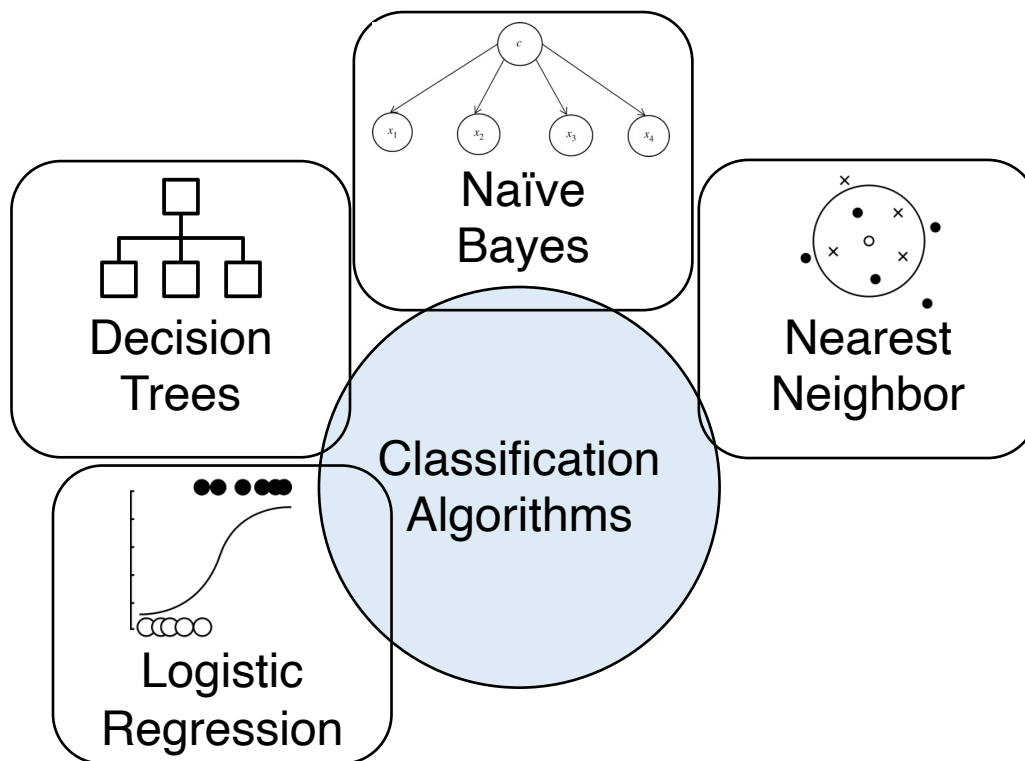# Machine Learning

## Lecture 7: Support Vector Machine

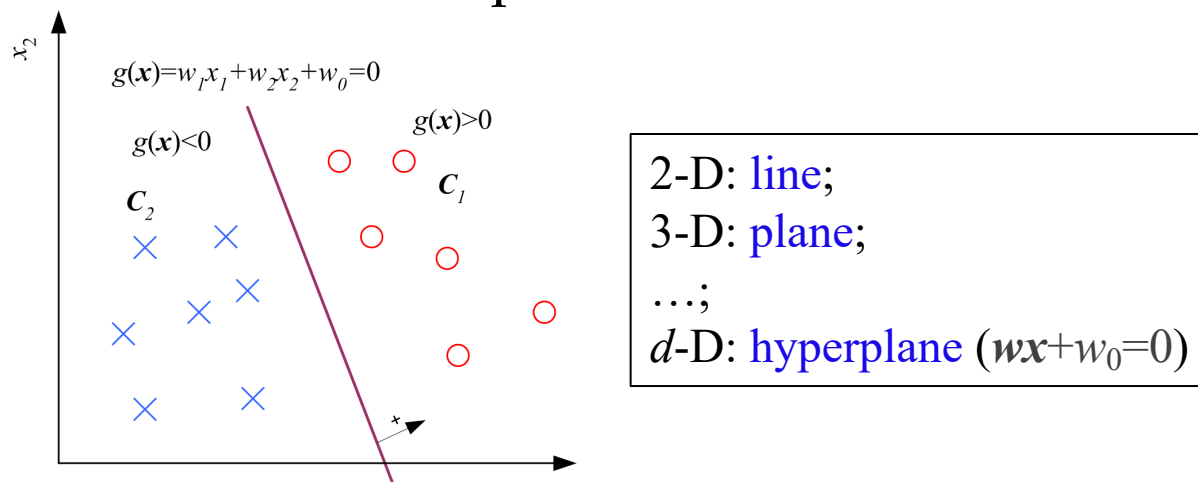### Fall 2022

**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 <span style="color:blue">linearly separable</span>: there exists a <span style="color:red">linear surface</span> to separate the two classes
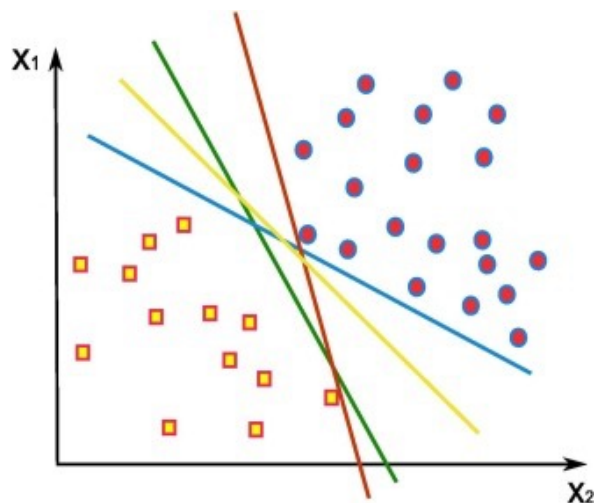


