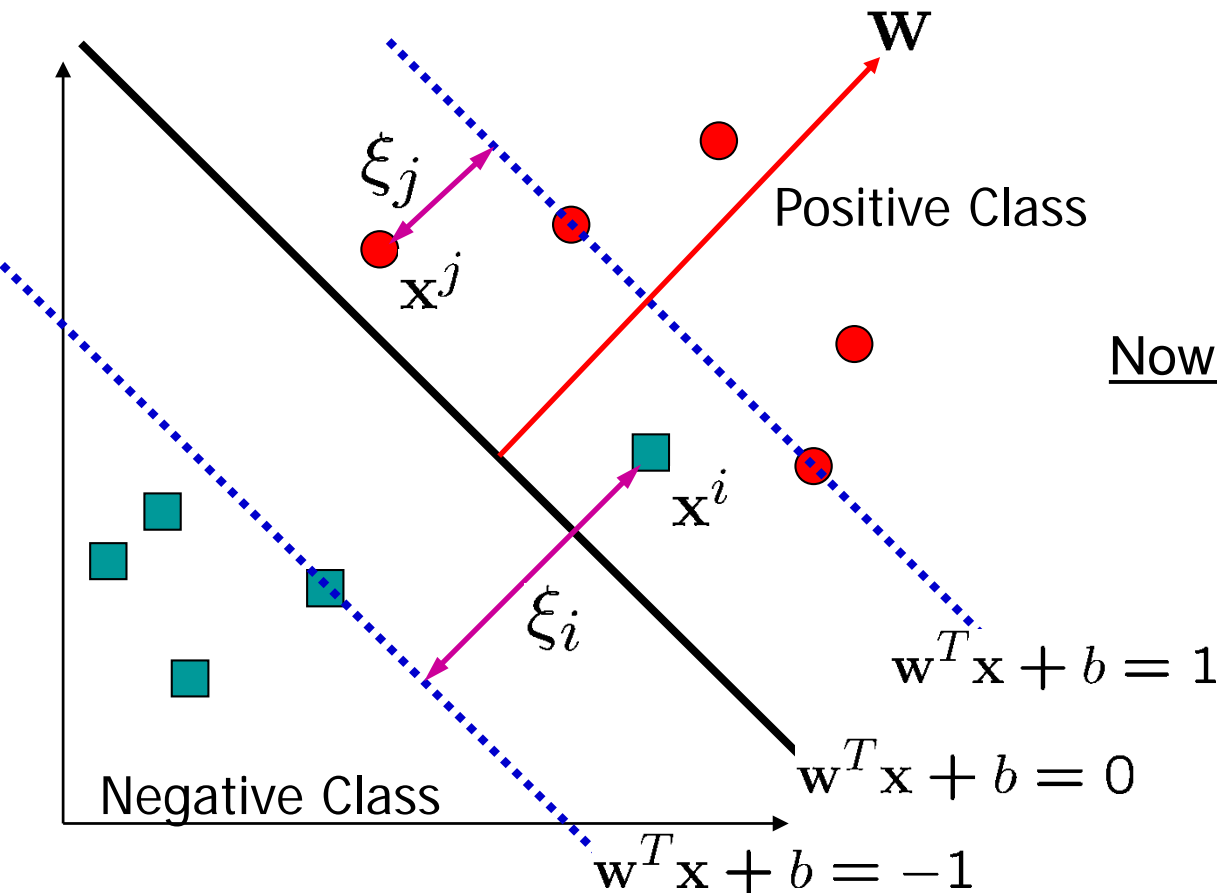
We would like to learn a boundary that ignores the outliers

But the margin will be defined by the outlier and we instead learn a boundary that overfit to the outliers

# Soft Margin

- Allow functional margins to be less than 1



**W**

$\xi_j$

$\mathbf{x}^j$

Positive Class

$\xi_i$

$\mathbf{x}^i$

$\mathbf{w}^T \mathbf{x} + b = 1$

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

Negative Class

$\mathbf{w}^T \mathbf{x} + b = -1$

Originally functional margins need to satisfy:

$$y^i(\mathbf{w}^T \mathbf{x}^i + b) \geq 1$$

Now we allow it to be less than 1:

$$y^i(\mathbf{w}^T \mathbf{x}^i + b) \geq 1 - \xi_i$$
$$\xi_i \geq 0$$

The objective changes to:

$$\min_{\mathbf{w}, b, \xi_i} \|\mathbf{w}\|^2 + c \sum_{i=1}^{N} \xi_i$$

# Soft-Margin Maximization

$$\min_{\mathbf{w},b} \|\mathbf{w}\|^2$$

$$\text{subject to}: \ y^i(\mathbf{w} \cdot \mathbf{x}^i + b) \geq 1, \ \ i = 1, \cdots, N$$

