

CSC 665: Support Vector Machines

Chicheng Zhang

October 3, 2019

1 Support vector machines - the maximum margin hyperplane problem

We consider linear classification, where examples $(x_i, y_i)_{i=1}^m$ are such that $x_i \in \mathbb{R}^d$ are features, and $y_i \in \{\pm 1\}$ are binary labels.

Suppose that the training set $S = (x_i, y_i)_{i=1}^m$ is linearly separable, i.e. there exists a linear classifier $(w, b) \in \mathbb{R}^{d+1}$, such that for all i ,

$$\begin{cases} \langle w, x_i \rangle + b > 0 & y_i = +1, \\ \langle w, x_i \rangle + b < 0 & y_i = -1. \end{cases} \quad (1)$$

One way to train a linear classifier would be to use the consistency algorithm, i.e. solving a linear program, that finds a (w, b) such that Equation (1) holds. However, note that not all consistent linear classifiers are created equal: some of them are closer to training examples than others. Formally, the distance of a point x in \mathbb{R}^d to a hyperplane $H_{w,b} = \{x_0 : w \cdot x_0 + b = 0\}$ is defined as the shortest distance of x to any of the points in $H_{w,b}$:

$$d(x, H_{w,b}) = \min \{ \|x - x_0\| : w \cdot x_0 + b = 0 \}. \quad (2)$$

Can we calculate this distance analytically? First, let us assume without loss of generality that $\|w\| = 1$, as any hyperplane $H_{w',b'}$ can be written as $H_{w,b}$ for $\|w\| = 1$ by letting $w = \frac{w'}{\|w'\|}$ and $b = \frac{b'}{\|w'\|}$. Now, consider a point $x_0 \in H_{w,b}$ such that $x_0 = x + \alpha w$ for some α . What is the value of α ? Note that

$$\langle w, x + \alpha w \rangle + b = 0,$$

which implies that $\alpha = -(\langle w, x \rangle + b)$.

Claim 1. For all x_1 in $H_{w,b}$,

$$\|x_1 - x\| \geq \|x_0 - x\|. \quad (3)$$

Consequently, $d(x, H_{w,b}) = |\langle w, x \rangle + b|$.

Proof. Note that x_0 and x_1 are both in $H_{w,b}$, $\langle w, x_0 \rangle + b = \langle w, x_1 \rangle + b = 0$. Therefore, $\langle x_1 - x_0, w \rangle = 0$. In other words,

$$\langle x_1 - x_0, x_0 - x \rangle = 0.$$

Now, by Pythagorean theorem,

$$\|x - x_1\|^2 = \|x - x_0\|^2 + \|x_0 - x_1\|^2 \geq \|x - x_0\|^2,$$

which proves Equation (3). This implies that

$$d(x, H_{w,b}) = \|x - x_0\| = |\langle w, x \rangle + b|.$$

□

Here is a proposal:

Find the linear classifier (w, b) that not only separates the examples but also maximizes the minimum distances to all examples.

Why is the proposal sensible? One observation is that this classifier is the most “robust”. For example, if test examples happen to be just a little distance away from training examples (with the same labels), then this classifier would still classify such examples correctly.

Formally, we can describe the proposal as an optimization problem:

$$\begin{aligned}
& \underset{w, b, A}{\text{maximize}} && A \\
& \text{s. t.} && A > 0, \quad \|w\| = 1, \\
& && y_i(\langle w, x_i \rangle + b) > 0, && \forall i \in \{1, \dots, n\}, \\
& && |\langle w, x_i \rangle + b| \geq A, && \forall i \in \{1, \dots, n\}.
\end{aligned} \tag{4}$$

The above program is not a convex program, and is difficult to optimize directly. Let’s make a few transformations to make it a convex program - i.e. finding a convex optimization problem whose solution is related to that of the above optimization problem.

Let’s consider the following optimization problem:

$$\begin{aligned}
& \underset{w, b, A}{\text{maximize}} && A \\
& \text{s. t.} && A > 0, \quad \|w\| = 1, \\
& && y_i(\langle w, x_i \rangle + b) \geq A, && \forall i \in \{1, \dots, n\},
\end{aligned} \tag{5}$$

Our claim is that the above two optimization problems have the same solutions. Why? Because under $A > 0$, constraints $y_i(\langle w, x_i \rangle + b) > 0$ and $|\langle w, x_i \rangle + b| \geq A$, together, are equivalent to $y_i(\langle w, x_i \rangle + b) \geq A$, as $y_i \in \{\pm 1\}$. For every i , the quantity $y_i(\langle w, x_i \rangle + b)$ is the *margin* of halfspace $H_{w, b}$ on example (x_i, y_i) . Therefore the above is also called the “maximum margin hyperplane” problem.

