



Machine Learning

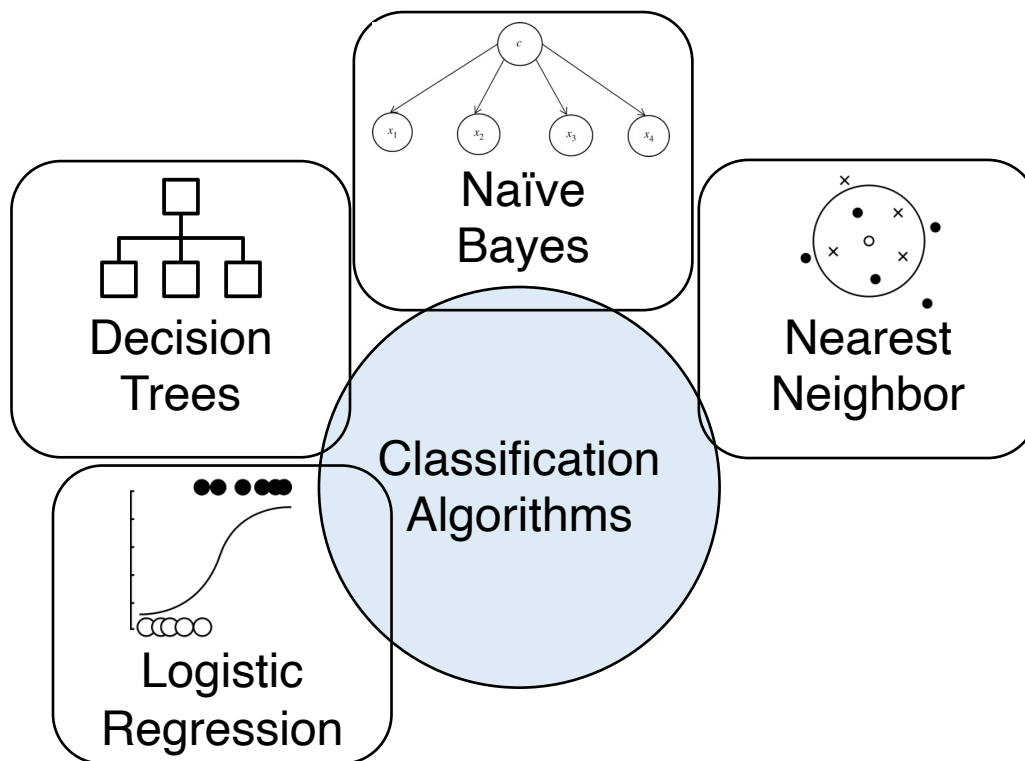
Lecture 7: Support Vector Machine

Fall 2023

Instructor: Xiaodong Gu



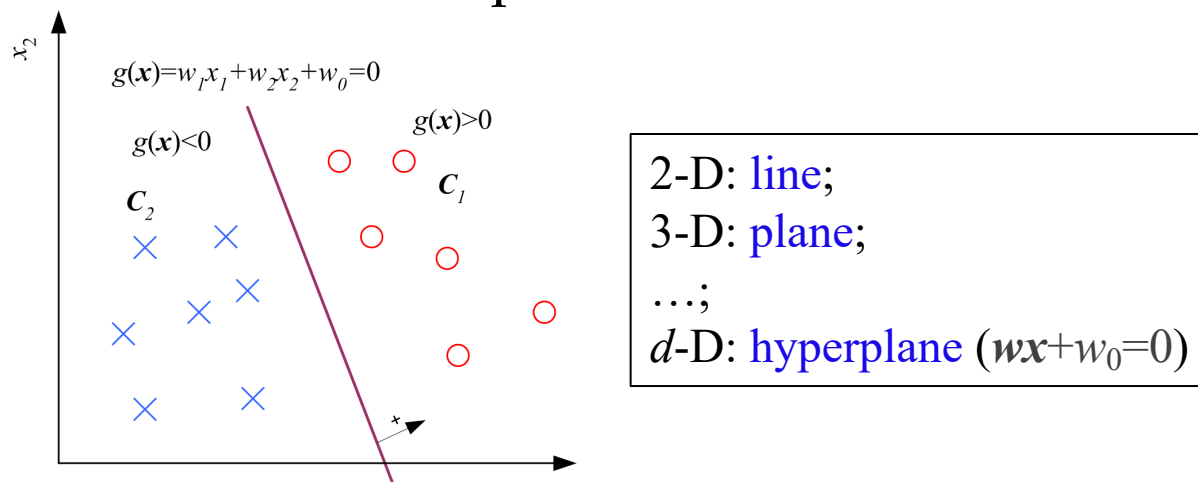
The family of classification





Recall: Linear Classifiers

- **Given:** a training set $\{(x^{(1)}, y^{(1)}), \dots, (x^{(N)}, y^{(N)})\}$ of N samples
 $x^{(\ell)} \in \mathbf{R}^d$: input, $y^{(\ell)} \in \{-1, 1\}$: target (label)
- Assume that the problem is **linearly separable**: there exists a **linear surface** to separate the two classes

