

# UVA CS 4774: Machine Learning

## S4: Lecture 21 Extra: (SVM Optimization and Dual basic)

Dr. Yanjun Qi

University of Virginia  
Department of Computer Science

Module IV  
Extra

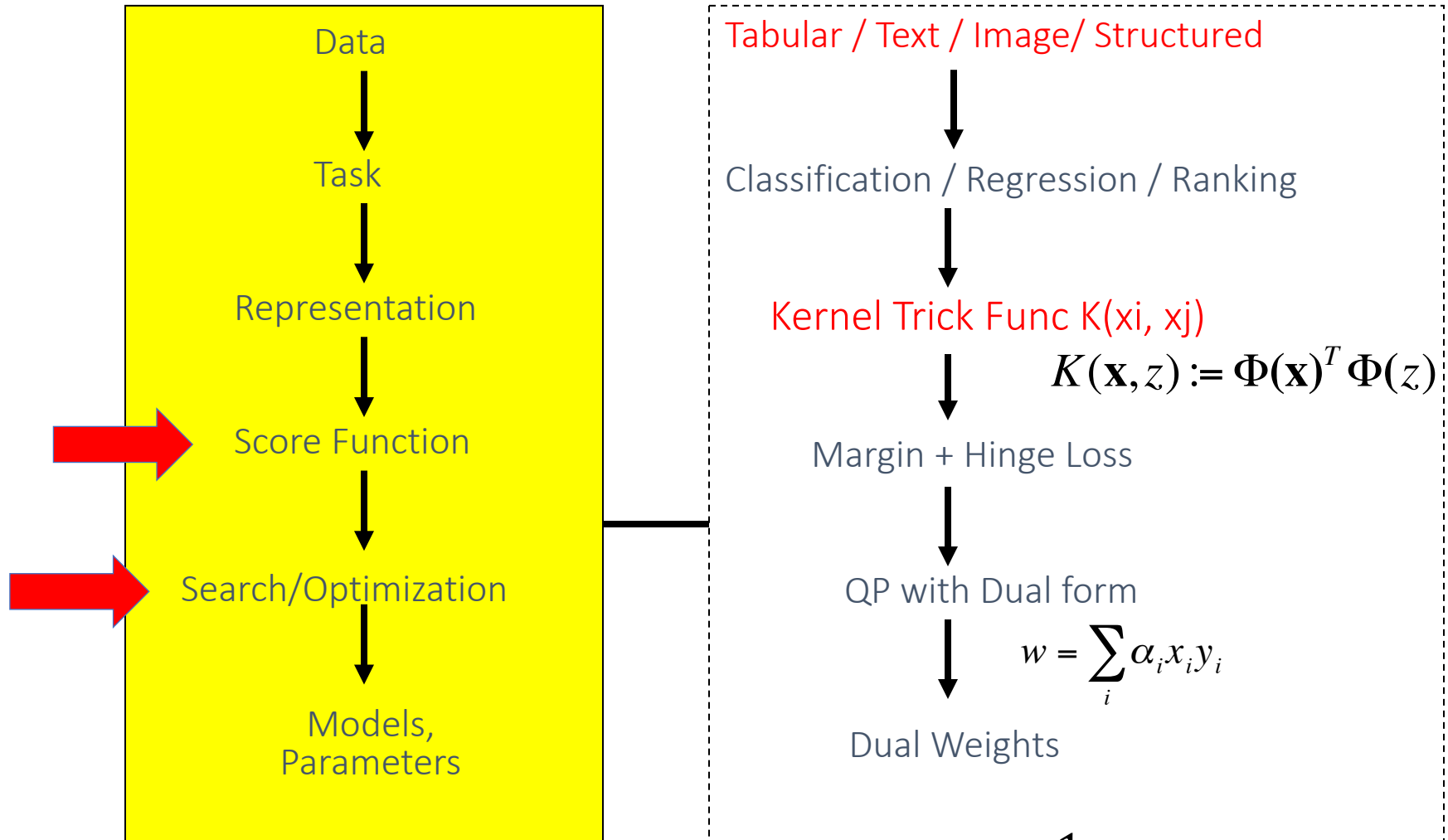
# What Left in SVM?

## ❑ Support Vector Machine (SVM)

- ✓ History of SVM
- ✓ Large Margin Linear Classifier
- ✓ Define Margin (M) in terms of model parameter
- ✓ Optimization to learn model parameters (w, b)
- ➔ ✓ Linearly Non-separable case (soft SVM)
- ✓ Optimization with dual form
- ✓ Nonlinear decision boundary
- ✓ Practical Guide

$K(x, z)$

# This: Kernel Support Vector Machine



$$\operatorname{argmin}_{\mathbf{w}, b} \sum_{i=1}^p w_i^2 + C \sum_{i=1}^n \varepsilon_i$$

$$\text{subject to } \forall \mathbf{x}_i \in D_{\text{train}} : y_i (\mathbf{x}_i \cdot \mathbf{w} + b) \geq 1 - \varepsilon_i$$

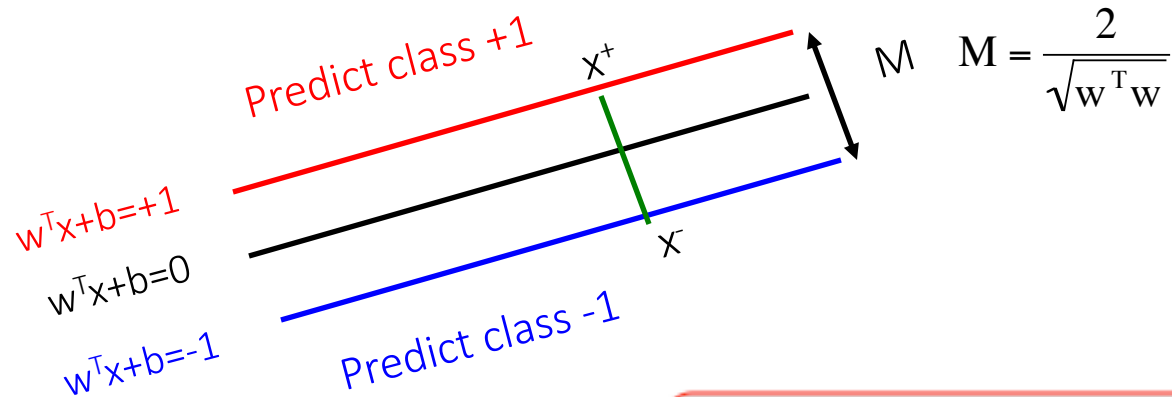


$$\max_{\alpha} \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j$$

$$\sum_i \alpha_i y_i = 0, \quad \alpha_i \geq 0 \quad \forall i$$

# Optimization Step

i.e. learning optimal parameter for SVM



1. Correctly classifies all points
2. Maximizes the margin (or equivalently minimizes  $w^T w$ )

# Optimization Reformulation

1. Correctly classifies all points
2. Maximizes the margin (or equivalently minimizes  $w^T w$ )

Min  $(w^T w)/2$

subject to the following constraints:

For all  $x$  in class + 1

$$w^T x + b \geq 1$$

$$y_i = 1$$

A total of  $n$  constraints if we have  $n$  input samples

For all  $x$  in class - 1

$$w^T x + b \leq -1$$

$$y_i = -1$$

$$\rightarrow \text{pos } y_i = 1, w^T x_i + b \geq 1$$

$$y_i (w^T x_i + b) \geq 1$$

$$\rightarrow \text{neg } y_i = -1, w^T x_i + b \leq -1$$

$$y_i (w^T x_i + b) \geq 1$$

# Optimization Reformulation

1. Correctly classifies all points
2. Maximizes the margin (or equivalently minimizes  $w^T w$ )

Min  $(w^T w)/2$

subject to the following constraints:

For all  $x$  in class + 1

$$w^T x + b \geq 1$$

For all  $x$  in class - 1

$$w^T x + b \leq -1$$



A total of  $n$  constraints if we have  $n$  input samples



$$\operatorname{argmin}_{w, b} \sum_{i=1}^p w_i^2 / 2$$

$$\text{subject to } \forall \mathbf{x}_i \in D_{\text{train}} : y_i (\mathbf{x}_i \cdot \mathbf{w} + b) \geq 1$$

# Linearly Non separable case

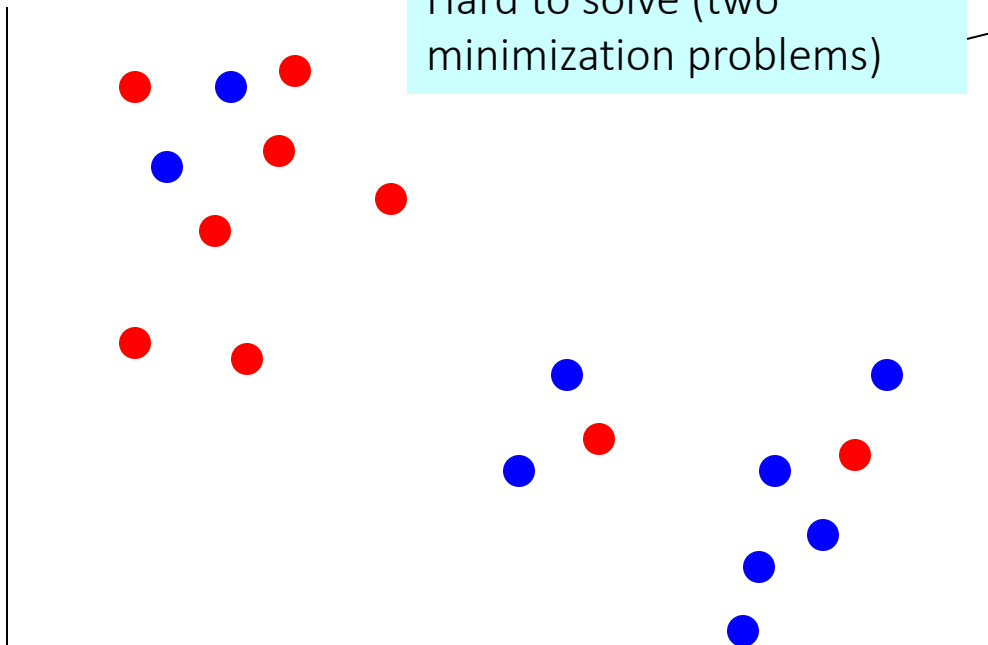
- So far we assumed that a linear hyperplane can perfectly separate the points
- But this is not usually the case
  - noise, outliers

How can we convert this to a QP problem?

