# Machine Learning
# CSCE 5215

## Support Vector Machines

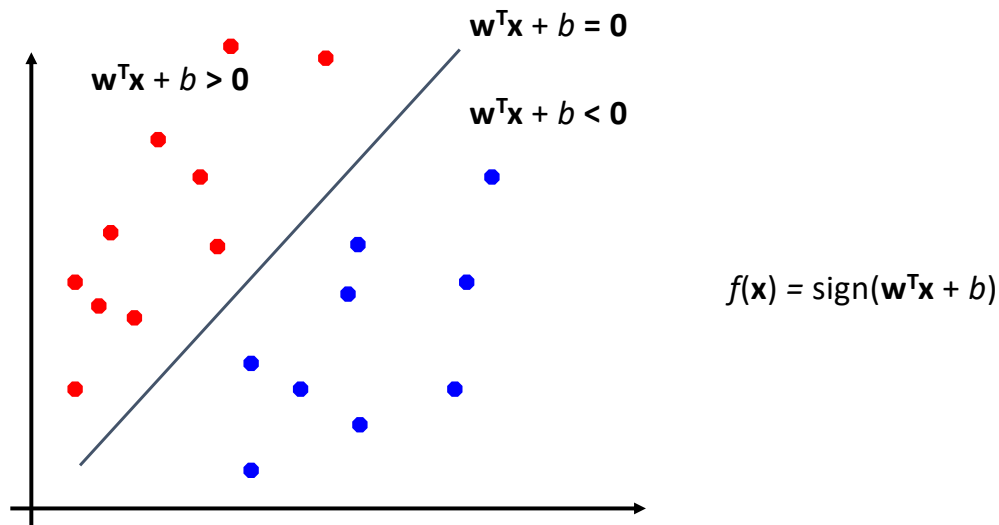Instructor: Zeenat Tariq

# Main Idea

- Max Margin Classifier: Formalize notion of the best linear separator

- Lagrangian Multipliers: Way to convert a constrained optimization problem to one that is easier to solve

- Kernel: Projecting data into higher-dimensional space makes it linearly separable

- Complexity: Depends only on the number of training examples, not on dimensionality of the kernel space!
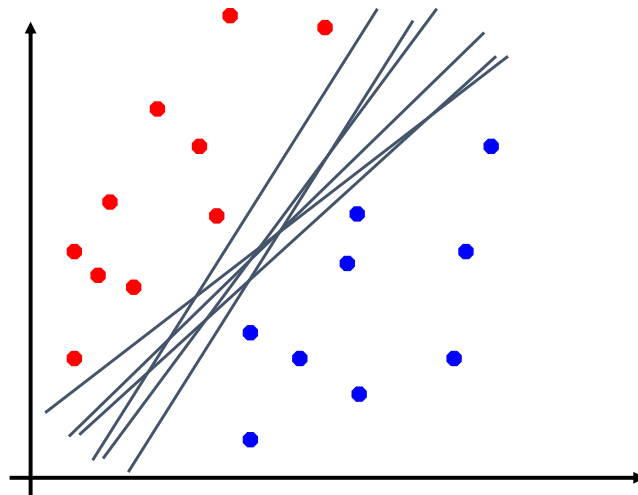
# Perceptron Revisited: Linear Separators

- Binary classification can be viewed as the task of separating classes in feature space:



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

$\mathbf{w}^T\mathbf{x} + b > 0$

$\mathbf{w}^T\mathbf{x} + b < 0$

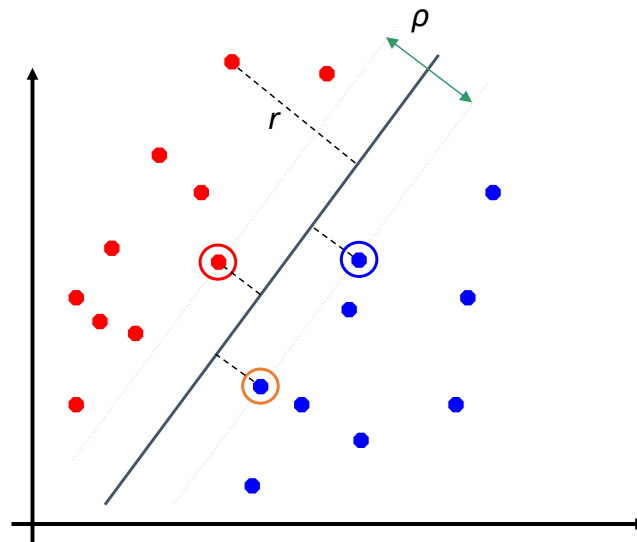$f(\mathbf{x}) = \text{sign}(\mathbf{w}^T\mathbf{x} + b)$

