

Support Vector Machines

Support vector machines (SVMs) are one of the central concepts in all of machine learning. They are simply a combination of two ideas: linear classification via maximum (or optimal soft) margin hyperplanes, and kernels.

Conceptually, given training data (\mathbf{x}_i, y_i) , $i = 1, \dots, n$, with $\mathbf{x}_i \in \mathbb{R}^d$ and $y_i \in \{-1, +1\}$, we are going to

1. Put the feature vectors through a nonlinear map: $\mathbf{x}_i \rightarrow \Phi(\mathbf{x}_i)$. This takes them from \mathbb{R}^d to some abstract Hilbert space \mathcal{H} .
2. Fit a linear classifier in \mathcal{H} . Since \mathcal{H} has a valid inner product, there is a notion of inner product (hence projections onto subspaces) and distance in \mathcal{H} , making this problem well-defined.
3. Translate the rule back to a function (classifier) in \mathbb{R}^d using

$$h(\mathbf{x}) = \text{sign}(\langle \mathbf{w}, \Phi(\mathbf{x}) \rangle + b)$$

The “trick” is that step 1 is actually implicit. By using a kernel, we avoid having to actually compute the map. This is also true for the inner product $\langle \mathbf{w}, \Phi(\mathbf{x}) \rangle$; we will see that this can also be computed using a kernel.

To see how this works, we turn to our recently acquired knowledge of constrained optimization and duality.

Dual of optimal soft margin

The optimal soft-margin classifier is found by solving

$$\begin{aligned} \text{(SM)} \quad & \underset{b, \mathbf{w}, \xi}{\text{minimize}} \quad \frac{1}{2} \|\mathbf{w}\|_2^2 + \frac{C}{n} \sum_{i=1}^n \xi_i \\ & \text{subject to} \quad y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 - \xi_i, \quad i = 1, \dots, n. \\ & \quad \quad \quad \xi_i \geq 0, \quad i = 1, \dots, n. \end{aligned}$$

If we try to solve this program after we map $\mathbf{x}_i \rightarrow \Phi(\mathbf{x}_i)$, then the optimization program is looking for a \mathbf{w} in the same space as $\Phi(\mathbf{x}_i)$. Since this space will in general be much higher dimensional than \mathbb{R}^d — and indeed could very well be infinite dimensional — this is problematic.

Fortunately, duality and the KKT conditions will come to our rescue.

Here are two facts, which we verify in the later sections of these notes.

1. The **dual** of program (SM) is

$$\begin{aligned} \text{(DSM)} \quad & \underset{\alpha \in \mathbb{R}^n}{\text{maximize}} \quad -\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i^\top \mathbf{x}_j + \sum_{i=1}^n \alpha_i \\ & \text{subject to} \quad \sum_{i=1}^n \alpha_i y_i = 0 \\ & \quad \quad \quad 0 \leq \alpha_i \leq C/n, \quad i = 1, \dots, n. \end{aligned}$$

Note that this is an optimization program in n variables. In fact, just like (SM), it is a (concave) quadratic program with linear equality and inequality constraints.

2. Given the solution $\boldsymbol{\alpha}^*$ to the dual, we can easily compute the primal solution \boldsymbol{w}^*, b^* . The weights \boldsymbol{w}^* are computed as

$$\boldsymbol{w}^* = \sum_{i=1}^n \alpha_i^* y_i \boldsymbol{x}_i.$$

Recall that the weights are used to define a linear functional which the classifier evaluates, and note that the expression above implies

$$\boldsymbol{w}^{*\top} \boldsymbol{x} = \sum_{i=1}^n \alpha_i^* y_i \boldsymbol{x}_i^\top \boldsymbol{x}.$$

The offset b^* is simply¹

$$b^* = y_{i_0} - \boldsymbol{w}^{*\top} \boldsymbol{x}_{i_0} = y_{i_0} - \sum_{i=1}^n \alpha_i^* y_i \boldsymbol{x}_i^\top \boldsymbol{x}_{i_0}, \quad (1)$$

where i_0 is any index where the inequality constraint for α_{i_0} is slack: $0 < \alpha_{i_0} < C/n$.