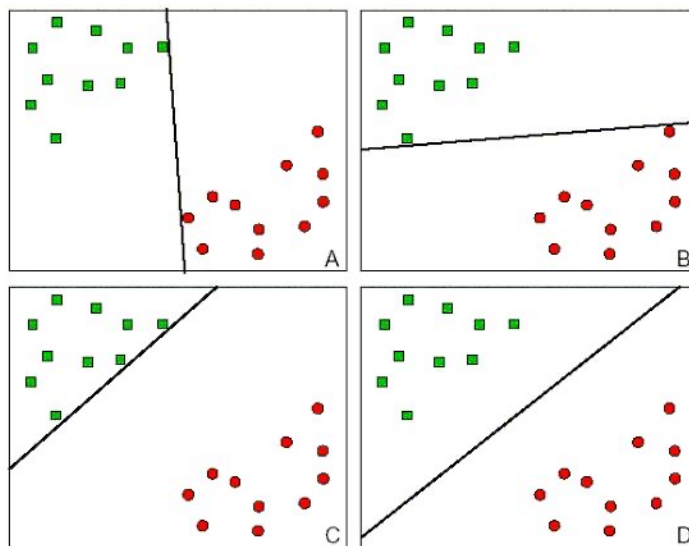
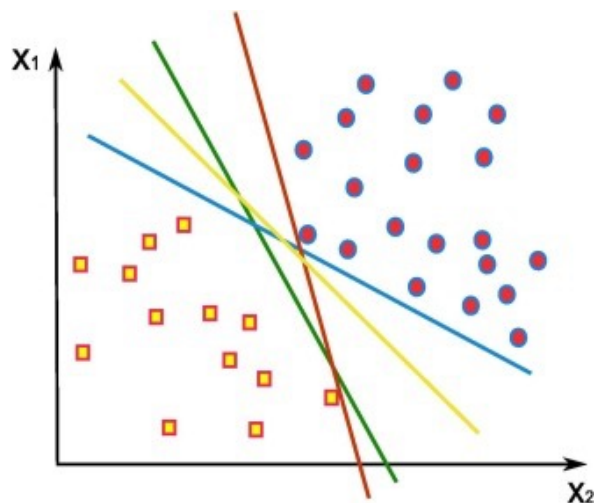


Goal:

Find $w\mathbf{x} + w_0 = 0$ that perfectly separates the two classes

The Problem

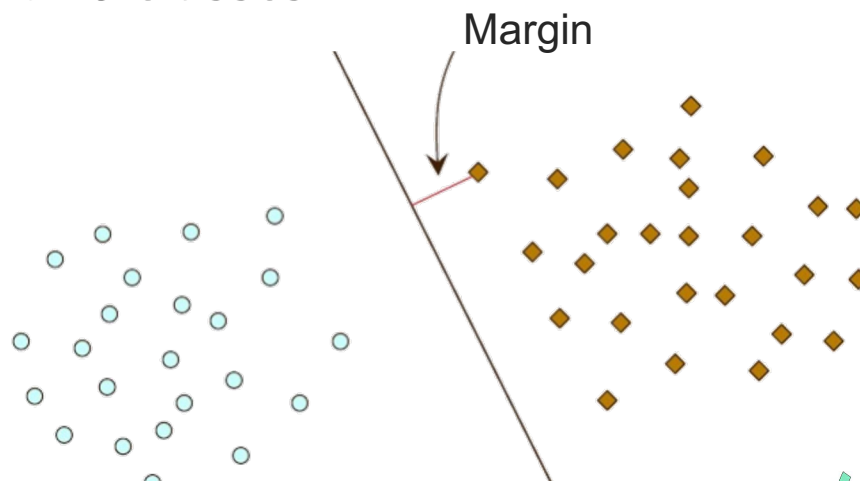
- There can be multiple separating hyperplanes.



Which one is the best?

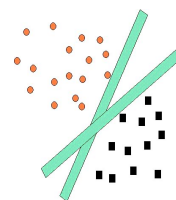
The Idea:

- Find a decision boundary that **maximizes the margin** between two classes.