# Linear Separators

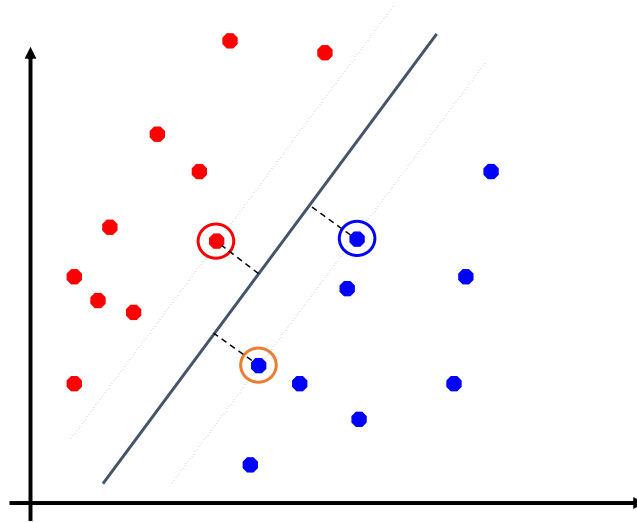- Which of the linear separators is optimal?

# Classification Margin

- Distance from example $\mathbf{x}_i$ to the separator is $r = \dfrac{\mathbf{w}^T\mathbf{x}_i + b}{\|\mathbf{w}\|}$
- Examples closest to the hyperplane are **support vectors**.
- **Margin** $\rho$ of the separator is the distance between support vectors.

# Maximum Margin Classification

- Maximizing the margin is good according to intuition and PAC theory.

- Implies that only support vectors matter; other training examples are ignorable.

# Linear SVM Mathematically

- Let training set $\{(\mathbf{x}_i, y_i)\}_{i=1..n}$, $\mathbf{x}_i \in \mathbf{R}^d$, $y_i \in \{-1, 1\}$ be separated by a hyperplane with margin $\rho$. Then for each training example $(\mathbf{x}_i, y_i)$:

$$\mathbf{w}^\mathsf{T}\mathbf{x}_i + b \leq -\rho/2 \quad \text{if } y_i = -1$$
$$\mathbf{w}^\mathsf{T}\mathbf{x}_i + b \geq \rho/2 \quad \text{if } y_i = 1$$

$\Longleftrightarrow \qquad y_i(\mathbf{w}^\mathsf{T}\mathbf{x}_i + b) \geq \rho/2$

- For every support vector $\mathbf{x}_s$ the above inequality is an equality. After rescaling $\mathbf{w}$ and $b$ by $\rho/2$ in the equality, we obtain that distance between each $\mathbf{x}_s$ and the hyperplane is

$$r = \frac{y_s(\mathbf{w}^T\mathbf{x}_s + b)}{\|\mathbf{w}\|} = \frac{1}{\|\mathbf{w}\|}$$

- Then the margin can be expressed through (rescaled) $\mathbf{w}$ and $b$ as:

$$\rho = 2r = \frac{2}{\|\mathbf{w}\|}$$

# Linear SVMs Mathematically (cont.)

Then we can formulate the *quadratic optimization problem:*

> Find **w** and $b$ such that
> $\rho = \dfrac{2}{\|\mathbf{w}\|}$ is maximized
> and for all ($\mathbf{x}_i$, $y_i$), $i=1..n$ :     $y_i(\mathbf{w}^\mathsf{T}\mathbf{x}_i + b) \geq 1$

Which can be reformulated as:

> Find **w** and $b$ such that
>
> $\Phi(\mathbf{w}) = ||\mathbf{w}||^2 = \mathbf{w}^\mathsf{T}\mathbf{w}$  is minimized
>
> and for all ($\mathbf{x}_i$, $y_i$), $i=1..n$ :     $y_i\,(\mathbf{w}^\mathsf{T}\mathbf{x}_i + b) \geq 1$

# Solving the Optimization Problem

> Find **w** and b such that
> **Φ(w)** =$\mathbf{w}^T\mathbf{w}$  is minimized
> and for all ($\mathbf{x}_i$, $y_i$), $i$=1..$n$ :     $y_i$ ($\mathbf{w}^T\mathbf{x}_i$ + $b$) ≥ 1

- Need to optimize a *quadratic* function subject to *linear* constraints.

