

Support Vector Machine II

Yanyan Lan

Institute of Computing Technology, Chinese Academy of Sciences

lanyanyan@ict.ac.cn

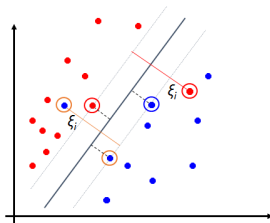
November 19, 2017

Support Vectors in Non-Separable Case

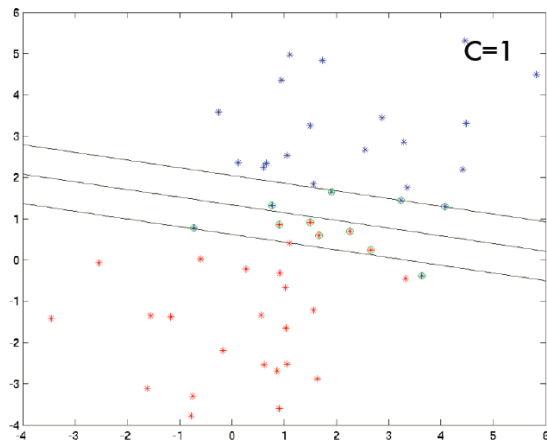
- We call the training examples $(x^{(i)}, y^{(i)})$ as the **support vectors**, if the corresponding $\alpha_i^* > 0$.
- Recall the KKT dual complementarity condition $\alpha_i^* g_i(\omega^*) = 0$, and $\eta_i^* \xi_i = 0$. That is:

$$\alpha_i^* [y^{(i)}(w^{*T} x^{(i)} + b) - 1 + \xi_i] = 0, \quad \eta_i^* \xi_i = 0.$$

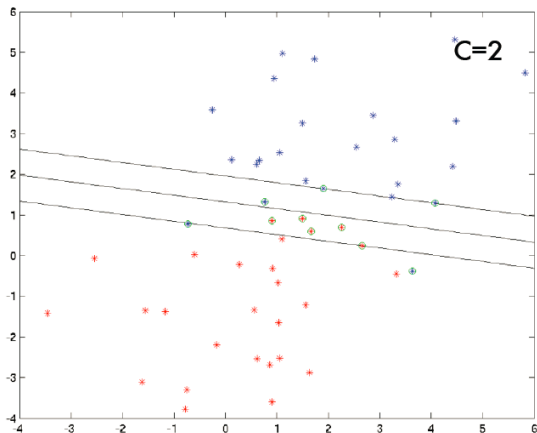
- 1 If $\alpha_i^* = 0$, we have $y^{(i)}(w^{*T} x^{(i)} + b) \geq 1$. (in the correct sides)
- 2 If $0 < \alpha_i^* < C$, we have $y^{(i)}(w^{*T} x^{(i)} + b) = 1$. (on the boundary)
- 3 If $\alpha_i^* = C$, we have $y^{(i)}(w^{*T} x^{(i)} + b) \leq 1$. (in the wrong sides)



