

**CSE/ISyE 6740**  
**Computational Data Analytics**

# **Support Vector Machine**

09/22/2025

Kai Wang, Assistant Professor in Computational Science and Engineering  
[kwang692@gatech.edu](mailto:kwang692@gatech.edu)

# Outline

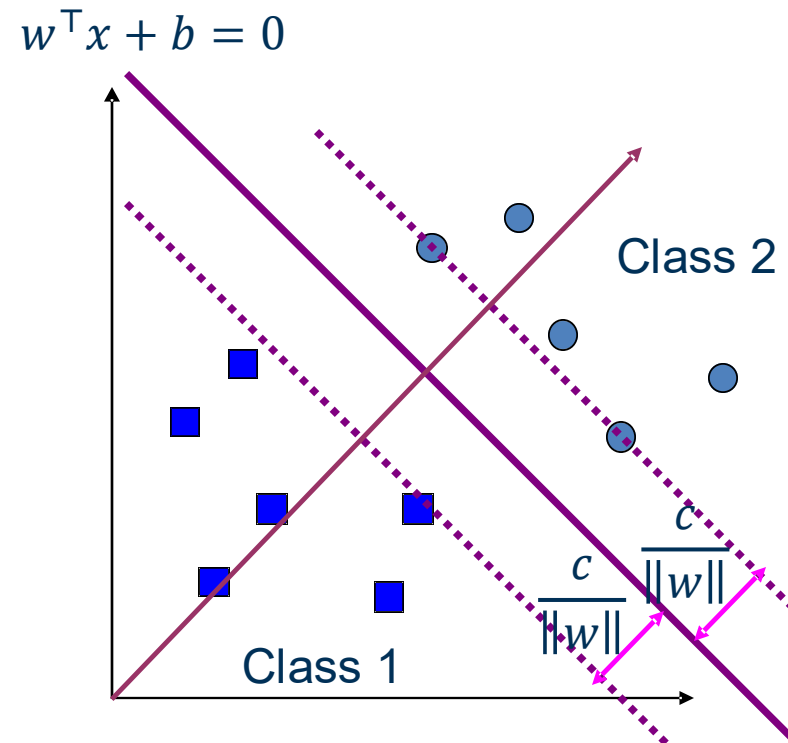
- **Supervised Learning**
  - **Support Vector Machine**
    - From hard-margin SVM to soft-margin SVM
    - Primal and dual problem of soft-margin SVM
    - Hinge loss
    - Support vectors
    - Comparison of SVM and logistic regression

# From Hard Margin to Soft Margin SVM

# Hard-margin SVM

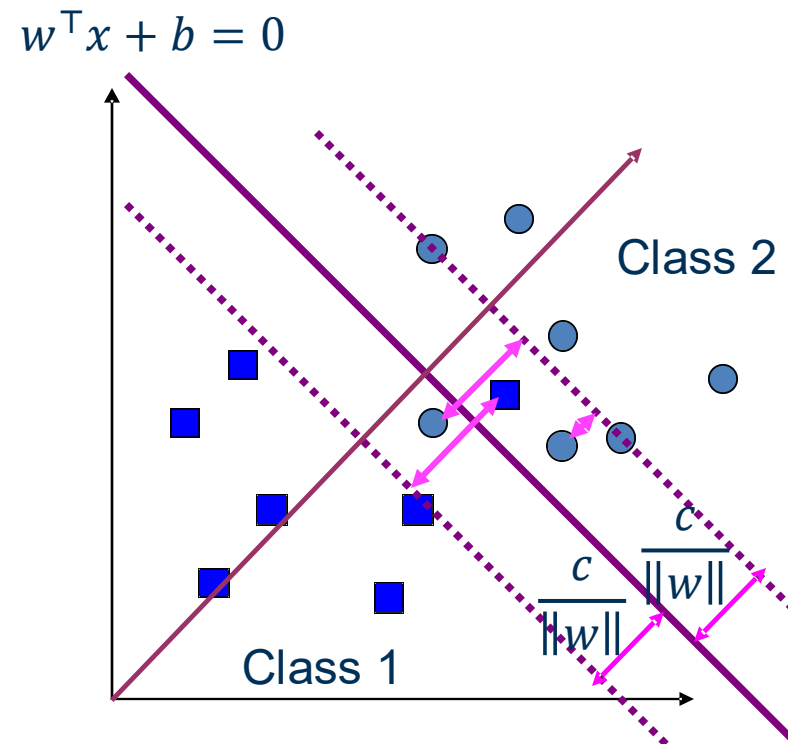
- Find decision boundary  $w$  as far from data point as possible

$$\min_{w,b} \frac{1}{2} \|w\|^2$$
$$s.t. y^i(w^\top x^i + b) \geq 1, \quad \forall i$$



# Soft-margin SVM

- What if the data is **NOT** linearly separable?
- We will allow points to violate the hard-margin constraint

