

CS 559 Machine Learning

Lecture 6: Support Vector Machines

Ping Wang

Department of Computer Science

Stevens Institute of Technology



Today's Lecture

- SVM for Linear Separable Case
- Non-separable Case, Penalties
- Non-linearity, Kernels

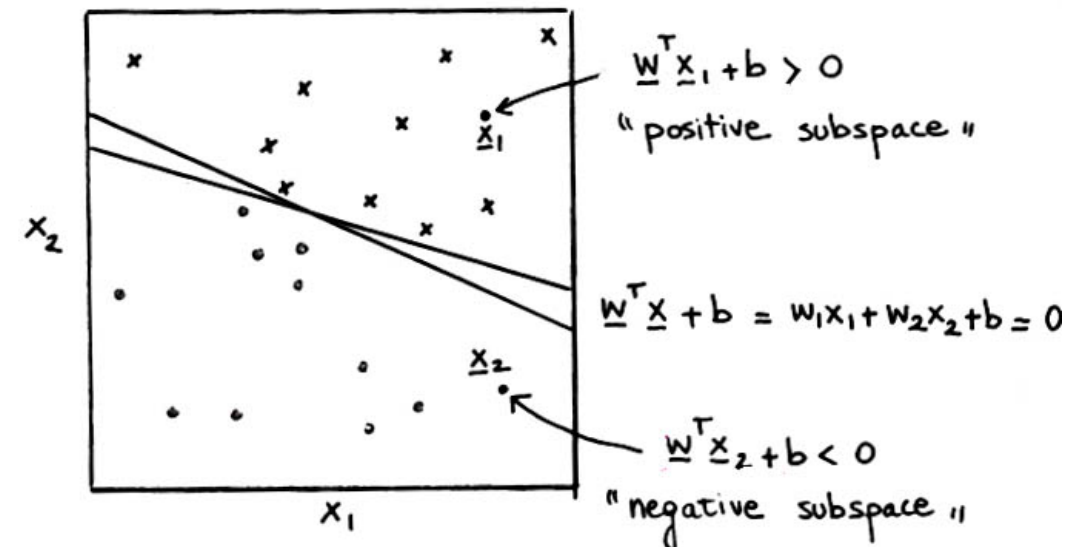
Linear Classifier

- Linear classifiers construct **linear decision boundaries** (hyperplanes) that try to separate the data into different classes as well as possible.
- Classification rule of the **Perceptron algorithm**:

Input: $x \in \mathbb{R}^d$

Output: $\text{sign}(w^T x + b)$

- The classifier computes a linear combination of the input features and return the sign.



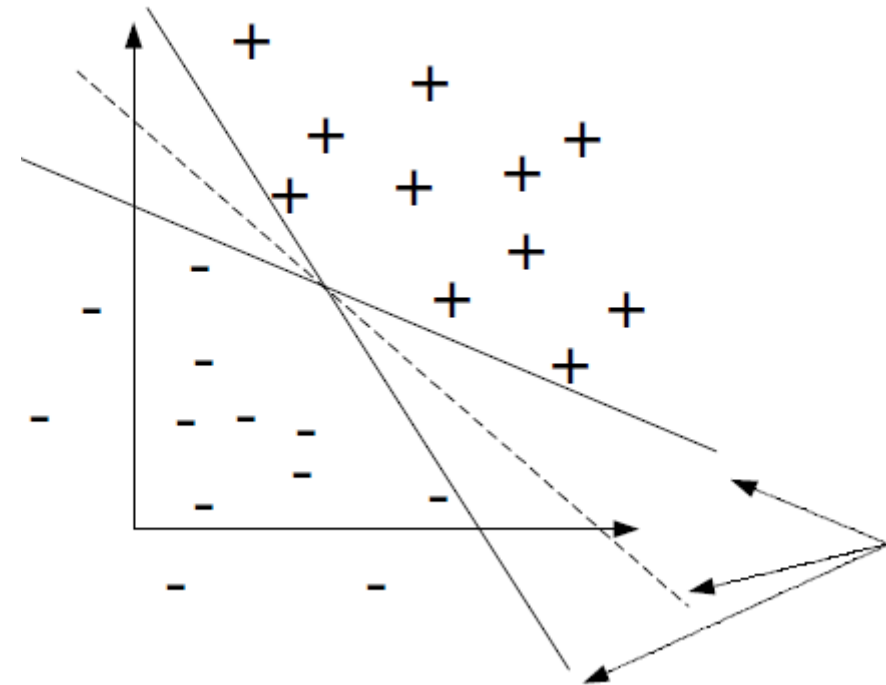
Limitations of Perceptron Algorithm

- When the data are **linearly separable**, there are many solutions.
 - The result depends on the starting values of the parameters and the order of data samples.
 - The finite number of steps to convergence can be very large: the smaller the gap between the two classes, the longer the time to find it.
- When the data are **not linearly separable**, the algorithm will not converge.