The kernelized dual

The dual program and the computation of the affine functional $\boldsymbol{w}^{*\top} \boldsymbol{x} + b^*$ rely on computations of inner products between feature vectors. We can kernelize the procedure above by replacing the inner products with a kernel evaluation,

$$\boldsymbol{x}_i^\top \boldsymbol{x}_j \rightarrow k(\boldsymbol{x}_i, \boldsymbol{x}_j),$$

¹In practice, you might average the associated values over all such indexes that have this property.

where $k(\cdot, \cdot)$ is a symmetric positive semidefinite kernel. As we have seen, $k(\mathbf{x}_i, \mathbf{x}_j) = \langle \Phi(\mathbf{x}_i), \Phi(\mathbf{x}_j) \rangle$ for some associated nonlinear mapping $\Phi(\cdot)$ and inner product $\langle \cdot, \cdot \rangle$ in an abstract Hilbert space.

We compute $\boldsymbol{\alpha}^*$ by solving

$$\begin{aligned}
 \text{(KDSM)} \quad & \underset{\boldsymbol{\alpha} \in \mathbb{R}^n}{\text{maximize}} && -\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j k(\mathbf{x}_i, \mathbf{x}_j) + \sum_{i=1}^n \alpha_i \\
 & \text{subject to} && \sum_{i=1}^n \alpha_i y_i = 0 \\
 & && 0 \leq \alpha_i \leq C/n, \quad i = 1, \dots, n.
 \end{aligned}$$

Notice that no matter what $k(\cdot, \cdot)$ is, this is an optimization program in \mathbb{R}^n . Given a solution $\boldsymbol{\alpha}^*$ to the above, we define our classifier as

$$h^*(\mathbf{x}) = \text{sign} \left(\sum_{i=1}^n \alpha_i^* y_i k(\mathbf{x}_i, \mathbf{x}) + b^* \right), \quad (2)$$

where b^* is defined as in (1). Notice that back in the feature space \mathbb{R}^d , this classifier is in general *nonlinear*, as the set

$$\left\{ \mathbf{x} \in \mathbb{R}^d : \sum_{i=1}^n \alpha_i^* y_i k(\mathbf{x}_i, \mathbf{x}) = b^* \right\},$$

is no longer a hyperplane. It is some smooth surface that depends on the properties of $k(\cdot, \cdot)$ (and of course the data).

The term *support vector machine* means solving the optimization program (KDSM) to implement the classifier (2). As we dig a little deeper into how the dual is derived, we will see where this term “support vector” comes from.

General approach to “kernelization”

Optimal soft margin is not the only linear machine learning algorithm that can be kernelized. It can be applied to other classification and regression algorithms as well (we will see more on the regression side of things a little later in the course).

There are three main steps to kernelizing an algorithm:

1. Show that the training process depends only on inner products between input feature vectors, $\mathbf{x}_i^T \mathbf{x}_j$.
2. Show that the applying the decision rule to a general \mathbf{x} only involves computing inner products (i.e. $\mathbf{w}^T \mathbf{x}$). (This means that your classification rule is linear.)
3. Replace all inner products with evaluations of the kernel function $k(\cdot, \cdot)$.

In short, this idea of kernleization has applications well beyond SVMs.

Technical details: Primals and duals in convex programming

For a general convex program² with inequality constraints,

$$\begin{aligned} \text{(P)} \quad & \underset{\mathbf{x}}{\text{minimize}} \quad f(\mathbf{x}) \\ & \text{subject to} \quad g_m(\mathbf{x}) \leq 0, \quad m = 1, \dots, M \end{aligned}$$

the Lagrangian is

$$L(\mathbf{x}, \boldsymbol{\lambda}) = f(\mathbf{x}) + \sum_{m=1}^M \lambda_m g_m(\mathbf{x}).$$

Recall that if \mathbf{x}_0 is feasible, meaning

$$g_m(\mathbf{x}_0) \leq 0, \quad \text{for all } m = 1, \dots, M,$$

then clearly

$$L(\mathbf{x}_0, \boldsymbol{\lambda}) \leq f(\mathbf{x}_0) \quad \text{for all } \boldsymbol{\lambda} \geq \mathbf{0}.$$

And in fact,

$$\max_{\boldsymbol{\lambda} \geq \mathbf{0}} L(\mathbf{x}_0, \boldsymbol{\lambda}) = f(\mathbf{x}_0).$$

Note that if \mathbf{x}' is not feasible ($g_m(\mathbf{x}') > 0$ for some m), then

$$\max_{\boldsymbol{\lambda} \geq \mathbf{0}} L(\mathbf{x}', \boldsymbol{\lambda}) = +\infty.$$

Thus solving (P) is the same as solving

$$\underset{\mathbf{x}}{\text{minimize}} \quad \underset{\boldsymbol{\lambda} \geq \mathbf{0}}{\text{maximize}} \quad L(\mathbf{x}, \boldsymbol{\lambda}).$$

²This means that $f(\cdot)$ and the $g_m(\cdot)$ are convex functions.