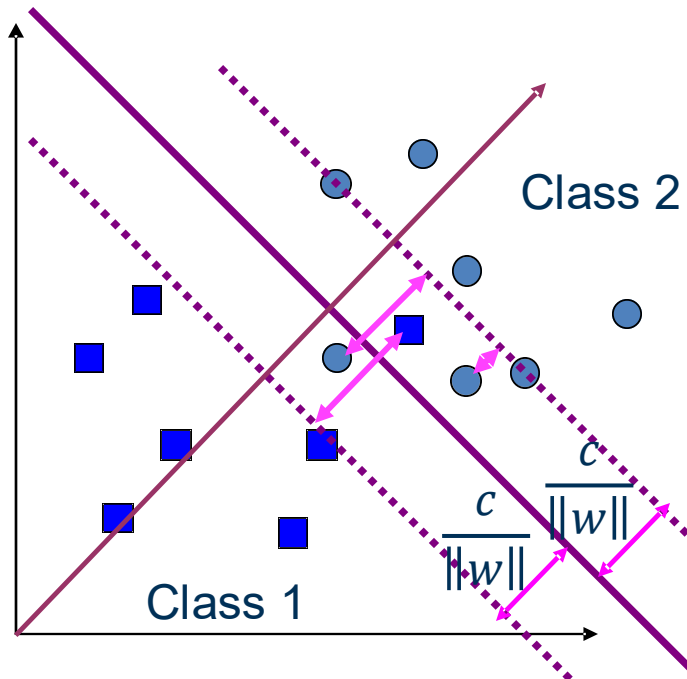


# Soft-margin SVM

- Find decision boundary  $w$  as far from data point as possible (with some exceptions  $\xi$ )

$$\min_{w, b, \xi} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi^i$$
$$s. t. y^i (w^\top x^i + b) \geq 1 - \xi^i, \quad \xi^i \geq 0 \quad \forall i$$

$$w^\top x + b = 0$$



penalty. should  
be positive

# Soft-margin SVM Primal Problem

- Express the optimization problem in standard form:

$$\begin{aligned} \min_{w,b,\xi} \quad & \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi^i \\ \text{s.t.} \quad & 1 - \xi^i - y^i(w^\top x^i + b) \leq 0 \quad \forall i, \\ & -\xi^i \leq 0 \quad \forall i. \end{aligned}$$

# Soft-Margin SVM Dual Problem



# Deriving the Dual Problem

- Define Lagrangian:

$$= \frac{1}{2} w^T w + \sum_{i=1}^n \overset{\substack{\uparrow \\ \text{constant number.}}}{C} \xi^i + \alpha_i (1 - \xi^i - y^i (w^T x^i + b)) - \beta_i \xi^i$$

*put C as C<sub>i</sub>*  
*=> there is diff weights on diff penalty.*

- Taking derivative and set to zero:

$$\frac{\partial L}{\partial w} = w^* - \sum_{i=1}^n \alpha_i y^i x^i = 0$$

$$\frac{\partial L}{\partial b} = \sum_{i=1}^n \alpha_i y^i = 0$$

$$\frac{\partial L}{\partial \xi^i} = C - \alpha_i - \beta_i = 0$$

# Plug Back Relation of $w$ , $b$ , and $\xi$

$$\begin{aligned}
 g(\alpha, \beta) &:= \inf_{w, b, \xi} L(w, b, \xi, \alpha, \beta) = L(w^*, b^*, \xi^*, \alpha, \beta) \\
 &= \frac{1}{2} \left( \sum_{i=1}^n \alpha_i y^i x^i \right)^\top \left( \sum_{i=1}^n \alpha_i y^i x^i \right) + \sum_{i=1}^n \cancel{C \xi^i} + \alpha_i \left( 1 - \cancel{\xi^i} - y^i \left( \left( \sum_{i=1}^n \alpha_i y^i x^i \right)^\top x^i + b \right) \right) - \cancel{\beta_i \xi^i} \\
 &= \frac{1}{2} \left( \sum_{i=1}^n \alpha_i y^i x^i \right)^\top \left( \sum_{j=1}^n \alpha_j y^j x^j \right) + \sum_{i=1}^n \alpha_i \left( 1 - y^i \left( \left( \sum_{j=1}^n \alpha_j y^j x^j \right)^\top x^i + b \right) \right)
 \end{aligned}$$

- After simplification

$$g(\alpha, \beta) := \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y^i y^j (x^{i^\top} x^j)$$

# The Dual Problem

$$\begin{aligned} \max_{\alpha, \beta} g(\alpha, \beta) &:= \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y^i y^j (x^i{}^\top x^j) \\ \text{s.t. } C - \alpha_i - \beta_i &= 0, \quad \alpha_i \geq 0, \quad \beta_i \geq 0 \quad \forall i \\ \sum_{i=1}^n \alpha_i y^i &= 0 \end{aligned}$$

- The constraint  $C - \alpha_i - \beta_i = 0$ ,  $\alpha_i \geq 0$ ,  $\beta_i \geq 0$  can be simplified to  $C \geq \alpha_i \geq 0$
- This is a constrained quadratic programming
- Nice and concave, and global maximum can be found

# The Dual Problem

soft margin. SVM

$$\begin{aligned} \max_{\alpha} g(\alpha) &:= \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y^i y^j (x^i{}^\top x^j) \\ \text{s.t. } C &\geq \alpha_i \geq 0 \quad \forall i \\ \sum_{i=1}^n \alpha_i y^i &= 0 \end{aligned}$$

- **Intuition:**  $\rightarrow$  no points violate margin.
  - When  $C \rightarrow \infty$ , then this becomes the hard margin SVM.
  - Nice and concave, and global maximum can be found

