# CS 559 Machine Learning
## Support Vector Machines

Yue Ning
Department of Computer Science
Stevens Institute of Technology

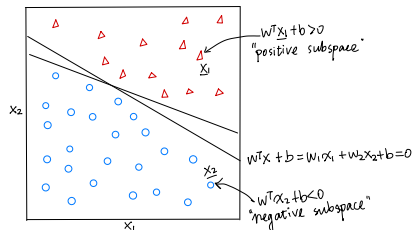Vector algebra, Formulation, Margin

Non-separable Case, Penalties

Non-linearity, Kernels

▶ Linear classifiers construct
<u>linear decision boundaries</u>(hyperplanes) that try to separate
the data into different classes as well as possible.



▶ <u>Classification Rule</u> (**the Perceptron model**):

$$\text{Input: } \mathbf{x} \in \mathbb{R}^d \qquad \text{Output: } \text{sign}(\mathbf{w}^\top \mathbf{x} + b)$$

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

# Perceptron Learning Algorithm as Gradient Descent

<u>Objective</u>: Find a separating hyperplane that correctly classifier all input patterns.

- ▶ There are two types of error:

$$y_i = 1 \text{ and } \mathbf{w}^\top \mathbf{x}_i + b < 0$$

$$y_i = -1 \text{ and } \mathbf{w}^\top \mathbf{x}_i + b > 0$$

- ▶ Thus, for all misclassified points:

$$y_i(\mathbf{w}^\top \mathbf{x}_i + b) < 0$$

- ▶ To reduce the number of misclassified points, the goal is to minimize:

$$D(\mathbf{w}, b) = -\sum_{i \in M} y_i(\mathbf{w}^\top \mathbf{x}_i + b)$$

  where M indexes the set of misclassified points.

# Perceptron Learning Algorithm as Gradient Descent

<u>Objective</u>: Find a separating hyperplane that correctly classifies all input patterns.

▶ To minimize $D(\mathbf{w}, b)$ we can perform gradient descent over the surface represented by $D(\mathbf{w}, b)$ in parameter space. We iteratively move along the opposite direction of the gradient till a minimum is reached.

▶ The gradient is given by:

$$\frac{\partial D(\mathbf{w}, b)}{\partial \mathbf{w}} = -\sum_{i \in M} y_i \mathbf{x}_i \qquad \frac{\partial D(\mathbf{w}, b)}{\partial b} = -\sum_{i \in M} y_i$$

▶ Update rule for parameters ("steepest descent"):

$$\mathbf{w}' = \mathbf{w} + \sum_{i \in M} y_i \mathbf{x}_i \qquad b' = b + \sum_{i \in M} y_i$$

▶ In practice, "Stochastic Gradient Descent (SGD)" is used: rather than computing the sum of the gradient contributions of each $x_i$, followed by a step in the negative gradient direction, a step is taken after each observation is visited. The resulting update rule for parameters is:

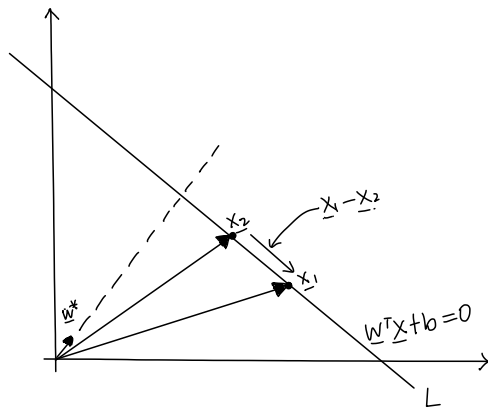$$\mathbf{w}' = \mathbf{w} + y_i \mathbf{x}_i$$

$$b' = b + y_i$$

▶ When the data are linearly separable, there are many solutions, and which one is found depends on the starting values of the parameters.

▶ 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 separable, the algorithm will **not converge**, and cycles develop. The cycles can be long, and therefore hard to detect. We can provide a solution to the first problem by adding additional constraints to the separating hyperplane we want to find!

▶ Hyperplane $L$: $f(\mathbf{x}) = (\mathbf{w}^\top\mathbf{x} + b) = 0$, when $\mathbf{x} \in \mathbb{R}^2$, $f(\mathbf{x})$ is a line.