Hard to solve (two minimization problems)

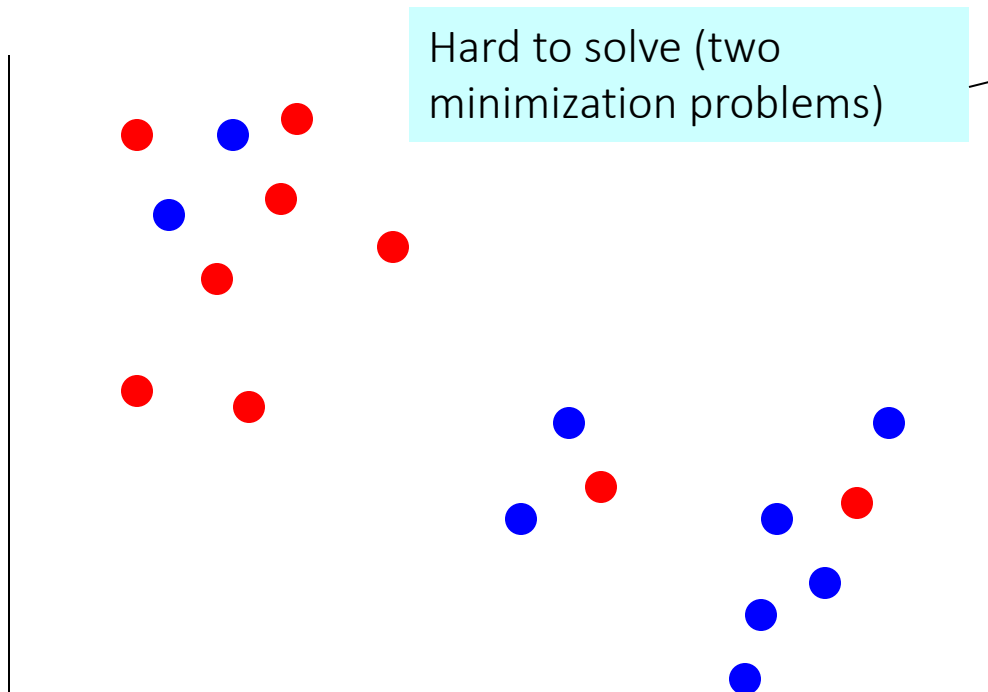
- Minimize training errors?

$\left\{ \begin{array}{l} \min w^T w / 2 \\ \min \text{\#errors} \end{array} \right.$



# Linearly Non separable case

- So far we assumed that a linear plane can perfectly separate the points
- But this is not usually the case
  - noise, outliers



How can we convert this to a QP problem?

- Minimize training errors?

$$\min w^T w / 2$$

$$\min \text{\textcolor{red}{\#errors}}$$

- Penalize training errors:

$$\min w^T w / 2 + C * (\text{\textcolor{red}{\#errors}})$$

Hard to encode in a QP problem



$C$   penalize errors more

$$\text{SVM: } \min_{\vec{w}, b} w^T w + \underbrace{C(\# \text{ errors})}$$

Ridge Regression:

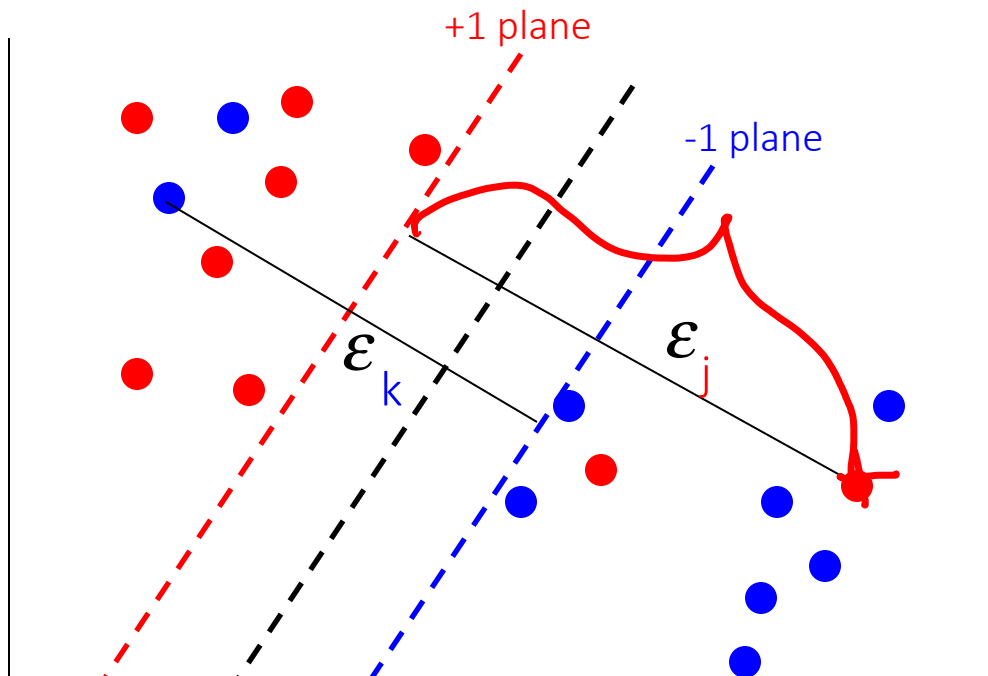
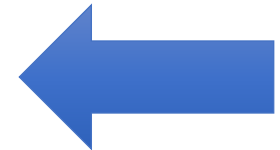
$$\min_{\vec{\theta}} \underbrace{\lambda \theta^T \theta} + J(\theta)$$

# Linearly Non separable case

- Instead of minimizing the number of misclassified points we can **minimize the distance between these points and their correct plane**

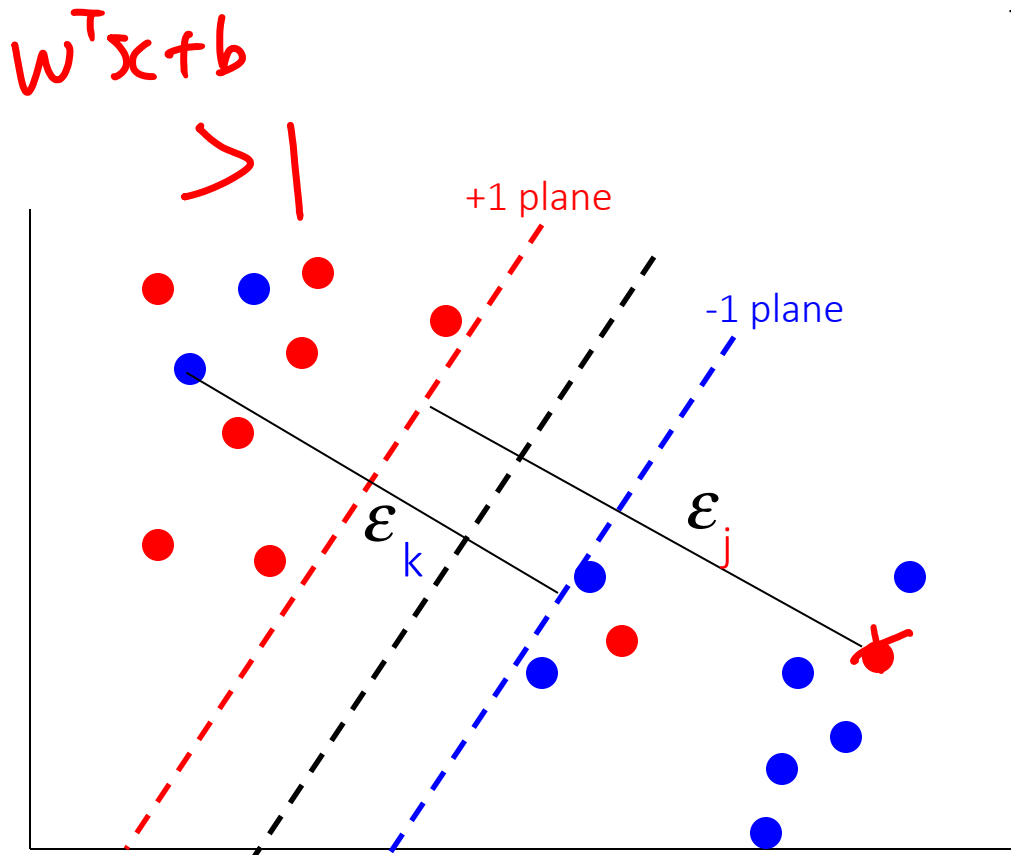
The new optimization problem is:

$$\min_w \frac{w^T w}{2} + C \sum_{i=1}^n \varepsilon_i$$



# Linearly Non separable case

- Instead of minimizing the number of misclassified points we can **minimize the distance between these points and their correct plane**



The new optimization problem is:

$$\min_w \frac{w^T w}{2} + C \sum_{i=1}^n \epsilon_i$$



subject to the following inequality constraint

For all  $x_i$  in class + 1

$$w^T x_i + b \geq 1 - \epsilon_i \quad \epsilon_i \geq 0$$

For all  $x_i$  in class - 1

$$w^T x_i + b \leq -1 + \epsilon_i$$

Wait. Are we missing something?

