

SDSU CS 549 Spring 2024 Supervised Machine Learning Lecture 5: Support Vector Machine (SVM)

IRFAN KHAN

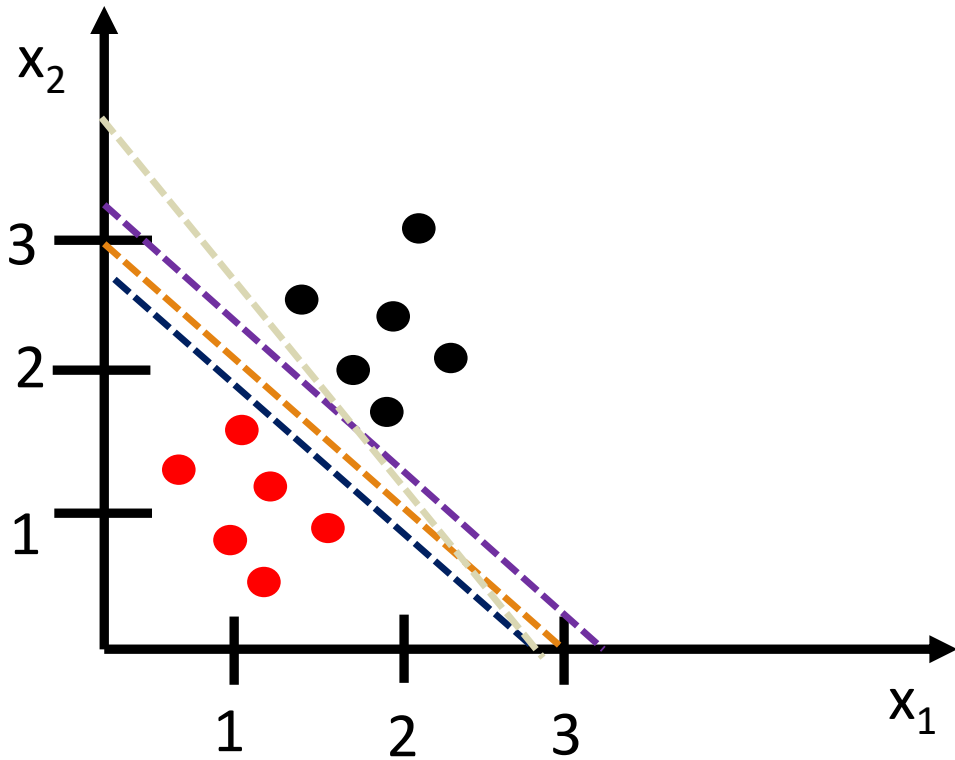
References

SDSU CS549 Lecture Notes by Prof Yang Xu, Spring 2023. Some updated slides used here

Introduction

- Linear regression and logistic regression solve unconstrained optimization problems.
- There are classes of problems that could benefit from constrained optimization, e.g.
 - A machine learning model to recognize handwritten digits (0-9) based on a set of input features extracted from images of the digits. The ***goal is to achieve high accuracy*** in classifying the digits into their corresponding classes
 - Spam vs non-spam email classifier. The ***goal is to achieve high accuracy*** in classifying emails to help prevent spam emails from reaching users' inboxes.
- Support Vector Machine (SVM) can be used here:
 - A supervised machine learning model
 - Developed by Vladimir N. Vapnik and his colleagues in the 1990s.

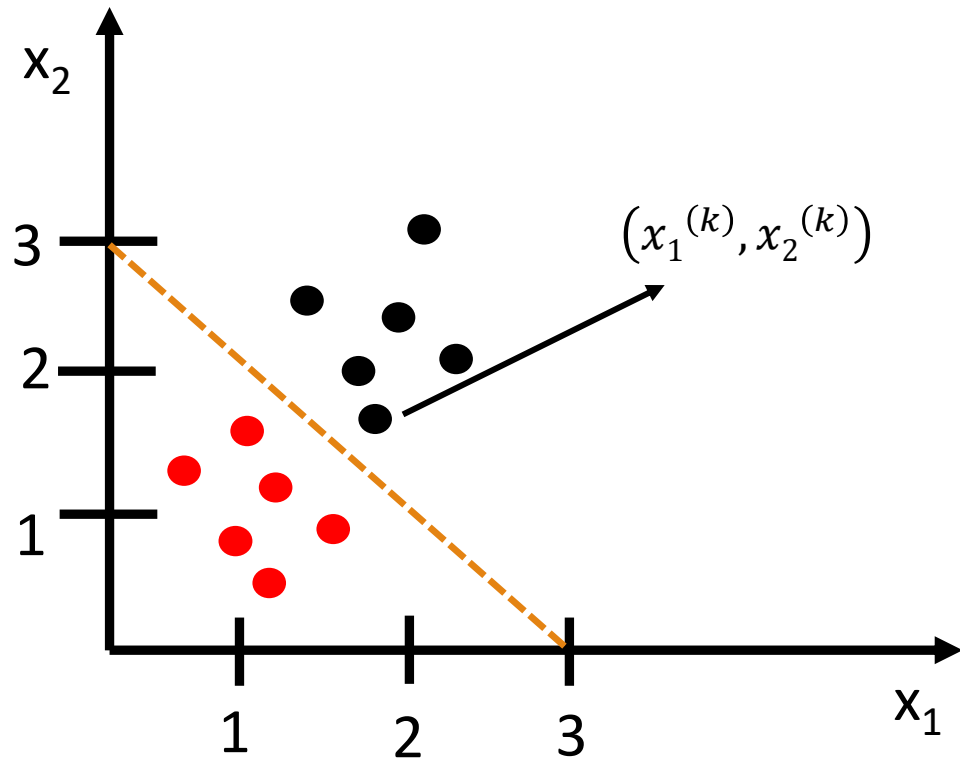
Recap - A Two-Feature example for Logistic Regression



- Intermediate variable $z = w_1x_1 + w_2x_2 + b$
- Pass it through a Sigmoid Function (to make decisions on 0/1 output)
 - $g(x) = \frac{1}{1+e^{-z}} = \frac{1}{1+e^{-(w_1x_1+w_2x_2+b)}}$
- What is the decision boundary?
 - $z = w_1x_1 + w_2x_2 + b = 0$
- For the example shown on the left, the following straight line seems to be a good decision boundary
 - $x_1 + x_2 = 3$
 - $w_1 = 1, w_2 = 1, b = -3$
- How about these new decision boundaries?
- Which is better?
- Maybe one with a wider margin since that will result in less misclassification.
 - This is the constraint!

Basic SVM

Constructing the Optimization Problem



- Decision Boundary is the straight line: $b + w_1x_1 + w_2x_2 = 0$
- For all data points $(x_1^{(i)}, x_2^{(i)})$, $i = 1, 2, \dots, m$: $|b + w_1x_1^{(i)} + w_2x_2^{(i)}| > 0$
- We define the classifier variable, $y^{(i)}$, for each data point $(x_1^{(i)}, x_2^{(i)})$ equal to 1 if the datapoint is over the decision line, or equal to -1 if under.
 - Using $(-1, 1)$ instead of $(0, 1)$ helps simplify the problem formulation
- Therefore, for all points, $y^{(i)}(b + w_1x_1 + w_2x_2) > 0$
- Consider the nearest point ("support vector") from the sample data points to the decision boundary, $(x_1^{(k)}, x_2^{(k)})$
- Since there can be infinite representations of the decision boundary line, let us place the normalization constraint on the decision boundary line:
 $|b + w_1x_1^{(k)} + w_2x_2^{(k)}| = 1$
- Our optimization problem finds w_1, w_2, b such that the distance from $(x_1^{(k)}, x_2^{(k)})$ to the decision boundary line is maximized.