- Quadratic optimization problems are a well-known class of mathematical programming problems for which several (non-trivial) algorithms exist.

- The solution involves constructing a *dual problem* where a *Lagrange multiplier $\alpha_i$* is associated with every inequality constraint in the primal (original) problem:

> Find $\alpha_1...\alpha_n$ such that
> **Q(α)** =$\Sigma\alpha_i$ - ½$\Sigma\Sigma\alpha_i\alpha_j y_i y_j \mathbf{x}_i^T\mathbf{x}_j$ is maximized and
> (1)  $\Sigma\alpha_i y_i$ = 0
> (2) $\alpha_i$ ≥ 0 for all $\alpha_i$

# The Optimization Problem Solution

- Given a solution $\alpha_1...\alpha_n$ to the dual problem, solution to the primal is:

$$\mathbf{w} = \Sigma \alpha_i y_i \mathbf{x}_i \qquad b = y_k - \Sigma \alpha_i y_i \mathbf{x}_i{}^\mathsf{T} \mathbf{x}_k \quad \text{for any } \alpha_k > 0$$
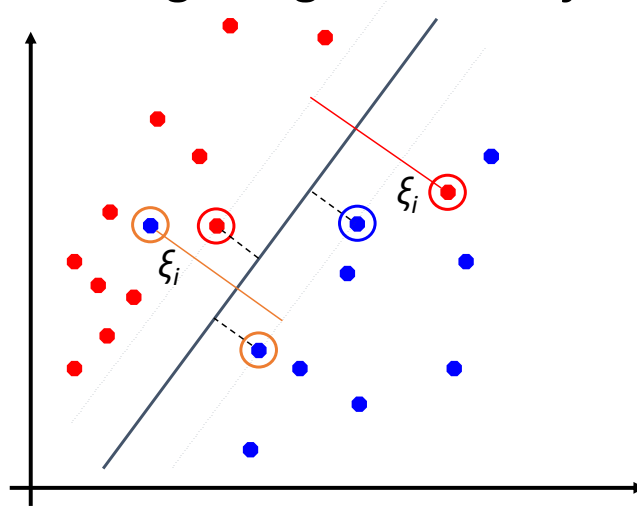
- Each non-zero $\alpha_i$ indicates that corresponding $\mathbf{x}_i$ is a support vector.
- Then the classifying function is (note that we don't need $\mathbf{w}$ explicitly):

$$f(\mathbf{x}) = \Sigma \alpha_i y_i \mathbf{x}_i{}^\mathsf{T} \mathbf{x} + b$$

- Notice that it relies on an *inner product* between the test point $\mathbf{x}$ and the support vectors $\mathbf{x}_i$
- Also keep in mind that solving the optimization problem involved computing the inner products $\mathbf{x}_i{}^\mathsf{T}\mathbf{x}_j$ between all training points.

# Soft Margin Classification

- What if the training set is not linearly separable?
- *Slack variables $\xi_i$* can be added to allow misclassification of difficult or noisy examples, resulting margin called *soft*.

# Soft Margin Classification Mathematically

- The old formulation:

  Find **w** and b such that
  $\Phi(\mathbf{w}) = \mathbf{w}^T\mathbf{w}$ is minimized
  and for all $(\mathbf{x}_i, y_i)$, $i=1..n$ :   $y_i(\mathbf{w^T x}_i + b) \geq 1$

- Modified formulation incorporates slack variables:

  Find **w** and b such that
  $\Phi(\mathbf{w}) = \mathbf{w}^T\mathbf{w} + C\Sigma\xi_i$ is minimized
  and for all $(\mathbf{x}_i, y_i)$, $i=1..n$ :   $y_i(\mathbf{w^T x}_i + b) \geq 1 - \xi_i$ ,   $\xi_i \geq 0$

- Parameter $C$ can be viewed as a way to control overfitting: it "trades off" the relative importance of maximizing the margin and fitting the training data.