Now let $w' = \frac{w}{A}$, $b' = \frac{b}{A}$. Note that the above optimization problem is equivalent to

$$\begin{aligned}
& \underset{w', b', A}{\text{maximize}} && A \\
& \text{s. t.} && A > 0, \quad \|w'\| = \frac{1}{A}, \\
& && y_i(\langle w', x_i \rangle + b') \geq 1, && \forall i \in \{1, \dots, n\},
\end{aligned}$$

Furthermore, this is equivalent to

$$\begin{aligned}
& \underset{w', b'}{\text{minimize}} && \|w'\| \\
& \text{s. t.} && y_i(\langle w', x_i \rangle + b') \geq 1, && \forall i \in \{1, \dots, n\},
\end{aligned}$$

As the function $x \mapsto \frac{1}{2}x^2$ is monotonically increasing for $x > 0$, we get that the above is equivalent to

$$\begin{aligned}
& \underset{w', b'}{\text{minimize}} && \frac{1}{2}\|w'\|^2 \\
& \text{s. t.} && y_i(\langle w', x_i \rangle + b') \geq 1, && \forall i \in \{1, \dots, n\},
\end{aligned} \tag{6}$$

Optimization problem (4) is called the *support vector machine* (SVM). Note that its constraints are all linear inequalities, which defines a convex constraint set. In addition, its optimization objective is a quadratic function over optimization variables, which is a convex function. This implies that it is a *convex optimization* problem.

Recovering the optimal solution of (4). Suppose we have a solution of (6), written as (w'^*, b'^*) . Note that the optimal A in (5) (thus, in (4)) is $1/\|w'^*\|$, which is the value of the minimum margin. This implies that in (5) (thus, in (4)), $w^* = A^* w'^* = \frac{w'^*}{\|w'^*\|}$, $b^* = A^* b'^* = \frac{b'^*}{\|w'^*\|}$. The optimal hyperplane is simply $H_{w^*, b^*} = H_{w'^*, b'^*}$.

1.1 Optimality condition

To avoid notation clutter, let us drop the apostrophes in optimization problem (6):

$$\begin{aligned} & \underset{w, b}{\text{minimize}} && \frac{1}{2} \|w\|^2 \\ & \text{s. t.} && y_i(\langle w, x_i \rangle + b) \geq 1, \quad \forall i \in \{1, \dots, n\}, \end{aligned} \tag{7}$$

What property does the optimal solution (w^*, b^*) have? We will take a detour and first discuss Lagrangian duality, a fundamental concept in constrained optimization. Let us first write (7) as an unconstrained optimization problem over a slightly more complicated objective:

$$\min_{w, b} \max_{\alpha \geq 0} L(w, b, \xi; \alpha), \tag{8}$$

where $L(w, b; \alpha) = \frac{\lambda}{2} \|w\|^2 + \sum_{i=1}^n \alpha_i (1 - y_i(\langle w, x_i \rangle + b))$.

Define $P(w, b) := \max_{\alpha \geq 0} L(w, b; \alpha)$. Observe that:

$$P(x, b) = \begin{cases} -\infty, & \exists i, 1 - y_i(\langle w, x_i \rangle + b) < 0 \\ \frac{\lambda}{2} \|w\|^2, & \forall i, 1 - y_i(\langle w, x_i \rangle + b) \geq 0 \end{cases}$$

Therefore, optimization problem (8) is equivalent to (7). Now consider switching the orders of min and max in (8):

$$\max_{\alpha \geq 0} \min_{w, b} L(w, b; \alpha).$$

This is called the dual problem of (8) ((8) is called the primal problem). Let's call the optimal primal value p^* and the optimal dual value d^* . What's the relationship between the primal and dual problems, and their respective optimal solutions?