The Linearly Separable Case

Limitations of Perceptron Algorithm

- When the data are **linearly separable**, there are **many solutions**, and which one is found depends on the starting values of the parameters and the order of data samples.
- There can be an infinite number of hyperplanes that achieve 100% accuracy on training data.
- Which hyperplane is the **optimal** with respect to the accuracy on test data?
- Solution: **adding additional constraints to the separating hyperplane.**

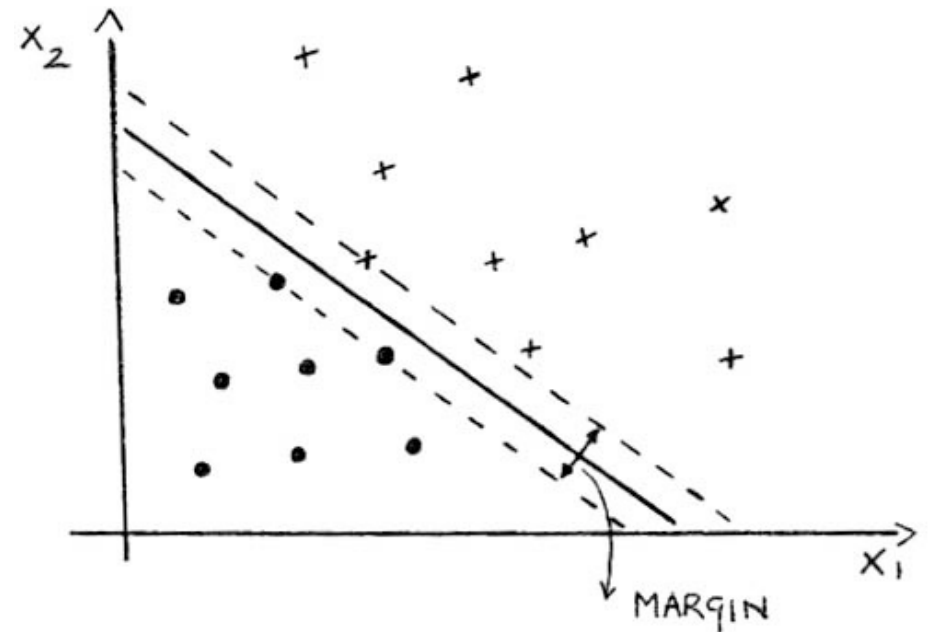


Three possible separations

Largest Margin Hyperplanes

Goal: Find the hyperplane that separates the two classes and maximizes the distance **to the closest points from each class**.

- Such distance is called **margin**.
- The added constraint:
 - Provide a unique solution to the separating hyperplane problem.
 - Maximizing the margin between the two classes on the training data gives better classification performance on test data.



Training Data

- For two classes:

$$(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$$

$$x_i \in R^d$$

$$y_i \in \{-1, +1\}$$

- We need to formalize **the largest margin criterion**.

Formulation

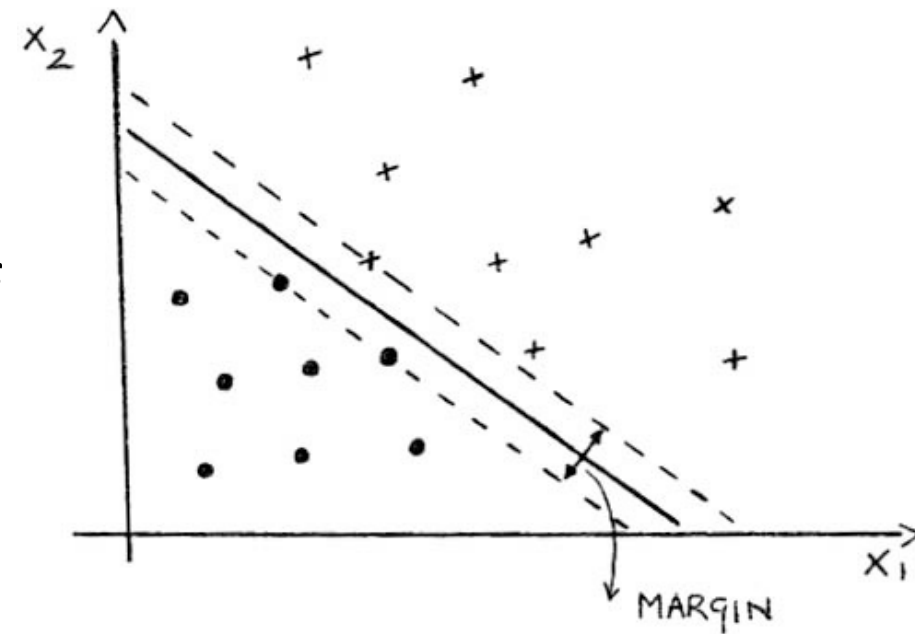
- Consider the following optimization problem:

$$\begin{aligned} & \max_{w,b} 2C \\ & \text{subject to } \frac{1}{\|w\|} y_i (w^T x_i + b) \geq C, i = 1, \dots, N \end{aligned}$$

- Remember a property: The signed distance of any point x to L is:

$$\frac{1}{\|w\|} (w^T x_i + b)$$

- Thus, the set of conditions above ensure that all the training data are **at least at distance C from the decision boundary**.



Formulation

- The optimization problem:

$$\max_{w,b} 2C \quad \text{subject to} \quad \frac{1}{\|w\|} y_i(w^T x_i + b) \geq C, i = 1, \dots, N$$

- We seek the largest C and associated parameters.
- We can rewrite the above conditions as:

$$y_i(w^T x_i + b) \geq C \|w\|$$

- Since $w^T x + b = 0$ and $c(w^T x + b) = 0$ define the same plane, we can arbitrarily normalize $\|w\| = \frac{1}{C}$.

Formulation

- The optimization problem:

$$\begin{aligned} & \max_{w,b} 2C \\ & \text{subject to } \frac{1}{\|w\|} y_i (w^T x_i + b) \geq C, i = 1, \dots, N \end{aligned}$$

- By **normalizing** $\|w\| = \frac{1}{C}$, the original maximization problem is equivalent to:

$$\begin{aligned} & \min_{w,b} \frac{1}{2} \|w\|^2 \\ & \text{subject to } y_i (w^T x_i + b) \geq 1, i = 1, \dots, N \end{aligned}$$

- The constraints define an empty margin around the linear decision boundary of **thickness** $\frac{2}{\|w\|}$. We choose w, b to maximize its thickness.
- This is a convex quadratic optimization problem subject to linear constraints and there is a unique minimum.