# Final optimization for linearly non-separable case

The new optimization problem is:

$$\min_w \frac{w^T w}{2} + C \sum_{i=1}^n \varepsilon_i$$

subject to the following inequality constraints:

For all  $x_i$  in class + 1

$$w^T x_i + b \geq 1 - \varepsilon_i$$

For all  $x_i$  in class - 1

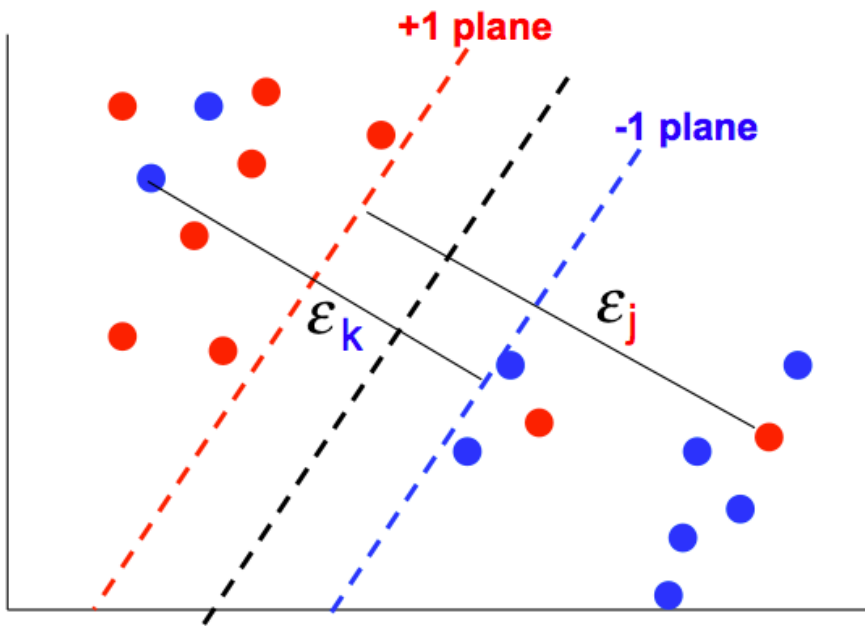
$$w^T x_i + b \leq -1 + \varepsilon_i$$

For all  $i$

$$\varepsilon_i \geq 0$$

A total of  $n$  constraints

Another  $n$  constraints



Two optimization problems:

For the separable and non separable cases

$$\min_w \frac{w^T w}{2}$$

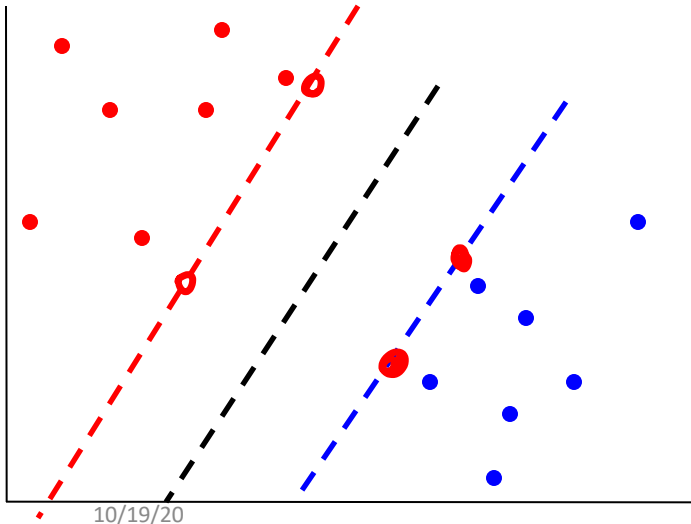
For all  $x$  in class + 1

$$w^T x + b \geq 1$$

For all  $x$  in class - 1

$$w^T x + b \leq -1$$

separable



$$\min_w \frac{w^T w}{2} + C \sum_{i=1}^n \varepsilon_i$$

For all  $x_i$  in class + 1

$$w^T x_i + b \geq 1 - \varepsilon_i$$

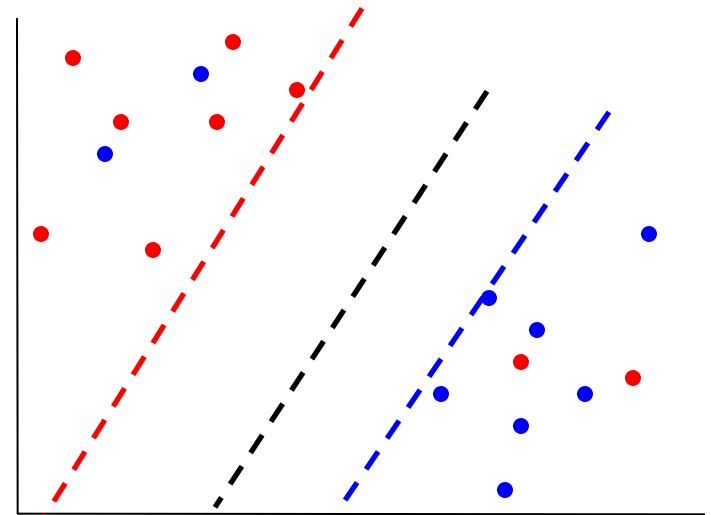
For all  $x_i$  in class - 1

$$w^T x_i + b \leq -1 + \varepsilon_i$$

For all  $i$

$$\varepsilon_i \geq 0$$

non separable



# Model Selection, find right C

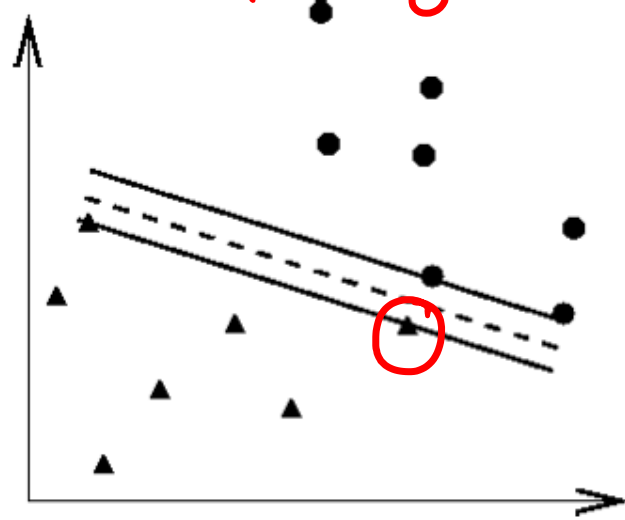
Training

Test

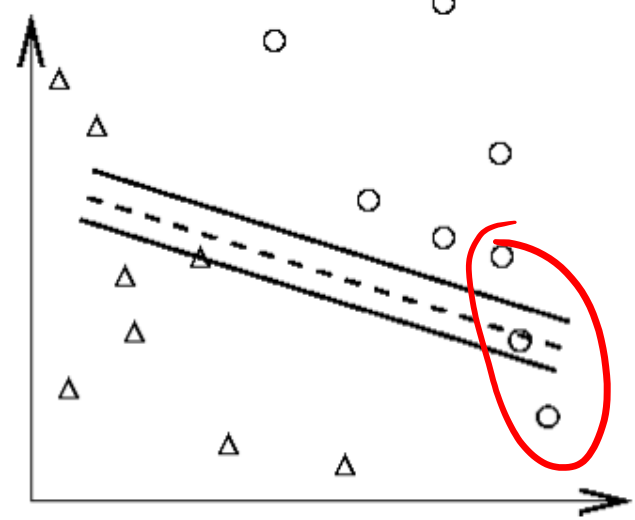
large C

Select the right penalty parameter C

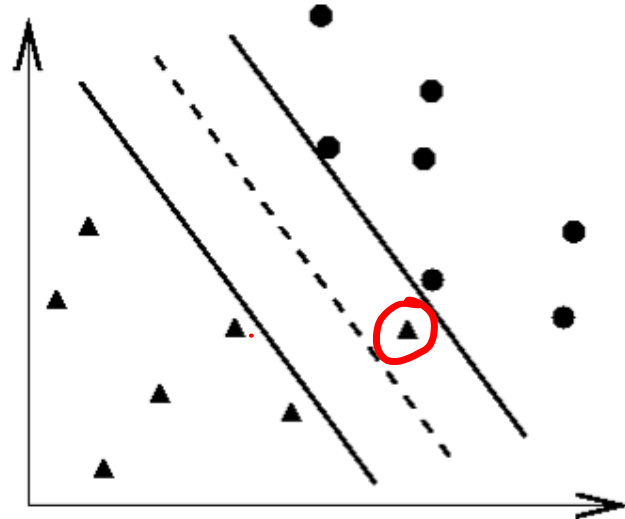
small C



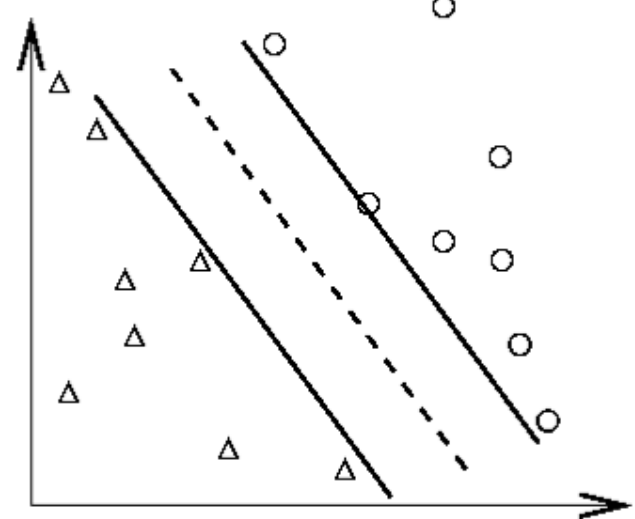
(a) Training data and an overfitting classifier



(b) Applying an overfitting classifier on testing data



(c) Training data and a better classifier



(d) Applying a better classifier on testing data

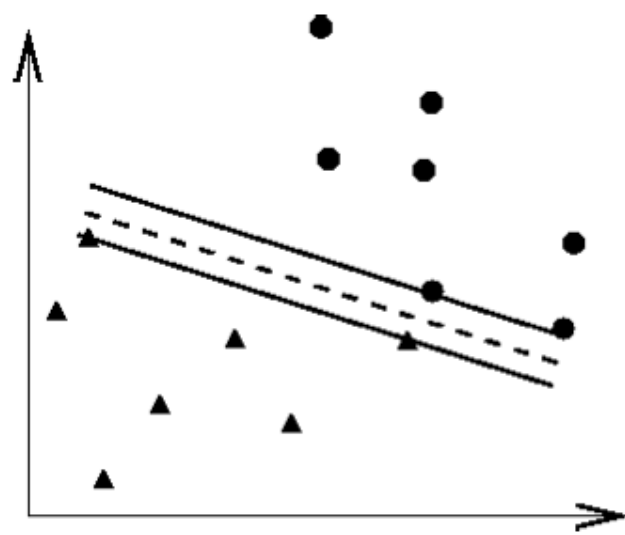
# Model Selection, find right C

large C

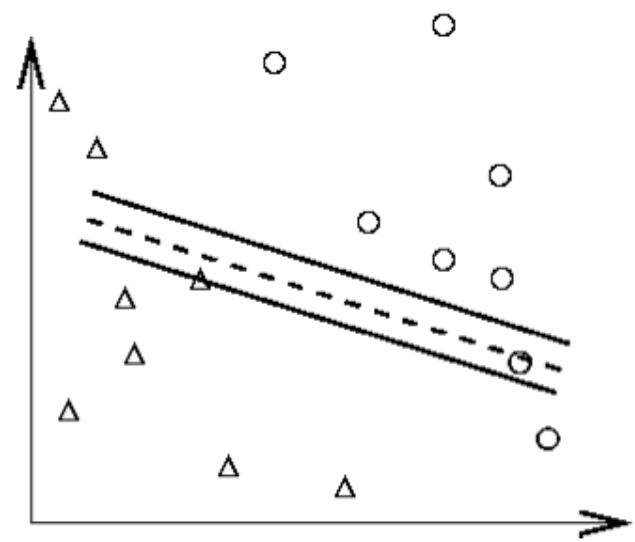
A large value of C means that misclassifications are bad - resulting in smaller margins and less training error (but more expected true error).