**Slack variables**

$$\min_{\mathbf{w},b} \|\mathbf{w}\|^2 + c\sum_{i=1}^{N} \xi_i$$

$$\text{subject to}: \ y^i(\mathbf{w} \cdot \mathbf{x}^i + b) \geq 1 - \xi_i, \ \ i = 1, \cdots, N$$

$$\xi_i \geq 0, \ \ i = 1, \cdots, N$$

- This allows some functional margins < 1 (could even be < 0)
- The $\xi_i$'s can be viewed as the "errors" of our *fat* decision boundary
- Adding $\xi_i$'s to the objective function to minimize errors
- We have a tradeoff between making the decision boundary fat and minimizing the error
- Parameter ***c*** controls the tradeoff:
  - Large c: $\xi_i$'s incur large penalty, so the optimal solution will try to avoid them
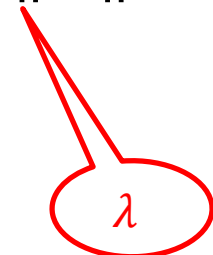  - Small c: small cost for $\xi_i$'s, we can sacrifice some training examples to have a large classifier margin

# Soft Margin SVM: Regularized Hinge loss

$$\min_{\mathbf{w},b} \|\mathbf{w}\|^2 + c \sum_{i=1}^{N} \xi_i$$

$$\text{subject to } y^i\left(\mathbf{w}^T\mathbf{x}^i + b\right) \geq 1 - \xi_i,$$

$$\xi_i \geq 0, \forall i = 1, \dots, N$$

Is equivalent to:

$$\min_{w,b} \|\mathbf{w}\|^2 + c \sum_{i}^{N} \max\left(0, 1 - y^i\left(w^T x^i + b\right)\right)$$

$\lambda$

$L_2$ Regularization

Hinge loss

# Different Loss functions

# Solutions to soft-margin SVM

$$w = \sum_{i=1}^{N} \alpha_i y^i x^i, \quad \textbf{s.t.} \sum_{i=1}^{N} \alpha_i y^i = 0$$

No soft margin

$$w = \sum_{i=1}^{N} \alpha_i y^i x^i, \quad \textbf{s.t.} \sum_{i=1}^{N} \alpha_i y^i = 0 \textbf{ and } 0 \le \alpha_i \le c$$

With soft margin

- *c* effectively puts a **box constraint** on $\alpha$, the weights of the support vectors
- It limits the influence of individual support vectors (maybe outliers)
- In practice, c is a parameter to be set, similar to *k* in k-nearest neighbor
- It can be set using cross-validation

# Kernel SVM with soft margin

$$\max L(\boldsymbol{\alpha}) = \sum_{i=1}^{N} \alpha_i - \frac{1}{2} \sum_{i=1}^{N} \sum_{j=1}^{N} \alpha_i \alpha_j \, y^i \, y^j \, K(\mathbf{x}^i, \mathbf{x}^j)$$

$$\text{subject to} \quad 0 \le \alpha_i \le c, \, i = 1 \dots, N; \qquad \sum_{i=1}^{N} \alpha_i y^i = 0$$

# Summary of SVM

- SVM aims to find the max margin linear separator

- Soft margin SVM can be interpreted as:
  - Introducing slack to the hard margin constraints – C-SVM, where C is the penalty weight for the accumulative slack
  - Minimizing $L2$ regularized hinge loss - $\lambda$-SVM, where $\lambda$ is the regularization parameter

- Large $C$ (or equivalently small $\lambda$): increased overfitting

- Small $C$ (or equivalently large $\lambda$): decreased overfitting

- By solving the dual problem with the kernel trick, we can learn max margin separator in the mapped nonlinear space