Lagrange Multipliers

- We introduce the **Lagrange multipliers** $\alpha_i \geq 0, i = 1, \dots, N$
- One for each of the inequality constraints.
- Recall the rule:
 - For constraints of the form $C_i \geq 0$, the constraint equations are **multiplied by** Lagrange multipliers and **subtracted from** the objective function, to form the Lagrangian.
- Lagrange multipliers allow us to take the constraints within the function to be minimized.

Lagrange Multipliers: Primal Form

- We then obtain the Lagrangian: (also called **primal form**):

$$L_p = \frac{1}{2} \|w\|^2 - \sum_{i=1}^N \alpha_i [y_i (w^T x_i + b) - 1]$$

- We now **minimize** L_p with respect to w and b :

$$\min_{w,b} \max_{\alpha_i \geq 0} L_p$$

- This indicates that this is the **primal form** of the optimization problem.
- We will actually solve the primal optimization problem by solving the **dual of the original problem**, since they provide the same solution.

Dual Form

- The solution to the dual form provides a lower bound to the solution of the primal form.
- What is the dual form?

$$\max_{\alpha_i \geq 0} \min_{w, b} L_p$$

- Setting the derivatives to zero gives:

$$\frac{\partial L_p}{\partial w} = w - \sum_{i=1}^N \alpha_i y_i x_i = 0 \quad \Rightarrow w = \sum_{i=1}^N \alpha_i y_i x_i \quad (1)$$

$$\frac{\partial L_p}{\partial b} = - \sum_{i=1}^N \alpha_i y_i = 0 \quad \Rightarrow \sum_{i=1}^N \alpha_i y_i = 0 \quad (2)$$

Dual Form

- Substituting Eq. (1) and (2) in L_p gives:

$$\begin{aligned} L_D &= \frac{1}{2} \left(\sum_{i=1}^N \alpha_i y_i x_i \right) \left(\sum_{k=1}^N \alpha_k y_k x_k \right) - \sum_{i=1}^N \alpha_i \left[y_i \left(x_i^T \left(\sum_{k=1}^N \alpha_k y_k x_k \right) + b \right) - 1 \right] \\ &= \frac{1}{2} \sum_{i=1}^N \sum_{k=1}^N \alpha_i \alpha_k y_i y_k x_i^T x_k - \sum_{i=1}^N \sum_{k=1}^N \alpha_i \alpha_k y_i y_k x_i^T x_k - b \sum_{i=1}^N \alpha_i y_i + \sum_{i=1}^N \alpha_i \\ &= \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{k=1}^N \alpha_i \alpha_k y_i y_k x_i^T x_k \end{aligned}$$

Subject to $\alpha_i \geq 0$

The Lagrangian Dual Form

- The solution is obtained by maximizing L_D with respect to the α_i .
- The solution must satisfy the conditions:

$$w = \sum_{i=1}^N \alpha_i y_i x_i$$

$$\sum_{i=1}^N \alpha_i y_i = 0$$

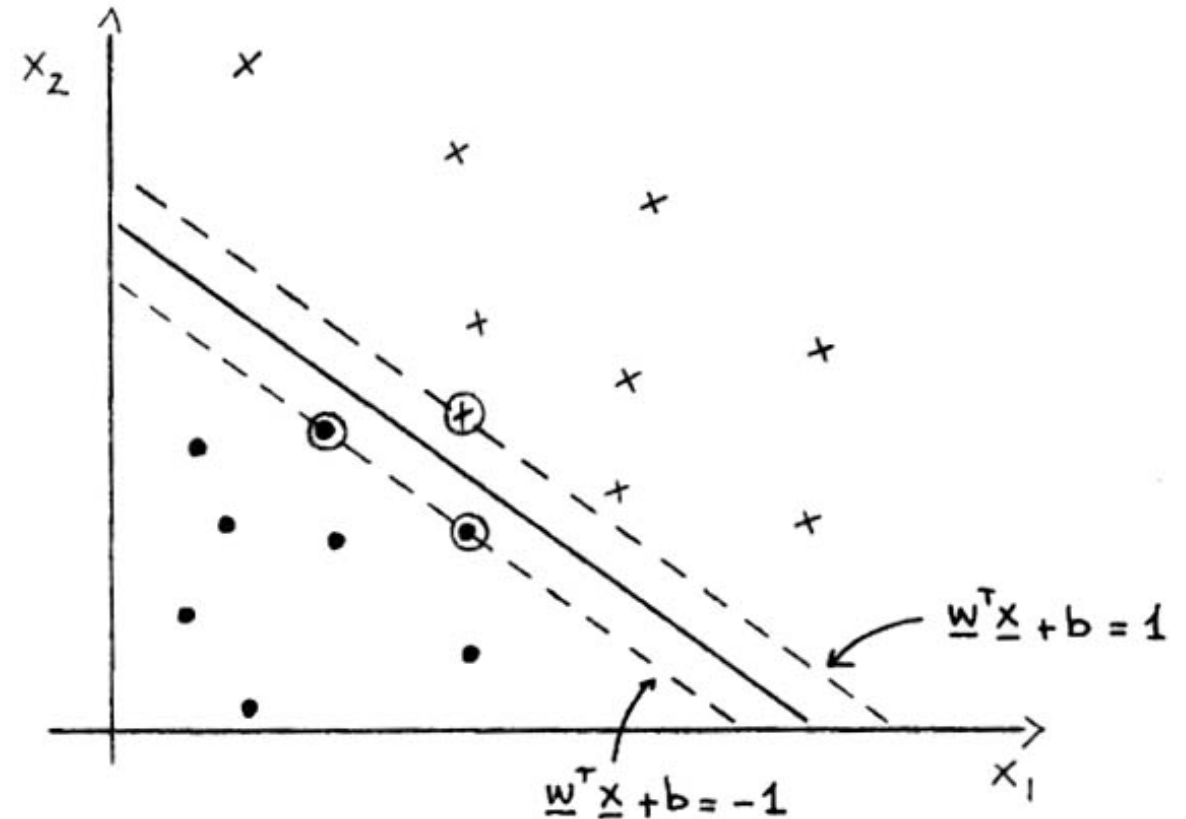
$$\alpha_i \geq 0$$

$$\alpha_i [y_i (w^T x_i + b) - 1] = 0 \quad \forall i = 1, \dots, N$$

Dual Form

$$\alpha_i [y_i(w^T x_i + b) - 1] = 0 \\ \forall i = 1, \dots, N$$

- If $\alpha_i > 0$, then $y_i(w^T x_i + b) = 1$, that is x_i is on the boundary of the margin.
- If $y_i(w^T x_i + b) > 1$, x_i is not on the boundary of the margin, and $\alpha_i = 0$.



Dual Form

- The solution vector w is: $w = \sum_{i=1}^N \alpha_i y_i x_i$. Thus: The solution is defined as a linear combination of those x_i for which $\alpha_i > 0$.
- Such x_i are the points **on the boundary of the margin**. They are called **SUPPORT VECTORS**. We have three support vectors in the above example.
- To obtain the value of b : solve $\alpha_i [y_i (w^T x_i + b) - 1] = 0$ for any of the support vectors.
- The largest margin hyperplane gives a function: $f(x) = w^T x + b$ for classifying new observations $\hat{y} = \text{sign}(f(x))$.

Observations

- The **support vectors** are the critical elements of the training set. They lie closest to the decision boundary.
- **Only the support vectors affect the solution** - If all other training points were removed (or moved around, but so as not to cross the margin), and training was repeated, the same separating hyperplane would be found.
- However, the identification of the support vectors **requires the use of all the training data**.
- Although none of the training observations fall within the margin (by construction), this will not necessarily be the case of test data. (The intuition is that a large margin on the training data indicates a good separation of the two classes and therefore a good separation on the test data as well).

The Non-separable Case

The Non-separable Case

- Suppose now the classes overlap. We can still maximize C , but allow for some points to be on the wrong side of the margin.

- We need to modify the constraints we had for the separable case:

$$\frac{1}{\|w\|} y_i (w^T x_i + b) \geq C, i = 1, \dots, N$$

- To achieve this goal, we define N **slack variables**:

$$\xi_1, \xi_2, \dots, \xi_N$$

- Then a natural way to modify the constraints above is:

$$\frac{1}{\|w\|} y_i (w^T x_i + b) \geq C(1 - \xi_i), i = 1, \dots, N$$

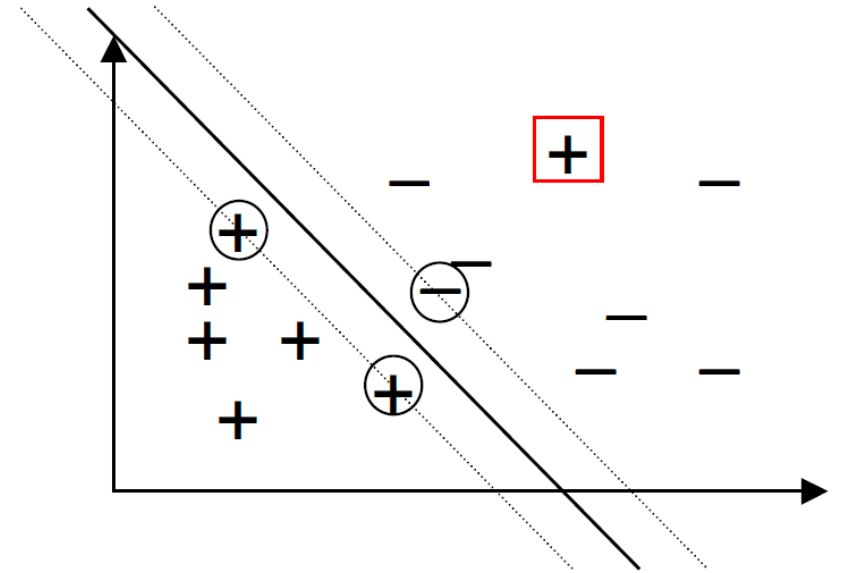
$$\text{with } \xi_i \geq 0, \forall i, \sum_{i=1}^N \xi_i \leq \text{Constant}$$