2-D: line;
3-D: plane;
…;
$d$-D: hyperplane ($\mathbf{w}x + w_0 = 0$)

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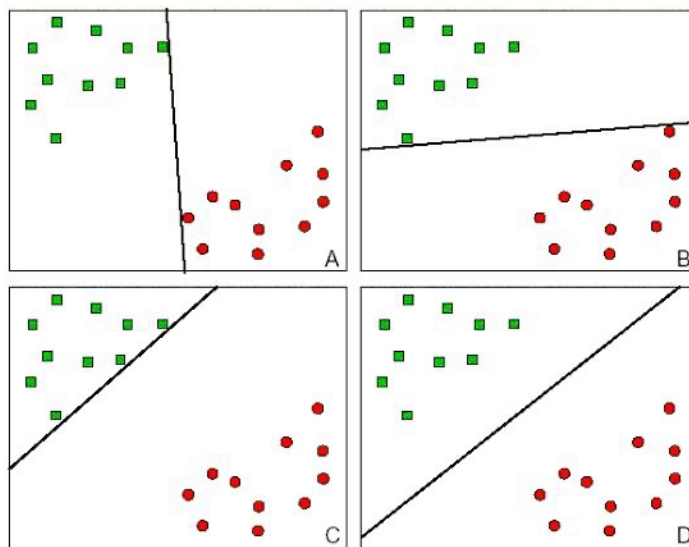
Find $\mathbf{w}x + w_0 = 0$ that perfectly separates the two classes

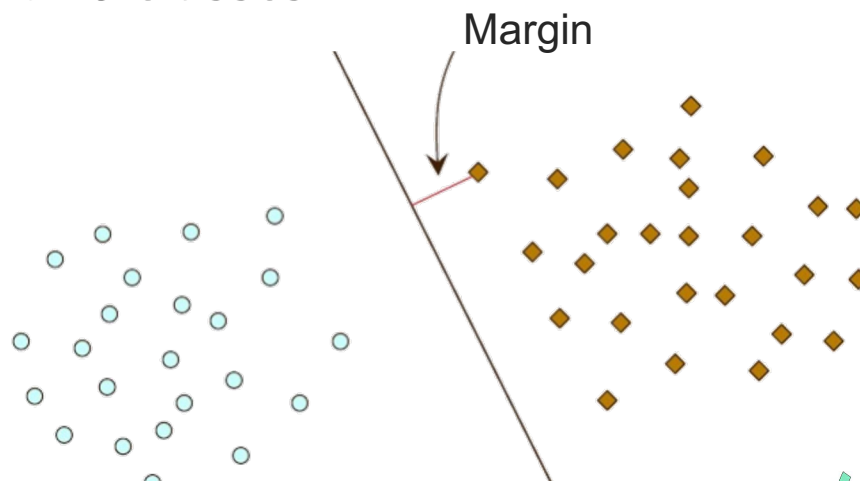# The Problem

- There can be multiple separating hyperplanes.