may lead

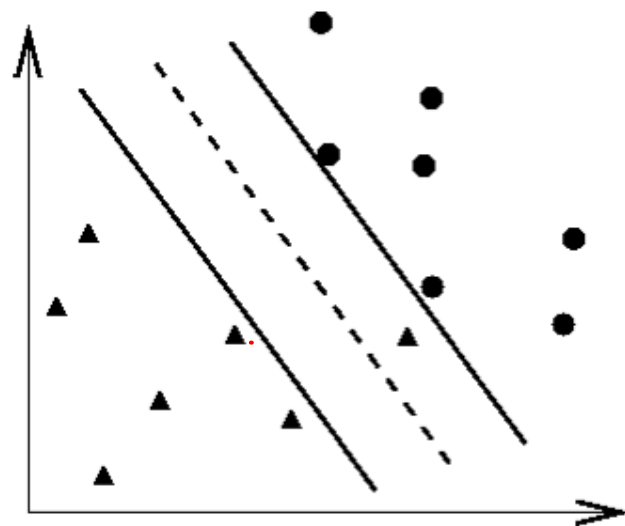
small C



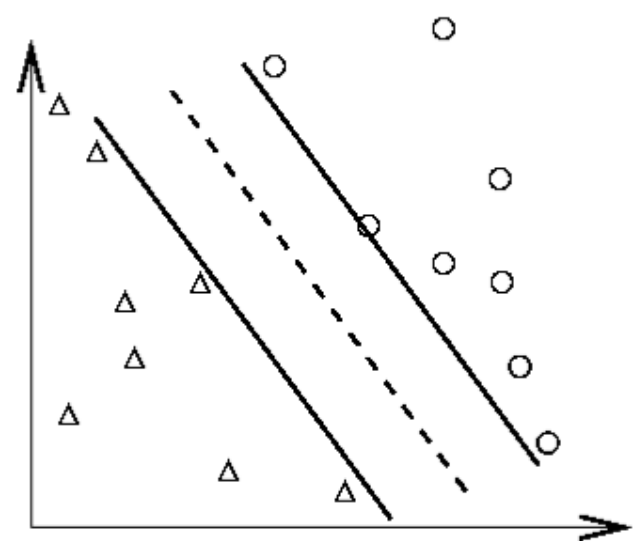
(a) Training data and an overfitting classifier



(b) Applying an overfitting classifier on testing data

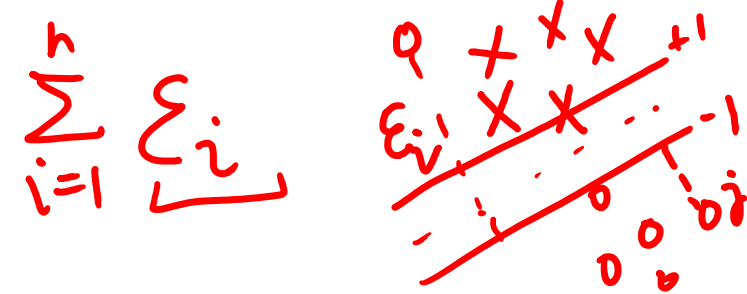


(c) Training data and a better classifier



(d) Applying a better classifier on testing data

# Hinge Loss for Soft SVM



$$\min_w \frac{\mathbf{w}^T \mathbf{w}}{2} + C \sum_{i=1}^n \epsilon_i$$

For all  $\mathbf{x}_i$  in class + 1

$$\mathbf{w}^T \mathbf{x}_i + b \geq 1 - \epsilon_i$$

For all  $\mathbf{x}_i$  in class - 1

$$\mathbf{w}^T \mathbf{x}_i + b \leq -1 + \epsilon_i$$

For all  $i$

$$\epsilon_i \geq 0$$



$$\sum_{i=1}^n \max(0, 1 - y_i f(\mathbf{x}_i))$$



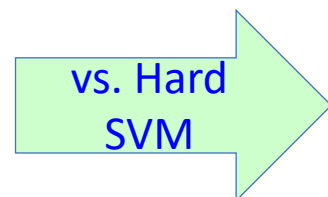
$$\operatorname{argmin}_{\mathbf{w}, b} \frac{\mathbf{w}^T \mathbf{w}}{2} + C \sum_{i=1}^n \max(0, 1 - y_i (\mathbf{w}^T \mathbf{x}_i + b))$$

subject to:

$$y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \epsilon_i$$

$$\epsilon_i \geq 0$$

$\geq 1$



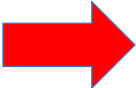
$$\operatorname{argmin}_{\mathbf{w}, b} \sum_{i=1}^p w_i^2 / 2$$

$$\text{subject to } \forall \mathbf{x}_i \in D_{\text{train}} : y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1$$

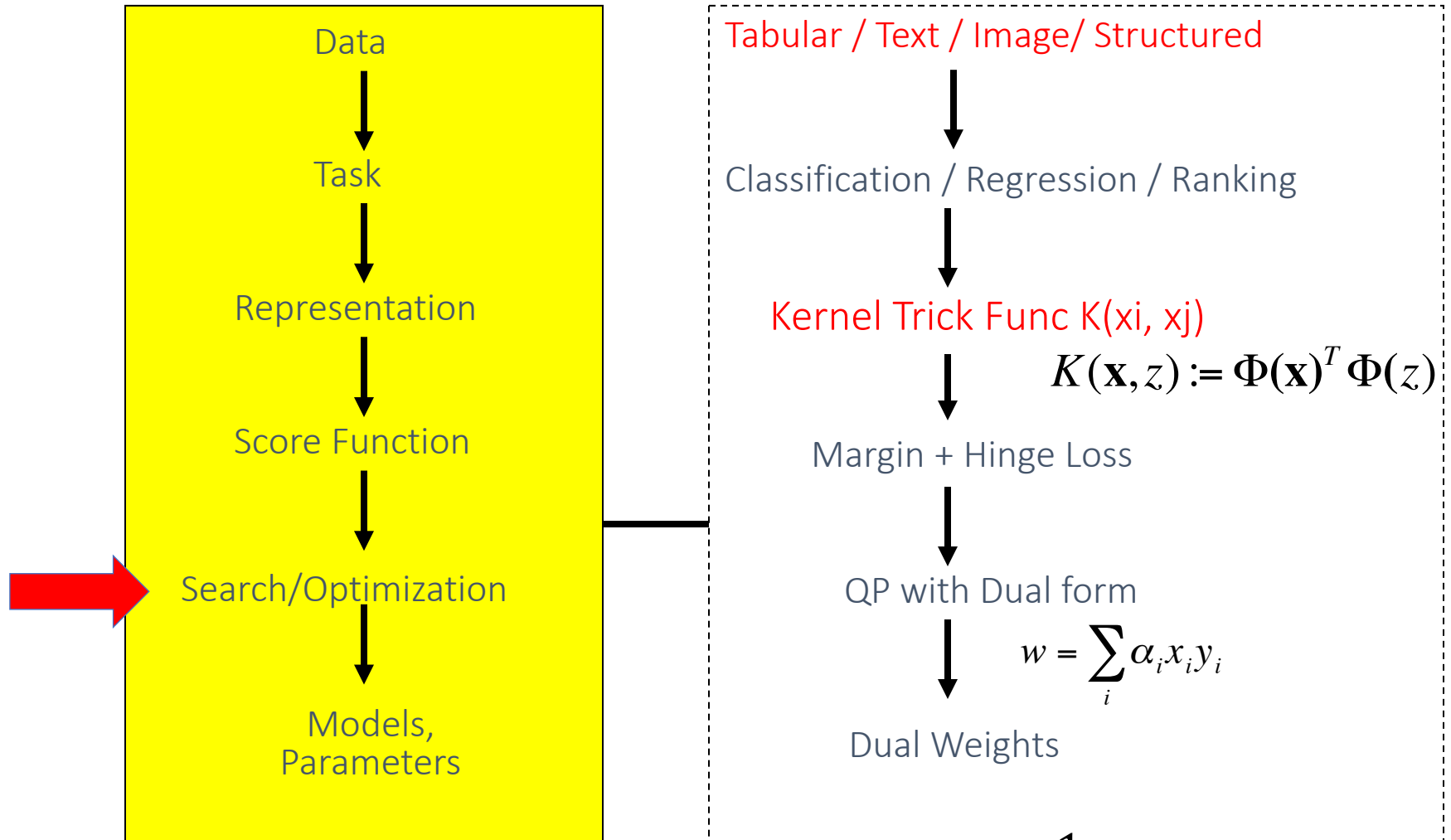


# What Left in SVM?

## Support Vector Machine (SVM)

- ✓ History of SVM
- ✓ Large Margin Linear Classifier
- ✓ Define Margin ( $M$ ) in terms of model parameter
- ✓ Optimization to learn model parameters ( $w, b$ )
- ✓ Linearly Non-separable case (soft SVM)
-  ✓ Optimization with dual form
- ✓ Nonlinear decision boundary
- ✓ Practical Guide

# This: Kernel Support Vector Machine



$$\operatorname{argmin}_{\mathbf{w}, b} \sum_{i=1}^p w_i^2 + C \sum_{i=1}^n \varepsilon_i$$

$$\text{subject to } \forall \mathbf{x}_i \in D_{\text{train}} : y_i (\mathbf{x}_i \cdot \mathbf{w} + b) \geq 1 - \varepsilon_i$$



$$\max_{\alpha} \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j$$

$$\sum_i \alpha_i y_i = 0,$$

$$\alpha_i \geq 0$$

$$\forall i$$

Two optimization problems: For the **separable** and **non separable** cases

$$\text{Min } (w^T w)/2$$

**For all  $x$  in class + 1**

$$w^T x + b \geq 1$$

**For all  $x$  in class - 1**

$$w^T x + b \leq -1$$

$$\min_w \frac{w^T w}{2} + C \sum_{i=1}^n \varepsilon_i$$

**For all  $x_i$  in class + 1**

$$w^T x_i + b \geq 1 - \varepsilon_i$$

**For all  $x_i$  in class - 1**

$$w^T x_i + b \leq -1 + \varepsilon_i$$

For all  $i$

$$\varepsilon_i \geq 0$$

- Instead of solving these QPs directly we will solve a **dual formulation of the SVM optimization problem**
- The main reason for switching to this type of representation is that it would allow us to use a **neat trick that will make our lives easier (and the run time faster)**