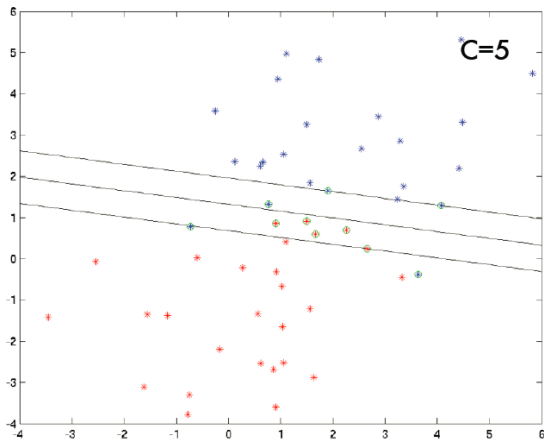
Parameter C



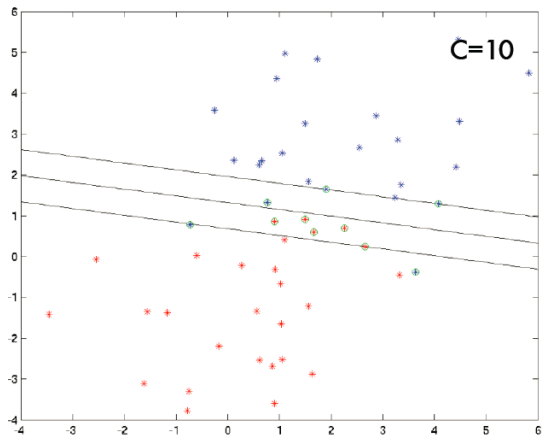
Parameter C



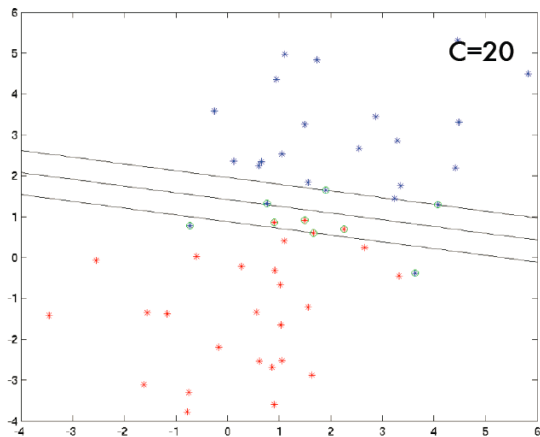
Parameter C



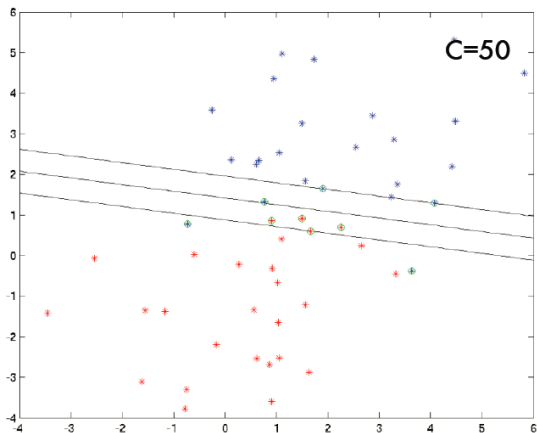
Parameter C



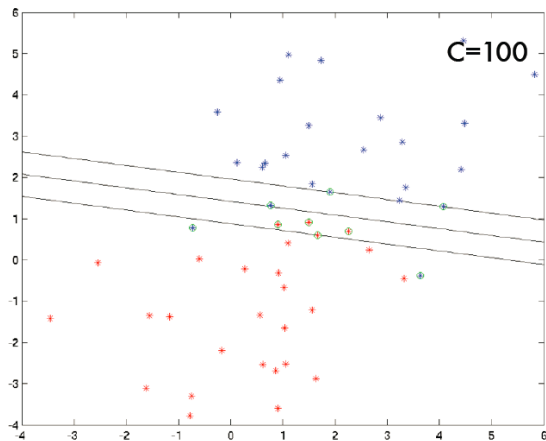
Parameter C



Parameter C



Parameter C

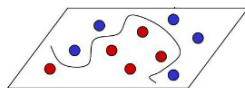


Intuitions of Kernel I: Non-Linear Support Vector Machine

- The **primal** optimization problem:

$$\begin{aligned} \min_{\omega, b} \quad & \frac{1}{2} \|\omega\|_2^2 + C \sum_{i=1}^n \xi_i \\ \text{s.t.} \quad & y^{(i)}(\omega^T x^{(i)} + b) \geq 1 - \xi_i, \quad i = 1, \dots, n, \\ & \xi_i \geq 0, \quad i = 1, \dots, n. \end{aligned}$$

- Nonlinear Separable: if we can use a hypersurface to correctly separate the positive and negative examples.
- But the nonlinear problem is usually not easy to solve. Therefore, a typical method is to transform the nonlinear problem to a linear problem. Then we can use linear techniques to solve it.



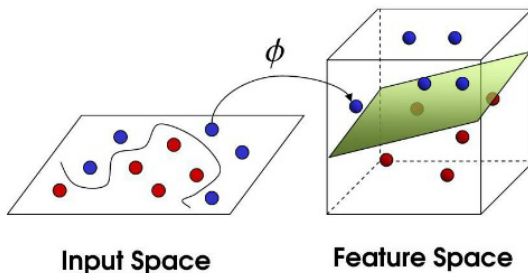
Input Space

Intuitions of Kernel I: Non-Linear Support Vector Machine

- Transformation $z = \phi(x)$.
- The **primal** optimization problem:

$$\min_{\omega, b} \frac{1}{2} \|\omega\|_2^2 + C \sum_{i=1}^n \xi_i$$

$$s.t. \quad y^{(i)}(\omega^T z^{(i)} + b) \geq 1 - \xi_i, \quad i = 1, \dots, n, \\ \xi_i \geq 0, \quad i = 1, \dots, n.$$



Intuitions of Kernel II: Dual Form of Support Vector Machine

- The **dual** optimization problem:

$$\begin{aligned} \max_{\alpha} \quad & \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n y^{(i)} y^{(j)} \alpha_i \alpha_j (x^{(i)})^T x^{(j)} \\ \text{s.t.} \quad & 0 \leq \alpha_i \leq C, i = 1, \dots, n, \\ & \sum_{i=1}^n \alpha_i y^{(i)} = 0. \end{aligned}$$

- Hyperplane:

$$\omega^* = \sum_{i=1}^n \alpha_i^* y^{(i)} x^{(i)}, \quad b^* = y^{(j)} - \sum_{i=1}^n \alpha_i^* y^{(i)} (x^{(i)})^T x^{(j)}.$$

- Largely dependent on a distance $\langle x^{(i)}, x^{(j)} \rangle$.

Intuitions of Kernel II: Dual Form of Regression

- Regularized sum of square error:

$$L(\omega) = \sum_{i=1}^n (\omega^T x^{(i)} - y^{(i)})^2 + \lambda \|\omega\|_2^2.$$

- Setting the derivatives of L with respect to ω to zero, yields:

$$\omega = -\frac{1}{\lambda} \sum_{i=1}^n (\omega^T x^{(i)} - y^{(i)}) x^{(i)}.$$

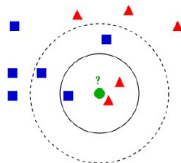
- Let $\alpha_i = \omega^T x^{(i)} - y^{(i)}$. Then we have $\omega = -\frac{1}{\lambda} \sum_{i=1}^n \alpha_i x^{(i)}$.
- Plug it back into $L(\omega)$, we get:

$$L(\alpha) = \sum_{i=1}^n \alpha_i^2 + \frac{1}{\lambda} \sum_{i,j=1}^n \alpha_i \alpha_j (x^{(i)})^T x^{(j)}.$$

- Largely dependent on a distance $\langle x^{(i)}, x^{(j)} \rangle$.

Intuitions of Kernel III: Nonparametric Methods

- Two kinds of methods for prediction:
 - Parametric models: During learning phase, we either get a maximum likelihood estimate of ω or a posterior distribution of ω . Training data is then discarded. We conduct our prediction based only on vector ω .
 - Nonparametric (Memory based) models: Training data points are used in prediction phase. Therefore, it largely depends on a defined distance (similarity), e.g. $\langle x, x^{(i)} \rangle$.
- Nearest neighbor classification: An object is classified by a majority vote of its neighbors, with the object being assigned to the class most common among its k nearest neighbors.



Definition: Kernel

Given the input space X and the feature space H (Hilbert Space), if there exists a mapping ϕ from X to \mathcal{H} , such that for any $x, z \in X$, the function $K(x, z)$ satisfies $K(x, z) = \langle \phi(x), \phi(z) \rangle$, then $K(x, z)$ is called a kernel, where ϕ is the mapping function.

- For a given kernel $K(x, z)$, the feature space \mathcal{H} and mapping function ϕ is usually not unique.
- Example: $X = R^2$, $K(x, z) = \langle x, z \rangle^2$.
 - Denote $x = (x_1, x_2)$ and $z = (z_1, z_2)$, we have:

$$\langle x, z \rangle^2 = (x_1 z_1 + x_2 z_2)^2 = (x_1 z_1)^2 + 2x_1 z_1 x_2 z_2 + (x_2 z_2)^2.$$

- ① $H = R^3$, $\phi(x) = (x_1^2, \sqrt{2}x_1 x_2, x_2^2)^T$.
- ② $H = R^3$, $\phi(x) = \frac{1}{\sqrt{2}}((x_1 - x_2)^2, 2x_1 x_2, (x_1 + x_2)^2)^T$.
- ③ $H = R^4$, $\phi(x) = (x_1^2, x_1 x_2, x_1 x_2, x_2^2)^T$.