# Soft Margin Classification – Solution

- Dual problem is identical to separable case (would *not* be identical if the penalty for slack variables $C\Sigma\xi_i^2$ was used in primal objective, we would need additional Lagrange multipliers for slack variables):

Find $\alpha_1...\alpha_N$ such that
$\mathbf{Q(\alpha)} = \Sigma\alpha_i - \frac{1}{2}\Sigma\Sigma\alpha_i\alpha_j y_i y_j \mathbf{x}_i^\mathsf{T}\mathbf{x}_j$ is maximized and
(1) $\Sigma\alpha_i y_i = 0$
(2) $0 \leq \alpha_i \leq C$ for all $\alpha_i$

- Again, $\mathbf{x}_i$ with non-zero $\alpha_i$ will be support vectors.

- Solution to the dual problem is:

$\mathbf{w} = \Sigma\alpha_i y_i \mathbf{x}_i$
$b = y_k(1 - \xi_k) - \Sigma\alpha_i y_i \mathbf{x}_i^\mathsf{T}\mathbf{x}_k$   for any $k$ s.t. $\alpha_k > 0$

Again, we don't need to compute **w** explicitly for classification:
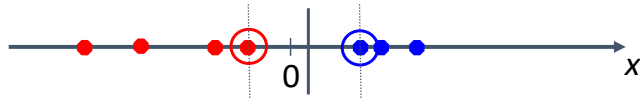
$f(\mathbf{x}) = \Sigma\alpha_i y_i \mathbf{x}_i^\mathsf{T}\mathbf{x} + b$

# Linear SVMs: Overview

- The classifier is a *separating hyperplane.*
- Most "important" training points are support vectors; they define the hyperplane.
- Quadratic optimization algorithms can identify which training points $x_i$ are support vectors with non-zero Lagrangian multipliers $\alpha_i$.
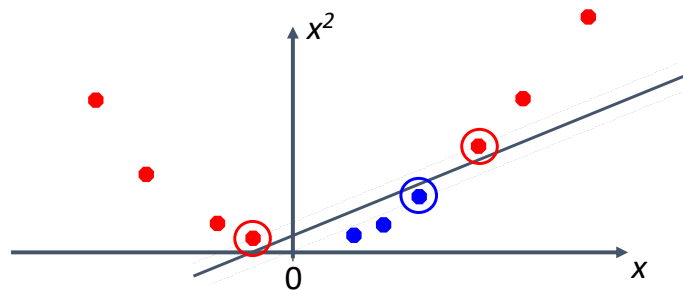
# Non-linear SVMs

- Datasets that are linearly separable with some noise work out great:

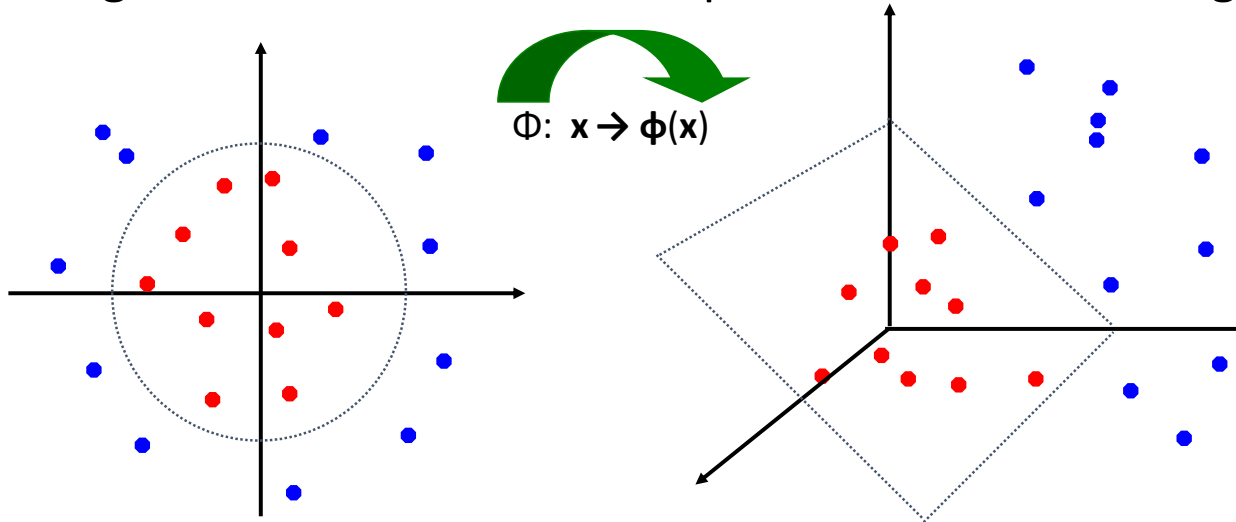- But what are we going to do if the dataset is just too hard?