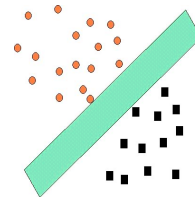


Margin and generalization:

Statistical learning theories have shown that the boundary with the **largest margin** generalizes best (i.e., has the **smallest generalization error**).



skinny margin is more flexible, thus more complex

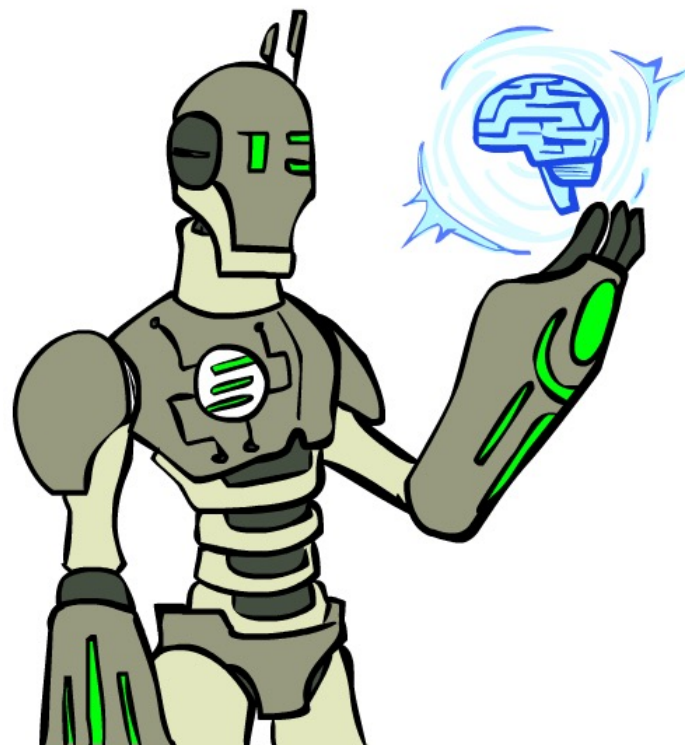


fat margin is less complex

Today

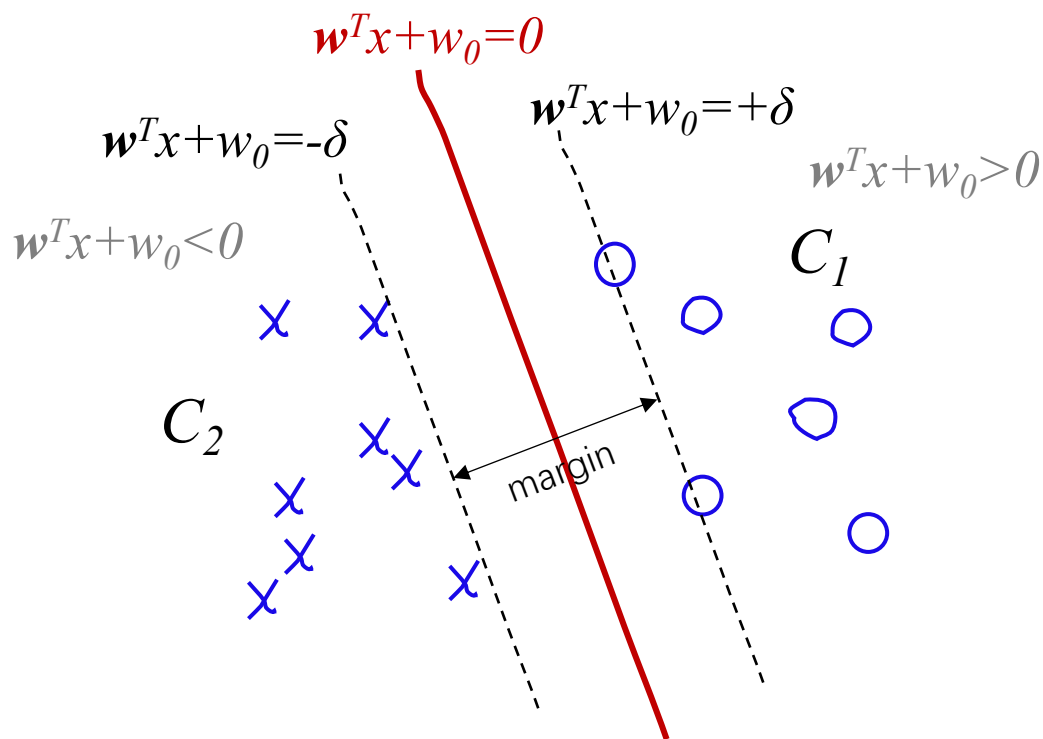
Support Vector Machine

- Problem Formulation
- Dual Problem