The Non-separable Case

$$\frac{1}{\|w\|} y_i (w^T x_i + b) \geq C(1 - \xi_i), i = 1, \dots, N$$

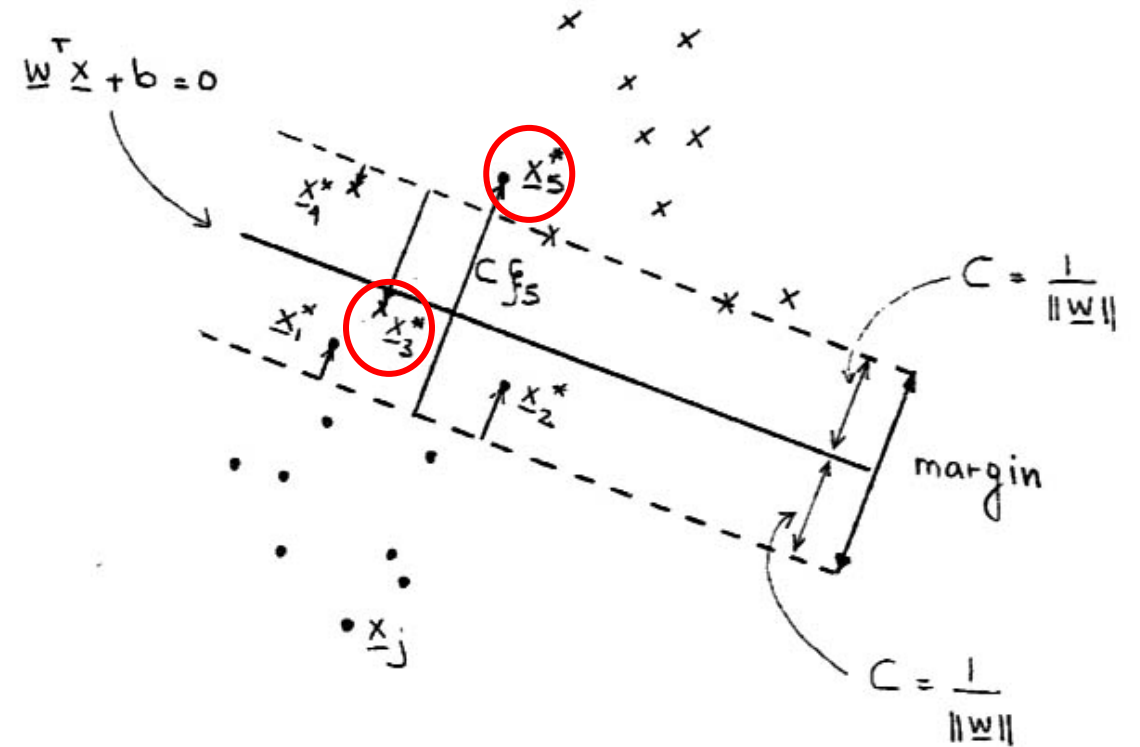
with $\xi_i \geq 0, \forall i, \sum_{i=1}^N \xi_i \leq \text{Constant}$

- Idea of the formulation: ξ_i is the proportional amount by which the prediction $f(x_i)$ is on **the wrong side** of the margin.



Slack Variables

- The points $(x_1^*, x_2^*, x_3^*, x_4^*, x_5^*)$ are on the wrong side of their margin.
- Point x_i^* is on the wrong side of its margin by an amount $C\xi_i$
- Point x_j^* on the correct side have $\xi_j = 0$
- Misclassification occurs when $\xi_i > 1 \Rightarrow C(1 - \xi_i) < 0$, e.g., points x_3^* and x_5^* are misclassified by the given boundary.



A geometric perspective

Slack Variables

- The condition $\sum_{i=1}^N \xi_i \leq \text{Constant}$ bounds the sum $\sum_{i=1}^N \xi_i$.
- Thus, it bounds the total proportional amount by which predictions fall on the wrong side of their margin.
- Since misclassification occur when $\xi_i > 1$ (in this case $y_i f(x_i) < 0$, bounding $\sum_{i=1}^N \xi_i < k$, bounds the **total number of training misclassifications** at k .
- So, for the non-separable case, we have the optimization problem:

$$\max_{w,b} 2C \quad \text{subject to} \quad \frac{1}{\|w\|} y_i (w^T x_i + b) \geq C(1 - \xi_i), i = 1, \dots, N$$
$$\text{with } \xi_i \geq 0, \forall i, \sum_{i=1}^N \xi_i \leq \text{Constant}$$

Slack Variables

- Similar to the separable case, we define $C = \frac{1}{\|w\|}$ and rewrite the above maximization problem in the equivalent form:

$$\min_{w,b} \frac{\|w\|^2}{2}$$

$$\text{subject to } y_i(w^T x_i + b) \geq 1 - \xi_i, i = 1, \dots, N$$

$$\text{with } \xi_i \geq 0, \forall i, \sum_{i=1}^N \xi_i \leq \text{Constant}$$

- We have obtained a quadratic optimization problem with linear constraints. We will solve it using Lagrange multipliers.

Lagrange Multipliers for Slack Variables

- First, one more step: we have seen that the condition $\sum_{i=1}^N \xi_i \leq \text{Constant}$, bounds the number of training misclassifications.
- We can incorporate this condition into the objective function by adding an extra cost for errors:

$$\min_{w,b} \frac{\|w\|^2}{2} + \gamma \sum_{i=1}^N \xi_i$$

subject to $y_i(w^T x_i + b) \geq 1 - \xi_i, i = 1, \dots, N$

with $\xi_i \geq 0, \forall i$

here, γ is a parameter to be chosen by the user. A larger γ corresponds to assigning a higher penalty to errors.

Lagrange Multipliers for Slack Variables

$$\begin{aligned} \min_{w,b} \quad & \frac{\|w\|^2}{2} + \gamma \sum_{i=1}^N \xi_i \\ \text{subject to} \quad & y_i(w^T x_i + b) \geq 1 - \xi_i, i = 1, \dots, N \\ & \text{with } \xi_i \geq 0, \forall i \end{aligned}$$

- Introducing the Lagrange multipliers α_i and μ_i (one for each constraint), gives the following Lagrange (primal) function:

$$L_p = \frac{1}{2} \|w\|^2 + \gamma \sum_{i=1}^N \xi_i - \sum_{i=1}^N \alpha_i [y_i(w^T x_i + b) - (1 - \xi_i)] - \sum_{i=1}^N \mu_i \xi_i$$

- Our objective is:

$$\min_{w,b,\xi_i} L_p$$

Lagrange Multipliers for Slack Variables

$$\min_{w,b} \frac{\|w\|^2}{2}$$

$$\text{subject to } y_i(w^T x_i + b) \geq 1 - \xi_i, i = 1, \dots, N$$

$$\xi_i \geq 0, \forall i$$

$$\sum_{i=1}^N \xi_i \leq \text{Constant}$$



$$L_p = \frac{1}{2} \|w\|^2 + \gamma \sum_{i=1}^N \xi_i - \sum_{i=1}^N \alpha_i [y_i(w^T x_i + b) - (1 - \xi_i)] - \sum_{i=1}^N \mu_i \xi_i$$