We state the following result from numerical optimization. Consider a constrained convex optimization problem that has both equality and inequality constraints:

$$\begin{aligned} & \underset{x}{\text{minimize}} && f(x) \\ & \text{s. t.} && g_i(x) \leq 0, \quad \forall i \in \{1, \dots, n\}, h_i(x) = 0, \quad \forall i \in \{1, \dots, m\}, \end{aligned}$$

Similar as before, we can define Lagrange function $L(x, \alpha, \beta) = f(x) + \sum_{i=1}^n \alpha_i g_i(x) + \sum_{i=1}^m \beta_i h_i(x)$. Define

$$\begin{aligned} P(x) &\triangleq \max_{\alpha \geq 0, \beta} L(x, \alpha, \beta), \\ D(\alpha, \beta) &= \min_x L(x, \alpha, \beta), \\ p^* &= \min_x P(x) = \min_x \max_{\alpha \geq 0, \beta} L(x, \alpha, \beta), \\ d^* &= \max_{\alpha \geq 0, \beta} D(\alpha, \beta) = \max_{\alpha \geq 0, \beta} \min_x L(x, \alpha, \beta), \end{aligned}$$

we have the following result.

Theorem 1. *Under mild assumptions¹, we have that there exists x^* , α^* , and β^* , such that*

¹specifically, f , g_i 's, h_i 's are convex, and there exists w, b, ξ such that all inequality constraints in are strictly satisfied, namely the Slater condition.

1. x^* is optimal solution of the primal problem and α^*, β^* is the optimal solution of the dual problem.
2. Strong duality holds:

$$p^* = L(x^*, \alpha^*, \beta^*) = d^*.$$

3. Karush-Kuhn-Tucker (KKT) condition holds:

$$\begin{array}{ll} \nabla_x L(x^*, \alpha^*, \beta^*) = 0, & \text{Stationarity} \\ \forall i, \quad g_i(x^*) \leq 0, h_i(x^*) \leq 0, & \text{Primal feasible} \\ \forall i, \quad \alpha_i \geq 0, & \text{Dual feasible} \\ \forall i, \quad \alpha_i g_i(x^*) = 0. & \text{Complementary slackness} \end{array}$$

Applying the theorem to SVM optimization, we can also recover the primal optimal solution (w^*, b^*) from dual solution α^* by invoking the KKT condition. To see why, recall that in SVM, $L(w, b; \alpha) = \frac{\lambda}{2} \|w\|^2 + \sum_{i=1}^n \alpha_i (1 - y_i (\langle w, x_i \rangle + b))$, hence by stationarity condition,

$$\nabla_w L(w, b; \alpha) = \lambda w - \sum_{i=1}^n \alpha_i y_i x_i = 0,$$

which implies that

$$w^* = \sum_{i=1}^n \alpha_i^* y_i x_i.$$

that is, the optimal solution is a linear combination of the feature vectors of training examples.

Furthermore, denote by $\mathcal{I} = \{i : y_i (\langle w^*, x_i \rangle + b^*) = 1\}$ the set of examples that has margin exactly equal to 1. Complementary slackness says that for all i ,

$$\alpha_i^* (1 - y_i (\langle w^*, x_i \rangle + b^*)) = 0.$$

This implies that for an $i \notin \mathcal{I}$, as $y_i (\langle w^*, x_i \rangle + b^*) > 1$, $\alpha_i = 0$. We call \mathcal{I} the set of *support vectors*, which are the vectors that “contribute” to the optimal solution w^* .

It can also be verified that there exists at least one i , $y_i (\langle w^*, x_i \rangle + b^*) = 1$. Pick one such i ; b^* can be recovered by the formula $b^* = y_i - \langle w^*, x_i \rangle$.

1.2 Coping with linear non-separability

Can we still train SVM if the data is not linearly separable? Note that optimization problem (6) will not find a solution, as now the constraint set become infeasible. Generally there are two ways to sidestep this problem: first, introduce nonlinear feature maps; second, relax the SVM formulation to allow for training examples to be classified incorrectly.

For the first approach, we can consider having a feature map $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^m$, so that every example (x_i, y_i) is transformed to $(\phi(x_i), y_i)$. Suppose $(\phi(x_i), y_i)$ is linearly separable, then we compute a SVM over these examples to get $w^* \in \mathbb{R}^{m+1}$ and $b \in \mathbb{R}$. Our output linear classifier is $\text{sign}(\langle w^*, \phi(x_i) \rangle + b)$. For example, suppose we have a distribution D over $\mathbb{R}^2 \times \{\pm 1\}$ such that for all examples (x, y) ’s on the support of D , $x_1^2 + x_2^2 \leq 1 \Leftrightarrow y = +1$. In this case, we can introduce feature map $\phi(x) = (x_1^2, x_2^2)$ to make the dataset linearly separable.