Kernel Trick and Its Application in SVM

- **Kernel Trick**: Only the kernel $K(x, z)$ is used in learning and prediction. We do not need to define the mapping ϕ explicitly.
- The **dual** optimization problem of SVM:

$$\begin{aligned} \max_{\alpha} \quad & \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n y^{(i)} y^{(j)} \alpha_i \alpha_j K(x^{(i)}, x^{(j)}) \\ \text{s.t.} \quad & 0 \leq \alpha_i \leq C, i = 1, \dots, n, \\ & \sum_{i=1}^n \alpha_i y^{(i)} = 0. \end{aligned}$$

- Hyperplane and **Optimal Marginal Classifier**:

$$\omega^* = \sum_{i=1}^n \alpha_i^* y^{(i)} x^{(i)}, \quad b^* = y^{(j)} - \sum_{i=1}^n \alpha_i^* y^{(i)} K(x^{(i)}, x^{(j)}).$$

$$f_{\omega^*, b^*}(x) = \omega^{*T} \phi(x) + b^* = \sum_{i=1}^n \alpha_i^* y^{(i)} K(x^{(i)}, x) + b^*.$$

- It can be viewed as we are mapping the original input space X to a feature space \mathcal{H} by the mapping function ϕ , and learning a linear SVM in the new feature space.
- When the mapping function is nonlinear, we can obtain a nonlinear classifier.
- Given the kernel $K(x, z)$, we can still use linear SVM methods to find the solution of the new nonlinear problem.
- Please note that the learning process is directly conducted in the feature space, we do not need to explicitly define the feature space \mathcal{H} and the mapping function ϕ . This is one advantage of kernel, which can directly use the linear classification methods to solve the nonlinear problem.

Valid Kernel

- **Question:** How can we directly define a kernel, not through ϕ ?
- **Question:** Given a kernel $K(x, z)$, how can we tell if it's a valid kernel, i.e., can we tell if there is some feature mapping ϕ so that $K(x, z) = \phi(x)^T \phi(z)$ for all x, z .

Mercer's Theorem

Let $K : R^n \times R^n \rightarrow R$ be given. Then for K to be a valid (Mercer) kernel, it is necessary and sufficient that for any $\{x^{(1)}, \dots, x^{(m)}\}, (m < \infty)$, the corresponding kernel matrix is symmetric positive semi-definite.

- **Kernel (Gram) Matrix:** consider some finite set of m points $\{x^{(1)}, \dots, x^{(m)}\}$, let a square, m by m matrix be defined so that its (i, j) -th entry is given by $K_{ij} = K(x^{(i)}, x^{(j)})$.

$$\mathbb{K} = \begin{pmatrix} K_{11} & K_{12} & \cdots & K_{1m} \\ K_{21} & K_{22} & \cdots & K_{2m} \\ \cdots & \cdots & \cdots & \cdots \\ K_{m1} & K_{m2} & \cdots & K_{mm} \end{pmatrix}$$

Reproducing Kernel Hilbert Space (RKHS)

- A reproducing kernel Hilbert space (RKHS) is a Hilbert space associated with a kernel that reproduces every function in the space.
- **Reproducing Kernel:** A kernel $K : X \times X \rightarrow R$ is called reproducing if it satisfies the following two properties:
 - ① For any $x_0 \in X$, $K(x, x_0)$ is a function of x in the space \mathcal{H} .
 - ② For any $x \in X$ and $f \in \mathcal{H}$, we have $f(x) = f(\cdot) \cdot K(\cdot, x)$. (**Reproducing Property**)
- The mapping function can be naturally be defined as $\phi(x) = K(\cdot, x)$.
- We can see that the following equation holds:

$$\phi(x)^T \phi(z) = K(\cdot, x) \cdot K(\cdot, z) = K(x, z).$$

- Note that for each $f(\cdot) = \sum_{i=1}^m \alpha_i K(\cdot, x^{(i)})$, we have:
 $f(\cdot) \cdot K(\cdot, x) = \sum_{i=1}^m \alpha_i K(x, x^{(i)}) = f(x)$. It defines an inner product in the space \mathcal{H} .

Some Common Kernel Functions

- Polynomial Kernel:

$$K(x, z) = (x^T z + 1)^p.$$

- p is the degree.
- If $p = 2$ and $x = (x_1, x_2)$.

$$K(x, z) = (x_1 z_1 + x_2 z_2 + 1)^2 = 1 + 2x_1 z_1 + 2x_2 z_2 + 2x_1 x_2 z_1 z_2 + x_1^2 z_1^2 + x_2^2 z_2^2.$$

Mapping Function: $\phi(x) = (1, \sqrt{2}x_1, \sqrt{2}x_2, \sqrt{2}x_1 x_2, x_1^2, x_2^2)^T$.

- The corresponding SVM is a polynomial classifier.

$$f(x) = \text{sgn}\left(\sum_{i=1}^n \alpha_i^* y^{(i)} ((x^{(i)})^T x + 1)^p + b^*\right).$$

Some Common Kernel Functions

- **Gaussian Kernel:**

$$K(x, z) = \exp\left\{-\frac{\|x - z\|^2}{2\sigma^2}\right\}.$$

- Sometimes it is also called **radial basis function (RBF)** kernel. Can be generalized to the following form:

$$K(x, z) = \exp\left\{-\frac{Dist(x, z)}{2\sigma^2}\right\}.$$

- The corresponding SVM is a gaussian radial basis function.

$$f(x) = \text{sgn}\left(\sum_{i=1}^n \alpha_i^* y^{(i)} \exp\left\{-\frac{\|x^{(i)} - z\|^2}{2\sigma^2}\right\} + b^*\right).$$

Representer Theorem

Let X be a nonempty set and K a positive-definite real-valued kernel on $X \times X$ with corresponding reproducing kernel Hilbert space \mathcal{H}_K . Given a training sample $\{(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)})\}$, a strictly monotonically increasing real-valued function $g : [0, \infty) \rightarrow R$, and an arbitrary empirical risk function \hat{E} , then for any $f^* \in \mathcal{H}_K$ satisfying:

$$f^* = \arg \min_{f \in \mathcal{H}_K} \{\hat{E} + g(\|f\|_K)\},$$

f^* admits a representation of the form:

$$f^*(\cdot) = \sum_{i=1}^n \alpha_i K(\cdot, x^{(i)}),$$

where $\alpha_i \in R$ for all $1 \leq i \leq n$.