- Loss and Optimization



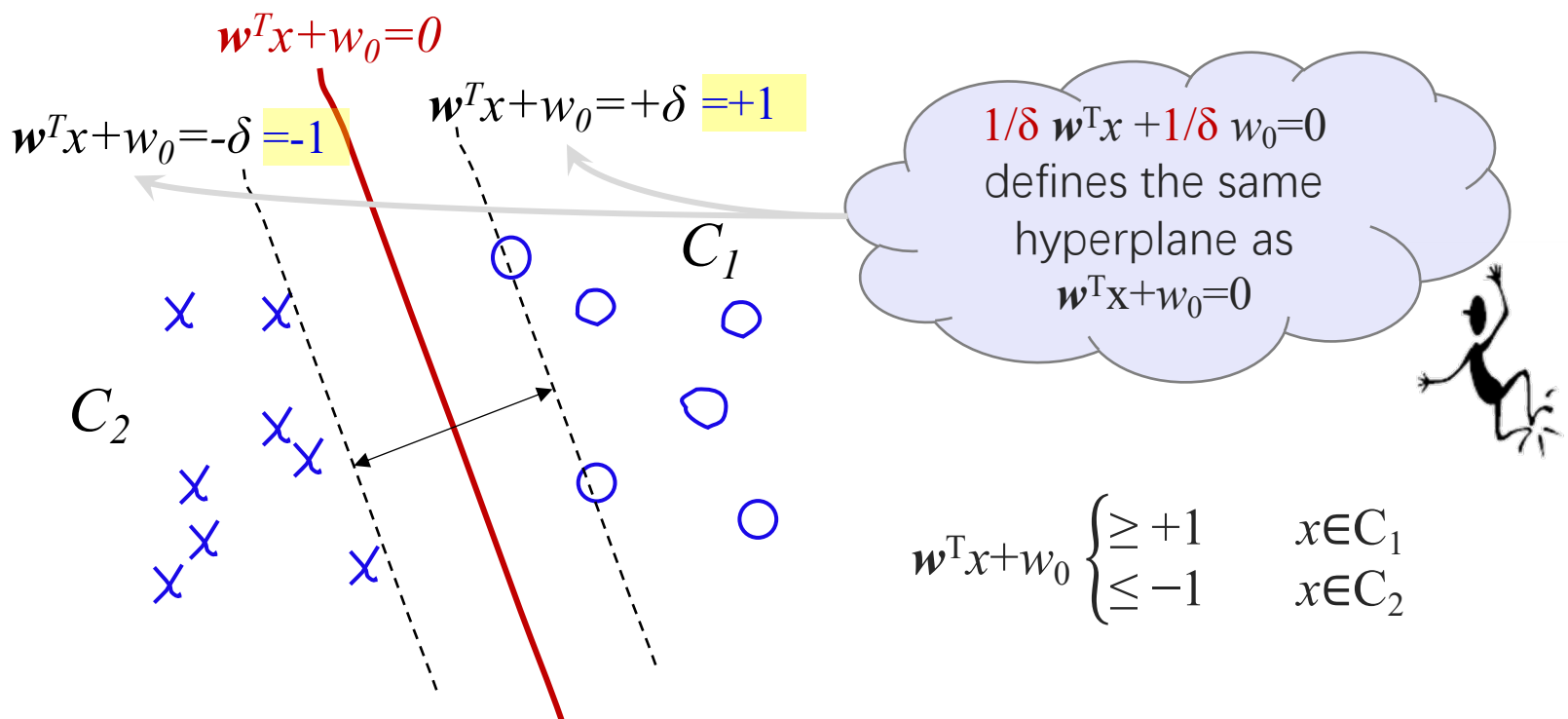
Problem Formulation

- A **linear** discriminant function models a **linear decision boundary** of two classes.



Problem Formulation

- A **linear** discriminant function models a **linear decision boundary** of two classes.





Maximizing the Margin

$$\mathbf{w}^T \mathbf{x} + w_0 = \begin{cases} 1 & \text{for the closest points on one side} \\ -1 & \text{for the closest points on the other} \end{cases}$$