Hard-margin SVM

$$\begin{aligned} \max_{\alpha} g(\alpha) &:= \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y^i y^j (x^i{}^\top x^j) \\ \text{s.t. } \alpha_i &\geq 0 \quad \forall i = 1, 2, \dots, m \\ \sum_{i=1}^n \alpha_i y^i &= 0 \end{aligned}$$

# Soft-Margin Kernel SVM

$$\begin{aligned} \max_{\alpha} g(\alpha) &:= \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y^i y^j K(x^i, x^j) \\ \text{s. t. } C &\geq \alpha_i \geq 0 \quad \forall i \\ \sum_{i=1}^n \alpha_i y^i &= 0 \end{aligned}$$

## Hard-margin kernel SVM

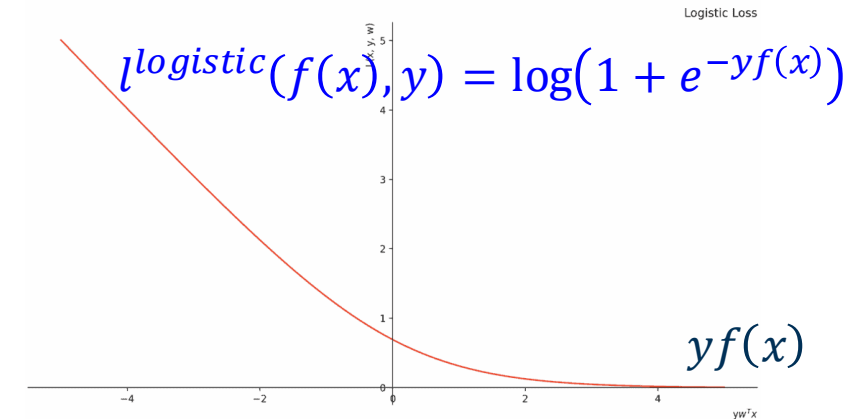
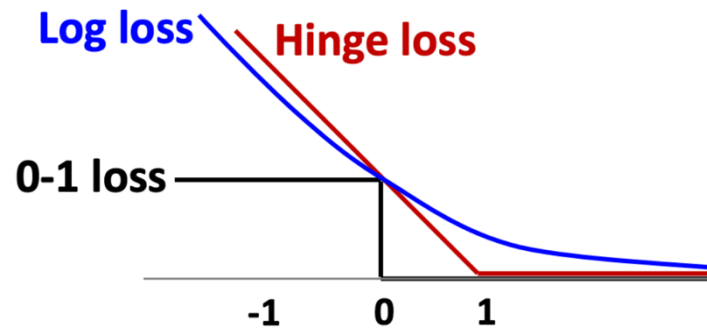
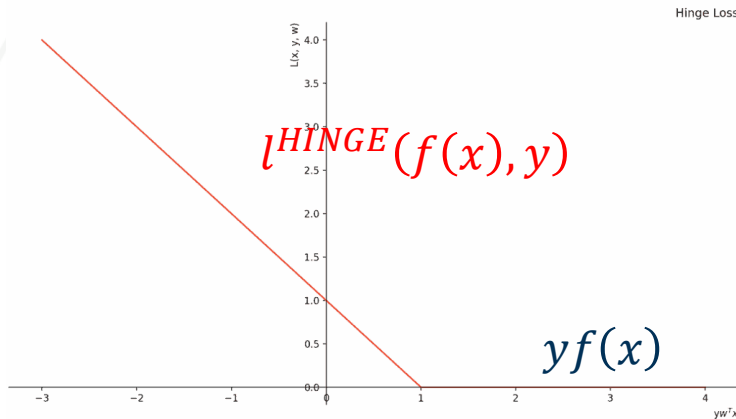
$$\begin{aligned} \max_{\alpha} g(\alpha) &:= \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y^i y^j K(x^i, x^j) \\ \text{s. t. } \alpha_i &\geq 0 \quad \forall i = 1, 2, \dots, m \\ \sum_{i=1}^n \alpha_i y^i &= 0 \end{aligned}$$

# Hinge Loss

# Hinge Loss

- **Definition:** Assuming the label  $y \in \{-1, 1\}$  and the decision rule is  $h(x) = \text{sign}(f(x))$  with  $f(x) = w^\top \phi(x) + b$

$$l^{\text{HINGE}}(f(x), y) = \begin{cases} 0 & \text{if } yf(x) \geq 1 \\ 1 - yf(x) & \text{otherwise} \end{cases}$$



- **Intuition:** penalize more if incorrectly classified (the left size of the curve)
- **Convenient shorthand:**

$$l^{\text{HINGE}}(f(x), y) = \max(0, 1 - yf(x))$$

# Equivalence of Hinge Loss to SVMs

- Minimizing the total hinge loss on all the training data

$$\min_{w,b} \sum_{i=1}^n \max(0, 1 - y^i [w^\top \phi(x^i) + b]) + \frac{\lambda}{2} \|w\|^2$$

Hinge loss

Regularization

- This is analogous to the regularized least square, which balances two terms (the loss and the regularizer). Conventionally, we rewrite the objective function as ( $C = \frac{1}{\lambda}$ ):

$$\min_{w,b} C \sum_{i=1}^n \max(0, 1 - y^i [w^\top \phi(x^i) + b]) + \frac{1}{2} \|w\|^2$$

- We further rewrite this into another equivalent form

$$\begin{aligned} \min_{w,b,\xi} C \sum_{i=1}^n \xi^i + \frac{1}{2} \|w\|^2 \\ \text{s.t. } \max(0, 1 - y^i [w^\top \phi(x^i) + b]) = \xi^i \quad \forall i \end{aligned}$$



# Equivalence of Hinge Loss to SVMs

- **Primal formulation:**

$$\begin{aligned} \min_{w,b,\xi} \quad & C \sum_{i=1}^n \xi^i + \frac{1}{2} \|w\|^2 \\ \text{s.t.} \quad & 1 - y^i [w^\top \phi(x^i) + b] \leq \xi^i \quad \forall i \\ & 0 \leq \xi^i \end{aligned}$$

Where all  $\xi^i$  are called slack variables.