SMO: Sequential Minimal Optimization

- We have mentioned that the primal SVM can be solved by traditional convex quadratic programming methods, which can guarantee the global optimal solution. However, these algorithm usually become slowly especially when the training data are large.
- SMO, proposed by John Platt in 1998, gives an efficient way of solving the dual problem arising from the derivation of the SVM.
- The **dual** optimization problem:

$$\begin{aligned} \max_{\alpha} \quad & \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n y^{(i)} y^{(j)} \alpha_i \alpha_j (x^{(i)})^T x^{(j)} \\ \text{s.t.} \quad & \sum_{i=1}^n \alpha_i y^{(i)} = 0, \\ & 0 \leq \alpha_i \leq C, i = 1, \dots, n. \end{aligned}$$

Motivation of SMO: Coordinate Ascent

- Consider the following unconstrained optimization problem:

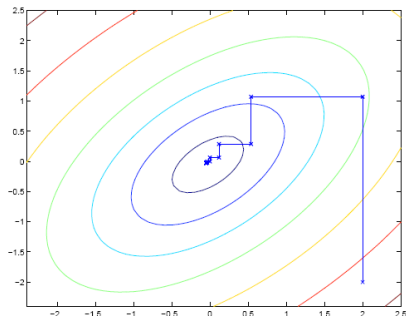
$$\max_{\alpha} W(\alpha_1, \dots, \alpha_m).$$

- Coordinate Ascent Optimization Algorithm:

Loop until convergence : {
 For $i = 1, \dots, m$, {
 $\alpha_i := \max_{\hat{\alpha}_i} W(\alpha_1, \dots, \alpha_{i-1}, \hat{\alpha}_i, \alpha_{i+1}, \dots, \alpha_m)$.
 }
}

- In the innermost loop of this algorithm, we will hold all the variables except for some α_i fixed, and re-optimize W with respect to just the parameter α_i .

Illustration of Coordinate Ascent



- The ellipses in the figure are the contours of a quadratic function that we want to optimize.
- Coordinate ascent was initialized at $(2, -2)$, and also plotted in the figure is the path that it took on its way to the global maximum.
- Notice that on each step, coordinate ascent takes a step that's parallel to one of the axes, since only one variable is being optimized at a time.

Coordinate Ascent for SVM?

- The **dual** optimization problem:

$$\begin{aligned} \max_{\alpha} \quad & \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n y^{(i)} y^{(j)} \alpha_i \alpha_j (x^{(i)})^T x^{(j)} \\ \text{s.t.} \quad & \sum_{i=1}^n \alpha_i y^{(i)} = 0, \\ & 0 \leq \alpha_i \leq C, i = 1, \dots, n. \end{aligned}$$

- Suppose we want to hold $\alpha_2, \dots, \alpha_n$ fixed, and take a coordinate ascent step and reoptimize the objective with respect to α_1 . Can we make any progress?
- The answer is **no**, because the constraint ensures that:

$$\alpha_1 = -y^{(1)} \sum_{i=2}^n \alpha_i y^{(i)}.$$

SMO: Update a Pair of α_i s

- If we want to update some subject of the α_i s, we must **update at least two of them** simultaneously in order to keep satisfying the constraint.
- This is exactly the motivation of SMO, which simply does the following:

Repeat until convergence : {

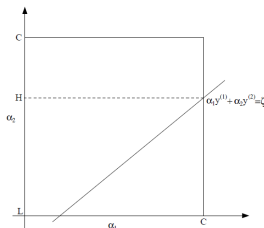
- (1) Select some pair α_i and α_j to update next (using a heuristic that tries to pick the two that will allow us to make the biggest progress towards the global maximum).
- (2) Re-optimize $W(\alpha)$ with respect to α_i and α_j , while holding all the other α_k s ($k = i, j$) fixed.

}

- **The key reason** that SMO is an efficient algorithm is that the update to α_i, α_j can be computed very efficiently.

SMO: Optimizing with respect to Two Variables

- Lets say we currently have some setting of the α_i s that satisfy the constraints, and suppose weve decided to hold $\alpha_3, \dots, \alpha_n$ fixed, and want to re-optimize $W(\alpha_1, \dots, \alpha_n)$ with respect to α_1 and α_2 .
- From the constraints, we require: $\alpha_1 y^{(1)} + \alpha_2 y^{(2)} = -\sum_{i=3}^n \alpha_i y^{(i)}$.
- Since the right hand side is fixed (as weve fixed $\alpha_3, \dots, \alpha_n$), we can just let it be denoted by some constant ζ : $\alpha_1 y^{(1)} + \alpha_2 y^{(2)} = \zeta$.
- We can thus picture the constraints on α_1 and α_2 as follows:



- α_1 and α_2 must lie within the box $[0, C] \times [0, C]$, on the line $\alpha_1 t^{(1)} + \alpha_2 y^{(2)} = \zeta$.
- $L \leq \alpha_2 \leq H$. In this example, $L = 0$; but more generally, there will be some lower-bound L and some upper-bound \mathcal{H} on the permissible values for α_2 that will ensure that α_1, α_2 lie within the box $[0, C] \times [0, C]$.

