

Classification with Support Vector Machines

In many situations we want our machine learning algorithm to predict one of a number of outcomes. For example an email client that sorts mail into personal mail and junk mail, which has two outcomes. Another example is a telescope that identifies whether an object in the night sky is a galaxy, star or planet. There are usually a small number of outcomes, and more importantly there is usually no additional structure on these outcomes. In this chapter, we consider predictors that output binary values, that is, there are only two possible outcomes. This is in contrast to Chapter 9 where we considered a prediction problem with continuous-valued outputs. This machine learning task is called *binary classification*. For binary classification the set of possible values that the label/output can attain is binary, and for this chapter we denote them as $\{+1, -1\}$. In other words, we consider predictors of the form

$$f : \mathbb{R}^D \rightarrow \{+1, -1\}. \quad (12.1)$$

Recall from Chapter 8 that we represent each example \mathbf{x}_n as a feature vector of D real numbers. The labels are often referred to as the positive and negative *classes*, respectively. One should be careful not to infer intuitive attributes of positiveness of the $+1$ class. For example, in a cancer detection task, a patient with cancer is often labelled $+1$. In principle, any two distinct values can be used, e.g., $\{\text{True}, \text{False}\}$, $\{0, 1\}$ or $\{\text{red}, \text{blue}\}$. The problem of binary classification is well studied, and we defer a survey of other approaches to Section 12.4.

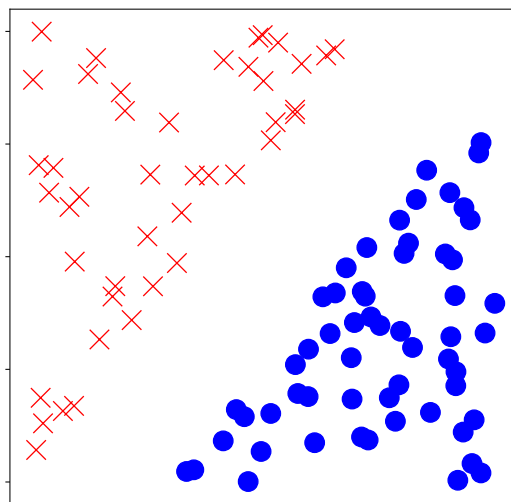
We present an approach known as the Support Vector Machine (SVM), which solves the binary classification task. Similar to regression, we have a supervised learning task, where we have a set of examples $\mathbf{x}_n \in \mathbb{R}^D$ along with their corresponding labels $y_n \in \{+1, -1\}$. Given the training data consisting of example-label pairs $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$, we would like to estimate parameters of the model that will give the best classification error. Similar to Chapter 9 we consider a linear model, and hide away the nonlinearity in a transformation ϕ of the examples (9.12). We will revisit ϕ later in this chapter in Section 12.3.3.

The SVM provides state of the art results in many applications, with sound theoretical guarantees (Steinwart and Christmann, 2008). In this book, the first reason we choose to discuss the SVM is to illustrate a

An example of structure is if the outcomes were ordered, like in the case of small, medium and large t-shirts.
binary classification

Input example \mathbf{x}_n may also be referred to as inputs, data points, features or instances.
classes
For probabilistic models, it is mathematically convenient to use $\{0, 1\}$ as a binary representation. See remark after Example 6.15.

Figure 12.1
Example 2D data,
illustrating the
intuition of data
where we can find a
linear classifier that
separates red
crosses from blue
dots.



geometric way to think about supervised machine learning. Whereas in Chapter 9 we considered the machine learning problem in terms of probabilistic models and attacked it using maximum likelihood estimation and Bayesian inference, here we will consider an alternative approach where we reason geometrically about the machine learning task. It relies heavily on concepts, such as inner products and projections, which we discussed in Chapter 3. In contrast to Chapter 9, the optimization problem for SVM does not admit an analytic solution. Hence, we resort to the optimization tools introduced in Chapter 7. This is the second reason for introducing the SVM – as an illustration of what to do when we cannot analytically derive a solution.