# Some Vector Algebra - property 1

▶ Consider any two points $x_1, x_2$, lying on hyperplane L:
$$\begin{aligned}\mathbf{w}\mathbf{x}_1 + b &= 0 \\ \mathbf{w}\mathbf{x}_2 + b &= 0\end{aligned} \to \mathbf{w}^\top(\mathbf{x}_1 - \mathbf{x}_2) = 0$$

▶ Since $\mathbf{w}^\top(\mathbf{x}_1 - \mathbf{x}_2) = \mathbf{w} \cdot (\mathbf{x}_1 - \mathbf{x}_2) = 0$, the two vectors $\mathbf{w}$ and $\mathbf{x}_1 - \mathbf{x}_2$ are orthogonal vectors.

$$\mathbf{w}^* = \frac{\mathbf{w}}{||\mathbf{w}||}$$

is the vector normal to the surface of L.

▶ <u>Note 1</u>: All vectors here are column vectors.

▶ <u>Note 2</u>: Dot product (inner product) of two vectors $\mathbf{a} \cdot \mathbf{b} = \mathbf{a}^\top \mathbf{b} = ||\mathbf{a}|| \times ||\mathbf{b}|| \times \cos\alpha$ where $\alpha$ is the angle between $a$ and $b$.
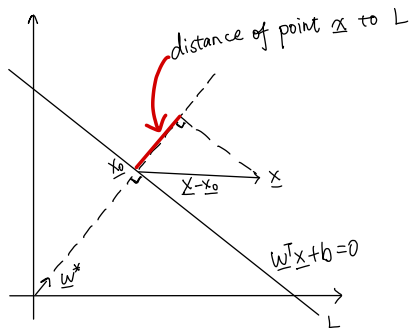
▶ For any point $\mathbf{x}_0$ on $L$:

$$\mathbf{w}^\top \mathbf{x}_0 + b = 0$$
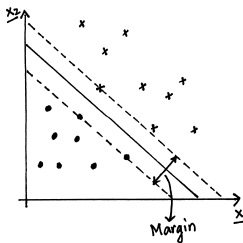
Thus:

$$\mathbf{w}^\top \mathbf{x}_0 = -b$$

The signed distance of any point $\mathbf{x}$ to $L$ is:

$$(\mathbf{w}^*)^\top(\mathbf{x} - \mathbf{x}_0) = \frac{\mathbf{w}^\top}{||\mathbf{w}||}(\mathbf{x} - \mathbf{x}_0) =$$

$$\frac{1}{||\mathbf{w}||}(\mathbf{w}^\top\mathbf{x} - \mathbf{w}^\top\mathbf{x}_0) = \frac{1}{||\mathbf{w}||}(\mathbf{w}^\top\mathbf{x} + b)$$

# Largest Margin Hyperplanes

<u>Goal</u>: Find the hyperplane that separates the two classes and maximizes the distance <u>to the closest points from either class</u>.

▶ Such distance is called <span style="color:red">margin</span>.



▶ 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.

For two classes:

$$(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), ..., (\mathbf{x}_N, y_N)$$

$$\mathbf{x}_i \in \mathbf{R}^d$$
$$y_i = \{-1, +1\}$$

We need to formalize the largest margin criterion.

# Formulation

Consider the following optimization problem:

$$\max_{\mathbf{w}, b} 2C$$
$$\text{subject to } \frac{1}{||\mathbf{w}||} y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq C \quad i = 1, ..., N$$

▶ Remember Property 3:
  The signed distance of any point $\mathbf{x}$ to L is:

$$\frac{1}{||\mathbf{w}||}(\mathbf{w}^\top \mathbf{x}_i + b)$$

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

▶ We seek the largest such C and associated parameters.

▶ We can rewrite the above conditions as:

$$y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq C ||\mathbf{w}||$$

# Formulation

▶ We can rewrite the above conditions as:

$$y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq C||\mathbf{w}||$$

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

▶ The original maximization problem is equivalent to:

$$\min_{\mathbf{w},b} \frac{1}{2}||\mathbf{w}||^2$$

$$\text{subject to } y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 \quad i = 1, ..., N$$

▶ The constraints define an empty margin around the linear decision boundary of thickness $\frac{2}{||\mathbf{w}||}$. We choose $\mathbf{w}$, $b$ to maximize its thickness.