SMO: Optimizing with respect to One Variable

- Write α_1 as a function of α_2 : $\alpha_1 = (\zeta - \alpha_2 y^{(2)})y^{(1)}$.
- The objective $W(\alpha)$ can be written as:

$$W(\alpha_1, \dots, \alpha_n) = W((\zeta - \alpha_2 y^{(2)})y^{(1)}, \alpha_2, \dots, \alpha_n).$$

- Treating other α_i s as constants, this is just a quadratic function in α_2 .
- If we ignore the constraint of $L \leq \alpha_2 \leq H$, then we can easily maximize this quadratic function by setting its derivative to zero and solving. Let $\alpha_2^{new, unclipped}$ denote the resulting value of α_2 .
- If we had instead wanted to maximize W with respect to α_2 but subject to the box constraint, then we can find the resulting value optimal simply by taking $\alpha_2^{new, unclipped}$ and clipping it to lie in the $[L, H]$ interval, to get:

$$\alpha_2^{new} = \begin{cases} H & \text{if } \alpha_2^{new, unclipped} > H \\ \alpha_2^{new, unclipped} & \text{if } L \leq \alpha_2^{new, unclipped} \leq H \\ L & \text{if } \alpha_2^{new, unclipped} < L \end{cases}$$

SMO: Optimal Solution of α_1 and α_2

- Having found the α_2^{new} , we can use $\alpha_1 = (\zeta - \alpha_2 y^{(2)}) y^{(1)}$ to go back and find the optimal value of α_1^{new} as:

$$\alpha_1^{new} = \alpha_1^{old} + y^{(1)} y^{(2)} (\alpha_2^{old} - \alpha_2^{new}).$$

- Recall that the optimal value of α_2^{new} is:

$$\alpha_2^{new} = \begin{cases} H & \text{if } \alpha_2^{new, unclipped} > H \\ \alpha_2^{new, unclipped} & \text{if } L \leq \alpha_2^{new, unclipped} \leq H \\ L & \text{if } \alpha_2^{new, unclipped} < L \end{cases}$$

SMO: Variable Selection for the First One

- **Main Idea:** Select the training example which most violate KKT condition, and treat the corresponding α_i as the first variable.
- Check whether the i -th training example $(x^{(i)}, y^{(i)})$ satisfies the KKT condition:

$$\begin{aligned}\alpha_i = 0 &\iff y^{(i)}g(x^{(i)}) \geq 1, \\ 0 < \alpha_i < C &\iff y^{(i)}g(x^{(i)}) = 1, \\ \alpha_i = C &\iff y^{(i)}g(x^{(i)}) \leq 1.\end{aligned}$$

where $g(x^{(i)}) = \sum_{j=1}^n \alpha_j y^{(j)} K(x^{(i)}, x^{(j)}) + b$.

- We first check all of the examples satisfying $0 < \alpha_i < C$, i.e. support vectors on the boundary. If all these examples satisfy the KKT condition, then we will check all the training examples.

SMO: Variable Selection for the Second One

- Assume that we have found the first variable α_1 , now we are finding α_2 . The criteria is to make α_2 change much.

Theorem

$$\alpha_2^{new, unclipped} = \alpha_2^{old} + \frac{y^{(2)}(E_1 - E_2)}{\eta},$$

where $\eta = K_{11} + K_{22} - 2K_{12} = \|\phi(x^{(1)}) - \phi(x^{(2)})\|^2, E_i = g(x^{(i)}) - y^{(i)}$.

- A simple method is to directly choose α_2 to make the corresponding $|E_1 - E_2|$ largest.
- Since α_1 is fixed, E_1 is also fixed. If E_1 is positive, we can choose the smallest E_i as E_2 ; Otherwise, if E_1 is negative, we can choose the largest E_i as E_2 .
- For efficiency, we can save all of the E_i in one table.

SMO: Bias b and D-Value E_i

- Each time we finish the optimization over two variables, we need to compute the new bias b .
- If $0 < \alpha_1^{new} < C$, from KKT condition: $\sum_{i=1}^n \alpha_i y^{(i)} K_{i1} + b = y^{(1)}$.
- Then: $b_1^{new} = y^{(1)} - \sum_{i=3}^n \alpha_i y^{(i)} K_{i1} - \alpha_1^{new} y^{(1)} K_{i1} - \alpha_2^{new} y^{(i)} K_{i1}$.
- By the definition of E_1 , we have:
$$E_1 = \sum_{i=3}^n \alpha_i y^{(i)} K_{i1} + \alpha_1^{old} y^{(i)} K_{i1} + \alpha_2^{old} y^{(i)} K_{i1} + b^{old} - y^{(1)}.$$
- Combine the above two equations, we have:
$$b_1^{new} = -E_1 - y^{(1)} K_{11} (\alpha_1^{new} - \alpha_1^{old}) - y^{(2)} K_{21} (\alpha_2^{new} - \alpha_2^{old}) + b^{old}.$$
- Similarly, we have:
$$b_2^{new} = -E_2 - y^{(1)} K_{12} (\alpha_1^{new} - \alpha_1^{old}) - y^{(2)} K_{22} (\alpha_2^{new} - \alpha_2^{old}) + b^{old}.$$
- If both α_1^{new} and α_2^{new} satisfy $0 < \alpha_i^{new} < C, i = 1, 2$, $b_1^{new} = b_2^{new}$.
- If $\alpha_1^{new}, \alpha_2^{new}$ is 0 or C , b_1^{new}, b_2^{new} and all the numbers between them satisfy the KKT condition, we choose the midpoint as b^{new} .
- Each time we finish the optimization over two variables, we need to update E_i : $E_i^{new} = \sum_S y^{(j)} \alpha_j K(x^{(i)}, x^{(j)}) + b^{new} - y^{(i)}$, where S is the set of all support vectors.

SMO Algorithm

- Input: a training set $S = \{(x^{(i)}, y^{(i)}), i = 1, \dots, n\}$, accuracy ϵ ;
- Output: $\hat{\alpha} = (\hat{\alpha}_1, \dots, \hat{\alpha}_n)$.
 - 1 Initialize $\alpha^{(0)} = 0, k = 0$.
 - 2 Select optimization variables $\alpha_1^{(k)}, \alpha_2^{(k)}$, find the solution of optimization problem, denoted as $\alpha_1^{(k+1)}, \alpha_2^{(k+1)}$, update α to $\alpha^{(k+1)}$.
 - 3 If with approximation error ϵ we satisfy the following stopping conditions, turn to step 4; otherwise let $k = k + 1$ and turn to step 2.

$$\begin{aligned} \sum_{i=1}^n \alpha_i y^{(i)} &= 0, \\ 0 \leq \alpha_i &\leq C, i = 1, \dots, n, \\ y^{(i)} g(x^{(i)}) &= \begin{cases} \geq 1 & \{x^{(i)} | \alpha_i = 0\} \\ = 1 & \{x^{(i)} | 0 < \alpha_i < C\} \\ \leq 1 & \{x^{(i)} | \alpha_i = C\} \end{cases} \end{aligned}$$

where $g(x^{(i)}) = \sum_{j=1}^n \alpha_j y^{(j)} K(x^{(j)}, x^{(i)}) + b$.