Distance of Support Vector to the Decision Boundary

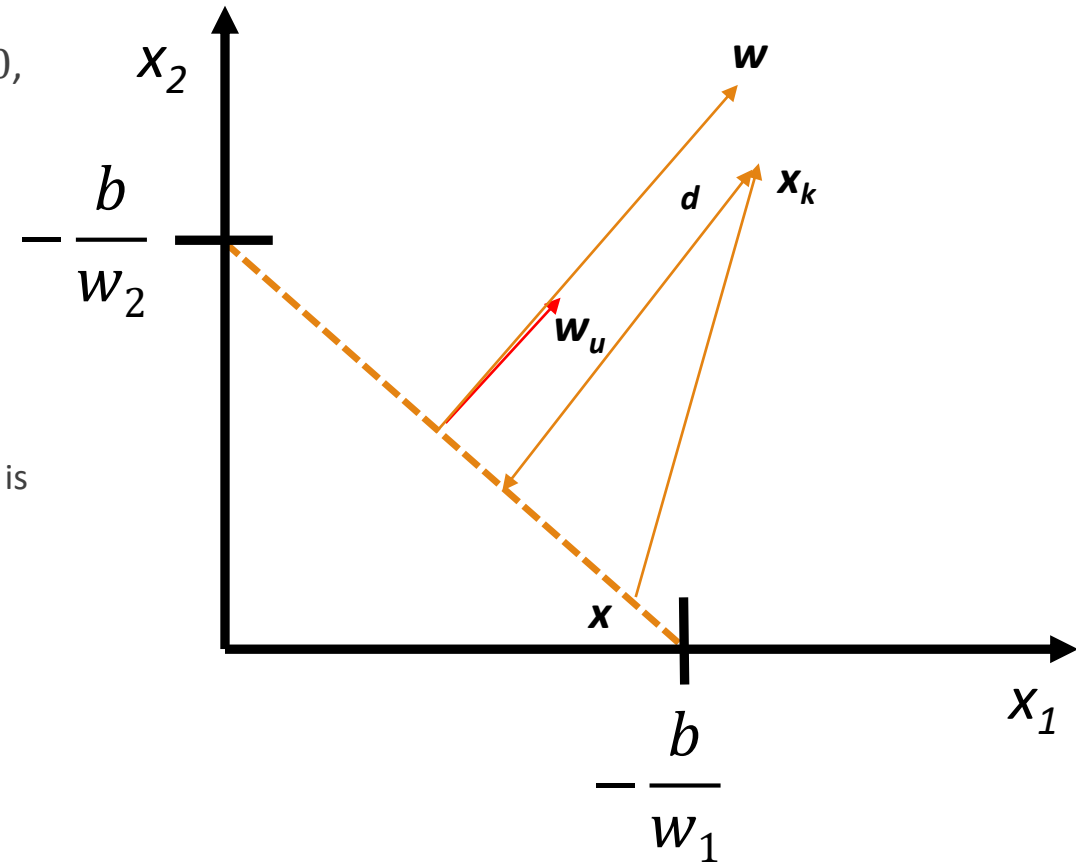
➤ For a straight-line decision boundary represented as $b + [w_1 \ w_2] \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = 0$,

in a 2D space, the vector $\mathbf{w} = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}$ is the Normal Vector that is perpendicular to the straight line.

➤ Proof: Take any two points on the straight line, say $\mathbf{x}_a = \begin{bmatrix} x_{a1} \\ x_{a2} \end{bmatrix}$, and $\mathbf{x}_b = \begin{bmatrix} x_{b1} \\ x_{b2} \end{bmatrix}$, then $\mathbf{x}_a - \mathbf{x}_b$ represents a vector on the straight line. A dot product of this vector with \mathbf{w} is always 0.

➤ The distance d from a point \mathbf{x}_k to the straight line decision boundary is the projection of the vector $(\mathbf{x}_k - \mathbf{x})$, for any \mathbf{x} on the decision boundary onto \mathbf{w}

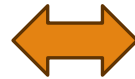
➤
$$d = \frac{1}{\|\mathbf{w}\|_2} |\mathbf{w}^T(\mathbf{x}_k - \mathbf{x})| = \frac{1}{\|\mathbf{w}\|_2}$$



SVM – Optimization Problem

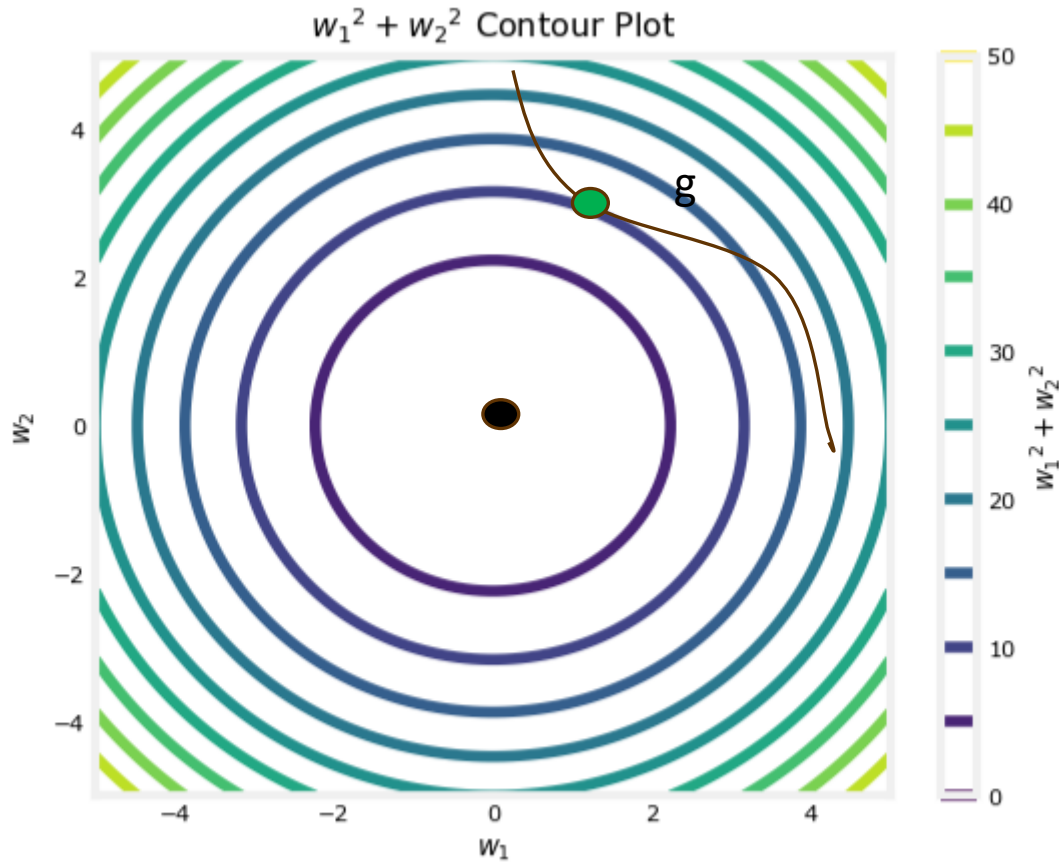
So, our optimization problem can be expressed as:

- Maximize $\frac{1}{\|w\|_2}$
- Subject to the constraint:
 - $y^{(k)}(b + w_1x_1^{(k)} + w_2x_2^{(k)}) = 1$ for support vector $(x_1^{(k)}, x_2^{(k)})$ and
 - $y^{(i)}(b + w_1x_1^{(i)} + w_2x_2^{(i)}) > 0$ for other data points $(x_1^{(i)}, x_2^{(i)})$



- Minimize $\|w\|_2 = \text{Minimize } \sqrt{w_1^2 + w_2^2} =$
 - Minimize $(w_1^2 + w_2^2) = \text{Minimize } \frac{1}{2}(w_1^2 + w_2^2)$
 - $= \text{Minimize } \frac{1}{2} w^T w$
- Subject to the constraint:
 - $y^{(k)}(b + w_1x_1^{(k)} + w_2x_2^{(k)}) = 1$ for support vector $(x_1^{(k)}, x_2^{(k)})$ and
 - $y^{(i)}(b + w_1x_1^{(i)} + w_2x_2^{(i)}) > 0$ for other data points $(x_1^{(i)}, x_2^{(i)})$

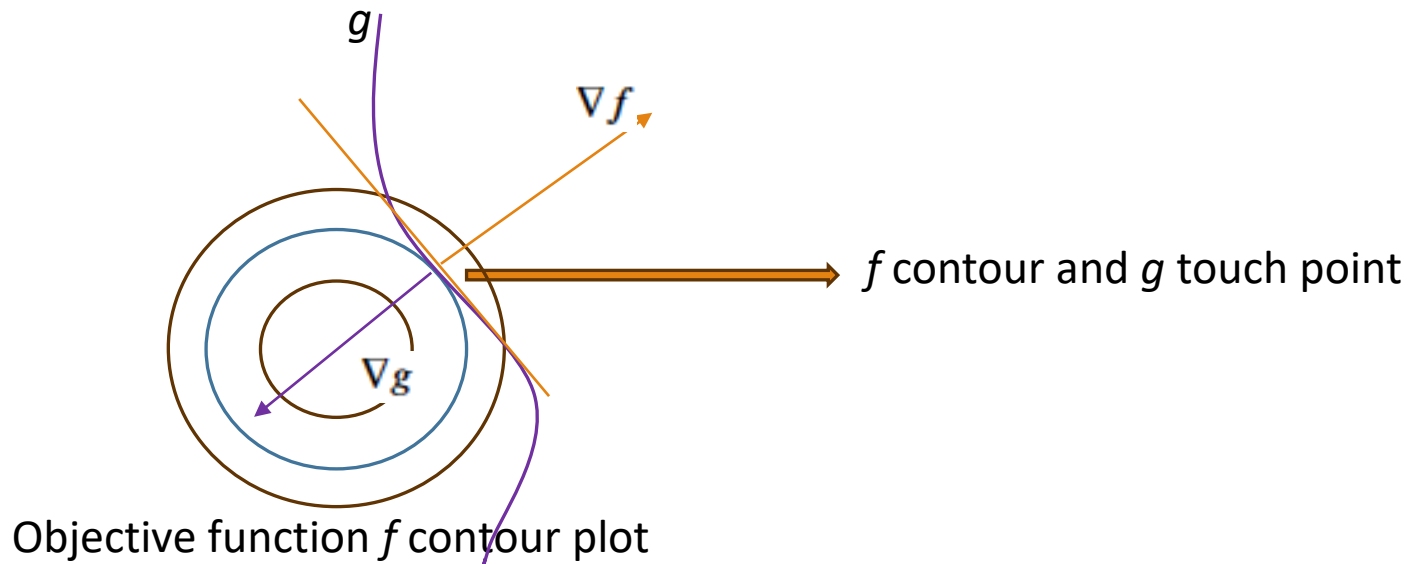
Solving the Constrained Optimization Problem



● Optimal Point

- In an unconstrained minimization problem, we would seek a point at the center of the innermost circle
- The constraint set is represented by a function g implementing:
 - $y^{(k)}(b + w_1x_1^{(k)} + w_2x_2^{(k)}) = 1$ for support vector $(x_1^{(k)}, x_2^{(k)})$ and
 - $y^{(i)}(b + w_1x_1^{(i)} + w_2x_2^{(i)}) > 0$ for other data points $(x_1^{(i)}, x_2^{(i)})$
 - It is a summation of all the constraints
- The optimal point for the constrained minimization problem is as shown

Solving Constrained Optimization Problem (Contd.)



- Objective function to be minimized $f(\mathbf{w}) = \frac{1}{2} \mathbf{w}^T \mathbf{w}$
- Constraints: $g^{(i)}(\mathbf{w}, b) = y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + b) - 1 \geq 0$
- At the f contour and g touch point, ∇f and ∇g point in opposite direction
- The two are related as:
 - $\nabla f(\mathbf{w}) = \sum_{i=1}^m \lambda^{(i)} \nabla g^{(i)}(\mathbf{w}, b)$
 - $\lambda^{(i)}$ are called Lagrange multipliers

Generalizing to any dimension

- A decision boundary in an n -dimensional space is a hyper-plane of $(n-1)$ dimensions
 - Previous slides explored $n=2$ because it is easier to visualize
- All the concepts carry over to an n -dimensional problem
 - Intermediate variable $z = \sum_j w_j x_j^{(i)} + b, j = 1, 2, \dots, n$
 - Instead of classifying datapoints as 0/1, we classify them as -1/1
 - Hence $y^{(i)} \sum_j w_j x_j^{(i)} + b$ is always >0 for all non-support vectors
 - We apply the normalizing constraint for the support vectors $|\sum_j w_j x_j^{(k)} + b| = 1$
 - As before, the objective function to be minimized is $\frac{1}{2} \mathbf{w}^T \mathbf{w}$, subject to $g^{(i)}(\mathbf{w}, b) = y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + b) - 1 \geq 0$

The Lagrangian and its Optimization

- We define the Lagrangian as:
 - $L(\mathbf{w}, b, \lambda) = \frac{1}{2} \mathbf{w}^T \mathbf{w} - \sum_{i=1}^m \lambda^{(i)} g^{(i)}(\mathbf{w}, b), \lambda^{(i)} \geq 0$
 - Satisfies the gradient intuition we developed earlier
 - $= \frac{1}{2} \mathbf{w}^T \mathbf{w} - \sum_{i=1}^m \lambda^{(i)} (y^{(i)} (\mathbf{w}^T \mathbf{x}^{(i)} + b) - 1)$
- We minimize L wrt \mathbf{w} and b and maximize L wrt $\lambda^{(i)}$
- Minimize L wrt \mathbf{w} and b
 - $\frac{\partial L}{\partial \mathbf{w}} = \mathbf{w} - \sum_{i=1}^m \lambda^{(i)} y^{(i)} \mathbf{x}^{(i)} = 0$
 - $\frac{\partial L}{\partial b} = - \sum_{i=1}^m \lambda^{(i)} y^{(i)} = 0$
 - $\Rightarrow \lambda \cdot \mathbf{y} = 0$
- In $L(\mathbf{w}, b, \lambda)$, substituting $\mathbf{w} = \sum_{i=1}^m \lambda^{(i)} y^{(i)} \mathbf{x}^{(i)}$ in L and utilizing $\sum_{i=1}^m \lambda^{(i)} y^{(i)} = 0$
- $L(\mathbf{w}, b, \lambda) = \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \lambda^{(i)} \lambda^{(j)} y^{(i)} y^{(j)} \mathbf{x}^{(i)} \cdot \mathbf{x}^{(j)} - \sum_{i=1}^m \sum_{j=1}^m \lambda^{(i)} \lambda^{(j)} y^{(i)} y^{(j)} \mathbf{x}^{(i)} \cdot \mathbf{x}^{(j)} + \sum_{i=1}^m \lambda^{(i)}$
 - $-\frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \lambda^{(i)} \lambda^{(j)} y^{(i)} y^{(j)} \mathbf{x}^{(i)} \cdot \mathbf{x}^{(j)} + \sum_{i=1}^m \lambda^{(i)}$
- L thus simplifies to $L(\lambda)$ a quadratic in only one variable λ

The Lagrangian and its Optimization (Contd.)