Lagrange Multipliers Solution

- Setting the respective derivatives to zero gives:

$$\frac{\partial L_p}{\partial w} = w - \sum_{i=1}^N \alpha_i y_i x_i = 0 \quad \Rightarrow w = \sum_{i=1}^N \alpha_i y_i x_i \quad (3)$$

$$\frac{\partial L_p}{\partial b} = - \sum_{i=1}^N \alpha_i y_i = 0 \quad \Rightarrow \sum_{i=1}^N \alpha_i y_i = 0 \quad (4)$$

$$\frac{\partial L_p}{\partial \xi_i} = \gamma - \alpha_i - \mu_i, \forall i \quad \Rightarrow \alpha_i = \gamma - \mu_i, \forall i \quad (5)$$

along with the positivity constraints $\alpha_i, \mu_i, \xi_i \geq 0, \forall i$

- Substituting Eq. (3), (4), (5) in L_p , we obtain the so called **dual objective function**:

$$L_D = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j x_i^T x_j$$

where L_D gives a lower bound on the objective function $\frac{1}{2} \|w\|^2 + \gamma \sum_{i=1}^N \xi_i$

Deriving the Dual Form

$$\begin{aligned}
 & \frac{1}{2} \left(\sum_{i=1}^N \alpha_i y_i x_i \right)^\top \left(\sum_{j=1}^N \alpha_j y_j x_j \right) + \gamma \sum_{i=1}^N \xi_i - \sum_{i=1}^N \alpha_i y_i x_i^\top w^\top - \sum_{i=1}^N \alpha_i y_i b + \sum_i \alpha_i (1 - \xi_i) - \sum_{i=1}^N \mu_i \xi_i \\
 &= \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j x_i x_j - \sum_i \sum_j \alpha_i \alpha_j y_i y_j x_i x_j + \gamma \sum_{i=1}^N \xi_i - \underbrace{b \sum_{i=1}^N \alpha_i y_i}_{=0} + \sum_i \alpha_i (1 - \xi_i) - \sum_{i=1}^N \mu_i \xi_i \\
 &= -\frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j x_i x_j + \gamma \sum_{i=1}^N \xi_i + \sum_{i=1}^N \alpha_i - \sum_i \alpha_i \xi_i - \sum_i \mu_i \xi_i \\
 &= \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j x_i x_j + \underbrace{\sum_i (\overbrace{\gamma - \mu_i}^{=\alpha_j}) \xi_i}_{=0} - \sum_i \alpha_i \xi_i \\
 &= \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j x_i x_j
 \end{aligned}$$

Lagrange Multipliers Solution

- Thus: the solution is obtained by maximizing L_D w.r.t the α_i , subject to:

$$\sum_{i=1}^N \alpha_i y_i = 0, \quad 0 \leq \alpha_i \leq \gamma$$

- The solution must satisfy the conditions:

- $w = \sum_{i=1}^N \alpha_i y_i x_i$ (6)

- $\sum_{i=1}^N \alpha_i y_i = 0$ (7)

- $\alpha_i = \gamma - \mu_i, \forall i$ (8)

- $\alpha_i [y_i (w^T x_i + b) - (1 - \xi_i)] = 0, \forall i$ (9)

- $\mu_i \xi_i = 0, \forall i$ (10)

- $y_i (w^T x_i + b) - (1 - \xi_i) \geq 0, \forall i$ (11)

Lagrange Multipliers Solution

- From (6), the solution is $w = \sum_{i=1}^N \alpha_i y_i x_i$.
- From (9), $\alpha_i > 0$ when constraint (11) is exactly met.
- The points (x_i) with $\alpha_i > 0$ are the **SUPPORT VECTORS**. Two types:
 - Those for which $\xi_i = 0$: they lie on the edge of the margin. From (8) and (10): $0 < \alpha_i < \gamma$
 - Those for which $\xi_i > 0$: they have $\alpha_i = \gamma$ and they lie on the wrong side of their margin.
- To estimate b , we can use (9) with any of the support vectors with $\xi_i = 0$.

$$w = \sum_{i=1}^N \alpha_i y_i x_i \quad (6)$$

$$\sum_{i=1}^N \alpha_i y_i = 0 \quad (7)$$

$$\alpha_i = \gamma - \mu_i, \forall i \quad (8)$$

$$\alpha_i [y_i (w^T x_i + b) - (1 - \xi_i)] = 0, \forall i \quad (9)$$

$$\mu_i \xi_i = 0, \forall i \quad (10)$$

$$y_i (w^T x_i + b) - (1 - \xi_i) \geq 0, \forall i \quad (11)$$

Lagrange Multipliers Solution

- Once we have w and b , the decision function can be written as:

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

- The tuning parameter of this procedure is γ . Its optimal value can be estimated via cross validation.