The SVM view of machine learning is also subtly different from the maximum likelihood view of Chapter 9. The maximum likelihood view proposes a model based on a probabilistic view of the data distribution, from which an optimization problem is derived. In contrast, the SVM view starts by designing a particular function that is to be optimized during training, based on geometric intuitions. In other words, we start by designing an objective function that is to be minimized on training data, following the principles of empirical risk minimization 8.1. This can also be understood as designing a particular loss function.

Let us derive the optimization problem corresponding to training an SVM on example-label pairs. Intuitively, we imagine binary classification data which can be separated by a hyperplane as illustrated in Figure 12.1, where the example (a vector of dimension 2) is used to indicate the location and the label is represented as different symbols (and colours). Hyperplane is a word that is commonly used in machine learning, and we saw them in Section 2.8 introduced as an affine subspace, which is the phrase used in linear algebra. The examples consists of two classes that

have features (the components of the vector representing the example) arranged in such a way as to allow us to separate/classify them by drawing a straight line.

In the following, we start by formalizing this idea of finding a linear separator. We introduce the idea of the margin and then extend linear separators to allow for examples to fall on the wrong side. We present two equivalent ways of formalizing the SVM: the geometric view (Section 12.2.4) and the loss function view (Section 12.2.5). We derive the dual version of the SVM using Lagrange multipliers (Section 7.2). The dual SVM allows us to observe a third way of formalizing the SVM: in terms of the convex hulls of the examples of each class (Section 12.3.2). We conclude by briefly describing kernels and how to numerically solve the nonlinear kernel-SVM optimization problem.

12.1 Separating Hyperplanes

Given two examples represented as vectors \mathbf{x}_i and \mathbf{x}_j , one way to compute the similarity between them is using an inner product $\langle \mathbf{x}_i, \mathbf{x}_j \rangle$. Recall from Section 3.2 that inner products measure the angle between two vectors. The value of the inner product also depends on the length (norm) of each vector. Furthermore, inner products allow us to rigorously define geometric concepts such as orthogonality and projections.

The main idea behind many classification algorithms is to represent data in \mathbb{R}^D and then partition this space. In the case of binary classification, the space would be split into two parts corresponding to the positive and negative classes, respectively. We consider a particularly convenient partition, which is to split the space into two halves using a hyperplane. Let example $\mathbf{x} \in \mathbb{R}^D$ be an element of the data space. Consider a function $f : \mathbb{R}^D \rightarrow \mathbb{R}$ parametrized by $\mathbf{w} \in \mathbb{R}^D$ and $b \in \mathbb{R}$ as follows

$$f(\mathbf{x}) = \langle \mathbf{w}, \mathbf{x} \rangle + b. \quad (12.2)$$

Recall from Section 2.8 that hyperplanes are affine subspaces. Therefore we define the hyperplane that separates the two classes in our binary classification problem as

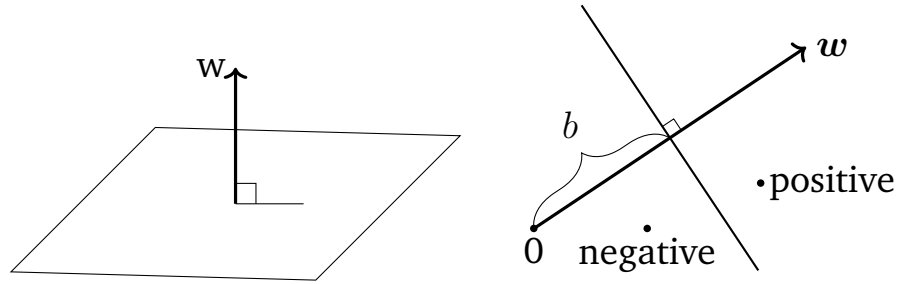
$$\{\mathbf{x} \in \mathbb{R}^D : f(\mathbf{x}) = 0\}. \quad (12.3)$$