The *dual* optimization program switches the order of the min and max, solving

$$(D) \quad \begin{aligned} & \underset{\boldsymbol{\lambda}}{\text{maximize}} && d(\boldsymbol{\lambda}) \\ & \text{subject to} && \boldsymbol{\lambda} \geq \mathbf{0}, \end{aligned}$$

where

$$d(\boldsymbol{\lambda}) = \min_{\boldsymbol{x}} L(\boldsymbol{x}, \boldsymbol{\lambda}).$$

If p^* is the optimal value of the primal program (P) and d^* is the optimal value of the dual (D), then it is always true that $d^* \leq p^*$. Under certain conditions on the constraint functions $g_m(\boldsymbol{x})$, we have strong duality, meaning $d^* = p^*$. One example of these conditions is that the $g_m(\boldsymbol{x})$ are affine³, which is the case for optimal soft margin.

In addition, when the $g_m(\boldsymbol{x})$ are affine, the KKT conditions are necessary and sufficient. That is, any solution \boldsymbol{x}^* to the primal (P) and any solution $\boldsymbol{\lambda}^*$ to the dual (D) will obey:

$$\begin{aligned} g_m(\boldsymbol{x}^*) &\leq 0, & m = 1, \dots, M, \\ \lambda_m &\geq 0, & m = 1, \dots, M, \\ \lambda_m g_m(\boldsymbol{x}^*) &= 0, & m = 1, \dots, M, \\ \nabla f(\boldsymbol{x}^*) + \sum_{m=1}^M \lambda_m^* \nabla g_m(\boldsymbol{x}^*) &= \mathbf{0}. \end{aligned}$$

³This means they can be written as $g_m(\boldsymbol{x}) = \boldsymbol{a}_m^T \boldsymbol{x} + b_m$ for some vector \boldsymbol{a}_m and scalar b_m .

Example: The dual of a QP

We know that optimal soft margin is a quadratic program with linear inequality constraints (QP). So for a quick example of how to derive a dual function, let's look at the general QP:

$$\begin{aligned} \text{(QP)} \quad & \underset{\mathbf{z} \in \mathbb{R}^N}{\text{minimize}} && \frac{1}{2} \mathbf{z}^T \mathbf{P} \mathbf{z} + \mathbf{c}^T \mathbf{z} \\ & \text{subject to} && \mathbf{A} \mathbf{z} \leq \mathbf{q}. \end{aligned}$$

Above, \mathbf{A} is an $M \times N$ matrix, and $\mathbf{q} \in \mathbb{R}^M$. The expression $\mathbf{A} \mathbf{z} \leq \mathbf{q}$ is a compact way to write the M affine inequality constraints

$$g_m(\mathbf{z}) = \mathbf{a}_m^T \mathbf{z} - q_m \leq 0,$$

where \mathbf{a}_m^T is m th row of \mathbf{A} , and q_m is the m th entry in \mathbf{q} . For (QP) to be convex, the $N \times N$ matrix \mathbf{P} needs to be symmetric positive semidefinite (i.e. all of its eigenvalues are ≥ 0). In this example, we will also assume \mathbf{P} is invertible — the discussion below can be extended to the general case, but this assumption just makes the derivation cleaner.

The Lagrangian is

$$L(\mathbf{z}, \boldsymbol{\lambda}) = \frac{1}{2} \mathbf{z}^T \mathbf{P} \mathbf{z} + \mathbf{c}^T \mathbf{z} + \sum_{m=1}^M \lambda_m (\mathbf{a}_m^T \mathbf{z} - q_m).$$

Notice that we can write the last term above more compactly as

$$\sum_{m=1}^M \lambda_m (\mathbf{a}_m^T \mathbf{z} - q_m) = \boldsymbol{\lambda}^T (\mathbf{A} \mathbf{z} - \mathbf{q}).$$

This is a convex quadratic function, so finding its unconstrained minimum amounts to finding a point where the gradient is equal to zero.