▶ This is a quadratic (convex) optimization problem subject to linear constraints and there is a unique minimum

# Lagrange Multipliers

▶ We introduce Lagrange multipliers, which allow to take the constraints within the function to be minimized. Two reasons for doing this:

1. The constraints will be replaced by constraints on the Lagrange multipliers themselves, which are easier to handle.
2. Training data will only appear in the form of dot products between vectors (crucial for nonlinear case).

▶ We introduce the Lagrange multipliers $\alpha_i \geq 0$, $i = 1, ..., N$, one for each of the inequality constraints.

▶ Recall the rule: for constraints of the form $F_i \geq 0$, the constraint equations are multiplied by Lagrange multipliers and subtracted from the objective function, to form the Lagrangian.

# Lagrange Multipliers - Primal Form

▶ We then obtain the Lagrangian: (also called "primal form"):

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

▶ We must now <u>minimize $L_p$</u> with respect to $\mathbf{w}$ and $b$:

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

▶ This indicates that this is the <u>primal form</u> of the optimization problem.

▶ We will actually solve the optimization problem by solving for the dual of the original problem

# Dual Form

- The solution to the dual form provides a **lower bound** to the solution of the primal form
- What is the <u>dual form</u>?

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

- Setting the derivatives to zero gives:

$$\frac{\partial L_p}{\partial \mathbf{w}} = \mathbf{w} - \sum_{i=1}^{N} \alpha_i y_i \mathbf{x}_i = 0 \Rightarrow \mathbf{w} = \sum_{i=1}^{N} \alpha_i y_i \mathbf{x}_i \qquad (1)$$

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

# Dual Form

Substituting eq. 1 and 2 in $L_p$ gives:

$$L_D = \frac{1}{2}\Big(\sum_{i=1}^{N}\alpha_i y_i \mathbf{x}_i\Big)\Big(\sum_{k=1}^{N}\alpha_k y_k \mathbf{x}_k\Big) - \sum_{i=1}^{N}\alpha_i\Big[y_i\big(\mathbf{x}_i^\top\big(\sum_{k=1}^{N}\alpha_k y_k \mathbf{x}_k\big) + b\big) - 1\Big]$$

$$= \frac{1}{2}\sum_{i=1}^{N}\sum_{k=1}^{N}\alpha_i\alpha_k y_i y_k \mathbf{x}_i^\top \mathbf{x}_k - \sum_{i=1}^{N}\sum_{k=1}^{N}\alpha_i\alpha_k y_i y_k \mathbf{x}_i^\top \mathbf{x}_k - b\underbrace{\sum_{i=1}^{N}\alpha_i y_i}_{=0} + \sum_{i=1}^{N}\alpha_i$$

$$= \sum_{i=1}^{N}\alpha_i - \frac{1}{2}\sum_{i=1}^{N}\sum_{k=1}^{N}\alpha_i\alpha_k y_i y_k \mathbf{x}_i^\top \mathbf{x}_k$$

subject to $\alpha_i \geq 0$

▶ Solution is obtained by maximizing $L_D$ with respect to the $\alpha_i$.

▶ Solution must satisfy the conditions:

$$\mathbf{w} = \sum_{i=1}^{N} \alpha_i y_i \mathbf{x}_i$$

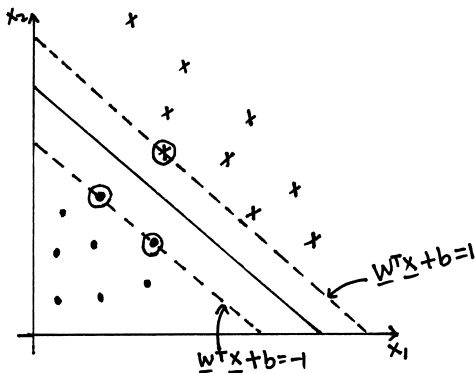$$\sum_{i=1}^{N} \alpha_i y_i = 0$$

$$\alpha_i \geq 0$$

$$\alpha_i \left[ y_i (\mathbf{w}^\top \mathbf{x}_i + b) - 1 \right] = 0 \quad \forall i = 1, ..., N$$

# Dual Form

- If $\alpha_i > 0$, then $y_i(\mathbf{w}^\top \mathbf{x}_i + b) = 1$, that is $\mathbf{x}_i$ is on the boundary of the margin.