# Optimization Review: Ingredients

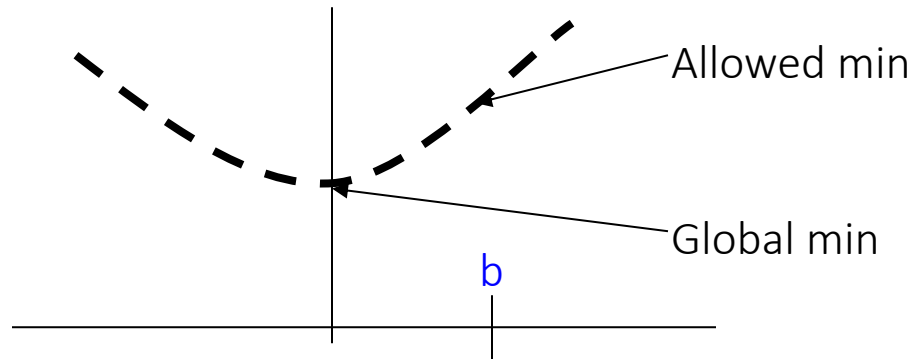
- Objective function
- Variables
- Constraints

**Find values of the variables  
that minimize or maximize the objective function  
while satisfying the constraints**

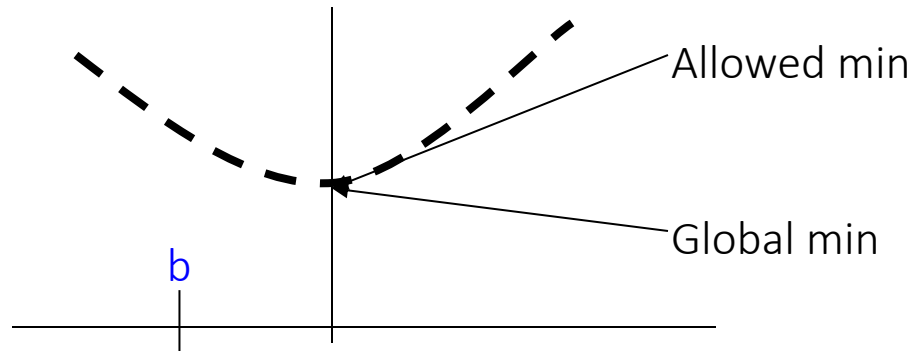
# Optimization Review: Constrained Optimization

$$\begin{array}{ll} \min_u & u^2 \\ \text{s.t.} & u \geq b \end{array}$$

Case 1:



Case 2:

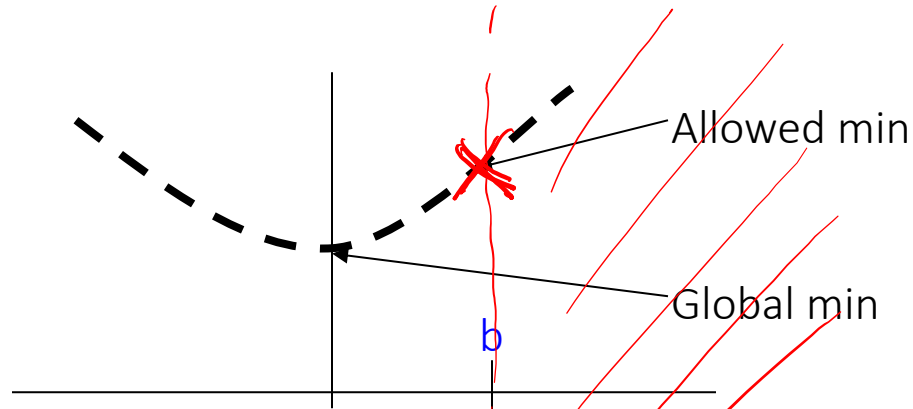


# Optimization Review: Constrained Optimization

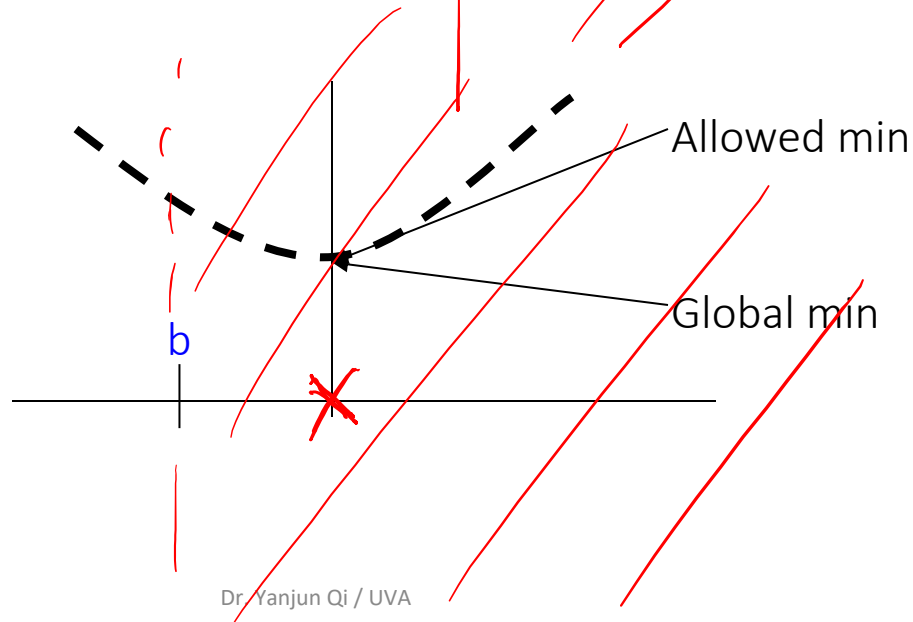
$$f(u) = u^2$$

$$\begin{array}{ll} \min_u & u^2 \\ \text{s.t.} & u \geq b \end{array}$$

Case 1:



Case 2:



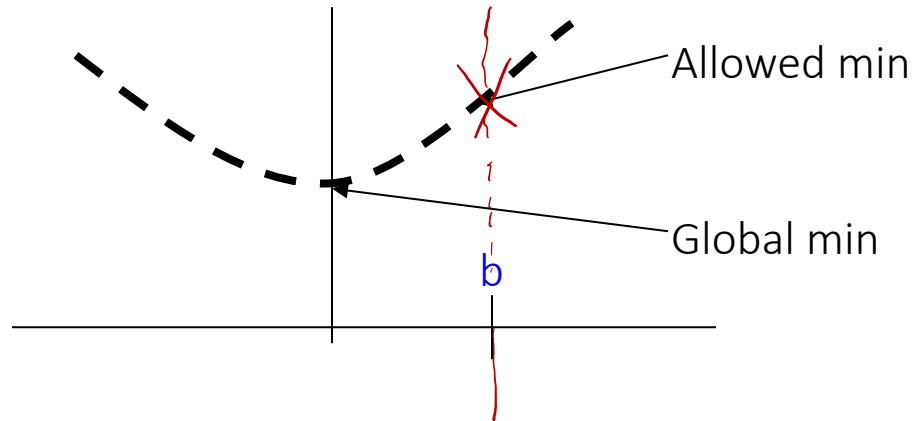
# Optimization Review: Constrained Optimization

$f(u)$

$$\begin{array}{ll} \min_u & u^2 \\ \text{s.t.} & u \geq b \end{array}$$

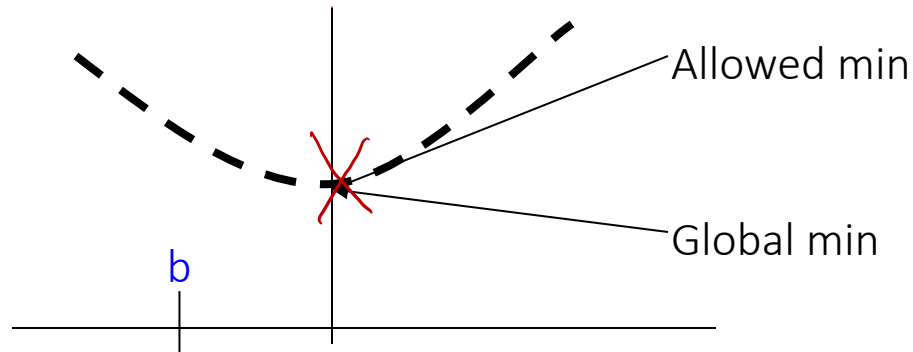
[Subject to]

Case 1:



$$\begin{cases} b > 0 \\ f(u^*) = b^2 \\ u^* = b \end{cases}$$

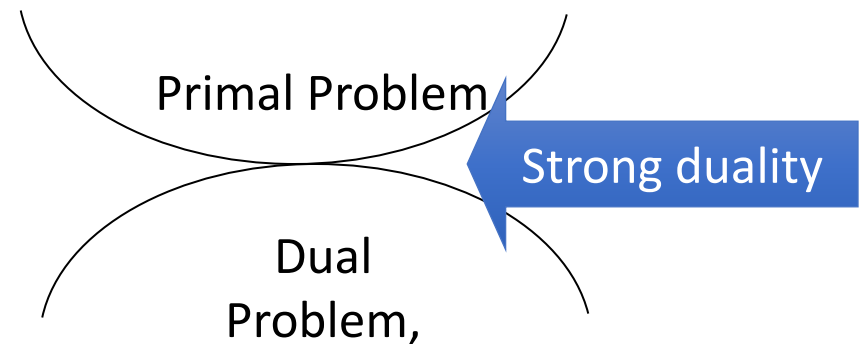
Case 2:



$$\begin{cases} b < 0 \\ f(u^*) = 0 \\ u^* = 0 \end{cases}$$

# Optimization Review: Dual Problem (Extra)

- Solving dual problem if the dual form is easier than primal form
- Need to change primal **minimization** to dual **maximization** (OR → Need to change primal **maximization** to dual **minimization**)
- Only valid when the original optimization problem is convex/concave (strong duality)





$$f(u): \begin{cases} \min u^2 \\ \text{s.t. } u \geq b \end{cases}$$

$$g(\alpha): \begin{cases} \max -\frac{\alpha^2}{4} + b\alpha = \max \left\{ -\underbrace{\left(\frac{\alpha}{2} - b\right)^2}_{\text{red}} + \underbrace{b^2}_{\text{red}} \right\} \\ \text{s.t. } \alpha \geq 0 \end{cases}$$

$$\begin{cases} \text{if } \underset{\text{red}}{b} \geq 0, & u^* = b, \quad g^* = b^2 \\ \text{if } \underset{\text{red}}{b} < 0, & \alpha^* = 0, \quad g^* = 0 \end{cases}$$

$$\Rightarrow \alpha (b - u) = 0 \quad \text{KKT condition}$$

# Optimization Review:

## Lagrangian Duality (Extra)

- The Primal Problem

$$\begin{array}{ll} \min_w & f_0(w) \\ \text{Primal:} & \text{s.t.} \quad f_i(w) \leq 0, \quad i = 1, \dots, k \end{array}$$