Recast as Unconstrained Optimization Problem

- A constrained optimization problem over w and ξ

$$\min_{w,b} \frac{\|w\|^2}{2} + \gamma \sum_{i=1}^N \xi_i$$

Subject to: $y_i(w^T x_i + b) \geq 1 - \xi_i, i = 1, \dots, N$

- The constraint can be written more concisely as: $y_i f(x_i) \geq 1 - \xi_i$ together with $\xi_i \geq 0$ is equivalent to

$$\xi_i = \max(0, 1 - y_i f(x_i))$$

- Hence the learning problem is equivalent to the **unconstrained** optimization problem over w :

$$\min_{w,b} \frac{\|w\|^2}{2} + \gamma \sum_{i=1}^N \max(0, 1 - y_i f(x_i))$$

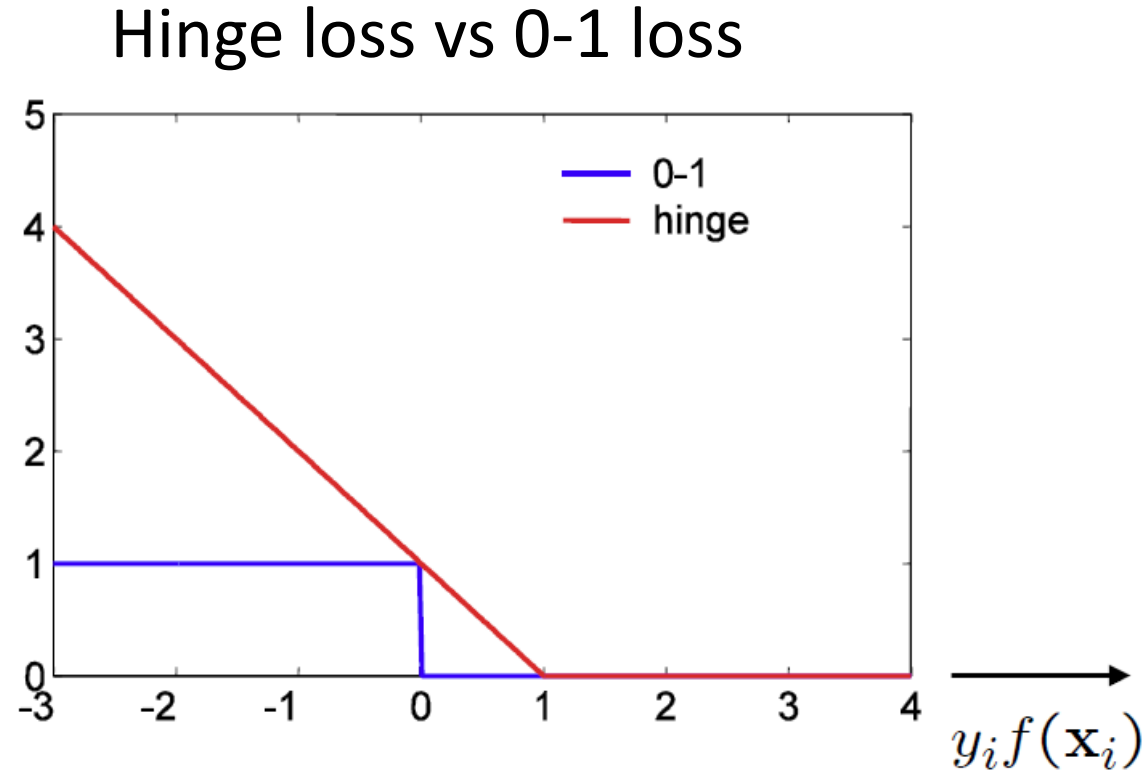
Loss Function

We can consider it as the minimization of a regularized error function.

$$\min_{w,b} \frac{\|w\|^2}{2} + \gamma \sum_{i=1}^N \underbrace{\max(0, 1 - y_i f(x_i))}_{\text{Hinge Loss}}$$

- $y_i f(x_i) > 1$: points outside margin. No contribution to loss.
- $y_i f(x_i) = 1$: points on margin. No contribution to loss (hard margin case)
- $y_i f(x_i) < 1$: points violates margin constraints. Contribute to loss.
- Optimization: solving the Quadratic Programming Problem
 - Sequential minimal optimization (SMO) algorithm
 - Stochastic sub-gradient descent algorithms

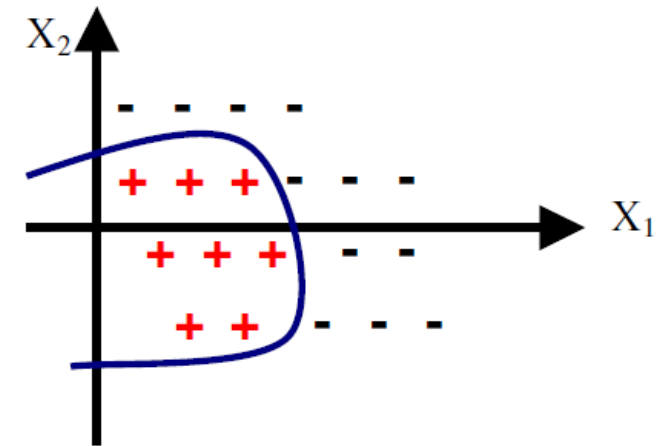
Hinge Loss



- A convex approximation to the 0-1 loss.
- It is an upper bound on the 0-1 loss.
- Not differentiable, we need to compute the sub-gradient.

Non-linear SVM

- **Problem:** SVM represented with **a linear function** have very limited representational power, and could not be very useful in practical classification problems.
- How to generalize the above methods to solve the case where the decision function is non-linear?
- **Good news:** With a slight modification, SVM could solve highly nonlinear classification problems!!
- **Assumption:** Suppose that dataset D is nonlinearly separable in the original attribute space. The attribute space can be **transformed** into a new attribute space where D is linearly separable!



Non-linear SVM

- It turns out that the generalization to a nonlinear boundary can be accomplished in a straightforward way using a **simple mathematical trick!**

- One major observation on the dual objective function:

$$\begin{aligned} L_D &= \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j x_i^T x_j \\ &= \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \text{ < } x_i, x_j \text{ >} \end{aligned}$$

Dot product

- The only way that the data appear in the training problem is in the form of dot products.