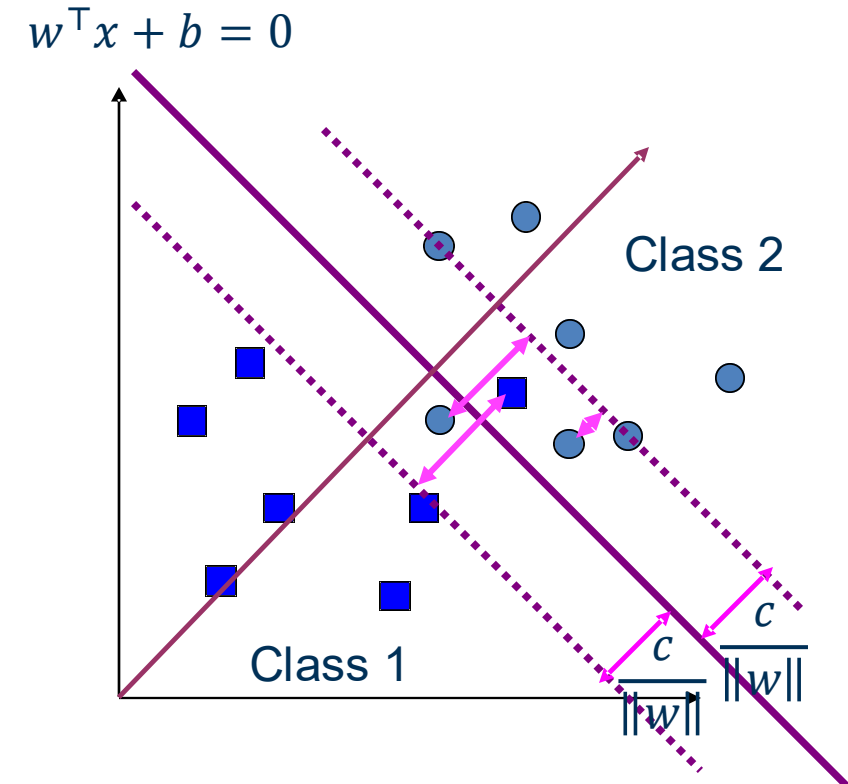
## Soft-margin SVM

$$\begin{aligned} \min_{w,b,\xi} \quad & \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi^i \\ \text{s.t.} \quad & y^i (w^\top x^i + b) \geq 1 - \xi^i, \\ & \xi^i \geq 0 \quad \forall i \end{aligned}$$

# Support Vectors

# Complementary Slackness

- What is the geometric meaning of  $\alpha_i, \beta_i$ ?



## Primal problem

$$\begin{aligned} \min_{w, b, \xi} \quad & \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi^i \\ \text{s.t.} \quad & 1 - \xi^i - y^i (w^\top x^i + b) \leq 0, \\ & -\xi^i \leq 0 \quad \forall i \end{aligned}$$