**Which one is the best?**

# The Idea:

- Find a decision boundary that <mark>**maximizes the margin**</mark> between two classes.
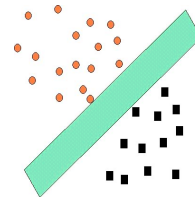
Margin



**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.



$$\mathbf{w}^T x + w_0 = 0$$

$$\mathbf{w}^T x + w_0 = -\delta = -1$$

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

$$C_1$$

$$C_2$$

$$1/\delta\ \mathbf{w}^T x + 1/\delta\ w_0 = 0$$ defines the same hyperplane as $$\mathbf{w}^T x + w_0 = 0$$

$$\mathbf{w}^T x + w_0 \begin{cases} \geq +1 & x \in C_1 \\ \leq -1 & x \in C_2 \end{cases}$$

# Maximizing the Margin

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

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

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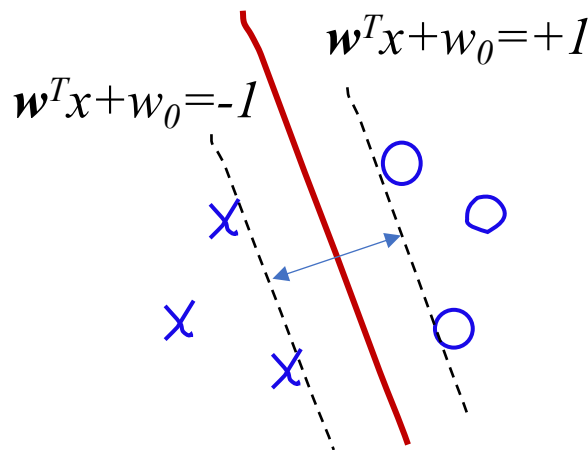$$w^\mathrm{T} x^{(2)} + w_0 = -1$$

which imply

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

Hence, the margin can be given by

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

$w^T x + w_0 = +1$

$w^T x + w_0 = -1$

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

# Inequality Constraints

- Given: $D = \{(x^{(\ell)}, y^{(\ell)}), \ldots, (x^{(N)}, y^{(N)})\}$, we want $\boldsymbol{w}$ and $w_0$ to satisfy

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

- Or, equivalently,

$$y^{(\ell)}\left(\boldsymbol{w}^{\mathrm{T}}x^{(\ell)}+w_0\right) \geq 1$$

$\boldsymbol{w}^Tx+w_0=+1$

$\boldsymbol{w}^Tx+w_0=-1$

# SVM = Solving an Optimization Problem

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

$$\text{minimize} \quad \frac{1}{2}\|\boldsymbol{w}\|^2$$
$$\text{subject to} \quad y^{(\ell)}(\boldsymbol{w}^\mathrm{T}x^{(\ell)}+w_0)\geq 1, \; \ell=1,\ldots, N$$

- This is a quadratic programming (QP) problem, which is one type of convex optimization problem.  （                 ）

- The complexity depends on the dimensionality $d$ of inputs
  （                    w        d
  x        N    ）

# Model Architecture

- The same architecture as Perceptron.



input — $w_0$ — $x_0=1$ ... $x_d$ — $w_d$ → Classifier $g(x)=\boldsymbol{w}^T x + w_0$ → predict $g(x)$ — loss — target $y$

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

# Loss Function

- For a given input $x$, the model outputs a score $g(x) = \boldsymbol{w}^\text{T}x+w_0$. 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 $y(\boldsymbol{w}^\text{T}x+w_0) < 1$: we aim to maximize $y(\boldsymbol{w}^\text{T}x+w_0)$ until reaching 1, cost is
    $$1-y(\boldsymbol{w}^\text{T}x+w_0)$$
  - if $y(\boldsymbol{w}^\text{T}x+w_0) > 1$: outlier points, no need to optimize, cost is **0**

- Can writhe this succinctly as a **Hinge** loss:

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

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

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

# Optimization – Gradient Descend

$$\boxed{w_{\text{new}} = w_{\text{old}} + \eta \frac{\partial L}{\partial w}}$$

What is $\dfrac{\partial L}{\partial w}$ ?