**The generalized Lagrangian:**

“Method of Lagrange multipliers”  
convert to a higher-dimensional problem

$$\mathcal{L}(w, \alpha) = f_0(w) + \sum_{i=1}^k \alpha_i f_i(w)$$

the  $\alpha$ 's ( $\alpha_i \geq 0$ ) are called the Lagrangian multipliers

**Lemma:**

$$\max_{\alpha, \alpha_i \geq 0} \mathcal{L}(w, \alpha) = \begin{cases} f_0(w) & \text{if } w \text{ satisfies primal constraints} \\ \infty & \text{o/w} \end{cases}$$

**A re-written Primal:**

$$\min_w \max_{\alpha, \alpha_i \geq 0} \mathcal{L}(w, \alpha)$$

# Optimization Review:

## Lagrangian Duality, cont. (Extra)

- Recall the Primal Problem:

$$\min_w \max_{\alpha, \alpha_i \geq 0} \mathcal{L}(w, \alpha)$$

- The Dual Problem:

$$\max_{\alpha, \alpha_i \geq 0} \min_w \mathcal{L}(w, \alpha)$$

- **Theorem (weak duality):**

$$d^* = \max_{\alpha, \alpha_i \geq 0} \min_w \mathcal{L}(w, \alpha) \leq \min_w \max_{\alpha, \alpha_i \geq 0} \mathcal{L}(w, \alpha) = p^*$$

- **Theorem (strong duality):**

Iff there exist a saddle point of  $\mathcal{L}(w, \alpha)$

we have

$$d^* = p^*$$

# Dual representation of the hard SVM QP

- We will start with the linearly separable case
- Instead of encoding the correct classification rule and constraint we will use Lagrange multipliers to encode it as part of the our minimization problem

$$\text{Min } (\mathbf{w}^T \mathbf{w})/2$$

s.t.

$$(\mathbf{w}^T \mathbf{x}_i + b) y_i \geq 1$$

Recall that Lagrange multipliers can be applied to turn the following problem:

$$L_{\text{primal}}(\mathbf{w}, b, \alpha) = \frac{1}{2} \mathbf{w} \cdot \mathbf{w} - \sum_{i=1}^N \alpha_i (y_i (\mathbf{w} \cdot \mathbf{x}_i + b) - 1)$$

# The Dual Problem (Extra)

$$\max_{\alpha_i \geq 0} \min_{w, b} \mathcal{L}(w, b, \alpha)$$

Dual formulation

- We minimize  $\mathcal{L}$  with respect to  $w$  and  $b$  first:

$$\nabla_w \mathcal{L}(w, b, \alpha) = w - \sum_{i=1}^{\text{train}} \alpha_i y_i x_i = 0, \quad (*)$$

$$\nabla_b \mathcal{L}(w, b, \alpha) = \sum_{i=1}^{\text{train}} \alpha_i y_i = 0, \quad (**)$$

Note that  $(*)$  implies:

$$w = \sum_{i=1}^{\text{train}} \alpha_i y_i x_i \quad (***)$$

- Plus  $(***)$  back to  $\mathcal{L}$ , and using  $(**)$ , we have:

$$\max_{\alpha_i} \mathcal{L}(w, b, \alpha) = \sum_{i=1} \alpha_i - \frac{1}{2} \sum_{i,j=1} \alpha_i \alpha_j y_i y_j (\mathbf{x}_i^T \mathbf{x}_j)$$

## Summary: Dual for hard SVM (Extra)

Solving for  $\mathbf{w}$  that gives maximum margin:

1. Combine objective function and constraints into new objective function, using **Lagrange multipliers**  $\alpha_i$

$$L_{primal} = \frac{1}{2} \mathbf{w} \cdot \mathbf{w} - \sum_{i=1}^N \alpha_i (y_i (\mathbf{w} \cdot \mathbf{x}_i + b) - 1)$$

2. To minimize this **Lagrangian**, we take derivatives of  $\mathbf{w}$  and  $b$  and set them to 0:

## Summary: Dual for hard SVM (Extra)

3. Substituting and rearranging gives the **dual** of the Lagrangian:

$$L_{dual} = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \mathbf{x}_i \cdot \mathbf{x}_j$$

which we try to maximize (not minimize).

4. Once we have the  $\alpha_i$ , we can substitute into previous equations to get  $\mathbf{w}$  and  $b$ .
5. This defines  $\mathbf{w}$  and  $b$  as **linear combinations of the training data**.

$$\mathbf{w} = \sum_{i=1}^{train} \alpha_i y_i \mathbf{x}_i$$

# Summary: Dual SVM for linearly separable case

Dual formulation

$$\begin{aligned} \max_{\alpha} \quad & \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j \\ \sum_i \alpha_i y_i &= 0 \\ \alpha_i &\geq 0 \quad \forall i \end{aligned}$$

$n \alpha_i$

$O(n)$   
#para

$w, b$  #para  
 $O(p)$

Min  $(\underline{w^T w})/2$

subject to the following inequality constraints:

For all  $x$  in class + 1

$$\underline{w^T x + b} \geq 1$$

For all  $x$  in class - 1

$$\underline{w^T x + b} \leq -1$$

A total of  $n$  constraints if we have  $n$  input samples



Easier than original QP, more efficient algorithms exist to find  $\alpha_i$ , e.g. SMO (see extra slides)



# Dual formulation for linearly non-separable case

Dual target function:

$$\max_{\alpha} \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j$$

$$\sum_i \alpha_i y_i = 0$$

$$C \geq \alpha_i \geq 0, \forall i$$

Hyperparameter  $C$   
should be tuned  
through k-folds CV

The only difference is that  
the  $\alpha$  are now  
bounded

upper

$O(n)$   
#para

primal linear soft  
margin + Hinge  $\left\{ \begin{array}{l} \epsilon_i \forall i \\ W, b \end{array} \right.$

$O(n + p + 1)$

This is very similar to the  
optimization problem in the linear  
separable case, except that there  
is an upper bound  $C$  on  $\alpha_i$  now

Once again, efficient algorithm  
exist to find  $\alpha_i$

# Prediction via Dual Weights for linear case

$$f(x_{ts}) = \text{Sign}(w^T x_{ts} + b)$$

$$w = \sum_{i=1}^n \alpha_i y_i x_i$$

Dual target function:

$$\max_{\alpha} \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j x_i^T x_j$$

$$\sum_i \alpha_i y_i = 0$$

$$C \geq \alpha_i \geq 0, \forall i$$

Hyperparameter C  
should be tuned  
through k-folds CV

The only difference is that  
the  $\alpha$  are now  
bounded

To evaluate a new sample  $x_{ts}$  we  
need to compute:

$$w^T x_{ts} + b = \sum_{i \in \text{supportV}} \alpha_i y_i x_i^T x_{ts} + b$$

This is very similar to the  
optimization problem in the linear  
separable case, except that there is  
an upper bound C on  $\alpha_i$  now

Once again, efficient algorithm exist  
to find  $\alpha_i$

# Dual SVM – Training using Kernel Matrix

Our dual target function:  $\max_{\alpha} \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j$