➤ $L(\lambda) = -\frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \lambda^{(i)} \lambda^{(j)} y^{(i)} y^{(j)} \mathbf{x}^{(i)} \cdot \mathbf{x}^{(j)} + \sum_{i=1}^m \lambda^{(i)}$

➤ Maximizing $L \Leftrightarrow$ Minimizing L given by:

➤ $L(\lambda) = \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \lambda^{(i)} \lambda^{(j)} y^{(i)} y^{(j)} \mathbf{x}^{(i)} \cdot \mathbf{x}^{(j)} - \sum_{i=1}^m \lambda^{(i)}$

➤ $= \frac{1}{2} \lambda^T \begin{bmatrix} y^{(1)} y^{(1)} \mathbf{x}^{(1)} \cdot \mathbf{x}^{(1)} & y^{(1)} y^{(m)} \mathbf{x}^{(1)} \cdot \mathbf{x}^{(m)} & \dots \\ y^{(m)} y^{(1)} \mathbf{x}^{(m)} \cdot \mathbf{x}^{(1)} & y^{(m)} y^{(m)} \mathbf{x}^{(m)} \cdot \mathbf{x}^{(m)} & \dots \end{bmatrix} \lambda - (\mathbf{1}^T) \lambda$

➤ Denote the matrix as \mathbf{Q}

➤ We can write our optimization problem in a compact form as:

➤ Minimize $\frac{1}{2} \lambda^T \mathbf{Q} \lambda - (\mathbf{1}^T) \lambda$ wrt λ

➤ $\mathbf{y}^T \lambda = \mathbf{0}, \lambda^i \geq 0$

➤ Known as Quadratic programming optimization problem

Obtaining \mathbf{w} and b

- Once we obtain $\boldsymbol{\lambda} = \begin{bmatrix} \lambda^{(1)} \\ \dots \\ \lambda^{(m)} \end{bmatrix}$
- We can obtain $\mathbf{w} = \begin{bmatrix} w_1 \\ \dots \\ w_n \end{bmatrix} = \sum_{i=1}^m \lambda^{(i)} y^{(i)} \mathbf{x}^{(i)}$
- Majority of $\lambda^{(i)}$ are 0, the few that are >0 correspond to the support vectors
- So, $\mathbf{w} = \begin{bmatrix} w_1 \\ \dots \\ w_n \end{bmatrix} = \sum_{\mathbf{x}^{(i)} \in SV} \lambda^{(i)} y^{(i)} \mathbf{x}^{(i)}$
- b can be obtained using one of the support vectors $\mathbf{x}^{(k)}$ from:
 - $|b + w_1 x_1^{(k)} + w_2 x_2^{(k)}| = 1$

Recap of the Steps

➤ We want to find the decision boundary $\sum_j w_j x_j^{(i)} + b = 0, j = 1, 2, \dots, n$

➤ Such that the distance to the support vectors is large

➤ Instead of 0/1 classifier variable, we use -1/1 and we apply normalization constraint for the support vector ($\mathbf{x}^{(k)}$)

➤ $y^{(k)}(b + \sum_j w_j x_j^{(k)}) = 1$ for support vector ($\mathbf{x}^{(k)}$) and

➤ $y^{(i)}(b + \sum_j w_j x_j^{(i)}) \geq 0$ for other data points ($\mathbf{x}^{(i)}$)

➤ The objective function that maximizes distance of the decision boundary to the support vector(s) simplifies to minimizing $\frac{1}{2} \mathbf{w}^T \mathbf{w}$ subject to the constraints $y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + b) - 1 \geq 0$

➤ Then we construct a Lagrangian $L(\mathbf{w}, b, \boldsymbol{\lambda}) = \frac{1}{2} \mathbf{w}^T \mathbf{w} - \sum_{i=1}^m \lambda^{(i)} g^{(i)}(\mathbf{w}, b), \lambda^{(i)} \geq 0$ to minimize

➤ Which results in a Quadratic optimization problem which can be solved using off-the-shelf libraries:

➤ Minimize $\frac{1}{2} \boldsymbol{\lambda}^T \mathbf{Q} \boldsymbol{\lambda} - (\mathbf{1}^T) \boldsymbol{\lambda}$ wrt $\boldsymbol{\lambda}$

➤ $\mathbf{y}^T \boldsymbol{\lambda} = 0, \lambda^i \geq 0$

➤ Majority of $\lambda^{(i)}$ are 0, the few that are >0 correspond to the support vectors

➤ So, $\mathbf{w} = \begin{bmatrix} w_1 \\ \dots \\ w_n \end{bmatrix} = \sum_{\mathbf{x}^{(i)} \in SV} \lambda^{(i)} y^{(i)} \mathbf{x}^{(i)}$

➤ b can be obtained using one of the support vectors $\mathbf{x}^{(k)}$ from:

➤ $|b + \sum_j w_j x_j^{(k)}| = 1$

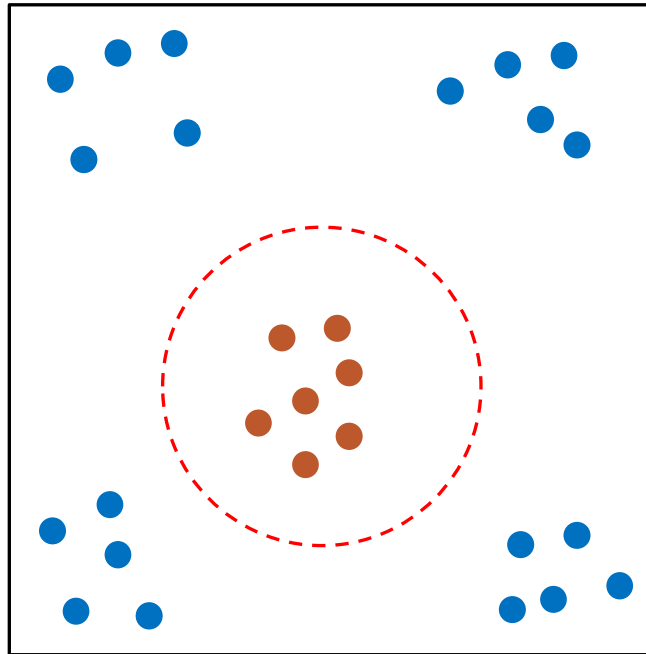
➤ Hence, we have the parameters for the decision boundary