$$L\ (w,\ w_0\ |\ D) = \begin{cases} \dfrac{\lambda}{2} \|w\|^2 & \textbf{if}\quad y^{(\ell)}(w^{\mathrm{T}}x^{(\ell)}+w_0) \geq 1 \\[3mm] \dfrac{1}{N}\sum_{\ell=1}^{N}\ 1-y^{(\ell)}(w^{\mathrm{T}}x^{(\ell)}+w_0)\ +\dfrac{\lambda}{2}\|w\|^2 & \text{otherwise} \end{cases}$$

For each $w_j$ ($j = 0,\dots,d$):

$$\frac{\partial L}{\partial w_j} = \begin{cases} \lambda w_j & \text{if } y^{(\ell)}(w^{\mathrm{T}}x^{(\ell)}+w_0) \geq 1 \\[2mm] \lambda w_j - y^{(\ell)}x^{(\ell)} & \text{o.w.} \end{cases}$$

$$\frac{\partial L}{\partial w_0} = \begin{cases} 0 & \text{if } y^{(\ell)}(w^{\mathrm{T}}x^{(\ell)}+w_0) \geq 1 \\[2mm] -y^{(\ell)} & \text{o.w.} \end{cases}$$

w_0
y=wx+w_0   w_0

# The Algorithm

## Gradient Descend for SVM

**Input:** $D = \{(\mathrm{x}^{(l)}, \mathrm{y}^{(l)})\}$ $(l = 1{:}N)$

**for** $j = 0,\ldots,d$

    $w_j \leftarrow rand\ (-0.01,\ 0.01)$

**repeat**

    **for** $j = 0,\ldots,d$

        $\Delta w_j \leftarrow 0$

    **for** $l = 1,\ldots,N$

        $a \leftarrow 0$

        **for** $j = 0, \ldots, d$

            $a \leftarrow a + w_j x_j^{(l)}$

        **for** $j = 0, \ldots, d$

            $\Delta w_j \leftarrow \Delta w_j + \lambda w_j$

            $\Delta w_0 \leftarrow 0$

            **if** $y^{(l)} a < 1$: $\Delta w_j \leftarrow \Delta w_j - y^{(l)} x^{(l)}$

    $\Delta w_j = \Delta w_j / \mathrm{N}$

    **for** $j = 0,\ldots,d$

        $w_j \leftarrow w_j + \eta \Delta w_j$

**until** convergence

# Lagrangian

- The primal optimization problem:

$$\text{minimize } \tfrac{1}{2}\|w\|^2$$
$$\text{subject to } y^{(\ell)}(\boldsymbol{w}x^{(\ell)}+w_0) \geq 1, \forall \ell$$

- **Lagrangian**: $\mathcal{L} = \tfrac{1}{2}\|\boldsymbol{w}\|^2 - \Sigma_\ell \alpha_\ell(y^{(\ell)}(\boldsymbol{w}x^{(\ell)}+w_0) -1)$

$$\nabla_{w,w_0}\mathcal{L}=0 \implies \begin{cases} \dfrac{\partial \mathcal{L}}{\partial \boldsymbol{w}} = 0 \implies \boldsymbol{w} = \sum_{\ell=1}^{N} \alpha_\ell y^{(\ell)} x^{(\ell)} \\[2mm] \dfrac{\partial \mathcal{L}}{\partial w_0} = 0 \implies \sum_{\ell=1}^{N} \alpha_\ell y^{(\ell)} = 0 \end{cases}$$

- Substitute back in the primal to get the <span style="color:blue">dual</span>

$$\textbf{maximize } \mathcal{L}(\alpha) = \Sigma_\ell \alpha_\ell - \tfrac{1}{2}\sum_{l=1}^{N}\sum_{l'=1}^{N} \alpha_l \alpha_{l'} y^{(l)} y^{(l')} (x^{(l)})^{\mathrm{T}} x^{(l')}$$
$$\textbf{subject to } \alpha_\ell \geq 0, \sum_{\ell=1}^{N} \alpha_\ell y^{(\ell)} = 0$$