- Let $x^{(1)}$ and $x^{(2)}$ be two closest points on each side of the hyperplane.
- Note that

$$\mathbf{w}^T x^{(1)} + w_0 = +1$$

$$\mathbf{w}^T x^{(2)} + w_0 = -1$$

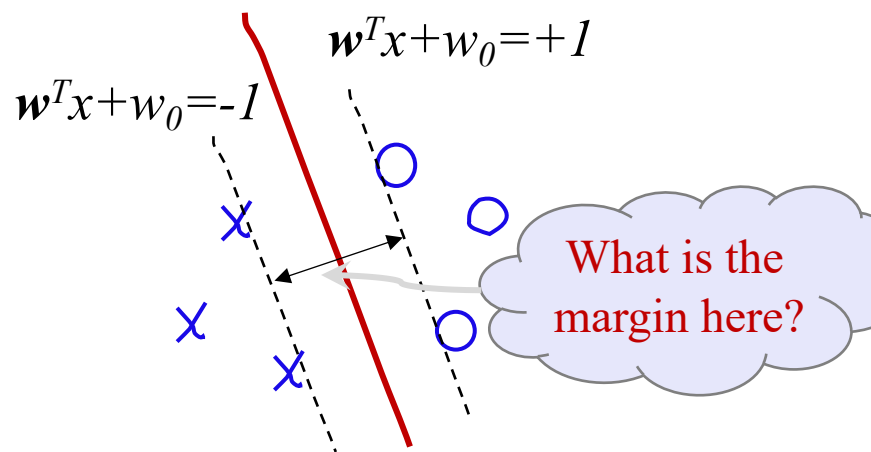
which imply

$$\mathbf{w}^T (x^{(1)} - x^{(2)}) = 2.$$

Hence, the margin can be given by

$$\text{margin} = \frac{\mathbf{w}^T (x^{(1)} - x^{(2)})}{\|\mathbf{w}\|} = \frac{2}{\|\mathbf{w}\|}$$

Maximizing the margin is equivalent to minimizing $\frac{1}{2} \|\mathbf{w}\|^2$





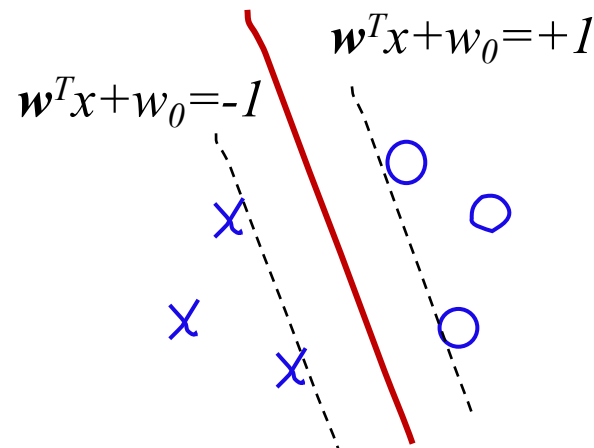
Inequality Constraints

- Given: $D = \{(x^{(1)}, y^{(1)}), \dots, (x^{(N)}, y^{(N)})\}$, we want \mathbf{w} and w_0 to satisfy

$$\mathbf{w}^T x^{(\ell)} + w_0 \begin{cases} \geq +1 & \text{if } y^{(\ell)} = +1 \\ \leq -1 & \text{if } y^{(\ell)} = -1 \end{cases}$$

- Or, equivalently,

$$y^{(\ell)} (\mathbf{w}^T x^{(\ell)} + w_0) \geq 1$$



SVM = Solving an Optimization Problem



In summary, SVM aims to solve a constrained optimization problem:

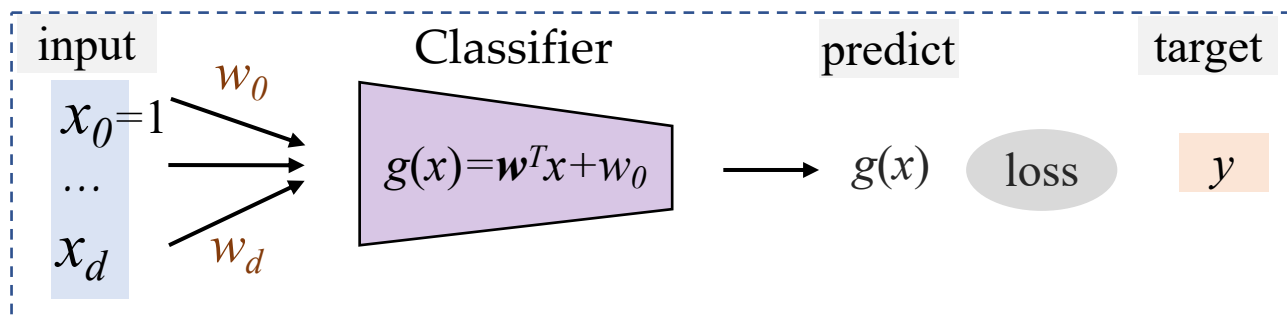
$$\begin{aligned} &\text{minimize} && \frac{1}{2} \|\mathbf{w}\|^2 \\ &\text{subject to} && y^{(\ell)}(\mathbf{w}^T \mathbf{x}^{(\ell)} + w_0) \geq 1, \ell = 1, \dots, N \end{aligned}$$

- This is a [quadratic programming](#) (QP) problem, which is one type of convex optimization problem.
- The complexity depends on the dimensionality d of inputs



Model Architecture

- The same as Perceptron.



