

Lecture 32 - Hard-Margin SVM; Slack Variables, Soft-Margin SVM

Hard-Margin Support Vector Machine (SVM): Separable Classes

Let's start with the two-class linearly separable task and then we will extend the method to more general cases where data are not separable. Let $\phi(x_i)$, $i = 1, 2, \dots, N$, be the feature vectors of the training set, X , and corresponding target values t_1, t_2, \dots, t_N where $t_n \in \{-1, 1\}$. These belong to either of two classes, C_1, C_2 , which are *assumed to be linearly separable*.

The goal, once more, is to design a hyperplane

$$y(x) = w^T \phi(x) + b = 0$$

that classifies correctly all the training vectors.

Because the training data is linearly separable in the feature space, by definition there exists at least one choice of the parameters w and b such that $y(x)$ satisfies $y(x_n) > 0$ for points having $t_n = +1$ and $y(x_n) < 0$ for points having $t_n = -1$, so that $t_n y(x_n) > 0$ for all training data points.

Such a hyperplane is not unique. The perceptron algorithm may converge to any one of the possible solutions. Having gained in experience, this time we will be more demanding.

- Which hyperplane would any sensible engineer choose as the classifier for operation in practice, where data outside the training set will be fed to it?
- Once again, the hyperplane that leaves more "room" on either side, so that data in both classes can move a bit more freely, with less risk of causing an error.
- A sensible choice for the hyperplane classifier would be the one that leaves the maximum margin from both classes.
- Thus such a hyperplane can be trusted more, when it is faced with the challenge of operating with unknown data.
 - It has a higher generalization performance.

Recall that the perpendicular distance of a point x from a hyperplane defined by $y(x) = 0$ where $y(x)$ takes the form $y(x) = w^T \phi(x) + b$ is given by $\frac{|y(x)|}{\|w\|}$.

Furthermore, we are only interested in solutions for which all data points are correctly classified, so that $t_n y(x_n) > 0, \forall n$. Thus the distance of a point x_n to the decision surface is given by

$$\frac{t_n y(x_n)}{\|w\|} = \frac{t_n (w^T \phi(x_n) + b)}{\|w\|}$$

The margin is given by the perpendicular distance to the closest point x_n from the data set, and we wish to optimize the parameters w and b in order to maximize this distance. Thus the maximum margin solution is found by solving

$$\arg_{w,b} \max \left\{ \frac{1}{\|w\|} \min_n [t_n(w^T \phi(x_n) + b)] \right\}$$

where the factor $\frac{1}{\|w\|}$ is taken outside the optimization over n because w does not depend on n .

- Direct solution of this optimization problem would be very complex, and so we shall convert it into an equivalent problem that is much easier to solve.
- To do this we note that if we make the rescaling $w \rightarrow \kappa w$ and $b \rightarrow \kappa b$, then the distance from any point x_n to the decision surface, given by $\frac{t_n y(x_n)}{\|w\|}$, is unchanged. We can use this freedom to set

$$t_n(w^T \phi(x_n) + b) = 1$$

for the point that is closest to the surface, called the **support vectors**. In this case, all data points will satisfy the constraints

$$t_n(w^T \phi(x_n) + b) \geq 1, n = 1, 2, \dots, N$$

- This is known as the canonical representation of the decision hyperplane.
- In the case of data points for which the equality holds (support vectors), the constraints are said to be active, whereas for the remainder they are said to be inactive.
- By definition, there will always be at least one active constraint, because there will always be a closest point, and once the margin has been maximized there will be at least two active constraints.

The optimization problem then simply requires that we maximize $\|w\|^{-1}$, which is equivalent to minimizing $\|w\|^2$, and so we have to solve the optimization problem

$$\arg_{w,b} \min \frac{1}{2} \|w\|^2 \tag{1}$$

$$\text{subject to } t_n(w^T \phi(x_n) + b) \geq 1 \tag{2}$$

In order to solve this constrained optimization problem, we introduce *Lagrange multipliers* $a_n \geq 0$, with one multiplier a_n for each of the constraints, giving the Lagrangian function

$$L(w, b, a) = \frac{1}{2} \|w\|^2 - \sum_{n=1}^N a_n (t_n(w^T \phi(x_n) + b) - 1)$$

- Note the minus sign in front of the Lagrange multiplier term, because we are minimizing with respect to w and b , and maximizing with respect to a .

Setting the derivatives of $L(w, b, a)$ with respect to w and b equal to zero, we obtain the following two conditions:

$$w = \sum_{n=1}^N a_n t_n \phi(x_n)$$

$$0 = \sum_{n=1}^N a_n t_n$$

Eliminating w and b from $L(w, b, a)$ using these conditions then gives the *dual representation* of the maximum margin problem in which we maximize

$$\tilde{L}(a) = \sum_{n=1}^N a_n - \sum_{n=1}^N \sum_{m=1}^N a_n a_m t_n t_m k(x_n, x_m)$$

with respect to a subject to the constraints

$$a_n \geq 0, n = 1, 2, \dots, N$$

$$\sum_{n=1}^N a_n t_n = 0$$

Here the kernel function is defined by $k(x, y) = \phi(x)^T \phi(y)$.