For the second approach, we introduce slack variables $\xi_i \geq 0$ for every example i , to allow some example to be misclassified. In addition, we introduce a regularization parameter $\lambda > 0$ that trades off misclassification and margin on correct examples:

$$\begin{array}{ll} \underset{w, b, \xi}{\text{minimize}} & \frac{\lambda}{2} \|w\|^2 + \sum_{i=1}^n \xi_i \\ \text{s. t.} & y_i (\langle w, x_i \rangle + b) \geq 1 - \xi_i, \quad \forall i \in \{1, \dots, n\}, \\ & \xi_i \geq 0, \quad \forall i \in \{1, \dots, n\}, \end{array} \tag{9}$$

Intuitively, when λ is larger, it focuses more on enforcing large margin on correct examples; when λ is smaller, it forces more on reducing misclassification. Notice that the last two lines can be summarized by: $\forall i \in \{1, \dots, n\}, \xi_i \geq \max(0, 1 - y_i(\langle w, x_i \rangle + b))$. Therefore, the optimal choice of ξ_i equals $\max(0, 1 - y_i(\langle w, x_i \rangle + b))$. Let $\phi(z) = \max(0, 1 - z)$ and $R(w) = \frac{\lambda}{2} \|w\|^2$. We thus can rewrite optimization problem (9) as:

$$\underset{w, b}{\text{minimize}} \quad \lambda R(w) + \sum_{i=1}^n \phi(y_i(\langle w, x_i \rangle + b)). \quad (10)$$

As a convention, we call $\phi(y(\langle w, x \rangle + b))$ the *hinge loss* of linear classifier (w, b) on example (x, y) , written as $\ell_{\text{hinge}}((w, b), (x, y))$. When the margin $y(\langle w, x \rangle + b)$ is larger, the hinge loss is smaller. The above form is also called a *regularized loss minimization* formulation, which captures a wide range of optimization problems in machine learning (by changing loss function ϕ and regularizer R), such as logistic regression, ridge regression, lasso, etc.

Both approaches has its own advantages and drawbacks. For the feature transformation approach, it is unclear if a ϕ will guarantee that the transformed dataset satisfies linear separability. For the soft margin approach, if the dataset is highly linearly nonseparable (e.g. the unit circle example), then as it is still learning a linear classifier, it will not perform well. It might be a good idea to combine nonlinear feature map with soft margin in practice.

2 The dual of SVM

Sometimes looking at the dual problem will yield unexpected insights about the original (primal) problem. Indeed, SVM is a canonical example for this statement - we have already seen that the KKT condition implies that we can write the optimal solution w^* in terms of dual optimal solution α^* . We have discussed the dual problem in an abstract way so far. But what exactly is the dual problem for SVM?

Let us first calculate the dual objective function $D(\alpha) = \min_{w, b} L(w, b; \alpha)$, where

$$L(w, b; \alpha) = \frac{\lambda}{2} \|w\|^2 + \sum_{i=1}^n \alpha_i (1 - y_i(\langle w, x_i \rangle + b)).$$

We can write $D(\alpha)$ as follows:

$$D(\alpha) = \sum_{i=1}^n \alpha_i + \min_w \left(\frac{\lambda}{2} \|w\|^2 - \left\langle w, \sum_{i=1}^n \alpha_i y_i x_i \right\rangle \right) + \min_b \left(\sum_{i=1}^n \alpha_i y_i b \right).$$

Define $g(z) = \begin{cases} -\infty & z = 0 \\ 0 & z \neq 0 \end{cases}$, then

$$D(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2\lambda} \left\| \sum_{i=1}^n \alpha_i y_i x_i \right\|^2 + g\left(\sum_{i=1}^n \alpha_i y_i\right).$$

Therefore, $\max_{\alpha \geq 0} D(\alpha)$ is equivalent to

$$\begin{aligned} & \underset{\alpha}{\text{maximize}} \quad \sum_{i=1}^n \alpha_i - \frac{1}{2\lambda} \left\| \sum_{i=1}^n \alpha_i y_i x_i \right\|^2 \\ & \text{s. t.} \quad \alpha_i \geq 0, \quad \forall i \in \{1, \dots, n\}, \end{aligned} \quad (11)$$

writing the objective more explicitly, it is

$$\begin{aligned} & \underset{\alpha}{\text{maximize}} \quad \sum_{i=1}^n \alpha_i - \frac{1}{2\lambda} \sum_{i=1}^n \sum_{j=1}^n y_i y_j \langle x_i, x_j \rangle \alpha_i \alpha_j \\ & \text{s. t.} \quad \alpha_i \geq 0, \quad \forall i \in \{1, \dots, n\}, \end{aligned} \quad (12)$$

Compared to (6), this is also a quadratic program, however, its objective becomes a complicated quadratic function, and its constraints says that α lies in the positive orthant of \mathbb{R}^n , which is simpler than the linear inequality constraints in (6).