- **Train:**
 - optimize the parameters \mathbf{w} and w_0 using data
- **Test:**
 - calculate $g(x) = \mathbf{w}^T \mathbf{x} + w_0$ and choose C_1 if $g(x) > 1$ or choose C_2 if $g(x) < -1$.

Loss Function

$$\begin{aligned} \min \quad & \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{s.t.} \quad & y^{(\ell)}(\mathbf{w}^T \mathbf{x}^{(\ell)} + w_0) \geq 1, \ell = 1, \dots, N \end{aligned}$$



- For a given input x , the model outputs a score $g(x) = \mathbf{w}^T x + w_0$. Let $y \in \{-1, +1\}$ be the label of the real class ($y = +1: x \in C_1, y = -1: x \in C_2$):
 - if $y(\mathbf{w}^T x + w_0) < 1$: we aim to maximize $y(\mathbf{w}^T x + w_0)$ until reaching 1, cost is $1 - y(\mathbf{w}^T x + w_0)$
 - if $y(\mathbf{w}^T x + w_0) > 1$: outlier points, no need to optimize, cost is 0
- Can write this succinctly as a **Hinge** loss:

$$\ell(\mathbf{w}, w_0 | x, y) = \max(0, 1 - y(\mathbf{w}^T x + w_0))$$

- Given: $D = \{(x^{(1)}, r^{(1)}), \dots, (x^{(N)}, r^{(N)})\}$, the loss over the dataset is defined as:

$$L(\mathbf{w}, w_0 | D) = \frac{1}{N} \sum_{\ell=1}^N \max(0, 1 - y^{(\ell)}(\mathbf{w}^T \mathbf{x}^{(\ell)} + w_0)) + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

Optimization – Gradient Descend



$$\mathbf{w}_{\text{new}} = \mathbf{w}_{\text{old}} - \eta \frac{\partial L}{\partial \mathbf{w}}$$

What is $\frac{\partial L}{\partial \mathbf{w}}$?

$$L(\mathbf{w}, w_0 | D) = \begin{cases} \frac{\lambda}{2} \|\mathbf{w}\|^2 & \text{if } y^{(\ell)}(\mathbf{w}^T \mathbf{x}^{(\ell)} + w_0) \geq 1 \\ \frac{1}{N} \sum_{\ell=1}^N 1 - y^{(\ell)}(\mathbf{w}^T \mathbf{x}^{(\ell)} + w_0) + \frac{\lambda}{2} \|\mathbf{w}\|^2 & \text{otherwise} \end{cases}$$

For each w_j ($j=0, \dots, d$): $\frac{\partial L}{\partial w_j} = \begin{cases} \lambda w_j & \text{if } y^{(\ell)}(\mathbf{w}^T \mathbf{x}^{(\ell)} + w_0) \geq 1 \\ \lambda w_j - y^{(\ell)} x_j^{(\ell)} & \text{o.w.} \end{cases}$

$$\frac{\partial L}{\partial w_0} = \begin{cases} 0 & \text{if } y^{(\ell)}(\mathbf{w}^T \mathbf{x}^{(\ell)} + w_0) \geq 1 \\ -y^{(\ell)} & \text{o.w.} \end{cases}$$



The Algorithm

Gradient Descend for SVM

```
Input:  $D = \{(x^{(\ell)}, y^{(\ell)})\} (\ell=1:N)$   
for  $j = 0, \dots, d$   
     $w_j \leftarrow \text{rand}(-0.01, 0.01)$   
repeat  
    for  $\ell = 1, \dots, N$   
         $a \leftarrow 0$   
        for  $j = 0, \dots, d$   
             $a \leftarrow a + w_j x_j^{(\ell)}$   
        for  $j = 0, \dots, d$   
             $\Delta w_j \leftarrow \lambda w_j$   
         $\Delta w_0 \leftarrow 0$   
        if  $y^{(\ell)} a < 1$ :  
            for  $j = 0, \dots, d$   
                 $\Delta w_j \leftarrow \Delta w_j - y^{(\ell)} x_j^{(\ell)}$   
             $\Delta w_0 \leftarrow \Delta w_0 - y^{(\ell)}$   
        for  $j = 0, \dots, d$   
             $w_j \leftarrow w_j - \eta \Delta w_j$   
until convergence
```