Non-linear SVM

- How about the solution function?
- From $w = \sum_{i=1}^N \alpha_i y_i x_i$, the solution function can be written as:

$$\begin{aligned} f(x) &= w^T x + b \\ &= \sum_i^{N_s} \alpha_i y_i x_i^T x + b \\ &= \sum_i^{N_s} \alpha_i y_i \langle x_i, x \rangle + b \end{aligned}$$

where N_s is the number of support vectors.

- In the solution function, the data also appear in the form of dot products where the (x_i) s are the support vectors.

Non-linear SVM

- Now, suppose we first map the data to some high dimension Euclidean space using a mapping Φ (usually $h > d$):
$$\Phi: \mathbb{R}^d \rightarrow \mathbb{R}^h$$
- The idea is to enlarge the input space to achieve better training class separation.
- In general, linear boundaries in the **enlarged space** translate to **nonlinear boundaries** in the original space (true for any nonlinear mapping).

Mapping

- Then, we **compute the largest margin hyperplane in the new space** \mathbb{R}^h .
- Of course, the training algorithm would only depend on the data through dot products in \mathbb{R}^h , i.e., $\langle \Phi(x_i), \Phi(x_j) \rangle$, where $\Phi(x_i) \in \mathbb{R}^h$.
- Suppose we have a function (called **kernel function**) K that computes such dot products in the **transformed space**:
$$K(x_i, x_j) = \langle \Phi(x_i), \Phi(x_j) \rangle$$
- Then: all we need in the training algorithm is K , and we would never need to explicitly even know what Φ is.
- **Resulting procedure**: replace $\langle \Phi(x_i), \Phi(x_j) \rangle$ with $K(x_i, x_j)$ everywhere in the training algorithm.

Mapping

- The algorithm constructs a **linear** support vector machine in \mathbb{R}^h .
- It achieves the objective in roughly the same amount of time it would take to train on the original data.

- How can we use such a machine? In test phase, given the test points x :

$$f(x) = \sum_i^{N_s} \alpha_i y_i \langle x_i, x \rangle + b = \sum_i^{N_s} \alpha_i y_i K(x_i, x) + b$$

where x_i are the support vectors and N_s is the number of support vectors.

- **Kernel trick**: with the kernel function K , we can work with vectors in input space, without even knowing the mapping function Φ .

Example: Kernel Functions

Example: an allowed kernel for which we can construct the mapping Φ :

- Training data are vectors in \mathbb{R}^2 .
- Suppose we choose $K(x_i, x_j) = (\langle x_i, x_j \rangle)^2$. We can find a mapping $\Phi: \mathbb{R}^2 \rightarrow \mathbb{R}^h$, such that $(\langle x_i, x_j \rangle)^2 = \langle \Phi(x_i), \Phi(x_j) \rangle$
- One such mapping is: $\Phi: \mathbb{R}^2 \rightarrow \mathbb{R}^3$ defined as

$$\Phi(x) = \begin{pmatrix} x_1^2 \\ \sqrt{2}x_1x_2 \\ x_2^2 \end{pmatrix}$$

where $x = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$

Example: Kernel Functions

- We can verify that this is indeed the case:

$$\begin{aligned} K(x, y) &= (\langle x, y \rangle)^2 = (x_1 y_1 + x_2 y_2)^2 \\ &= x_1^2 y_1^2 + x_2^2 y_2^2 + 2x_1 y_1 x_2 y_2 \end{aligned}$$

$$\langle \Phi(x), \Phi(y) \rangle = \Phi(x)^T \Phi(y)$$

$$\begin{aligned} &= (x_1^2, \sqrt{2}x_1 x_2, x_2^2) \begin{pmatrix} y_1^2 \\ \sqrt{2}y_1 y_2 \\ y_2^2 \end{pmatrix} \\ &= x_1^2 y_1^2 + x_2^2 y_2^2 + 2x_1 y_1 x_2 y_2 \end{aligned}$$

- Note: in general, the mapping Φ and the space \mathbb{R}^h are **not unique** for a given kernel.

Example: Kernel Functions

- You can verify the following two mappings Φ also satisfy $K(x, y) = \langle \Phi(x), \Phi(y) \rangle$ for kernel given above.
- Example 1: $\Phi: \mathbb{R}^2 \rightarrow \mathbb{R}^3$

$$\Phi(x) = \frac{1}{\sqrt{2}} \begin{pmatrix} x_1^2 - x_2^2 \\ 2x_1x_2 \\ x_1^2 + x_2^2 \end{pmatrix}$$

- Example 2: $\Phi: \mathbb{R}^2 \rightarrow \mathbb{R}^4$

$$\Phi(x) = \frac{1}{\sqrt{2}} \begin{pmatrix} x_1^2 \\ x_1x_2 \\ x_1x_2 \\ x_2^2 \end{pmatrix}$$

Properties of a **Valid** Kernel

- To ensure that a mapping Φ and an expansion $K(x, y) = \langle \Phi(x), \Phi(y) \rangle$ exist, the mathematical properties that a function K have been studied and are well known as **Mercer theorem**.
- Two popular choices for K are:
 - d^{th} degree polynomial: $K(x, y) = (1 + \langle x, y \rangle)^d$
 - Radial Basis Function (RBF) Kernel: $K(x, y) = e^{\frac{-\|x-y\|^2}{2\sigma^2}}$
- The best choice of a kernel for a given problem is still a research issue; (e.g., Latent Semantic Kernel for document classification)

Readings

1. <https://faculty.sites.iastate.edu/jia/files/inline-files/lagrange-multiplier.pdf>
2. <http://cs229.stanford.edu/summer2020/cs229-notes3.pdf>

Summary of Today's Lecture

- SVM for Linear Separable Case
- Non-separable Case, Penalties
- Non-linearity, Kernels