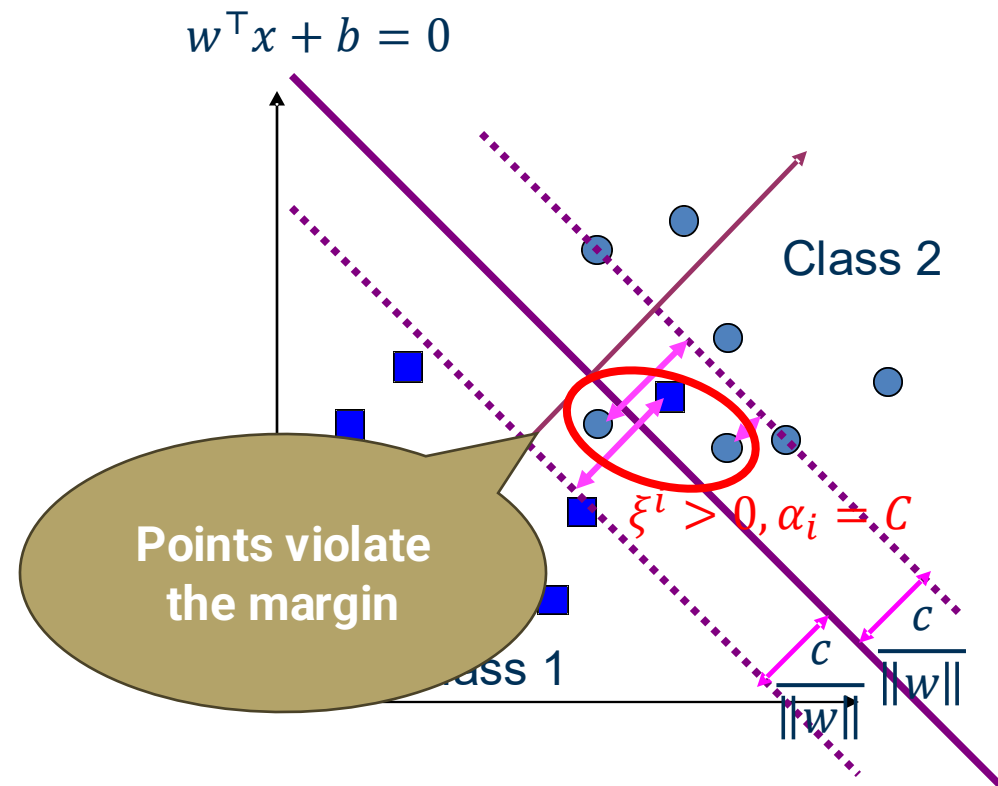
## Dual problem

$$\begin{aligned} \max_{\alpha} \quad & g(\alpha) := \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y^i y^j (x^i{}^\top x^j) \\ \text{s.t.} \quad & C \geq \alpha_i \geq 0 \quad \forall i \\ & \sum_{i=1}^n \alpha_i y^i = 0 \end{aligned}$$

# Complementary Slackness

- **Non-margin support vectors  $\alpha_i = C \neq 0$ :**
  - This means  $\beta_i = C - \alpha_i = 0$
  - Therefore,  $\xi_i > 0$ , and  $1 - y^i(w^\top x^i + b) = \xi_i > 0$

all the  $\alpha$  will be  $C$   
 $\alpha = C$



## Primal problem

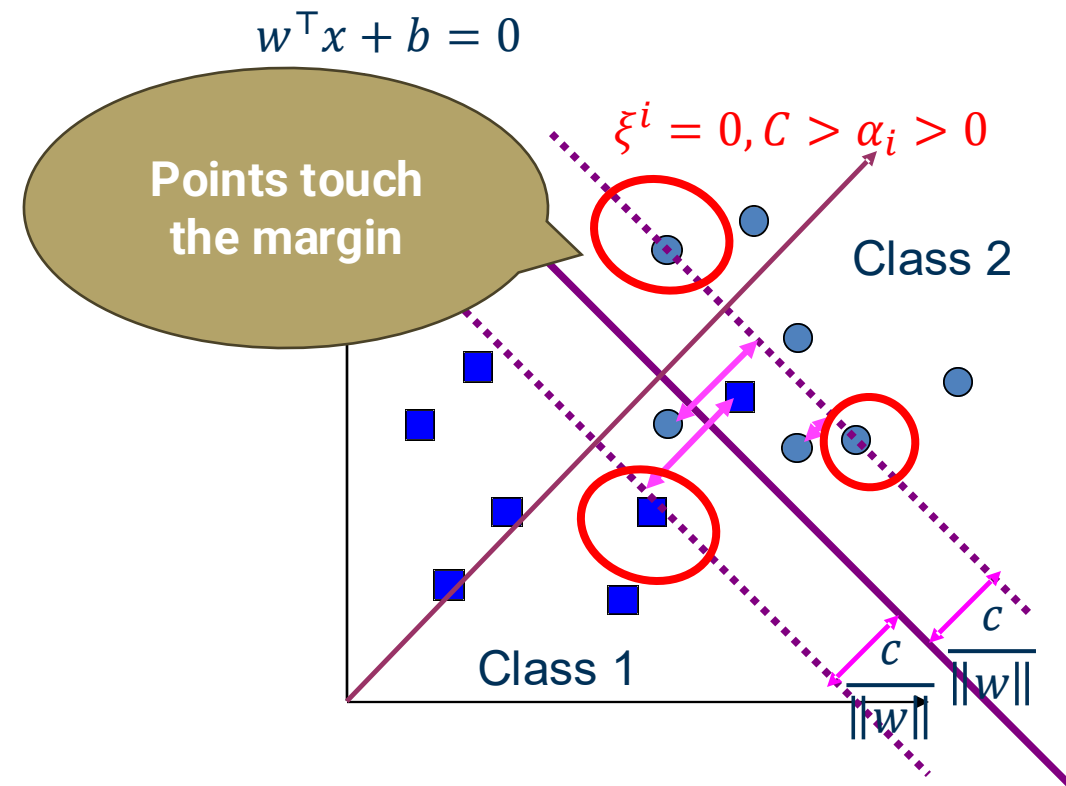
$$\begin{aligned} \min_{w, b, \xi} \quad & \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi^i \\ \text{s.t.} \quad & 1 - \xi^i - y^i(w^\top x^i + b) \leq 0, \\ & -\xi^i \leq 0 \quad \forall i \end{aligned}$$

## Dual problem

$$\begin{aligned} \max_{\alpha} \quad & g(\alpha) := \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y^i y^j (x^{i\top} x^j) \\ \text{s.t.} \quad & C \geq \alpha_i \geq 0 \quad \forall i \\ & \sum_{i=1}^n \alpha_i y^i = 0 \end{aligned}$$

# Complementary Slackness

- **Non-margin support vectors  $\alpha_i = C \neq 0$ :**
  - This means  $\beta_i = C - \alpha_i = 0$
  - Therefore,  $\xi_i > 0$ , and  $1 - y^i(w^\top x^i + b) = \xi_i > 0$
- **Margin support vectors  $0 < \alpha_i < C$ :**
  - This means  $\beta_i = C - \alpha_i > 0$
  - Therefore,  $\xi_i = 0$ , and  $1 - y^i(w^\top x^i + b) = \xi_i = 0$