- 4 let $\hat{\alpha} = \alpha^{(k+1)}$.

Summary of SMO

- SMO is a heuristic algorithm.
- The **basic idea** is that: If all variables satisfy the KKT condition, then we have obtained the optimal solution of the optimization problem. Otherwise, we can choose two variables, while keeping the other variables fixed, to construct a quadratic programming problem. The solution of the subproblem is nearer to the primal optimization problem, since it will make the objective function much smaller. More importantly, the subproblem can usually be efficiently solved by its closed form solution. The two variables are chosen as follows: the first variable is chosen to be the one who most violates the KKT condition, while the second one is automatically determined by the constraints. With the above strategies, SMO continuously divides the primal optimization problem to several subproblems, and it can find the optimal solution by iteratively solve these subproblems.

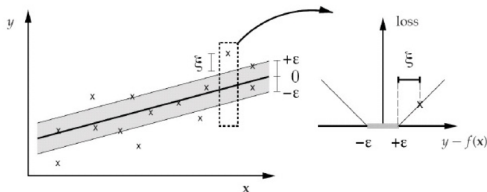
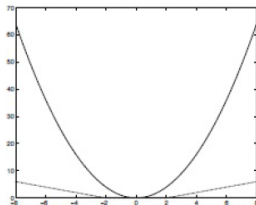


Support Vector Regression: ϵ -Insensitive Loss Function

- In order for the sparseness property of support vectors in SVM for classification to carry over to support vector regression (SVR), we do not use the squared loss but the ϵ -insensitive loss function:

$$\epsilon_i(y^{(i)}, f(x^{(i)})) = \begin{cases} 0 & \text{if } |y^{(i)} - f(x^{(i)})| \leq \epsilon, \\ |y^{(i)} - f(x^{(i)})| - \epsilon & \text{otherwise} \end{cases}$$

- Two characteristics:
 - Errors are tolerated up to a threshold of ϵ .
 - Errors beyond ϵ have linear (rather than quadratic) effect so that the models is more robust against noise.



Support Vector Regression: Primal Optimization Problem

- Primal Optimization Problem:

$$\begin{aligned} \min_{\omega, b} \quad & \frac{1}{2} \|\omega\|_2^2 + C \sum_{i=1}^n (\xi_i^+ + \xi_i^-) \\ \text{s.t.} \quad & y^{(i)} - (\omega^T x^{(i)} + b) \leq \epsilon + \xi_i^+, \quad i = 1, \dots, n, \\ & \omega^T x^{(i)} + b - y^{(i)} \leq \epsilon + \xi_i^-, \quad i = 1, \dots, n, \\ & \xi_i^+, \xi_i^- \geq 0, \quad i = 1, \dots, n. \end{aligned}$$

- Two Types of Slack Variables:

- ξ_i^+ : for positive deviation such that $y^{(i)} - (\omega^T x^{(i)} + b) > \epsilon$,
 - ξ_i^- : for negative deviation such that $\omega^T x^{(i)} + b - y^{(i)} > \epsilon$.
- If $y^{(i)} - (\omega^T x^{(i)} + b) \leq \epsilon$ and $\omega^T x^{(i)} + b - y^{(i)} \leq \epsilon$, then $\xi_i^+ = \xi_i^- = 0$, contributing no cost to the objective function.

Support Vector Regression: Dual and Kernel

- Similar to SVM for classification, the optimization problem for SVR can also be rewritten in the **dual form**.
- **Nonlinear kernel extension** is possible by introducing appropriate kernel functions.
- Due to the **sparseness** property of the ϵ -insensitive loss function, only a small fraction of the training instances are **support vectors** which are used in defining the regression function (like the discriminant function for classification).

Multi-Class SVM: One Versus All

- Approach: view as multi-class classification task, with every complex output y is one class.
- Training example: $\{(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)})\}$, $x^{(i)} \in R^D$, $y^{(i)} \in \{1, \dots, K\}$.
- **Optimization Problem:**

$$\begin{aligned} \min_{\omega, \xi} \quad & \sum_{k=1}^K \|\omega_k\|_2^2 + C \sum_{i=1}^n \xi_i, \\ \text{s.t.} \quad & \forall j \neq y^{(1)} : \quad \omega_{y^{(1)}}^T x^{(1)} \geq \omega_j^T x^{(1)} + 1 - \xi_1, \\ & \dots \\ & \forall j \neq y^{(n)} : \quad \omega_{y^{(n)}}^T x^{(n)} \geq \omega_j^T x^{(n)} + 1 - \xi_n, \\ & \xi_i \geq 0, \quad i = 1, \dots, n. \end{aligned}$$

Learning Theory: Notations

- **Generalization Error:** $\varepsilon(h) = P_{(x,y) \sim \mathcal{D}}(h(x) \neq y)$.
- **Empirical Risk (Error):** $\hat{\varepsilon}(h) = \frac{1}{n} \sum_{i=1}^n I_{\{h(x^{(i)}) \neq y^{(i)}\}}$.
- **PAC (IID) Assumption:** $(x^{(i)}, y^{(i)}), i = 1, \dots, n$ are drawn independently from the same distribution of \mathcal{D} .
- Consider the setting of linear classification, and let $h_{\theta}(x) = I_{\{\theta^T x \geq 0\}}$.
- What's a reasonable way of fitting the parameters θ ?
- **Empirical Risk Minimization (ERM):** $\hat{\theta} = \arg \min_{\theta} \hat{\varepsilon}(h_{\theta})$.
- **Question:** Why should doing well on the training set tell us anything about generalization error? Specifically, can we relate error on the training set to generalization error?
- Define the **hypothesis class** \mathcal{H} used by a learning algorithm to be the set of all classifiers considered by it. For linear classification, $\mathcal{H} = \{h_{\theta} : h_{\theta}(x) = I_{\{\theta^T x \geq 0\}}\}$ is thus the set of all classifiers over X where the decision boundary is linear.
- ERM can now be written as: $\hat{h} = \arg \min_{h \in \mathcal{H}} \hat{\varepsilon}(h)$.