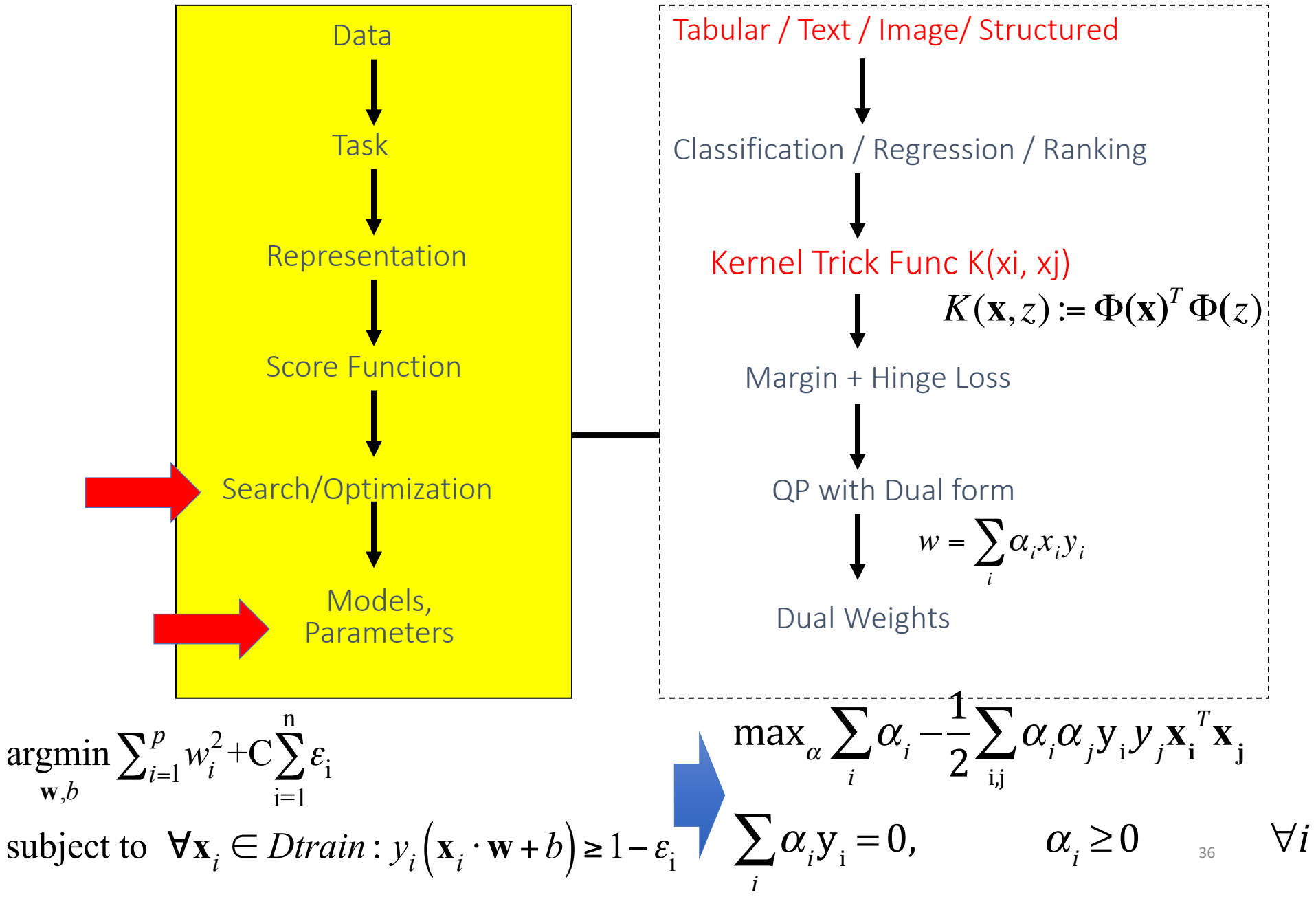
$$\sum_i \alpha_i y_i = 0$$

$$C \geq \alpha_i \geq 0, \forall i$$

Dot product among all training samples

Handwritten diagram illustrating the kernel matrix structure. The matrix is  $n \times n$  and contains dot products  $\mathbf{x}_i^T \mathbf{x}_j$  between training samples  $\mathbf{x}_i$  and  $\mathbf{x}_j$ . The matrix is labeled "matrix" in red.

# This: Kernel Support Vector Machine



Support vectors: non-zero  $\alpha_i$

- only a few  $\alpha_i$  can be nonzero!!

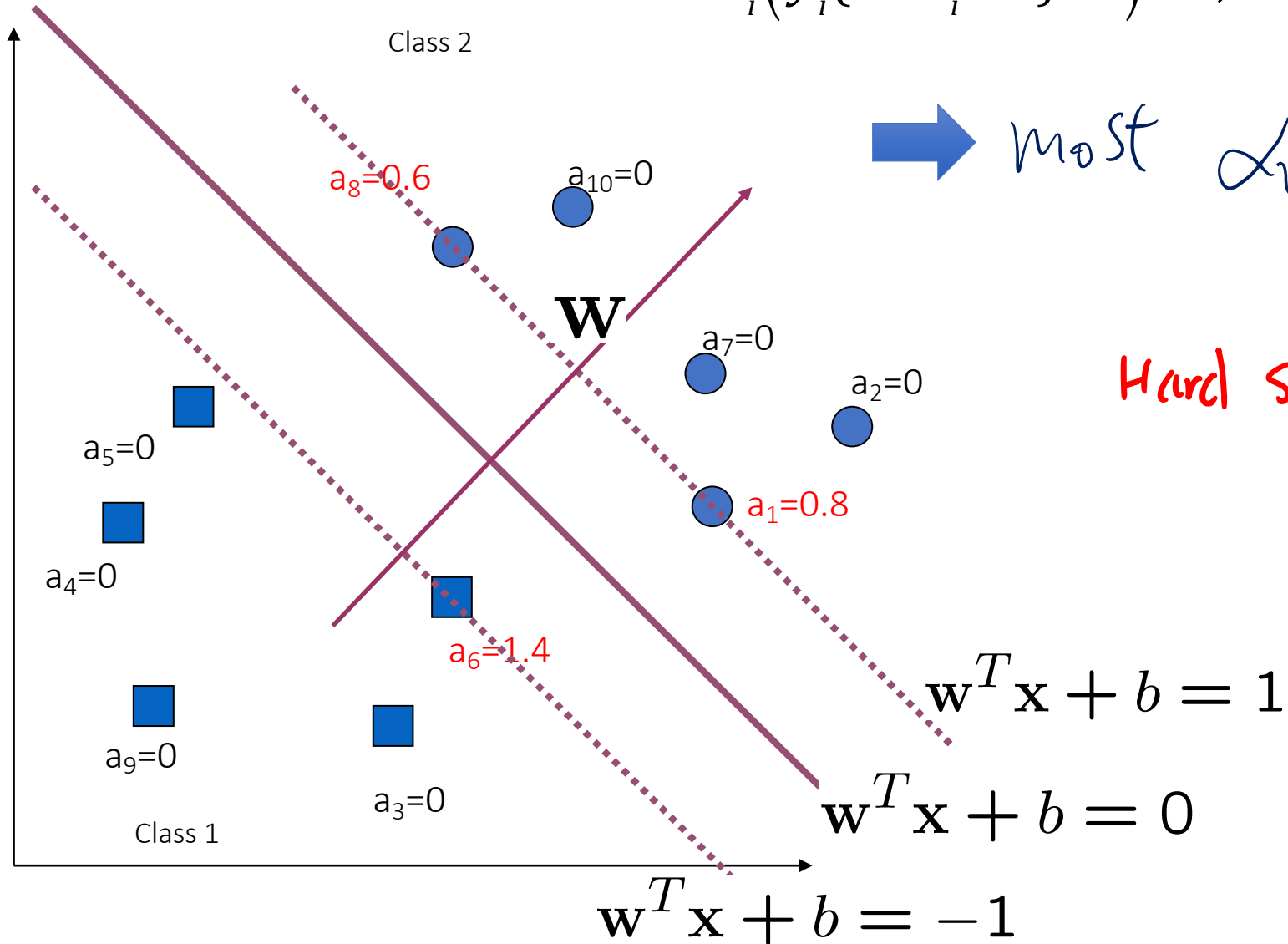
$$\forall i \Rightarrow \alpha_i (y_i (\mathbf{w}^T \mathbf{x}_i + b) - 1) = 0, \quad i = 1, \dots, n$$

Handwritten diagram illustrating the SVM margin and slack variables. It shows three parallel lines representing the decision boundary and margins. The top solid line is labeled  $y_i (\mathbf{w}^T \mathbf{x}_i + b) > 1 \Rightarrow \alpha_i = 0$ . The middle dashed line is labeled  $+1 \quad \mathbf{w}^T \mathbf{x} + b = 1$  and  $- \alpha_i > 0$ . The bottom solid line is labeled  $< -1 \quad \alpha_i = 0$  and  $-1 \quad \mathbf{w}^T \mathbf{x} + b = -1$ .

$$\alpha_i (y_i (\mathbf{w} \cdot \mathbf{x}_i + b) - 1) = 0, \quad i = 1, \dots, n$$

→ most  $\alpha_i = 0$

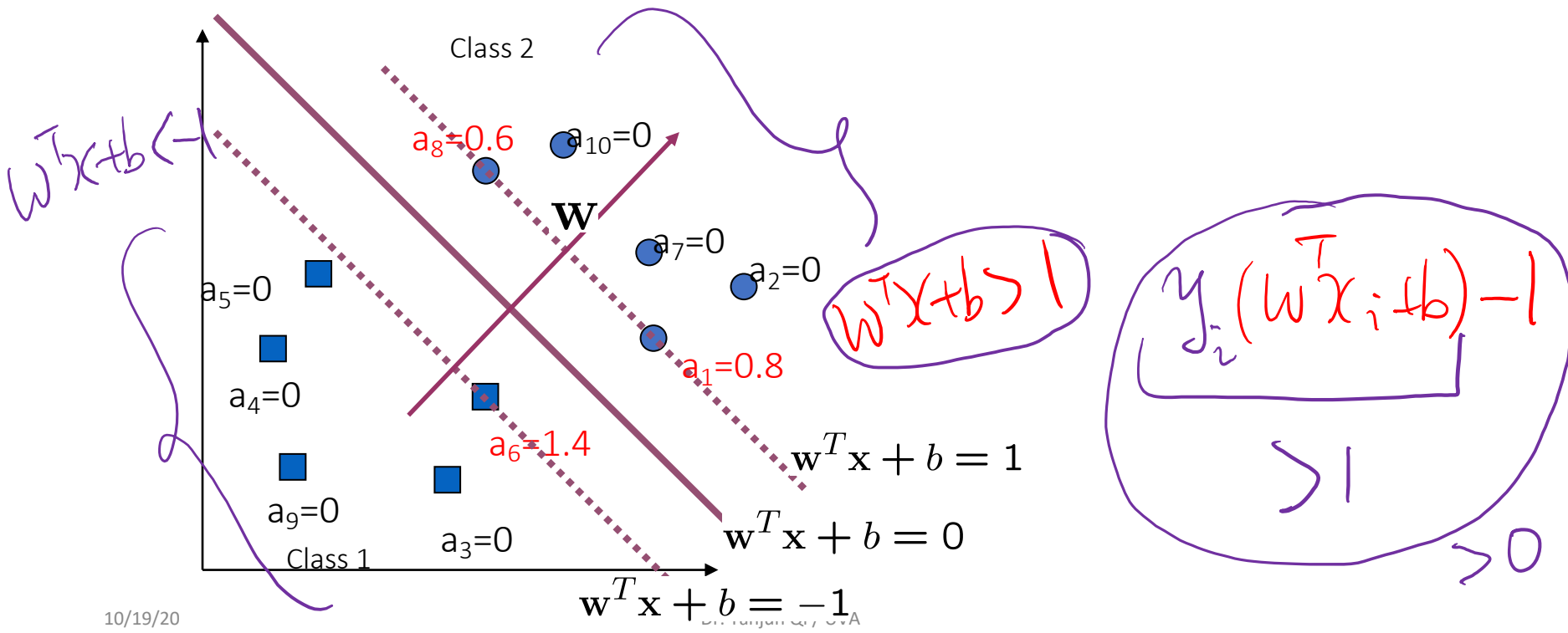
Hard SVM



- only a few  $\alpha_i$  can be nonzero! i.e.

$$y_i(\mathbf{w} \cdot \mathbf{x}_i + b) = 1$$

$$\alpha_i(y_i(\mathbf{w} \cdot \mathbf{x}_i + b) - 1) = 0, \quad i = 1, \dots, n \quad \rightarrow \quad \text{for most } \alpha_i \Rightarrow \alpha_i = 0$$



# Support vectors: non-zero $a_i$

- only a few  $a_i$  can be nonzero!!