## Primal problem

$$\begin{aligned} \min_{w, b, \xi} \quad & \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi^i \\ \text{s.t.} \quad & 1 - \xi^i - y^i(w^\top x^i + b) \leq 0, \\ & -\xi^i \leq 0 \quad \forall i \end{aligned}$$

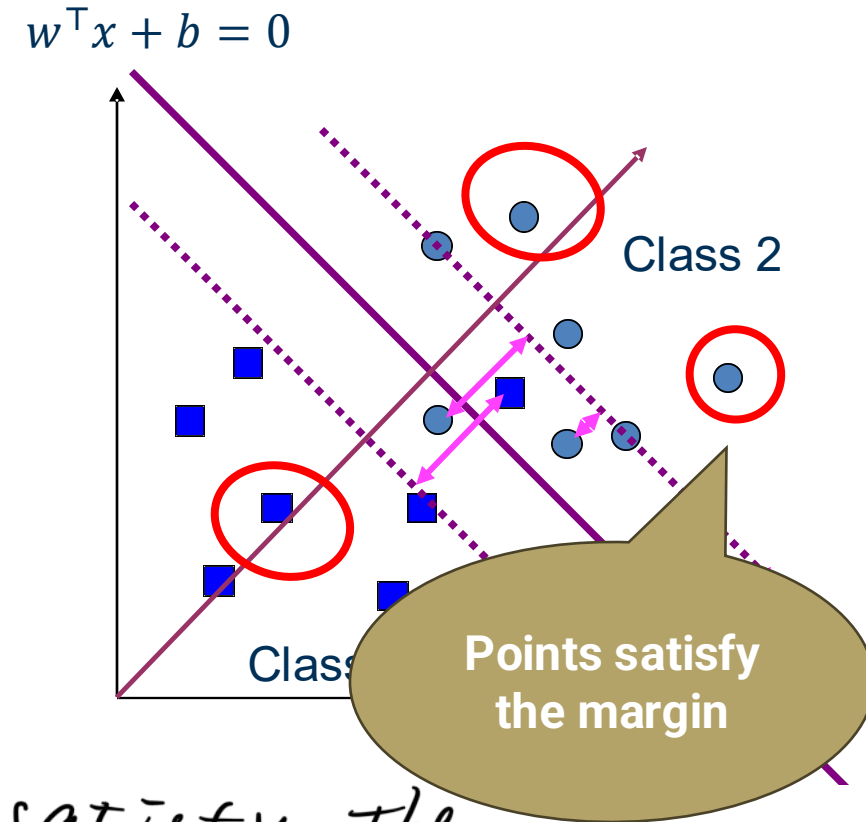
## Dual problem

$$\begin{aligned} \max_{\alpha} \quad & g(\alpha) := \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y^i y^j (x^{i\top} x^j) \\ \text{s.t.} \quad & C \geq \alpha_i \geq 0 \quad \forall i \\ & \sum_{i=1}^n \alpha_i y^i = 0 \end{aligned}$$

# Complementary Slackness

- **Non-margin support vectors**  $\alpha_i = C \neq 0$ :
  - This means  $\beta_i = C - \alpha_i = 0$
  - Therefore,  $\xi_i > 0$ , and  $1 - y^i(w^\top x^i + b) = \xi_i > 0$
- **Margin support vectors**  $0 < \alpha_i < C$ :
  - This means  $\beta_i = C - \alpha_i > 0$
  - Therefore,  $\xi_i = 0$ , and  $1 - y^i(w^\top x^i + b) = \xi_i = 0$
- **Non-support vectors**  $0 = \alpha_i$ :
  - This means  $\beta_i = C - \alpha_i = C > 0$
  - Therefore,  $\xi_i = 0$ , i.e.,  $1 - y^i(w^\top x^i + b) < \xi_i = 0$

*→ perfectly satisfy the margin*



# Inference of Support Vector Machines

- The inference is the same as hard-margin SVM. We make decisions by comparing each new example  $z$  with only the support vectors:

$$y^* = \text{sign}(w^\top z + b) = \text{sign} \left( \left( \sum_{i \in \text{support vectors}} \alpha_i y^i (x^i{}^\top z) \right) + b \right)$$

## Primal

- Primal solution  $w^*, b^*$
- Inference:  $\text{sign}(w^\top z + b)$
- Computation cost:  $O(d)$
- Kernel formulation: no