# The Dual Problem

**Dual optimization problem:**

$$\text{maximize} \quad \sum_{\ell=1}^{N} \alpha_\ell - \frac{1}{2} \sum_{\ell=1}^{N} \sum_{\ell'=1}^{N} \alpha_\ell \alpha_{\ell'} y^{(\ell)} y^{(\ell')} (x^{(\ell)})^{\mathrm{T}} x^{(\ell')}$$

$$\text{subject to} \quad \sum_{\ell=1}^{N} \alpha_\ell y^{(\ell)} = 0$$

$$\alpha_\ell \geq 0, \; \ell = 1 \ldots N$$

- 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 | Dual |
|---|---|

**Primal**

$$\min \ \frac{1}{2}\|w\|^2$$
$$\text{s.t.} \quad y^{(l)}(w^{\mathrm{T}}x^{(l)}+w_0)\geq 1, \ \forall l$$

**Dual**

$$\max \ \Sigma_l \alpha_l - \tfrac{1}{2}\Sigma_l \Sigma_{l'} \alpha_l \alpha_{l'} y^{(l)} y^{(l')} (x^{(l)})^{\mathrm{T}} x^{(l')}$$
$$\text{s.t.} \quad \Sigma_l \alpha_l y^{(l)}=0$$
$$\alpha_l \geq 0, \ l=1\ldots N$$

The complexity depends on the dimensionality $d$ of inputs

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 $W$ in **Primal**

# Training (Dual)

- Given: $D = \{(x^{(\ell)}, y^{(\ell)}), \ldots, (x^{(N)}, y^{(N)})\}$

- 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:**

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

- **Iteration:**
  Update two variables each time until convergence {
  1. Select an $\alpha_1$ that violates the KKT condition
  2. Pick a second multiplier $\alpha_2$ and optimize $L(\alpha)$ w.r.t. $\alpha_1$ and $\alpha_2$
  }

---

# Support Vectors

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

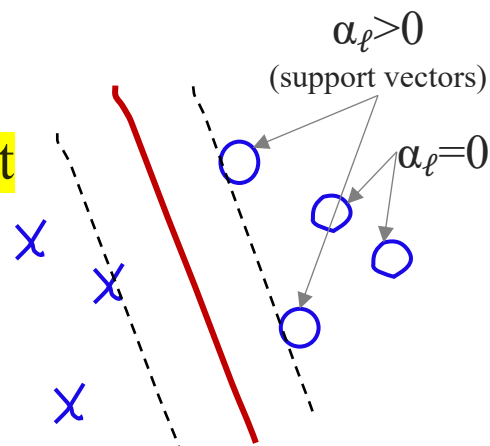- Patterns for which $y^{(\ell)}(\boldsymbol{w}x^{(\ell)}+w_0) > 1$

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

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

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

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

$\alpha_\ell>0$
(support vectors)

$\alpha_\ell=0$

- 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).** support vectors

# Computation of Primal Variables

- Having obtained the optimal $\boldsymbol{\alpha}$, we can obtain $w$:

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

  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 $w_0$:

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

# Prediction

- **Discriminant Function:**

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

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

- **Decision Rule:**

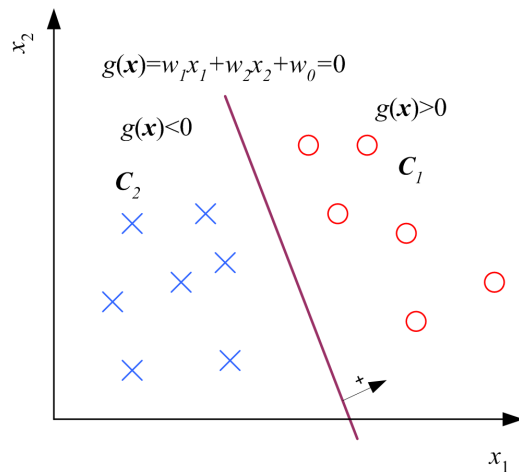$$\text{Choose} \begin{cases} C_1 & \text{if } g(\mathbf{x}) > 0 \\ C_2 & \text{otherwise} \end{cases}$$

# A Python Tutorial

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

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

# What's Next



$g(x)=w_1x_1+w_2x_2+w_0=0$

$g(x)<0$   $g(x)>0$

$C_2$   $C_1$

Perceprton

Logistic Regression

SVM

I can approximate **any** non-linear functions !

Neural Networks

The curse of nonlinearity