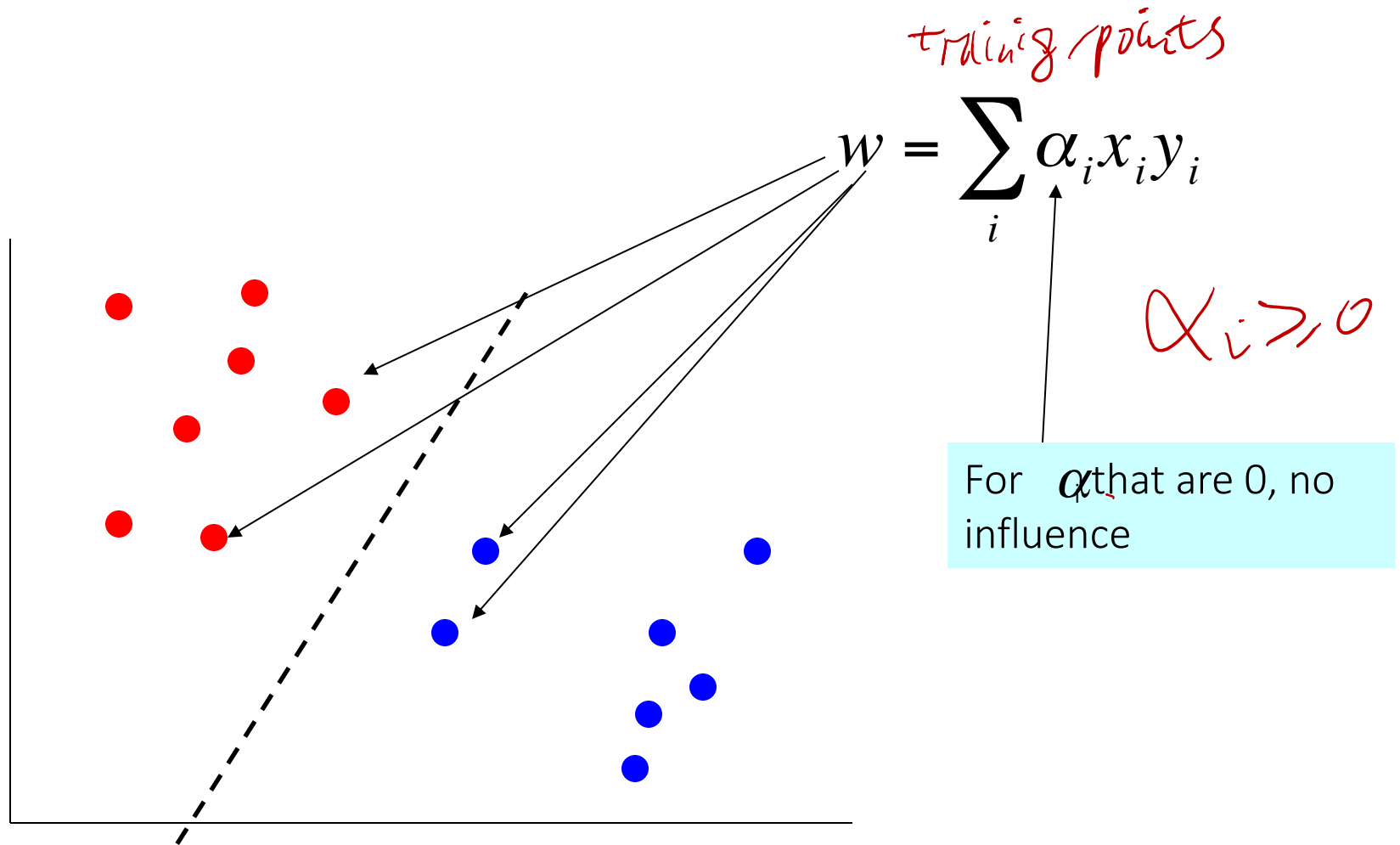
$$\alpha_i (y_i (\mathbf{w} \cdot \mathbf{x}_i + b) - 1) = 0, \quad i = 1, \dots, n$$

for most  $i \Rightarrow \alpha_i = 0$

We call the training data points whose  $a_i$ 's are nonzero the **support vectors** (SV)



# Dual SVM - interpretation




# Dual SVM– Testing

To evaluate a new sample  $\mathbf{x}_{ts}$  we need to compute:

$$\hat{y}_{ts} = \text{sign}(w^T \mathbf{x}_{ts} + b) = \text{sign}\left(\sum_{i=1..n} \alpha_i y_i \mathbf{x}_i^T \mathbf{x}_{ts} + b\right)$$

Dot product with (“all” ??)  
training samples


$$\hat{y}_{ts} = \text{sign}\left(\sum_{i \in \text{Support Vectors}} \alpha_i y_i (\mathbf{x}_i^T \mathbf{x}_{ts}) + b\right)$$

↓  
 $\alpha_i > 0$

For  $\alpha_i$  that are 0,  
no influence

# Support Vectors for the Soft-SVM

$$\alpha_i \left( y_i (\mathbf{w} \cdot \mathbf{x}_i + b) - 1 + \varepsilon_i \right) = 0, \quad i = 1, \dots, n$$

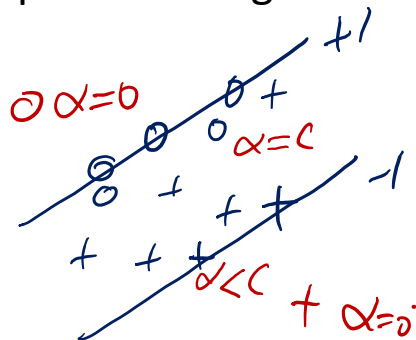
- Support vectors are

$$y_i (\mathbf{x}_i \cdot \mathbf{w} + b) = 1,$$

- Samples on the margin:

$$0 < \alpha_i < C$$

- Samples violating the margin (mostly inside the margin area):



$$y_i (\mathbf{x}_i \cdot \mathbf{w} + b) < 1,$$

$$\alpha_i = C$$

More in L11Extra-SVMoptimDual

# Value C and Number of Support Vectors (no clear relation!!!)

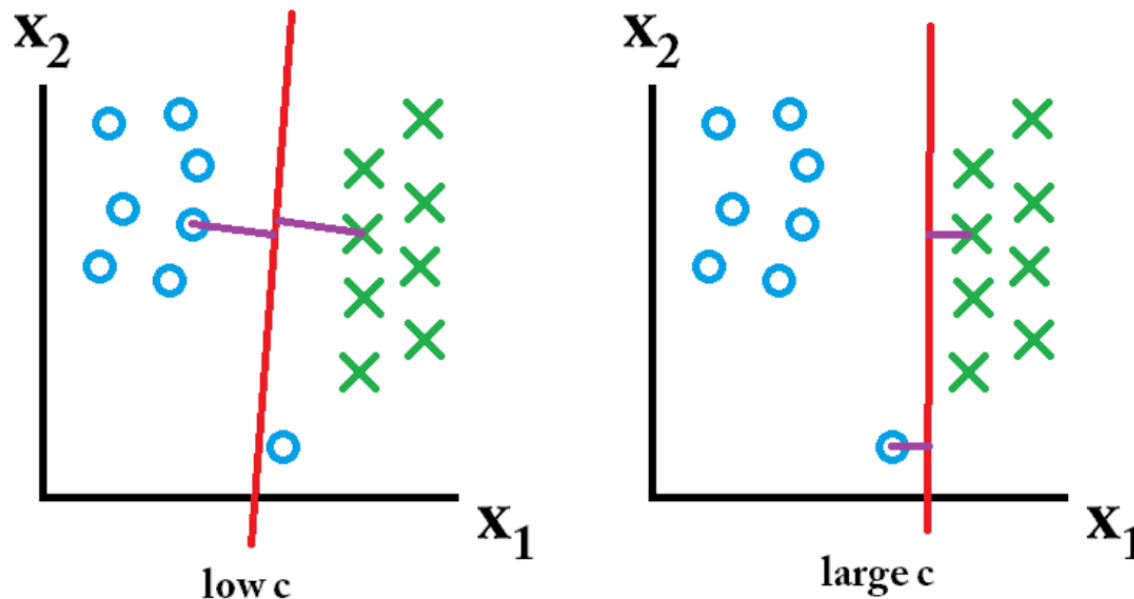
C is how hard we want to punish misclassified examples in the training set.

So for large values of C, it would make sense that if misclassified examples are severely punished, then it will choose a small margin with not many support vectors.

For small values of C, it would broaden the margin, and as a result, end up getting more points inside of the margin (if not easily separable), so there would be possibly more support vectors.

I thought that this [StackExchange post](#) described it pretty well.

As in this example, the small c has more "support vectors" than the large c



<https://github.com/scikit-learn/scikit-learn/issues/7955>

<https://stats.stackexchange.com/questions/31066/what-is-the-influence-of-c-in-svms-with-linear-kernel>

(Recap)

KNN:  $\hat{y}_{ts} = \frac{1}{k} \sum_{i \in k \text{ Neighbors of } X_{ts}} y_i$

find  $k$  neighbor of  $X_{ts} \sim O(n \times k)$

SVM:  $\hat{y}_{ts} = \sum_{i \in SV} \alpha_i y_i k(\vec{x}_i, \vec{x}_{ts}) + b$

# para  $\sim O(n)$

Logistic Regression / Linear Classifier  
 $\hat{y}_{ts} = \sigma(W^T X_{ts} + b)$

# para  $\sim O(p)$

# Summary of SVM

## □ Support Vector Machine (SVM)

- ✓ History of SVM
- ✓ Large Margin Linear Classifier
- ✓ Define Margin ( $M$ ) in terms of model parameter
- ✓ Optimization to learn model parameters ( $w, b$ )
- ✓ Non linearly separable case
- ✓ Optimization with dual form
- ✓ Nonlinear decision boundary
- ✓ Practical Guide
  - ✓ File format / LIBSVM
  - ✓ Feature preprocsssing
  - ✓ Model selection
  - ✓ Pipeline procedure

# References

- Big thanks to Prof. Ziv Bar-Joseph and Prof. Eric Xing @ CMU for allowing me to reuse some of his slides
- Elements of Statistical Learning, by Hastie, Tibshirani and Friedman
- Prof. Andrew Moore @ CMU's slides
- Tutorial slides from Dr. Tie-Yan Liu, MSR Asia
- A Practical Guide to Support Vector Classification Chih-Wei Hsu, Chih-Chung Chang, and Chih-Jen Lin, 2003-2010
- Tutorial slides from Stanford "Convex Optimization I — Boyd & Vandenberghe