3 The kernel trick

The dual of SVM (12) uncovers an interesting fact: if we would like to compute the optimal solution of (6), it suffices to solve the dual optimization problem, whose objective function only depends on the pairwise inner product between training examples (as opposed to the original feature vectors of training examples).

This opens up a new opportunity: suppose we have a feature map that is extremely high dimensional (say has dimensionality M) but has succinct representation on pairwise inner product $\langle \phi(x), \phi(x') \rangle$ (say can be evaluated with time m), then we may avoid paying a time complexity of M in learning the SVM classifier on the transformed examples. Here is the full proposal:

1. Define $k(x, x') = \langle \phi(x), \phi(x') \rangle$ be the *kernel function* associated with feature mapping ϕ .
2. Solve the dual optimization problem (6), get $(\alpha_i)_{i=1}^m$.
3. By KKT condition, we can recover

$$w^* = \sum_{i=1}^n \alpha_i y_i \phi(x_i),$$

but we only store w^* *implicitly*, i.e. storing the value of all α_i 's.

4. To recover b^* , find an j such that $\alpha_j > 0$, and let

$$b^* = y_j - \langle w^*, x_j \rangle = y_j - \sum_{i=1}^n \alpha_i^* y_i k(x_i, x_j),$$

where we directly evaluate $k(x_i, x_j)$ as opposed to calculating $\phi(x_i)$, $\phi(x_j)$ and take their inner product.

5. To make prediction on future example x , we compute

$$\langle w^*, \phi(x) \rangle + b^* = \sum_{i=1}^n \alpha_i^* k(x_i, x) + b^*.$$

Same as before, we directly evaluate the kernel function.

As discussed before, each feature map corresponds to a kernel function. Some feature map gives succinct kernel functions, whereas others may not. For example, for input domain \mathbb{R}^2 , define $\phi(x) = (x_1^2, \sqrt{2}x_1x_2, x_2^2)$. It can be checked that its associated kernel function has a succinct form:

$$\langle \phi(x), \phi(x') \rangle = (x_1x'_1 + x_2x'_2)^2 = (\langle x, x' \rangle)^2.$$

However, if we define $\phi(x) = (x_1^2, x_1x_2, x_2^2)$, then its corresponding $k(x, x')$ does not have a succinct form.

Basic properties of kernel functions:

1. if K is the kernel function of ϕ , then for positive c , cK is the kernel function of $\sqrt{c}\phi$.
2. if K_1 (resp. K_2) is the kernel function of ϕ_1 (resp. ϕ_2), then $K_1 + K_2$ is the kernel function of $\phi(x) = (\phi_1(x), \phi_2(x))$.

3. if K_1 (resp. K_2) is the kernel function of ϕ_1 (resp. ϕ_2), then $K_1 \cdot K_2$ is the kernel function of $\phi(x) = \phi_1(x) \otimes \phi_2(x)$, where the \otimes notation denotes the Kronecker product. Suppose $a = (a_1, \dots, a_n)$ and $b = (b_1, \dots, b_m)$. Then,

$$a \otimes b = \begin{bmatrix} a_1 b \\ \dots \\ a_n b \end{bmatrix} = \begin{bmatrix} a_1 b_1 \\ \dots \\ a_1 b_m \\ \dots \\ a_n b_1 \\ \dots \\ a_n b_m \end{bmatrix}.$$

The claim follow from a basic fact about Kronecker product:

$$\langle a \otimes b, c \otimes d \rangle = \langle a, c \rangle \cdot \langle b, d \rangle.$$

Popular choices of kernel functions:

1. Polynomial kernel $k(x, x') = (1 + \sum_{i=1}^d x_i x'_i)^s$.
2. Radial basis function kernel $k(x, x') = \exp\left(-\frac{\|x - x'\|^2}{2\sigma^2}\right)$.

It can be verified by the basic properties above that $k(x, x')$ are kernel functions. Can you find out their respective feature mappings?