Since

$$\nabla_{\mathbf{z}}L(\mathbf{z}, \boldsymbol{\lambda}) = \mathbf{P}\mathbf{z} + \mathbf{c} + \mathbf{A}^T\boldsymbol{\lambda},$$

we know that $\min_{\mathbf{z}} L(\mathbf{z}, \boldsymbol{\lambda})$ will be achieved when

$$\mathbf{P}\mathbf{z} + \mathbf{c} + \mathbf{A}^T\boldsymbol{\lambda} = \mathbf{0} \quad \Rightarrow \quad \mathbf{z} = -\mathbf{P}^{-1}(\mathbf{c} + \mathbf{A}^T\boldsymbol{\lambda}).$$

Thus the dual function is

$$\begin{aligned} d(\boldsymbol{\lambda}) &= \frac{1}{2}(\mathbf{c} + \mathbf{A}^T\boldsymbol{\lambda})^T \mathbf{P}^{-1}(\mathbf{c} + \mathbf{A}^T\boldsymbol{\lambda}) - \mathbf{c}^T \mathbf{P}^{-1}(\mathbf{c} + \mathbf{A}^T\boldsymbol{\lambda}) \\ &\quad + \boldsymbol{\lambda}^T(-\mathbf{A}\mathbf{P}^{-1}(\mathbf{c} + \mathbf{A}^T\boldsymbol{\lambda}) - \mathbf{q}) \\ &= -\frac{1}{2}(\mathbf{c} + \mathbf{A}^T\boldsymbol{\lambda})^T \mathbf{P}^{-1}(\mathbf{c} + \mathbf{A}^T\boldsymbol{\lambda}) - \mathbf{q}^T\boldsymbol{\lambda}, \end{aligned}$$

and the dual optimization program is

$$\begin{aligned} \text{(DQP)} \quad & \underset{\boldsymbol{\lambda} \in \mathbb{R}^M}{\text{maximize}} \quad -\frac{1}{2}(\mathbf{c} + \mathbf{A}^T\boldsymbol{\lambda})^T \mathbf{P}^{-1}(\mathbf{c} + \mathbf{A}^T\boldsymbol{\lambda}) - \mathbf{q}^T\boldsymbol{\lambda} \\ & \text{subject to} \quad \boldsymbol{\lambda} \geq \mathbf{0}. \end{aligned}$$

This is also a (concave) quadratic program with linear inequality constraints.

The dual of optimal soft margin

Now let's derive the dual of the optimal soft margin program

$$\begin{aligned} \text{(SM)} \quad & \underset{b, \mathbf{w}, \boldsymbol{\xi}}{\text{minimize}} \quad \frac{1}{2} \|\mathbf{w}\|_2^2 + \frac{C}{n} \sum_{i=1}^n \xi_i \\ & \text{subject to} \quad y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 - \xi_i, \quad i = 1, \dots, n. \\ & \quad \quad \quad \xi_i \geq 0, \quad i = 1, \dots, n. \end{aligned}$$

We could put this in the form of a QP and use our result above, but it is just as easy (and more enlightening) to compute the dual directly. We will use two sets of Lagrange multipliers, $\boldsymbol{\alpha}, \boldsymbol{\beta} \in \mathbb{R}^n$ for the two sets of linear inequality constraints above. The Lagrangian is

$$L = \frac{1}{2} \|\mathbf{w}\|_2^2 + \frac{C}{n} \sum_{i=1}^n \xi_i + \sum_{i=1}^n \alpha_i (1 - \xi_i - y_i \mathbf{w}^\top \mathbf{x}_i - y_i b) - \sum_{i=1}^n \beta_i \xi_i.$$

This function is smooth in all of its variables $(b, \mathbf{w}, \boldsymbol{\xi}, \boldsymbol{\alpha}, \boldsymbol{\beta})$; it will be minimized in the primal variables $b, \mathbf{w}, \boldsymbol{\xi}$ when⁴

$$\begin{aligned} \nabla_b L &= - \sum_{i=1}^n \alpha_i y_i = 0, \\ \nabla_{\mathbf{w}} L &= \mathbf{w} - \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i = \mathbf{0}, \\ \nabla_{\boldsymbol{\xi}} L &= \frac{C}{n} \mathbf{1} - \boldsymbol{\alpha} - \boldsymbol{\beta} = \mathbf{0}. \end{aligned}$$