Implementation of QP and SVM

- We use **cvxopt** for the assignment
- **cvxopt** is a Python package for solving optimization problems including quadratic programming problems.
- You can install it using the following command:
 - `conda install -c conda-forge cvxopt`
- **cvxopt** requires the matrices to be **cvxopt matrix**
- From **cvxopt** we use **matrix** to convert our usual numpy matrices to the cvxopt matrices and the quadratic programming optimization function **solvers**
- The scikit-learn library is a widely used machine learning library that includes a built-in implementation of Support Vector Machines (SVM) for classification
 - It can be used directly to solve SVM problems

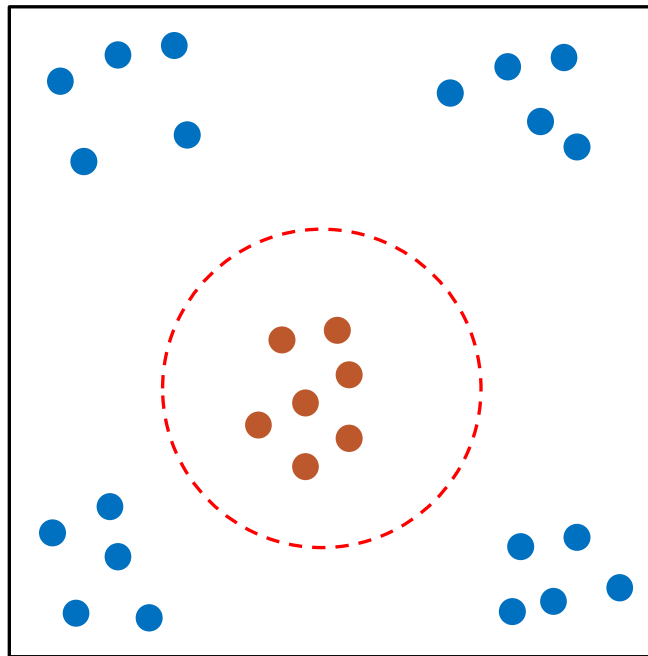
Kernel Methods

Data may not lend itself to a linear boundary for classification

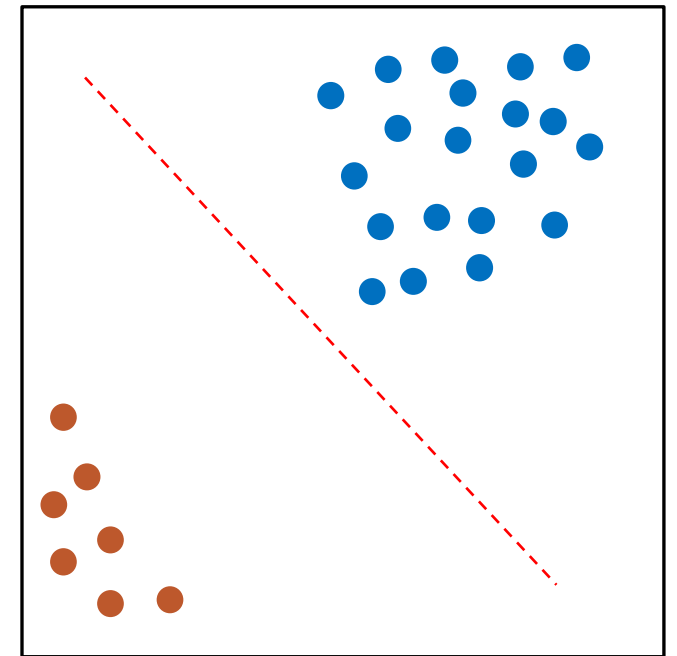
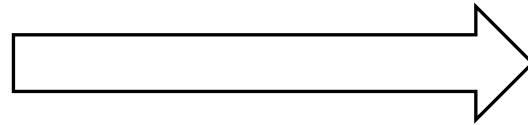


- Image Classification
- Text Classification:
- Bioinformatics and Genomics
- Speech Recognition
- Anomaly Detection
- Chemoinformatics
- Social Network Analysis
- Financial Data Analysis