An illustration of the hyperplane is shown in Figure 12.2 where the vector \mathbf{w} is a vector normal to the hyperplane and b the intercept. We can derive that \mathbf{w} is a normal vector to the hyperplane in (12.3) by choosing any two examples \mathbf{x}_a and \mathbf{x}_b on the hyperplane and showing that the vector between them is orthogonal to \mathbf{w} . In the form of an equation,

$$f(\mathbf{x}_a) - f(\mathbf{x}_b) = \langle \mathbf{w}, \mathbf{x}_a \rangle + b - (\langle \mathbf{w}, \mathbf{x}_b \rangle + b) \quad (12.4)$$

$$= \langle \mathbf{w}, \mathbf{x}_a - \mathbf{x}_b \rangle, \quad (12.5)$$

Figure 12.2
Equation of a
separating
hyperplane (12.3).
(left) The standard
way of representing
the equation in 3D.
(right) For ease of
drawing, we look at
the hyperplane edge
on.



where the second line is obtained by the linearity of the inner product (Section 3.2). Since we have chosen x_a and x_b to be on the hyperplane, this implies that $f(x_a) = 0$ and $f(x_b) = 0$ and hence $\langle w, x_a - x_b \rangle = 0$. Recall that two vectors are orthogonal when their inner product is zero, therefore we obtain that w is orthogonal to any vector on the hyperplane.

Remark. Recall from Chapter 1 that we can think of vectors in different ways. In this chapter, we think of the parameter vector w as an arrow indicating a direction. That is we consider w to be a geometric vector. In contrast, we think of the example vector x as a point (as indicated by its coordinates). That is we consider x to be the coordinates of a vector with respect to the standard basis. \diamond

When presented with a test example, we classify the example as positive or negative by deciding on which side of the hyperplane it occurs. Note that (12.3) not only defines a hyperplane, it additionally defines a direction. In other words, it defines the positive and negative side of the hyperplane. Therefore, to classify a test example x_{test} , we calculate the value of the function $f(x_{\text{test}})$ and classify the example as $+1$ if $f(x_{\text{test}}) \geq 0$ and -1 otherwise. Thinking geometrically, the positive examples lie “above” the hyperplane and the negative examples “below” the hyperplane.

When training the classifier, we want to ensure that the examples with positive labels are on the positive side of the hyperplane, i.e.,

$$\langle w, x_n \rangle + b \geq 0 \quad \text{when} \quad y_n = +1 \quad (12.6)$$

and the examples with negative labels are on the negative side,

$$\langle w, x_n \rangle + b < 0 \quad \text{when} \quad y_n = -1. \quad (12.7)$$

Refer to Figure 12.2 for a geometric intuition of positive and negative examples. These two conditions are often presented in a single equation,

$$y_n(\langle w, x_n \rangle + b) \geq 0. \quad (12.8)$$

The equation (12.8) above is equivalent to (12.6) and (12.7) when we multiply both sides of (12.6) and (12.7) with $y_n = 1$ and $y_n = -1$, respectively.

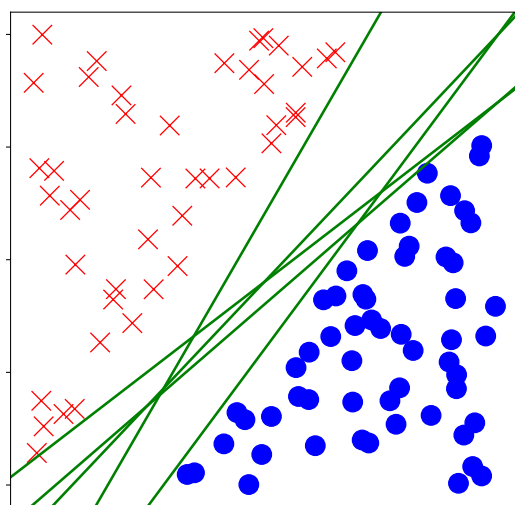


Figure 12.3
Possible separating hyperplanes. There are many linear classifiers (green lines) that separate red crosses from blue dots.