## Dual

- Dual solution  $\alpha$  with  $n'$  support vectors (nonzero  $\alpha_i$ )
- Inference:  $\text{sign} \left( \left( \sum_{i \in \text{support vectors}} \alpha_i y^i (x^i{}^\top z) \right) + b \right)$ 
  - Need to pre-compute the offset  $b \in \mathbb{R}$  with  $O(1)$  cost
- Computation cost:  $O(dn')$
- Kernel formulation: yes

biggest  
diff

Can get better mapping

# Comparison of SVM and Logistic Regression



# SVMs v.s. Logistic Regression

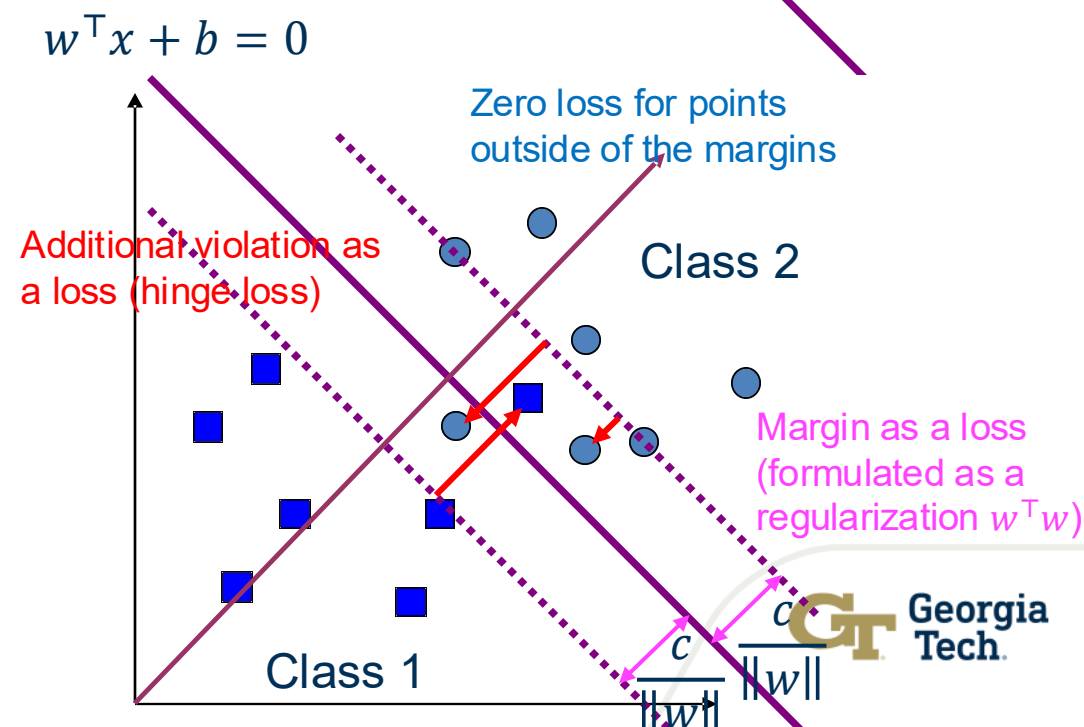
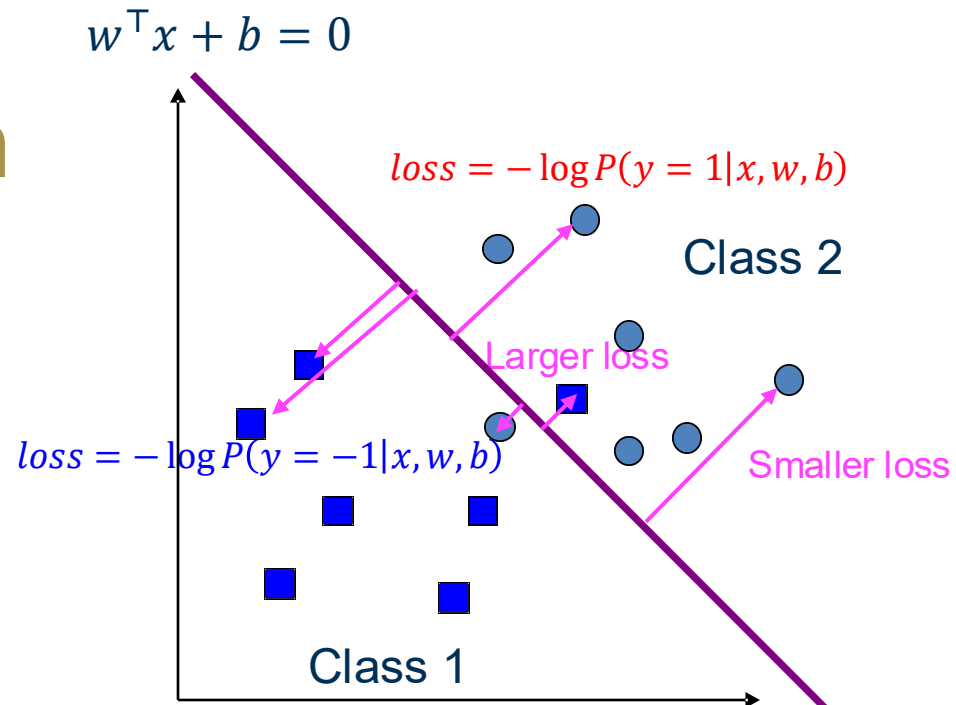
- **Logistic regression** focuses on maximizing the probability of the data. The farther the data lies from the separating hyperplane (on the correct side), the happier LR is.

All points have loss.

- **SVM** tries to find the separating hyperplane that maximizes the distance of the closest points to the margin (the support vectors). If a point is not a support vector, it doesn't really matter.