Nonlinear Transformation to New Space



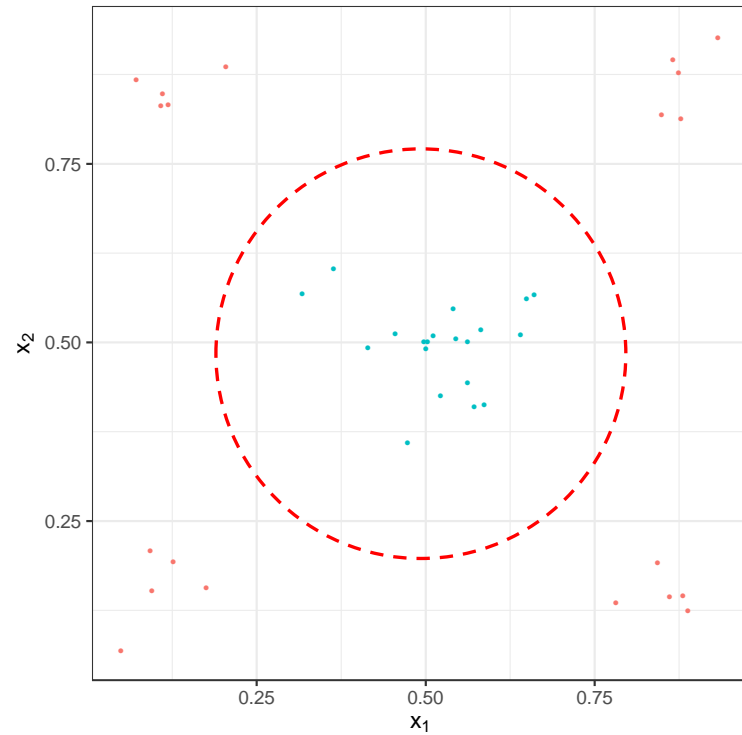
Apply nonlinear transform to data



Obtain a decision boundary in a new space

Example of a Non-Linear Transform

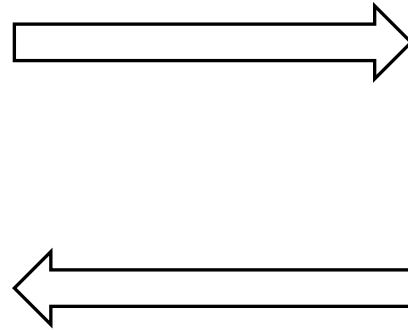
X space: x_1, x_2



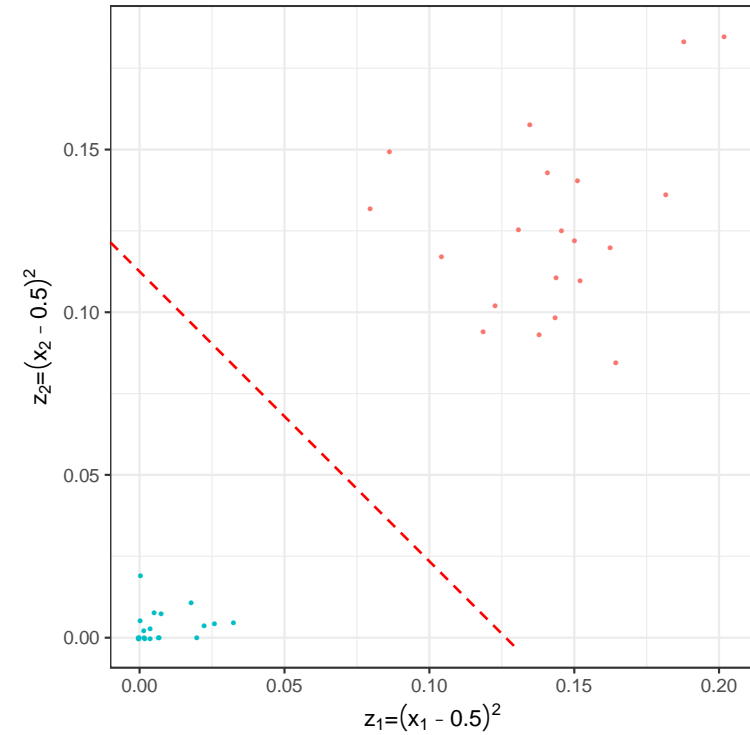
Non-linear **transform**

$$z_1 = (x_1 - 0.5)^2$$

$$z_2 = (x_2 - 0.5)^2$$



Z space: z_1, z_2



Modification to the Basic SVM Lagrangian

- $L(\lambda) = -\frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \lambda^{(i)} \lambda^{(j)} y^{(i)} y^{(j)} \mathbf{z}^{(i)} \cdot \mathbf{z}^{(j)} + \sum_{i=1}^m \lambda^{(i)}$ subject to the constraints:
 - $\lambda^{(i)} \geq 0$
 - $\lambda \cdot \mathbf{y} = 0$
- We can solve this in the \mathbf{z} -space as before
- Dimension of \mathbf{z} does not matter
- We just need the inner product term
- Quadratic programming gives us $\lambda^{(i)}$
- Most of these will be 0 except for support vectors where $\lambda^{(k)} > 0$
 - So, $\mathbf{w} = \begin{bmatrix} w_1 \\ \dots \\ w_n \end{bmatrix} = \sum_{\mathbf{z}^{(i)} \in SV} \lambda^{(i)} y^{(i)} \mathbf{z}^{(i)}$
 - b can be obtained using one of the support vectors $\mathbf{z}^{(k)}$ from:
 - $|b + \sum_j w_j z_j^{(k)}| = 1$
 - Hence, we have the parameters for the decision boundary in the \mathbf{z} -space: $\mathbf{w}^T \mathbf{z} + b = 0$

What do we need from \mathbf{Z} space?

$$\mathcal{L}(\boldsymbol{\lambda}) = \sum_{i=1}^m \lambda^{(i)} - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m y^{(i)} y^{(j)} \lambda^{(i)} \lambda^{(j)} \mathbf{z}^{(i)T} \mathbf{z}^{(j)}$$

Constraints:

1. $\lambda^{(i)} \geq 0$ for $i = 1, \dots, m$
 2. $\sum_{i=1}^m \alpha^{(i)} y^{(i)} = 0$
- } No requirement of \mathbf{z}

Solve this in **\mathbf{Z} space** using quadratic programming

$$\boldsymbol{\lambda} = (\lambda^{(1)}, \lambda^{(2)}, \dots, \lambda^{(m)})$$

$$\mathbf{w} = \sum_{\mathbf{z}^{(i)} \in \text{SV}} \lambda^{(i)} y^{(i)} \mathbf{z}^{(i)}$$

b can be solved by taking any SV $\mathbf{z}^{(k)}$: $y^{(k)} (\mathbf{w}^T \mathbf{z}^{(k)} + b) = 1$

The decision boundary is: $\mathbf{w}^T \mathbf{z} + b = 0$

Only need to know $\mathbf{z}^{(i)T} \mathbf{z}^{(k)}$, i.e., inner product

Do not need to know exactly what the transform $\mathbf{x} \rightarrow \mathbf{z}$ is

Existence of \mathbf{z} space is sufficient:

- We do not to know *what* it is.
- We just need the inner product

Inner product in \mathbf{z} space

Given any two points in X space: \mathbf{x} and \mathbf{x}'

We need $\mathbf{z}^T \mathbf{z}'$

Let $\mathbf{z}^T \mathbf{z}' = K(\mathbf{x}, \mathbf{x}') \longrightarrow$ **Kernel** a **kernel** is a function that computes the dot product between the transformed data points in a higher-dimensional space

Example:

$$\begin{array}{ccc} \mathbf{x} = (x_1, x_2) & \xrightarrow{\text{2nd order polynomial transformation } \Phi} & \mathbf{z} = \Phi(\mathbf{x}) = (1, x_1, x_2, x_1^2, x_2^2, x_1 x_2) \\ \mathbf{x}' = (x'_1, x'_2) & & \mathbf{z}' = \Phi(\mathbf{x}') = (1, x'_1, x'_2, x'^2_1, x'^2_2, x'_1 x'_2) \end{array}$$

$$K(\mathbf{x}, \mathbf{x}') = \mathbf{z}^T \mathbf{z}' = 1 + x_1 x'_1 + x_2 x'_2 + x_1^2 x'^2_1 + x_2^2 x'^2_2 + x_1 x_2 x'_1 x'_2$$

Can we compute $K(\mathbf{x}, \mathbf{x}')$ *without* transforming \mathbf{x} and \mathbf{x}' , i.e., without explicit knowledge of Φ ?

Inner product in Z space

Consider the form $K(\mathbf{x}, \mathbf{x}') = (1 + \mathbf{x}^T \mathbf{x}')^2$ $\left\{ \begin{array}{l} \text{It's just a function, not inner product in } X \text{ space} \\ \text{We don't know if it is the inner product in any other space} \end{array} \right.$

$$= (1 + x_1 x'_1 + x_2 x'_2)^2$$

$$= 1 + x_1^2 x'^2_1 + x_2^2 x'^2_2 + \underline{2x_1 x'_1} + \underline{2x_2 x'_2} + \underline{2x_1 x_2 x'_1 x'_2}$$

This is still an inner product

If we use such a $\mathbf{x} \rightarrow \mathbf{z}$ transformation:

$$\mathbf{z} = (1, x_1^2, x_2^2, \sqrt{2}x_1, \sqrt{2}x_2, \sqrt{2}x_1x_2)$$

$$\mathbf{z}' = (1, x'^2_1, x'^2_2, \sqrt{2}x'_1, \sqrt{2}x'_2, \sqrt{2}x'_1x'_2)$$

Kernel becomes handy when raising to higher order space

Example: Polynomial kernel

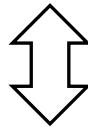
$\mathbf{x} \in \mathbb{R}^d$ $\Phi: \mathbf{x} \rightarrow \mathbf{z}$ is a polynomial transformation of order Q

The equivalent kernel is: $K(\mathbf{x}, \mathbf{x}') = (\mathbf{1} + \mathbf{x}^T \mathbf{x}')^Q$
 $= (1 + x_1 x'_1 + x_2 x'_2 + \cdots + x_d x'_d)^Q$



This is a rather simple computation:

sum of d products, and raise to the power of Q



Don't need to actually expand the polynomial

Scales can be adjusted: $K(\mathbf{x}, \mathbf{x}') = (\mathbf{a} \mathbf{x}^T \mathbf{x}' + \mathbf{b})^Q$

Example: RBF kernel

We only need \mathbf{z} to exist. Don't need to know what \mathbf{z} is.