The last relation above means we can eliminate the second set of lagrange multipliers by substituting $\beta_i = C/n - \alpha_i$. Plugging these

⁴In this expression, $\mathbf{1}$ is a vector with every entry equal to 1.

relations in the Lagrangian gives us the dual function

$$\begin{aligned}
 d(\boldsymbol{\alpha}) &= \frac{1}{2} \left\| \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i \right\|_2^2 + \sum_{i=1}^n \alpha_i - \sum_{i=1}^n \alpha_i y_i \left(\sum_{j=1}^n \alpha_j y_j \mathbf{x}_j \right)^T \mathbf{x}_i \\
 &= -\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j + \sum_{i=1}^n \alpha_i,
 \end{aligned}$$

where this simplification implicitly incorporates the constraint that $-\sum_{i=1}^n \alpha_i y_i = 0$.

Finally, this gives us the dual program

$$\begin{aligned}
 \text{(DSM)} \quad & \underset{\boldsymbol{\alpha} \in \mathbb{R}^n}{\text{maximize}} && -\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j + \sum_{i=1}^n \alpha_i \\
 & \text{subject to} && \sum_{i=1}^n \alpha_i y_i = 0 \\
 & && 0 \leq \alpha_i \leq C/n, \quad i = 1, \dots, n.
 \end{aligned}$$

Support vectors

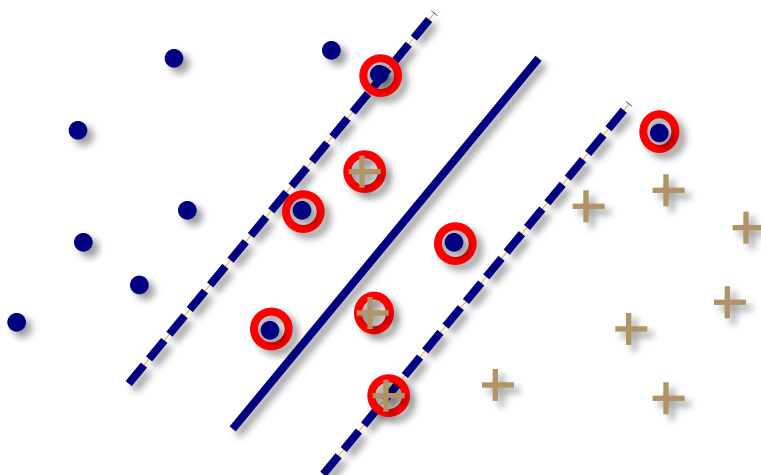
From the complementary slackness conditions ($\lambda_m g_m(\mathbf{x}^*) = 0$), we know that

$$\alpha_i^* (1 - \xi_i^* - y_i(\mathbf{w}^{*\top} \mathbf{x}_i + b^*)) = 0, \quad i = 1, \dots, n.$$

Thus⁵

$$\alpha_i^* > 0 \quad \Leftrightarrow \quad \mathbf{w}^{*\top} \mathbf{x}_i + b^* = y_i(1 - \xi_i^*).$$

This means the i for which $\alpha_i > 0$ correspond to data point \mathbf{x}_i that are on or inside the margin of separation.



These are called the **support vectors**, and in practice, there tend to be a small number of them. One consequence for this is that the weights can be calculated from just this small portion of the data. Recall that we derived the following expression:

$$\mathbf{w}^* = \sum_{i=1}^n \alpha_i^* y_i \mathbf{x}_i = \sum_{i \in \text{SV}} \alpha_i^* y_i \mathbf{x}_i, \quad \text{where SV} = \text{set of support vectors}$$

⁵Note that since $y_i \in \{-1, +1\}$, $y_i^{-1} = y_i$.

If the set of support vectors is small, then computing \mathbf{w}^* is can be very inexpensive.

The other thing that the complementary slackness conditions tell us is that

$$\left(\frac{C}{n} - \alpha_i^*\right) \xi_i^* = 0 \quad \text{for all } i = 1, \dots, n.$$

This means that $\alpha_i^* < C/n \Rightarrow \xi_i^* = 0$. Combining this with the condition above means that

$$0 < \alpha_i^* < C/n \quad \Rightarrow \quad b^* = y_i - \mathbf{w}^{*\top} \mathbf{x}_i.$$

This gives us the relation that allows us to recover the affine offset from the dual solution.