Inside margins : loss ✓

outside margins : no loss



# SVMs v.s. Logistic Regression

	SVMs	Logistic regression
Loss function	Hinge loss	Log-loss

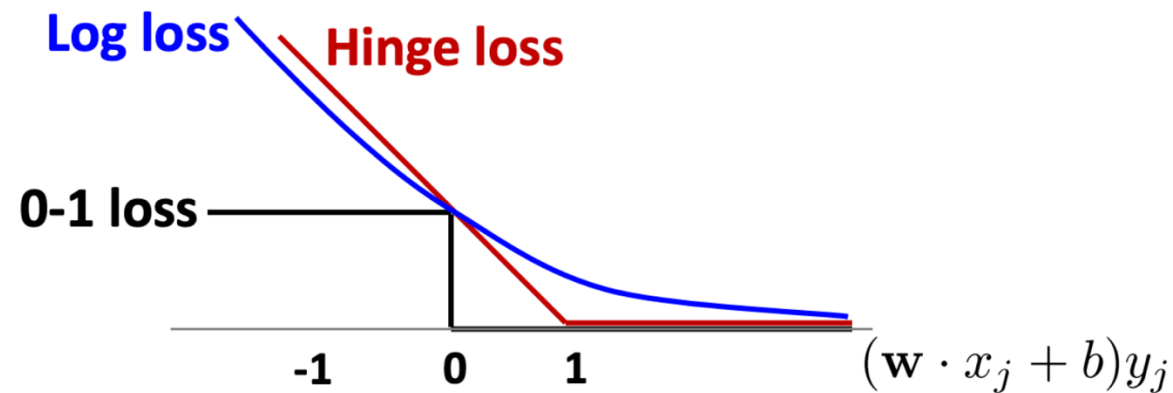
# SVMs v.s. Logistic Regression

- **SVM: Hinge loss**

$$\text{loss}(f(x^i), y^i) = (1 - (w^\top x^i + b)y^i)_+$$

- **Logistic regression: log loss** (negative log likelihood)

$$\text{loss}(f(x^i), y^i) = -\log P(y^i | x^i, w, b) = \log(1 + e^{-(w^\top x^i + b)y^i})$$



# SVMs v.s. Logistic Regression

	<b>SVMs</b>	<b>Logistic regression</b>
<b>Loss function</b>	Hinge loss	Log-loss
<b>High dimensional features</b>	Yes	Yes
<b>Dual solution sparse</b>	Often yes!	Almost always no!
<b>Semantics of output</b>	Margin	Probabilities

# SVMs v.s. Logistic Regression: Complexity

- $n$ : number of training examples
- $n'$ : number of support vectors
- $d$ : number of features
- $E$ : number of stochastic gradient descent (SGD) epochs
- **Logistic regression**
  - Train time complexity:  $O(ndE)$  (SGD)
  - Test time complexity:  $O(d)$
- **SVM**
  - Train time complexity:  $O(n^2d)$  (QP) or  $O(ndE)$  (SGD)
  - Test time complexity:  $O(d)$  (linear SVM) or  $O(dn')$  (kernel SVM)

# Which One to Use?

- **Logistic regression**

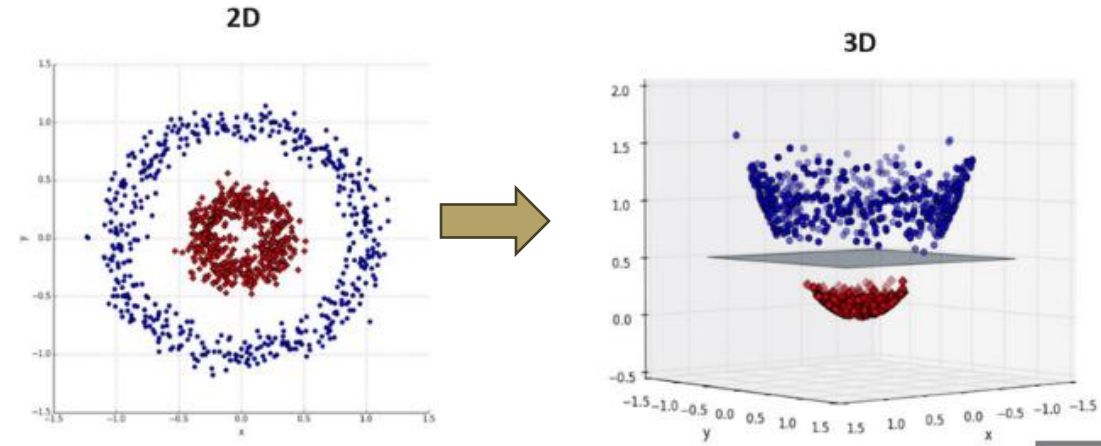
- LR gives calibrated probabilities that can be interpreted as confidence in a decision, while SVM does not have a direct probabilistic interpretation.
- LR can be (straightforwardly) used within Bayesian models.
- LR gives us an unconstrained, smooth objective.

- **SVMs**

- SVMs don't penalize examples for which the correct decision is made with sufficient confidence. This may be good for generalization.
- SVMs have a nice dual form, giving sparse solutions.
- SVM results are strongly dependent on a suitable choice for the softening parameter  $C$ .

# Next Lecture: Kernel Methods

- What is the kernel function  $K$ ?
- What exactly does the kernel method do?
- Why can the kernel method help extend to nonlinearity?



$$\begin{aligned} \max_{\alpha} g(\alpha) &:= \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y^i y^j K(x^i, x^j) \\ \text{s.t. } & \sum_{i=1}^n \alpha_i y^i = 0 \\ & \alpha_i \geq 0 \quad \forall i \end{aligned}$$