Radius basis function (**RBF**) kernel:

$$K(\mathbf{x}, \mathbf{x}') = e^{-\eta \|\mathbf{x} - \mathbf{x}'\|^2} \longrightarrow \text{Is this an inner product in *some* space? YES!}$$

An easier example: $x \in \mathbb{R}$, $\eta = 1$

$$K(x, x') = e^{-(x-x')^2}$$

$$= e^{-x^2} e^{x'^2} \boxed{e^{2xx'}}$$

Taylor expansion:

$$e^x = 1 + x + \frac{x^2}{2} + \frac{x^3}{3!} + \cdots + \frac{x^n}{n!} + \cdots$$

$$= \exp(-x^2) \exp(-x'^2) \sum_{k=0}^{\infty} \frac{2^k x^k x'^k}{k!}$$

This is an inner product within a \mathbf{z} space of **infinite dimensions!**

$$\left\{ \begin{aligned} &= \exp(-x^2) \exp(-x'^2) (1 + 2xx' + \frac{2^2}{2!} x^2 x'^2 + \frac{2^3}{3!} x^3 x'^3 + \cdots) \end{aligned} \right.$$

Quadratic programming with kernel function

$$\max_{\boldsymbol{\lambda}} \left(-\frac{1}{2} \boldsymbol{\lambda}^T \mathbf{Q} \boldsymbol{\lambda} + (\mathbf{1}^T) \boldsymbol{\lambda} \right) \quad \text{Subject to:} \quad \mathbf{y}^T \boldsymbol{\lambda} = 0; \boldsymbol{\lambda} > 0$$

$$\mathbf{Q} = \begin{bmatrix} y^{(1)}y^{(1)}\mathbf{K}(\mathbf{x}^{(1)}, \mathbf{x}^{(1)}) & \dots & y^{(1)}y^{(m)}\mathbf{K}(\mathbf{x}^{(1)}, \mathbf{x}^{(m)}) \\ \vdots & \ddots & \vdots \\ y^{(m)}y^{(1)}\mathbf{K}(\mathbf{x}^{(m)}, \mathbf{x}^{(1)}) & \dots & y^{(m)}y^{(m)}\mathbf{K}(\mathbf{x}^{(m)}, \mathbf{x}^{(m)}) \end{bmatrix}$$

Replace $\mathbf{x}^{(i)T} \mathbf{x}^{(j)}$ with $K(\mathbf{x}^{(i)}, \mathbf{x}^{(j)})$

Everything else is the same!

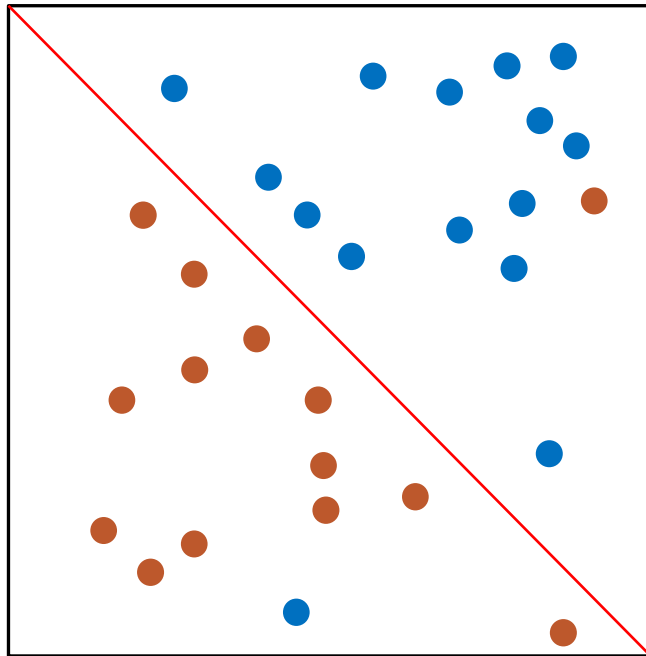
Soft Margin

Introduction

- The soft margin in Support Vector Machines (SVM) is beneficial in scenarios where the data may not be perfectly separable, or when there is a need to tolerate some level of misclassification. Here are some example use cases where the soft margin in SVM could be advantageous:
 - In medical diagnosis, patient data may not have clear boundaries between different classes.
 - Image classification tasks may involve noisy data, and perfect separation between classes may not be possible.
 - In text classification, features may not provide a clear separation between different categories.
 - Financial data which is often noisy and subject to market fluctuations.
 - Biological data, such as gene expression profiles, may exhibit overlapping patterns.
 - Fraud detection may involve imbalanced datasets and ambiguous patterns
 - Handwriting recognition tasks may encounter variations in writing styles.
 - Speech signals often include background noise, making perfect separation difficult.
 - Customer behavior data may not always lead to clear distinctions between market segments.
 - Network security data often contains anomalies that may not be clearly separated from normal behavior.
- In these use cases, the soft margin in SVM provides a balance between achieving a wider margin and allowing for misclassifications, making the model more adaptable to real-world scenarios with inherent uncertainty and complexity.

Use of Soft Margin vs Kernel Method

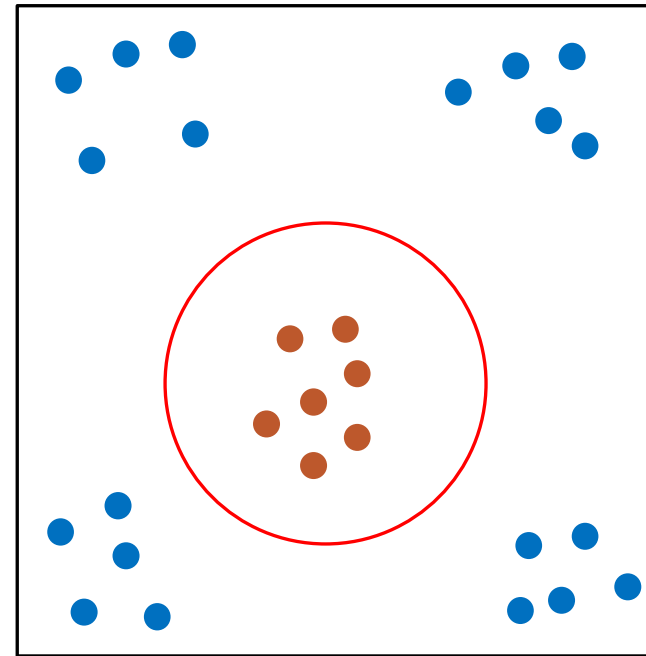
Linear Separation May be Tolerable



Soft-margin deals with this

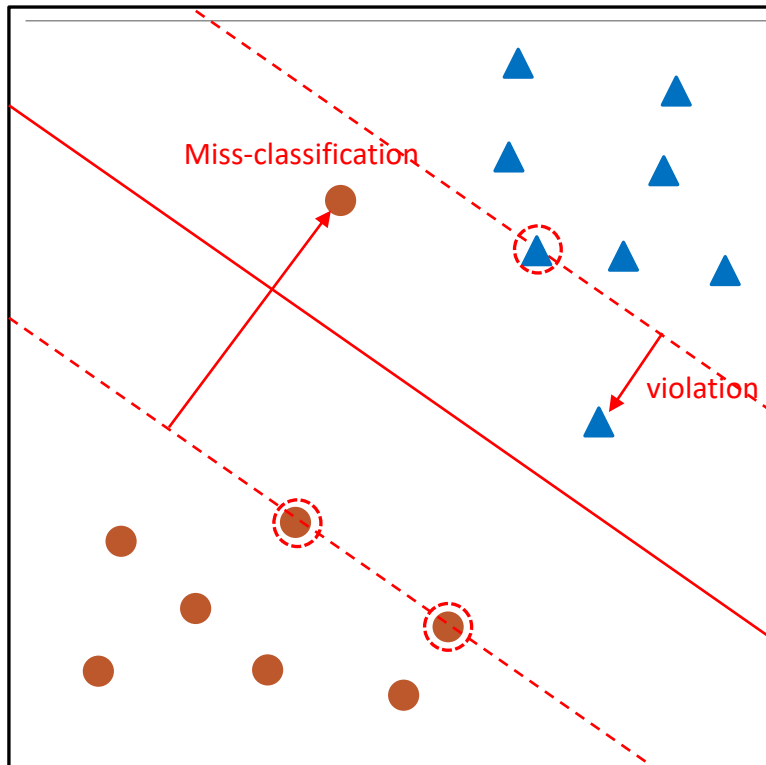
Combine

Linear Separation Won't Work



Kernels deal with this

Violation of margin



For SVs: $y^{(k)}(\mathbf{w}^T \mathbf{x}^{(k)} + b) = 1$

For any points outside the margin: $y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + b) > 0$