Lagrangian

- The primal optimization problem:

$$\begin{array}{ll} \text{minimize} & \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{subject to} & y^{(\ell)}(\mathbf{w}\mathbf{x}^{(\ell)} + w_0) \geq 1, \forall \ell \end{array}$$

- Lagrangian: $\mathcal{L} = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{\ell} \alpha_{\ell} (y^{(\ell)}(\mathbf{w}\mathbf{x}^{(\ell)} + w_0) - 1)$

$$\nabla_{\mathbf{w}, w_0} \mathcal{L} = 0 \Rightarrow \begin{cases} \frac{\partial \mathcal{L}}{\partial \mathbf{w}} = 0 \Rightarrow \mathbf{w} = \sum_{\ell=1}^N \alpha_{\ell} y^{(\ell)} \mathbf{x}^{(\ell)} \\ \frac{\partial \mathcal{L}}{\partial w_0} = 0 \Rightarrow \sum_{\ell=1}^N \alpha_{\ell} y^{(\ell)} = 0 \end{cases}$$

- Substitute back in the primal to get the dual

$$\begin{array}{ll} \text{maximize} & \mathcal{L}(\alpha) = \sum_{\ell} \alpha_{\ell} - \frac{1}{2} \sum_{\ell=1}^N \sum_{\ell'=1}^N \alpha_{\ell} \alpha_{\ell'} y^{(\ell)} y^{(\ell')} (\mathbf{x}^{(\ell)})^T \mathbf{x}^{(\ell')} \\ \text{subject to} & \alpha_{\ell} \geq 0, \sum_{\ell=1}^N \alpha_{\ell} y^{(\ell)} = 0 \end{array}$$



The Dual Problem

Dual optimization problem:

$$\begin{aligned} & \text{maximize} \quad \sum_{\ell=1}^N \alpha_{\ell} - \frac{1}{2} \sum_{\ell=1}^N \sum_{\ell'=1}^N \alpha_{\ell} \alpha_{\ell'} y^{(\ell)} y^{(\ell')} (x^{(\ell)})^T x^{(\ell')} \\ & \text{subject to} \quad \sum_{\ell=1}^N \alpha_{\ell} y^{(\ell)} = 0 \\ & \quad \quad \quad \alpha_{\ell} \geq 0, \ell=1 \dots N \end{aligned}$$

- This is also a QP problem, but its complexity depends on the sample size N (rather than the input dimensionality d)



Primal and Dual

Primal

$$\begin{aligned} \min \quad & \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{s.t.} \quad & y^{(\ell)}(\mathbf{w}^T \mathbf{x}^{(\ell)} + w_0) \geq 1, \forall \ell \end{aligned}$$

The complexity depends on the dimensionality d of inputs

Dual

$$\begin{aligned} \max \quad & \sum_{\ell} \alpha_{\ell} - \frac{1}{2} \sum_{\ell} \sum_{\ell'} \alpha_{\ell} \alpha_{\ell'} y^{(\ell)} y^{(\ell')} (\mathbf{x}^{(\ell)})^T \mathbf{x}^{(\ell')} \\ \text{s.t.} \quad & \sum_{\ell} \alpha_{\ell} y^{(\ell)} = 0 \\ & \alpha_{\ell} \geq 0, \ell = 1 \dots N \end{aligned}$$

The complexity depends on the sample size N

- It turns out to be more convenient to solve the dual problem than solving the primal problem ($N < d$)
- We can firstly solve **Dual** to obtain $\{\alpha_{\ell}\}$, and then obtain the \mathbf{W} in **Primal**



Training (Dual)

- Given: $D = \{(x^{(1)}, y^{(1)}), \dots, (x^{(N)}, y^{(N)})\}$
- minimize the **loss function** $L(\alpha)$ using any general purpose optimization toolkits (e.g., Matlab)
- But SVM is usually optimized using the **SMO** (sequential minimal optimization)

- **Goal:**

$$\min_{\alpha} L(\alpha)$$

- **Iteration:**

Update two variables each time until convergence {

1. Select an α_1 that violates the KKT condition
 2. Pick a second multiplier α_2 and optimize $L(\alpha)$ w.r.t. α_1 and α_2
- }

https://en.wikipedia.org/wiki/Sequential_minimal_optimization



Support Vectors

Suppose the optimal $\{\alpha_\ell\}$ have been obtained

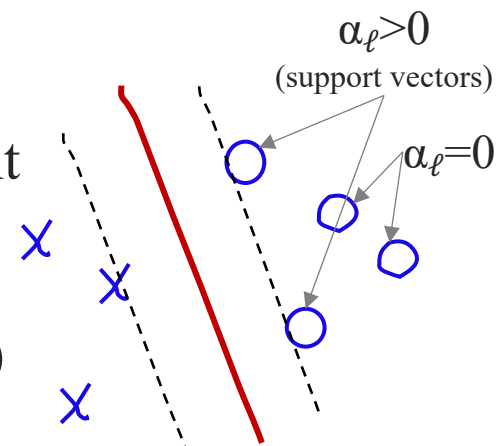
- Patterns for which $y^{(\ell)}(\mathbf{w}x^{(\ell)} + w_0) > 1$