- If $y_i(\mathbf{w}^\top \mathbf{x}_i + b) > 1$, $\mathbf{x}_i$ is not on the boundary of the margin, and $\alpha_i = 0$

# Dual Form

- The solution vector $\mathbf{w}$ is: $\mathbf{w} = \sum_{i=1}^{N} \alpha_i y_i \mathbf{x}_i$.
- $\mathbf{x}_i$ are <u>SUPPORT VECTORS</u> when $\alpha_i > 0$. We have three support vectors in the above example.
- To obtain the value of $b$: solve $\alpha_i \big[ y_i(\mathbf{w}^\top \mathbf{x}_i + b) - 1 \big] = 0$ for any of the support vectors.
- The largest margin hyperplane gives a function:
  $f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} + b$ for classifying new observations
  $\hat{y} = \text{sign}(f(\mathbf{x}))$

▶ The support vectors are the critical elements of the training set. They lie closet to the decision boundary.

▶ Only the support vectors affect the solution

▶ However, the identification of the support vectors requires the use of all the training data.

# Summary so far

- Training data: $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), ..., (\mathbf{x}_N, y_N)$,
  $\mathbf{x}_i \in \mathbb{R}^d, y_i \in \{-1, +1\}$
- When the two classes are linearly separable, we can find a function $f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} + b$ with $y_i f(\mathbf{x}_i) > 0 \ \forall i$
- In particular, we can find the hyperplane that creates the largest margin between the training points.
- The optimization problem captures this concept

$$\max_{\mathbf{w}, b} 2C$$
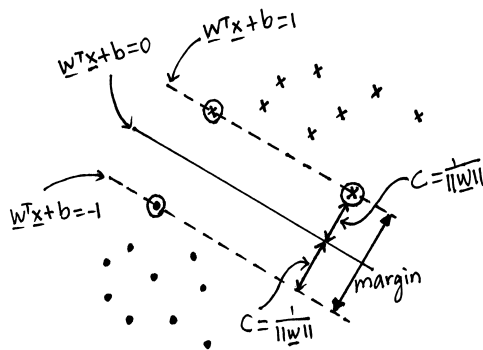$$\text{subject to } \frac{1}{||\mathbf{w}||} y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq C \quad i = 1, ..., N$$

- It can be more conveniently rewritten as below where $C = \frac{1}{||\mathbf{w}||}$

$$\min_{\mathbf{w}, b} \frac{1}{2} ||\mathbf{w}||^2$$
$$\text{subject to } y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 \quad i = 1, ..., N$$

# 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: $\frac{1}{||\mathbf{w}||} y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq C \quad i = 1, ..., N$.

▶ To achieve this goal, we define N <u>slack variables</u>:

$$\xi_1, \xi_2, ..., \xi_N$$

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

$$\frac{1}{||\mathbf{w}||} y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq C(1 - \xi_i) \quad i = 1, ..., N$$

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

▶ <u>Idea of the formulation</u>: $\xi_i$ is the proportional amount by which the prediction $f(\mathbf{x}_i)$ is on the wrong side of the margin.

▶ Note: $C(1 - \xi_i) = C - C\xi_i$

# Slack Variables

▶ A geometric perspective:



▶ The points $(\mathbf{x}_1^*, \mathbf{x}_2^*, \mathbf{x}_3^*, \mathbf{x}_4^*, \mathbf{x}_5^*)$ are on the wrong side of their margin.

▶ Point $\mathbf{x}_i^*$ is on the wrong side of its margin by an amount $C\xi_i$

▶ Point $\mathbf{x}_j^*$ on the correct side have $\xi_j = 0$

▶ Misclassification occurs when $\xi_i > 1 \Rightarrow C(1 - \xi_i) < 0$, e.g., points $\mathbf{x}_3^*$ and $\mathbf{x}_5^*$ are misclassified by the given boundary.