Learning Theory: the case of finite \mathcal{H}

- Consider a finite hypothesis class $\mathcal{H} = \{h_1, \dots, h_k\}$.
- We would like to give guarantees on the generalization error of \hat{h} :
- Take any one fixed $h_i \in \mathcal{H}$. Consider a Bernoulli random variable Z whose distribution is defined as follows. Sample $(x, y) \sim \mathcal{D}$ and set $Z = I_{\{h_i(x) \neq y\}}$. Similarly, we define $Z_j = I_{\{h_j(x^{(j)}) \neq y^{(j)}\}}$. Since our data was drawn iid from \mathcal{D} , Z and the Z_j s have the same distribution.
- The training error can be written as: $\hat{\varepsilon}(h_i) = \frac{1}{n} \sum_{i=1}^n Z_j$.
- Applying the **Hoeffding inequality**, and obtain:

$$P(|\varepsilon(h_i) - \varepsilon(\hat{h}_i)| > \gamma) \leq 2 \exp(-2\gamma^2 n).$$

Hoeffding Inequality (Chernoff bound)

Let Z_1, \dots, Z_n be n iid random variables drawn from a Bernoulli distribution. I.e., $P(Z_i = 1) = \phi$, and $P(Z_i = 0) = 1 - \phi$. Let $\hat{\phi} = \frac{1}{n} \sum_{i=1}^n Z_i$ be the mean of these random variables, and let any $\gamma > 0$ be fixed. Then:

$$P(|\phi - \hat{\phi}| > \gamma) \leq 2 \exp(-2\gamma^2 n).$$

Learning Theory: the case of finite \mathcal{H}

- This shows that, for our particular h_i , training error will be close to generalization error with high probability, assuming n is large. We want to prove that this will be true for simultaneously for all $h \in \mathcal{H}$.
- Let $A_i = \{|\varepsilon(h_i) - \varepsilon(\hat{h}_i)| > \gamma\}$, then we have already show that, for any particular A_i , it holds true that $P(A_i) \leq 2 \exp(-2\gamma^2 n)$.
- Using the union bound, we have that:

$$\begin{aligned} P(\exists h \in \mathcal{H}, |\varepsilon(h_i) - \varepsilon(\hat{h}_i)| > \gamma) &= P(A_1 \cup \dots \cup A_k) \\ &\leq \sum_{i=1}^k P(A_i) \leq \sum_{i=1}^k 2 \exp(-2\gamma^2 n) = 2k \exp(-2\gamma^2 n). \end{aligned}$$

- Then: $P(\forall h \in \mathcal{H}, |\varepsilon(h_i) - \varepsilon(\hat{h}_i)| \leq \gamma) \geq 1 - 2k \exp(-2\gamma^2 n)$.
- With probability at least $1 - 2k \exp(-2\gamma^2 n)$, we have that $\varepsilon(h)$ will be within of $\hat{\varepsilon}(h)$ for all $h \in \mathcal{H}$. This is called **uniform convergence** result, because this is a bound that holds simultaneously for all (as opposed to just one) $h \in \mathcal{H}$.

Learning Theory: the case of finite \mathcal{H}

- There are three quantities of interest here: n , γ and the probability of error; we can bound either one in terms of the other two.
- Given γ and some $\delta > 0$, how large must n be before we can guarantee that with probability at least $1 - \delta$, training error will be within of generalization error?
 - By setting $\delta = 2k \exp(-2\gamma^2 n)$ and solving for n , we find that if $n \geq \frac{1}{2\gamma^2} \log \frac{2k}{\delta}$, then with probability at least $1 - \delta$, we have that $|\varepsilon(h) - \hat{\varepsilon}(h)| > \gamma$ for some $h \in \mathcal{H}$. ((Equivalently, this shows that the probability that $|\varepsilon(h) - \hat{\varepsilon}(h)| > \gamma$ for some $h \in \mathcal{H}$ is at most δ .)
 - This bound tells us how many training examples we need in order make a guarantee. The training set size n that a certain method or algorithm requires in order to achieve a certain level of performance is also called the algorithms **sample complexity**.
- we can also hold n and δ fixed, and solve for γ in the previous equation, and show that with probability $1 - \delta$, we have that for all $h \in \mathcal{H}$,

$$|\hat{\varepsilon}(h) - \varepsilon(h)| \leq \sqrt{\frac{1}{2n} \log \frac{2k}{\delta}}.$$

Learning Theory: the case of finite \mathcal{H}

- Now, let's assume that uniform convergence holds, i.e., that $|\varepsilon(h) - \hat{\varepsilon}(h)| \leq \gamma$ for all $h \in \mathcal{H}$. What can we prove about the generalization of our learning algorithm that picked $\hat{h} = \arg \min_{h \in \mathcal{H}} \hat{\varepsilon}(h)$?
- Define $h^* = \arg \min_{h \in \mathcal{H}} \varepsilon(h)$ to be the best possible hypothesis in \mathcal{H} . Note that h^* is the best that we could possibly do given that we are using \mathcal{H} , so it makes sense to compare our performance to that of h^* . We have: $\varepsilon(\hat{h}) \leq \hat{\varepsilon}(\hat{h}) + \gamma \leq \hat{\varepsilon}(h^*) + \gamma \leq \varepsilon(h^*) + 2\gamma$.
- If uniform convergence occurs, then the generalization error of \hat{h} is at most 2γ worse than the best possible hypothesis in \mathcal{H} !