$\alpha_\ell = 0$ (inactive constraints) $\Rightarrow x^{(\ell)}$ irrelevant

recall complementary slackness: $\lambda_i^* g_i(x^*) = 0$

- Patterns that have $\alpha_\ell > 0$ (active constraints)

$y^{(\ell)}(\mathbf{w}x^{(\ell)} + w_0) = 1 \Rightarrow x^{(\ell)}$ lies on margin



- Most of the dual variables vanish with $\alpha_\ell = 0$. They are points lying beyond the margin with **no effect** on the hyperplane.
- Solution is only determined by the examples on the margin (**support vectors**), i.e., $x^{(\ell)}$ with $\alpha_\ell > 0$, hence the name **support vector machine (SVM)**.



Computation of Primal Variables

- Having obtained the optimal α , we can obtain \mathbf{w} :

$$\mathbf{w} = \sum_{\ell=1}^N \alpha_{\ell} y^{(\ell)} \mathbf{x}^{(\ell)} = \sum_{\mathbf{x}^{(\ell)} \in \mathcal{SV}} \alpha_{\ell} y^{(\ell)} \mathbf{x}^{(\ell)}$$

where \mathcal{SV} denotes the set of support vectors.

- The support vectors must lie on the margin, so they should satisfy $y^{(\ell)}(\mathbf{w}^T \mathbf{x}^{(\ell)} + w_0) = 1$ or $w_0 = y^{(\ell)} - \mathbf{w}^T \mathbf{x}^{(\ell)}$
- For numerical stability, all support vectors are used to compute w_0 :

$$w_0 = \frac{1}{|\mathcal{SV}|} \sum_{\mathbf{x}^{(\ell)} \in \mathcal{SV}} (y^{(\ell)} - \mathbf{w}^T \mathbf{x}^{(\ell)})$$



Prediction

- **Discriminant Function:**

$$\begin{aligned} g(\mathbf{x}) &= \mathbf{w}^T \mathbf{x} + w_0 \\ &= \left(\sum_{\mathbf{x}^{(\ell)} \in \mathcal{SV}} \alpha_{\ell} y^{(\ell)} \mathbf{x}^{(\ell)} \right)^T \mathbf{x} + \frac{1}{|\mathcal{SV}|} \sum_{\mathbf{x}^{(\ell)} \in \mathcal{SV}} (y^{(\ell)} - \mathbf{w}^T \mathbf{x}^{(\ell)}) \end{aligned}$$

- **Decision Rule:**

$$\text{Choose } \begin{cases} C_1 & \text{if } g(x) > +1 \\ C_2 & \text{if } g(x) < -1 \end{cases}$$

Programming Time



Tutorial: SVM from scratch with explanation (Python)

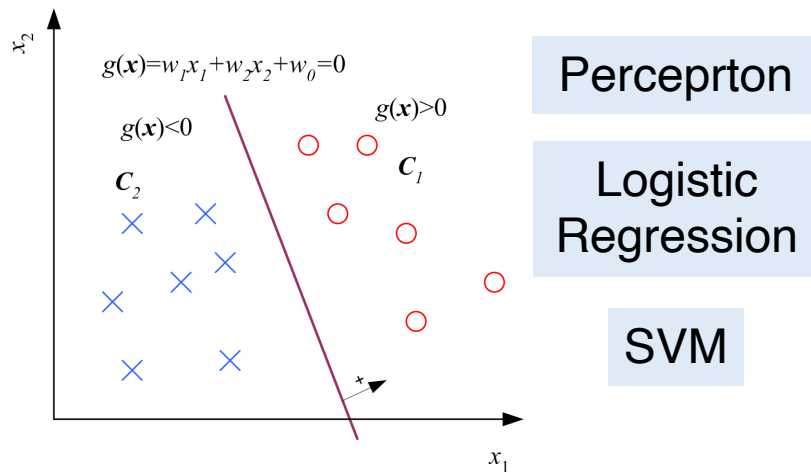
<https://towardsdatascience.com/implementing-svm-from-scratch-784e4ad0bc6a>

<https://www.kaggle.com/code/misbahbilgili/svm-from-scratch-with-explanation/notebook>



What's Next

WHAT'S
NEXT?



The curse of nonlinearity