# Slack Variables (2)

▶ The condition $\sum_{i=1}^{N} \xi_i \leq$ constant bounds the sum $\sum_i \xi_i$

▶ 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 $y_i f(\mathbf{x}_i) < 0$), bounding $\sum_{i=1}^{N} \xi_i \leq k$, bounds the total number of training misclassifications at $k$.

▶ So, for the non-separable case, we have the optimization problem:

$$\max_{\mathbf{w}, b} 2C$$

$$\text{subject to: } \frac{1}{||\mathbf{w}||} y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq C(1 - \xi_i) \quad i = 1, ..., N$$

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

▶ As for the separable case, we define $C = \frac{1}{||\mathbf{w}||}$ and rewrite the above maximization problem in the equivalent form:

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

$$\text{subject to: } y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 - \xi_i \quad i = 1, ..., N$$

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

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

▶ First, one more step: we have seen that the condition $\sum_i \xi_i \leq$ constant, bounds the number of training misclassifications.

▶ We can incorporate this condition into the objective function by adding an extra cost for errors:

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

subject to: $y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 - \xi_i \quad i = 1, ..., N$

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

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

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

$$\text{subject to: } y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 - \xi_i \quad i = 1, ..., N$$

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

▶ Introducing the <u>Lagrange multipliers</u> $\alpha_i$ and $\mu_i$ (one for each constraint), gives the following Lagrange (primal) function:

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

▶ Our objective is: $\min_{\mathbf{w},b,\xi_i} L_p$

▶ Setting the respective derivatives to zero gives:

$$\frac{\partial L_p}{\partial \mathbf{w}} = \mathbf{w} - \sum_{i=1}^{N} \alpha_i y_i \mathbf{x}_i = 0 \Rightarrow \mathbf{w} = \sum_{i=1}^{N} \alpha_i y_i \mathbf{x}_i \qquad (3)$$

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

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

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

▶ Substituting eq. 3, 4, 5 in $L_p$, we obtain the so called <u>dual objective function</u>:

$$L_D = \sum_{i=1}^{N} \alpha_i - \frac{1}{2} \sum_{i=1}^{N} \sum_{j=1}^{N} \alpha_i \alpha_j y_i y_j \mathbf{x}_i^\top \mathbf{x}_j$$

where $L_D$ gives a lower bound on the objective function

# Deriving the Dual Form

$$\frac{1}{2} \left( \sum_{i=1}^{N} \alpha_i y_i \mathbf{x}_i \right)^{\top} \left( \sum_{j=1}^{N} \alpha_j y_j \mathbf{x}_j \right) + \gamma \sum_{i=1}^{N} \xi_i - \sum_{i=1}^{N} \alpha_i y_i \mathbf{x}_i^{\top} \mathbf{w}^{\top} - \sum_{i=1}^{N} \alpha_i y_i b + \sum_i \alpha_i (1 - \xi_i) - \sum_{i=1}^{N} \mu_i \xi_i$$

$$= \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \mathbf{x}_i \mathbf{x}_j - \sum_i \sum_j \alpha_i \alpha_j y_i y_j \mathbf{x}_i \mathbf{x}_j + \gamma \sum_{i=1}^{N} \xi_i - b \underbrace{\sum_{i=1}^{N} \alpha_i y_i}_{=0} + \sum_i \alpha_i (1 - \xi_i) - \sum_{i=1}^{N} \mu_i \xi_i$$

$$= -\frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \mathbf{x}_i \mathbf{x}_j + \gamma \sum_{i=1}^{N} \xi_i + \sum_{i=1}^{N} \alpha_i - \sum_i^{N} \alpha_i \xi_i - \sum_i^{N} \mu_i \xi_i$$

$$= \sum_i^{N} \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \mathbf{x}_i \mathbf{x}_j + \underbrace{\sum_i \overbrace{(\gamma - \mu_i)}^{=\alpha_i} \xi_i - \sum_i \alpha_i \xi_i}_{=0}$$

$$= \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \mathbf{x}_i \mathbf{x}_j$$