- This takes the form of a quadratic programming problem in which we optimize a quadratic function of a subject to a set of inequality constraints.
- The solution to a quadratic programming problem in M variables in general has computational complexity that is $O(M^3)$. In going to the dual formulation we have turned the original optimization problem into the dual problem, which has N variables.
 - For a fixed set of basis functions whose number M is smaller than the number N of data points, the move to the dual problem appears disadvantageous. However, it allows the model to be reformulated using kernels, and so the maximum margin classifier can be applied efficiently to feature spaces whose dimensionality exceeds the number of data points, including infinite feature spaces.
 - The kernel formulation also makes clear the role of the constraint that the kernel function $k(x, y)$ be positive definite, because this ensures that the Lagrangian function $\tilde{L}(a)$ is bounded below, giving rise to a well defined optimization problem.

In order to classify new data points using the trained model, we evaluate the sign of $y(x) = w^T \phi(x) + b$. This can be expressed in terms of the parameters $\{a_n\}$ and the kernel function by substituting for w to give

$$y(x) = \sum_{n=1}^N a_n t_n k(x, x_n) + b$$

- A constrained optimization of this form satisfies the Karush-Kuhn-Tucker (KKT) conditions, which in this case require that the following three properties hold

$$a_n \geq 0 \quad (3)$$

$$t_n y(x_n) - 1 \geq 0 \quad (4)$$

$$a_n (t_n y(x_n) - 1) = 0 \quad (5)$$

Thus for every data point, either $a_n = 0$ or $t_n y(x_n) = 1$.

- Any data point for which $a_n = 0$ will not appear in the sum of $y(x)$ and hence plays no role in making predictions for new data points.
- The remaining data points are called **support vectors**, and because they satisfy $t_n y(x_n) = 1$, they correspond to points that lie on the maximum margin hyperplanes in feature space.
- Once the model is trained, a significant proportion of the data points can be discarded and only the support vectors retained!

Having solved the quadratic programming problem and found a value for a , we can then determine the value of the threshold parameter b by noting that any support vector x_n satisfies $t_n y(x_n) = 1$.

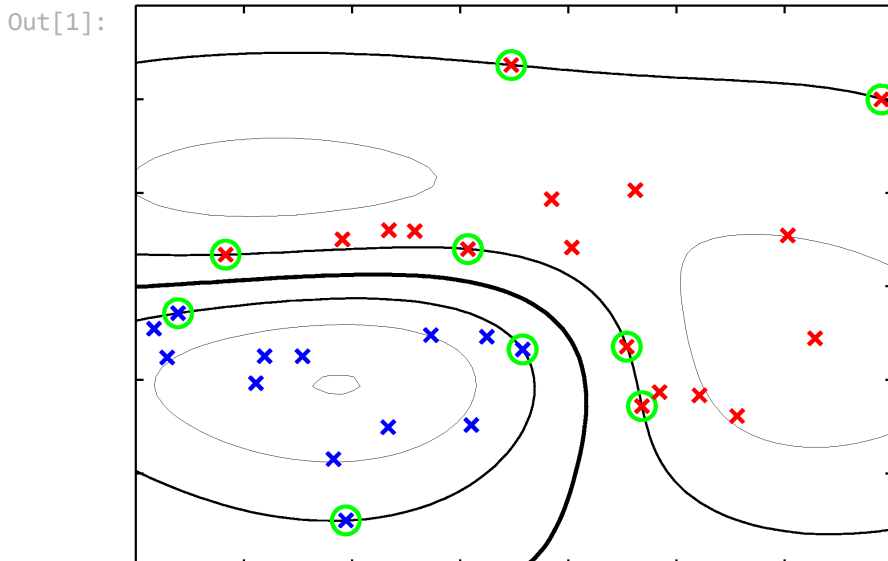
$$t_n \left(\sum_{m \in S} a_m t_m k(x_n, x_m) + b \right) = 1$$

where S denotes the set of indices of the support vectors. Although we can solve this equation for b using an arbitrarily chosen support vector x_n , a numerically more stable solution is obtained by first multiplying through by t_n , making use of $t_n^2 = 1$, and then averaging these equations over all support vectors and solving for b to give

$$b = \frac{1}{N_S} \sum_{n \in S} \left(t_n - \sum_{m \in S} a_m t_m k(x_n, x_m) \right)$$

where N_S is the total number of support vectors.

```
In [1]: from IPython.display import Image
        Image("figures/Figure7.2.png", width=400)
```



PAC Learning & VC Dimension

Historically, support vector machines have largely been motivated and analyzed using a theoretical framework known as *computational learning theory*, also sometimes called *statistical learning theory*. This has its origins with Valiant (1984) who formulated the **probably approximately correct**, or **PAC, learning** framework. The goal of the PAC framework is to understand how large a data set needs to be in order to give good generalization. It also gives bounds for the computational cost of learning.

A key quantity in PAC learning is the Vapnik-Chervonenkis dimension, or **VC dimension**, which provides a measure of the complexity of a space of functions, and which allows the PAC framework to be extended to spaces containing an infinite number of functions.