Allow violations:

$$\begin{cases} 0 < y^{(k)}(\mathbf{w}^T \mathbf{x}^{(k)} + b) < 1 & \text{Violate margin} \\ y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + b) < 0 & \text{Miss-classification} \end{cases}$$

Use a **slack** to measure the violation:

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

Penalize the *total violation*: $= \sum_{i=1}^m \xi^{(i)}$

Change the optimization goal by including this error term

New optimization goal

Minimize: $\frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^m \xi^{(i)}$

Subject to $y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + b) \geq 1 - \xi^{(i)}, \text{ for } i = 1, \dots, m$

and $\xi^{(i)} \geq 0, \text{ for } i = 1, \dots, m$

Hyper-parameter C
quantifies the relative
importance of avoiding
violations

or tolerance to violations
Like L2 regularization

If C is big:

w will be learned in a way that only small $\xi^{(i)}$ are allowed

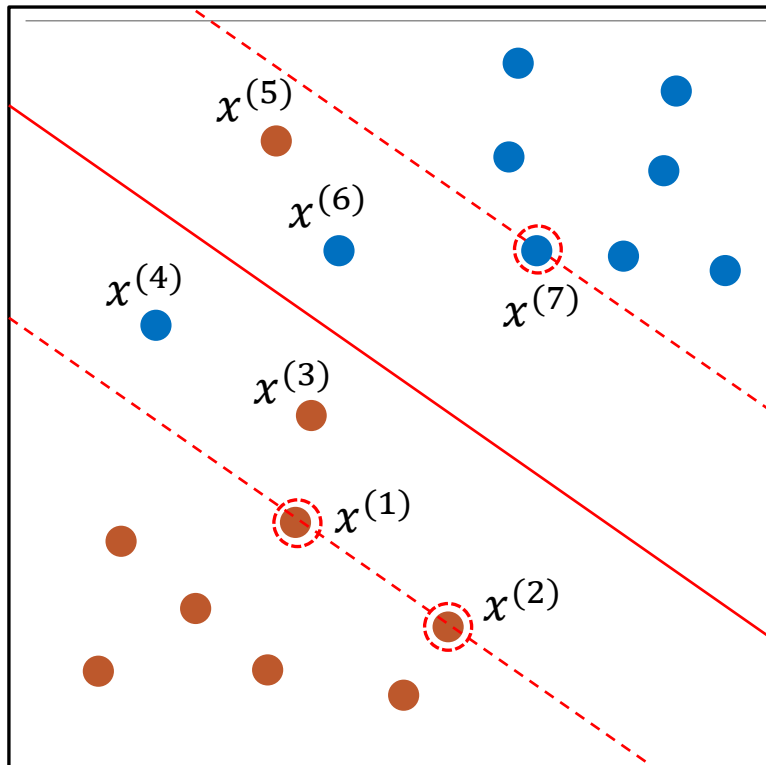
If C is small:

w will be learned in a way that big $\xi^{(i)}$ are allowed

Updated Lagrangian

- We define the Lagrangian as:
 - $L(\mathbf{w}, b, \boldsymbol{\xi}, \boldsymbol{\lambda}, \boldsymbol{\gamma}) = \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^m \xi^{(i)} - \sum_{i=1}^m \lambda^{(i)} (y^{(i)} (\mathbf{w}^T \mathbf{x}^{(i)} + b) - 1 + \xi^{(i)}) - \sum_{i=1}^m \gamma^{(i)} \xi^{(i)}$
- We minimize L wrt \mathbf{w} and b and $\boldsymbol{\xi}$ and maximize L wrt $\lambda^{(i)} \geq 0$ and $\gamma^{(i)} \geq 0$
- Minimize L wrt \mathbf{w} and b
 - $\frac{\partial L}{\partial \mathbf{w}} = \mathbf{w} - \sum_{i=1}^m \lambda^{(i)} y^{(i)} \mathbf{x}^{(i)} = 0$
 - $\frac{\partial L}{\partial b} = - \sum_{i=1}^m \lambda^{(i)} y^{(i)} = 0$
 - $\Rightarrow \boldsymbol{\lambda} \cdot \mathbf{y} = 0$
 - $\frac{\partial L}{\partial \xi^{(i)}} = C - \lambda^{(i)} - \gamma^{(i)} = 0$
- In $L(\mathbf{w}, b, \boldsymbol{\xi}, \boldsymbol{\lambda}, \boldsymbol{\gamma})$, substituting $\mathbf{w} = \sum_{i=1}^m \lambda^{(i)} y^{(i)} \mathbf{x}^{(i)}$ in L and utilizing $\sum_{i=1}^m \lambda^{(i)} y^{(i)} = 0$ and utilizing $C - \lambda^{(i)} - \gamma^{(i)} = 0$
- $L(\mathbf{w}, b, \boldsymbol{\lambda}) = \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \lambda^{(i)} \lambda^{(j)} y^{(i)} y^{(j)} \mathbf{x}^{(i)} \cdot \mathbf{x}^{(j)} - \sum_{i=1}^m \sum_{j=1}^m \lambda^{(i)} \lambda^{(j)} y^{(i)} y^{(j)} \mathbf{x}^{(i)} \cdot \mathbf{x}^{(j)} + \sum_{i=1}^m \lambda^{(i)}$
- Same Lagrangian as before with one additional constraint, $C \geq \lambda^{(i)} \geq 0$ because $\gamma^{(i)} \geq 0$
- We can utilize Quadratic programming with the additional constraint to solve for $\lambda^{(i)}$ and then for \mathbf{w} and b

Types of support vectors



Margin support vectors: $C > \lambda^{(i)} > 0$

$$y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + b) = 1 \quad \xi^{(i)} = 0 \quad \gamma^{(i)} > 0$$

Non-margin support vectors: $\lambda^{(i)} = C$

$$y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + b) < 1 \quad \xi^{(i)} > 0 \quad \gamma^{(i)} = 0$$

Violate the margin, but correctly classified

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

Violate the margin, and misclassified

$$y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + b) < 0$$

Choose proper hyperparameters

Soft margin is usually the default

Choose $C \in \{2^{-4}, 2^{-3}, \dots, 2^4\}$ Best C can be determined by cross-validation

Combined with kernel methods

E.g., RBF kernel $K(\mathbf{x}, \mathbf{x}') = \exp(-\eta \|\mathbf{x} - \mathbf{x}'\|^2)$

Choose best C and γ using grid search: $C, \eta \in \{2^{-4}, 2^{-3}, \dots, 2^4\}$