# Lagrange Multipliers Solution

▶ Thus: the solution is obtained by <u>maximizing $L_D$</u> w.r.t the $\alpha_i$, subject to:

$$\sum_i^N \alpha_i y_i = 0, \quad 0 \le \alpha_i \le \gamma$$

▶ The solution must satisfy the conditions:

$$\mathbf{w} = \sum_i^N \alpha_i y_i \mathbf{x}_i \tag{6}$$

$$\sum_i^N \alpha_i y_i = 0 \tag{7}$$

$$\alpha_i = \gamma - \mu_i, \quad \forall i \tag{8}$$

$$\alpha_i \big[ y_i(\mathbf{w}^\top \mathbf{x}_i + b) - (1 - \xi_i) \big] = 0, \quad \forall i \tag{9}$$

$$\mu_i \xi_i = 0 \quad \forall i \tag{10}$$

$$y_i(\mathbf{w}^\top \mathbf{x}_i + b) - (1 - \xi_i) \ge 0 \quad \forall i \tag{11}$$

# Lagrange Multipliers Solution

▶ From Eq.6, the solution is $\mathbf{w} = \sum_i^N \alpha_i y_i \mathbf{x}_i$. From Eq. 9, $\alpha_i > 0$ when constraint 11 are exactly met, i.e., $y_i(\mathbf{w}^\top \mathbf{x}_i + b) - (1 - \xi_i) = 0$

▶ The points ($\mathbf{x}_i$) with $\alpha_i > 0$ are the <u>SUPPORT VECTORS</u>.

▶ We have two kinds of support vectors:
  - ▶ Those for which $\xi_i = 0$: they lie on the edge of the margin. From Eqs. 10 and 8: $0 < \alpha_i < \gamma$
  - ▶ 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 Eq. 9 with any of the support vectors with $\xi_i = 0$.

▶ Once we have $\mathbf{w}$ and $b$, the decision function can be written:

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

▶ The tuning parameter of this procedure is $\gamma$. Its optimal value can be estimated via cross validation.

## Optimization

▶ A constrained optimization problem over **w** and $\xi$

$$\min_{\mathbf{w},b} \frac{||\mathbf{w}||^2}{2} + \gamma \sum_{i=1}^{N} \xi_i \qquad (12)$$

subject to: $y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 - \xi_i \quad i = 1, ..., N$

$$(13)$$

▶ The constraint can be written more concisely as:

$$y_i f(\mathbf{x}_i) \geq 1 - \xi_i \qquad (14)$$

together with $\xi_i > 0$ is equivalent to

$$\xi = max(0, 1 - y_i f(\mathbf{x}_i)) \qquad (15)$$

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

$$\min_{\mathbf{w},b} \frac{||\mathbf{w}||^2}{2} + \gamma \sum_{i=1}^{N} max(0, 1 - y_i f(\mathbf{x}_i)) \qquad (16)$$

$$\min_{\mathbf{w},b} \frac{||\mathbf{w}||^2}{2} + \gamma \sum_{i=1}^{N} \underbrace{\max(0, 1 - y_i f(\mathbf{x}_i))}_{\text{Hinge loss}} \tag{17}$$

- $y_i f(\mathbf{x}_i) > 1$: points outside margin. No contribution to loss
- $y_i f(\mathbf{x}_i) = 1$: points on margin. No contribution to loss (hard margin case)
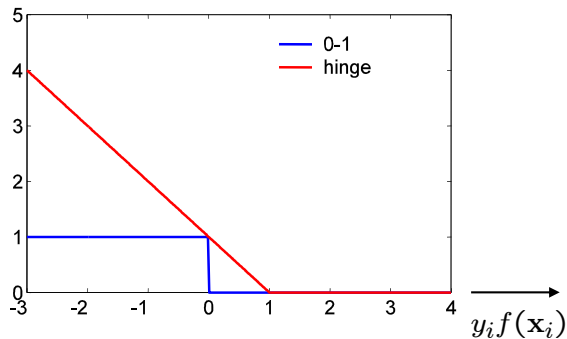- $y_i f(\mathbf{x}_i) < 1$: points violates margin constraints. Contribute to loss.

Figure: Hinge loss vs 0-1 loss

An approximation to the 0-1 loss. Convex?

- ▶ Solving the Quadratic Programming Problems (slow)
- ▶ Use an interior point method that uses Newton-like iterations to find a solution of the Karush–Kuhn–Tucker conditions of the primal and dual problems
- ▶ Platt's sequential minimal optimization (SMO) algorithm
- ▶ Stochastic sub-gradient descent algorithms.