- How about… mapping data to a higher-dimensional space:

# Non-linear SVMs: Feature spaces

- General idea: the original feature space can always be mapped to some higher-dimensional feature space where the training set is separable:

$$\Phi: \ \mathbf{x} \rightarrow \boldsymbol{\phi}(\mathbf{x})$$

# What Functions are Kernels?

- For some functions $K(\mathbf{x}_i, \mathbf{x}_j)$ checking that $K(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^\top \phi(\mathbf{x}_j)$ can be cumbersome.

- Mercer's theorem:

  ***Every semi-positive definite symmetric function is a kernel***

- Semi-positive definite symmetric functions correspond to a semi-positive definite symmetric Gram matrix:

K=

| $K(\mathbf{x}_1,\mathbf{x}_1)$ | $K(\mathbf{x}_1,\mathbf{x}_2)$ | $K(\mathbf{x}_1,\mathbf{x}_3)$ | … | $K(\mathbf{x}_1,\mathbf{x}_n)$ |
|---|---|---|---|---|
| $K(\mathbf{x}_2,\mathbf{x}_1)$ | $K(\mathbf{x}_2,\mathbf{x}_2)$ | $K(\mathbf{x}_2,\mathbf{x}_3)$ | | $K(\mathbf{x}_2,\mathbf{x}_n)$ |
| | | | | |
| … | … | … | … | … |
| $K(\mathbf{x}_n,\mathbf{x}_1)$ | $K(\mathbf{x}_n,\mathbf{x}_2)$ | $K(\mathbf{x}_n,\mathbf{x}_3)$ | … | $K(\mathbf{x}_n,\mathbf{x}_n)$ |

# Examples of Kernel Functions

- Linear: $K(\mathbf{x}_i,\mathbf{x}_j)= \mathbf{x}_i^\mathsf{T}\mathbf{x}_j$
  - Mapping $\Phi$:   $\mathbf{x} \to \boldsymbol{\phi}(\mathbf{x})$, where $\boldsymbol{\phi}(\mathbf{x})$ is $\mathbf{x}$ itself


- Polynomial of power $p$: $K(\mathbf{x}_i,\mathbf{x}_j)= (1+ \mathbf{x}_i^\mathsf{T}\mathbf{x}_j)^p$
  - Mapping $\Phi$:   $\mathbf{x} \to \boldsymbol{\phi}(\mathbf{x})$, where $\boldsymbol{\phi}(\mathbf{x})$ has $\binom{d+p}{p}$ dimensions


- Gaussian (radial-basis function): $K(\mathbf{x}_i,\mathbf{x}_j) = e^{-\frac{\left\|\mathbf{x}_i-\mathbf{x}_j\right\|^2}{2\sigma^2}}$

  - Mapping $\Phi$:  $\mathbf{x} \to \boldsymbol{\phi}(\mathbf{x})$, where $\boldsymbol{\phi}(\mathbf{x})$ is *infinite-dimensional*: every point is mapped to *a function* (a Gaussian); combination of functions for support vectors is the separator.

- Higher-dimensional space still has *intrinsic* dimensionality $d$ (the mapping is not *onto*), but linear separators in it correspond to *non-linear* separators in original space.

# Non-linear SVMs Mathematically

- Dual problem formulation:

Find $\alpha_1 \ldots \alpha_n$ such that
$\mathbf{Q(\alpha)} = \Sigma \alpha_i - \frac{1}{2} \Sigma \Sigma \alpha_i \alpha_j y_i y_j K(\mathbf{x}_i, \mathbf{x}_j)$ is maximized and
(1) $\Sigma \alpha_i y_i = 0$
(2) $\alpha_i \geq 0$ for all $\alpha_i$

- The solution is:

$f(\mathbf{x}) = \Sigma \alpha_i y_i K(\mathbf{x}_i, \mathbf{x}_j) + b$

- Optimization techniques for finding $\alpha_i$'s remain the same!