Theorem

Let $|\mathcal{H}| = k$, any n, δ be fixed. Then with probability at least $1 - \delta$:

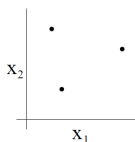
$$\varepsilon(\hat{h}) \leq \left(\min_{h \in \mathcal{H}} \varepsilon(h) \right) + 2\sqrt{\frac{1}{2n} \log \frac{2k}{\delta}}.$$

Learning Theory: the case of infinite \mathcal{H}

- We have proved some useful theorems for the case of finite hypothesis classes. But many hypothesis classes, including any parameterized by real numbers (as in linear classification) actually contain an infinite number of functions. **Can we prove similar results for this setting?**
- Given a set $S = \{x^{(1)}, \dots, x^{(d)}\}$ (no relation to the training set) of points $x^{(i)} \in X$, we say that **\mathcal{H} shatters S** if \mathcal{H} can realize any labeling on S . I.e., if for any set of labels $\{y^{(1)}, \dots, y^{(d)}\}$, there exists some $h \in \mathcal{H}$ so that $h(x^{(i)}) = y^{(i)}$ for all $i = 1, \dots, d$.
- Given a hypothesis class \mathcal{H} , we then define its **Vapnik-Chervonenkis dimension**, written $VC(\mathcal{H})$, to be the size of the largest set that is shattered by \mathcal{H} . (If \mathcal{H} can shatter arbitrarily large sets, then $VC(\mathcal{H}) = \infty$.)

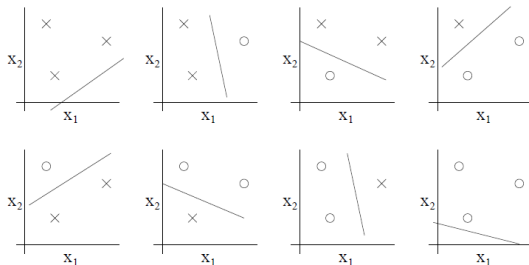
VC Dimension

- For instance, consider the following set of three points:

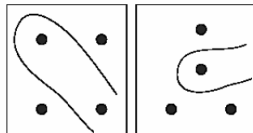


Can the set \mathcal{H} of linear classifiers in two dimensions ($h(x) = I_{\{\theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0\}}$) shatter the set above? The answer is yes.

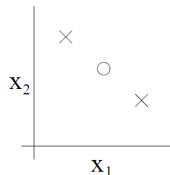
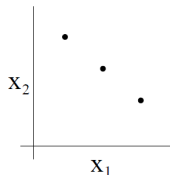
- For any of the eight possible labelings of these points, we can find a linear classifier that obtains zero training error on them:



VC Dimension



- There is no set of 4 points that this hypothesis class can shatter. Thus, the largest set that \mathcal{H} can shatter is of size 3, i.e., $VC(\mathcal{H}) = 3$.
- Note that the VC dimension of \mathcal{H} here is 3 even though there may be sets of size 3 that it cannot shatter. In other words, in order to prove that $VC(\mathcal{H})$ is at least d , we need to show only that there is at least one set of size d that \mathcal{H} can shatter.



Theorem

Let \mathcal{H} be given, and let $d = VC(\mathcal{H})$. Then with probability at least $1 - \delta$, we have that for all $h \in \mathcal{H}$,

$$|\varepsilon(h) - \hat{\varepsilon}(h)| \leq O\left(\sqrt{\frac{d}{n} \log \frac{n}{d} + \frac{1}{n} \log \frac{1}{\delta}}\right).$$

Thus with probability at least $1 - \delta$, we also have that:

$$\varepsilon(\hat{h}) \leq \varepsilon(h^*) + O\left(\sqrt{\frac{d}{n} \log \frac{n}{d} + \frac{1}{n} \log \frac{1}{\delta}}\right).$$

- In other words, if a hypothesis class has finite VC dimension, then uniform convergence occurs as n becomes large. As before, this allows us to give a bound on $\varepsilon(h)$ in terms of $\hat{\varepsilon}(h)$.

- We also have the following corollary:

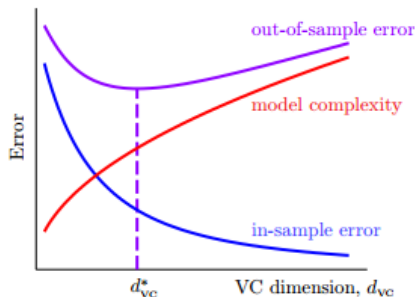
Corollary

For $|\varepsilon(h) - \hat{\varepsilon}(h)| \leq \gamma$ to hold for all $h \in \mathcal{H}$ (and hence $\varepsilon(\hat{h}) \leq \varepsilon(h^*) + 2\gamma$) with probability at least $1 - \delta$, it suffices that $n = O_{\gamma, \delta}(d)$.

- In other words, the number of training examples needed to learn well using H is linear in the VC dimension of \mathcal{H} . It turns out that, for most hypothesis classes, the VC dimension (assuming a reasonable parameterization) is also roughly linear in the number of parameters. Putting these together, we conclude that (for an algorithm that tries to minimize training error) the number of training examples needed is usually roughly linear in the number of parameters of \mathcal{H} .

Summary of VC Dimension

- VC dimension is a measure of model complexity.



- When n is large, we can accommodate models with large complexity (VC dimension).

Generalization Analysis of SVM

- The following Theorem by Vapnik (1982) provides the essential link between margin and realized classifier class complexity for SVMs.

Theorem: Vapnik 1982

The class of optimal linear separators has VC dimension h bounded from above as:

$$h \leq \min \left\{ \left\lceil \frac{4r^2}{\rho^2} \right\rceil, m \right\} + 1,$$

where ρ is the margin, r is the radius of the smallest sphere that can enclose all of the training examples, and m is the dimensionality of X .

- Intuitively, this implies that regardless of dimensionality m we can minimize the VC dimension by maximizing the margin ρ .
- Thus, complexity of the classifier is kept small regardless of dimensionality.

Summary of Supervised Learning

- Framework of Statistical Machine Learning
- Bias and Variance Decomposition
- Overfitting and Regularization
- Discriminative and Generative Approach
- Linear Regression
- Logistic Regression
- Naive Bayesian
- Stochastic Gradient Descent
- SVM and Dual
- Kernel and Learning Theory

Wish You All Good Luck!



- At the entrance to science, as at the entrance to hell, the demand must be made: Here all hesitations must be wiped out. All cowardice must here be dead". ([Marx, Preface to A Contribution to the Critique of Political Economy, Germany](#))