# Non-linear SVM

▶ How can the above methods be generalized to the case where the decision function is not a linear function of the data?

▶ 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 (look at the dual objective function obtained earlier):

$$L_D = \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \mathbf{x}_i \mathbf{x}_j \qquad (18)$$

$$= \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \underbrace{< \mathbf{x}_i, \mathbf{x}_j >}_{\text{dot product}} \qquad (19)$$

▶ The only way in which the data appear in the training problem is in the form of dot products.

▶ How about the solution function?

▶ From $\mathbf{w} = \sum_i^N \alpha_i y_i \mathbf{x}_i$, the solution function can be written as:

$$
\begin{aligned}
f(\mathbf{x}) &= \mathbf{w}^\top \mathbf{x} + b \\
&= \sum_i^{N_s} \alpha_i y_i \mathbf{x}_i^\top \mathbf{x} + b \\
&= \sum_i^{N_s} \alpha_i y_i < \mathbf{x}_i, \mathbf{x} > + b
\end{aligned}
$$

where $N_s$ is the number of support vectors.

$\Rightarrow$Also in the solution function, the data appear in the form of dot products where the ($\mathbf{x}_i$)s are the support vectors.

▶ Now, suppose we first map the data to some high dimension Euclidean space using a mapping $\Phi$:

$$\Phi : \mathbb{R}^d \to \mathbb{R}^h$$

usually $h > d$

▶ The idea is to enlarge the input space to achieve better training class separation.

▶ In general, linear boundaries in the enlarged space translate to nonliear boundaries in the original space (true for any nonlinear mapping $\Phi$)

# 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., $< \Phi(\mathbf{x}_i), \Phi(\mathbf{x}_j) >$ where $\Phi(\mathbf{x}_i) \in \mathbb{R}^h$.

▶ Suppose we have a function (called **kernel function**) K that computes such dot products in the transformed space:

$$K(\mathbf{x}_i, \mathbf{x}_j) = < \Phi(\mathbf{x}_i), \Phi(\mathbf{x}_j) >$$

▶ Then: all we need in the training algorithm is K, and we would never need to explicitly even know what $\Phi$ is.

▶ Resulting procedure: replace $< \Phi(\mathbf{x}_i), \Phi(\mathbf{x}_j) >$ with $K(\mathbf{x}_i, \mathbf{x}_j)$ everywhere in the training algorithm.

▶ Example of K: $K(\mathbf{x}_i, \mathbf{x}_j) = e^{\frac{-||\mathbf{x}_i - \mathbf{x}_j||^2}{2\sigma^2}}$ (Gaussian kernel)

# Mapping

▶ The algorithm constructs a <span style="color:red">linear</span> support vector machine in $\mathbb{R}^h$.

▶ It achieves this objective in roughly the same amount of time it would take to train on the un-mapped data.

▶ How can we use such a machine? In test phase, given the test points $\mathbf{x}$:

$$f(x) = \sum_i^{N_s} \alpha_i y_i <\mathbf{x}_i, \mathbf{x}> + b$$

$$= \sum_i^{N_s} \alpha_i y_i K(\mathbf{x}_i, \mathbf{x}) + b$$

where $\mathbf{x}_i$ are the support vectors and $N_s$ is the number of support vectors.

▶ The fact that, through the kernel function $K$, we can work with vectors in input space, without even knowing the mapping function $\Phi$ is known as the "**kernel trick**"

# Example - kernel functions

Example: an allowed kernel for which we can construct the mapping $\Phi$

- Training data are vectors in $\mathbb{R}^2$.
- Suppose we choose $K(\mathbf{x}_i, \mathbf{x}_j) = (<\mathbf{x}_i, \mathbf{x}_j>)^2$.
- We can find a mapping

$$\Phi : \mathbb{R}^2 \to \mathbb{R}^h$$

such that $(<\mathbf{x}_i, \mathbf{x}_j>)^2 = <\Phi(\mathbf{x}_i), \Phi(\mathbf{x}_j)>$
- One such mapping is:

$$\Phi : \mathbb{R}^2 \to \mathbb{R}^3$$

defined as

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

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

# Example - kernel functions

▶ We can verify that this is indeed the case:

$$K(\mathbf{x}, \mathbf{y}) = (\mathbf{x}, \mathbf{y})^2 = (x_1 y_1 + x_2 y_2)^2$$
$$= x_1^2 y_1^2 + x_2^2 y_2^2 + 2 x_1 y_1 x_2 y_2$$

$$< \Phi(\mathbf{x}), \Phi(\mathbf{y}) > = \Phi(\mathbf{x})^\top \Phi(\mathbf{y})$$
$$= (x_1^2, \sqrt{2} x_1 x_2, x_2^2) \begin{pmatrix} y_1^2 \\ \sqrt{2} y_1 y_2 \\ y_2^2 \end{pmatrix}$$
$$= x_1^2 y_1^2 + x_2^2 y_2^2 + 2 x_1 y_1 x_2 y_2$$

▶ Note: in general neither the mapping $\Phi$ nor the space $\mathbb{R}^h$ are unique for a given kernel.

# Exercise - kernel functions

▶ You can verify the following two mapping $\Phi$ also satisfy $K(\mathbf{x}, \mathbf{y}) = \ <\Phi(\mathbf{x}), \Phi(\mathbf{y}) >$ for kernel given above.

▶ Example 1

$$\Phi : \mathbb{R}^2 \to \mathbb{R}^3$$

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

▶ Example 2

$$\Phi : \mathbb{R}^2 \to \mathbb{R}^4$$

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

- ▶ The mathematical properties that a function $K$ must have so that a mapping $\Phi$ and an expansion $K(\mathbf{x}, \mathbf{y}) = <\Phi(\mathbf{x}), \Phi(\mathbf{y})>$ exist, have been studied and are well known (Mercer theorem)
- ▶ Two popular choices for $K$ are:
  - ▶ $d^{th}$ degree polynomial: $K(\mathbf{x}, \mathbf{y}) = (1 + <\mathbf{x}, \mathbf{y}>)^d$
  - ▶ Radial basis: $K(\mathbf{x}, \mathbf{y}) = e^{\frac{-||\mathbf{x}_i - \mathbf{x}_j||^2}{2\sigma^2}}$

# Importance of Parameter $\gamma$ for Non-linear SVMs

- Perfect separation of training data can in general be achieved in the enlarged space $\mathbb{R}^h$.

- Such perfect separation leads to the danger of <u>overfitting</u> the data. As a consequence, the classifier will generalize poorly.

- A proper setting of $\gamma$ allows to avoid overfitting.

- Lets look at the objective function minimized by an SVM: $\frac{1}{2}||\mathbf{w}||^2 + \gamma \sum_i^N \xi_i$ where $\gamma$ is the penalty factor for errors.

- Large $\gamma \rightarrow$ discourage any positive $\xi_i \rightarrow$ tendency to overfit the data $\rightarrow$ highly complicated decision boundary in input space.

- Small $\gamma \rightarrow$ encourage a small value of $||\mathbf{w}|| \rightarrow$ larger margin $\rightarrow$ more data on the wrong side of their margin $\rightarrow$ smoother decision boundary in input space.

- In practice: we need to tune $\gamma$ so to achieve best test error performance.

# Limitations and Open Issues

- Choice of the kernel
  - performance may depend on the chosen kernel, and on the values of associated parameters;
  - the best choice of a kernel for a given problem is still a research issue; (e.g., Latent Semantic Kernel for document classification)
- Speed and size for training and testing
  - Training: the evaluation of the dual objective function requires the computation of all dot products $< \mathbf{x}_i, \mathbf{x}_j > \Rightarrow$ time complexity $O(N^2 d)$ where $N$ is the number of training data and $d$ is the dimension.
  - Testing: need to evaluate: $f(\mathbf{x}) = \sum_i^{N_s} \alpha_i y_i \mathbf{x}_i^\top \mathbf{x} + b$.
    - $\Rightarrow$ time complexity $O(M \cdot N_s)$ where M is the number of operation required to evaluate the kernel.
    - $\Rightarrow$ time complexity $O(d \cdot N_s)$ when K is a RBF kernel, $M = O(d)$.
- Extension to multiple classes

## Acknowledgements