12.2 Primal Support Vector Machine

Based on the concept of distances from points to a hyperplane, we now are in a position to discuss the support vector machine. For a dataset $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$ that is linearly separable, we have many candidate hyperplanes (refer to Figure 12.3) that solve our classification problem without any (training) errors. In other words, for a given training set we have many candidate classifiers. One idea is to choose the separating hyperplane that maximizes the margin between the positive and negative examples. In the following, we use the concept of a hyperplane, see also Section 2.8, and derive the distance between an example and a hyperplane. Recall that the closest point on the hyperplane to a given point (example \mathbf{x}_n) is obtained by the orthogonal projection (Section 3.7). We will see in the next section how to use the orthogonal projection to derive the margin.

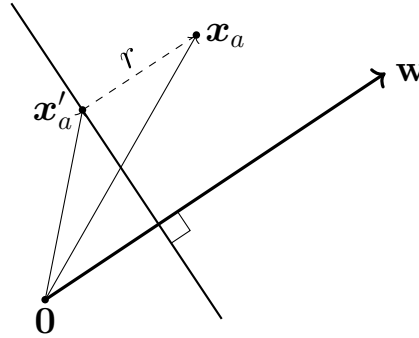
A classifier with large margin turns out to generalize well (Steinwart and Christmann, 2008).

12.2.1 Concept Of The Margin

The concept of the *margin* is intuitively simple: It is the distance of the separating hyperplane to the closest examples in the dataset, assuming that the dataset is linearly separable. However, when trying to formalize this distance, there is a technical wrinkle that is confusing. The technical wrinkle is that we need to define a scale at which to measure the distance. A potential scale is to consider the scale of the data, i.e., the raw values of \mathbf{x}_n . There are problems with this, as we could change the units of measurement of \mathbf{x}_n and change the values in \mathbf{x}_n , and, hence, change the distance to the hyperplane. As we will see shortly, we define the scale based on the equation of the hyperplane (12.3) itself.

margin
There could be two or more closest examples to a hyperplane.

Figure 12.4 Vector addition to express distance to hyperplane:
 $\mathbf{x}_a = \mathbf{x}'_a + r \frac{\mathbf{w}}{\|\mathbf{w}\|}$.



Consider a hyperplane $\langle \mathbf{w}, \mathbf{x} \rangle + b$, and an example \mathbf{x}_a as illustrated in Figure 12.4. Without loss of generality, we can consider the example \mathbf{x}_a to be on the positive side of the hyperplane, i.e., $\langle \mathbf{w}, \mathbf{x}_a \rangle + b > 0$. We would like to derive the distance $r > 0$ of \mathbf{x}_a from the hyperplane. We do so by considering the orthogonal projection (Section 3.7) of \mathbf{x}_a onto the hyperplane, which we denote by \mathbf{x}'_a . Since \mathbf{w} is orthogonal to the hyperplane, we know that the distance r is just a scaling of this vector \mathbf{w} . However, we need to use a vector of unit length (its norm must be 1), and obtain this by dividing \mathbf{w} by its norm, $\frac{\mathbf{w}}{\|\mathbf{w}\|}$. Using vector addition (Section 2.4) we obtain

$$\mathbf{x}_a = \mathbf{x}'_a + r \frac{\mathbf{w}}{\|\mathbf{w}\|}. \quad (12.9)$$

Another way of thinking about r is that it is the coordinate of \mathbf{x}_a in the subspace spanned by \mathbf{w} . We have now expressed the distance of \mathbf{x}_a from the hyperplane as r , and if we choose \mathbf{x}_a to be the point closest to the hyperplane, this distance r is the margin.

Recall that we would like the positive examples to be further than r from the hyperplane, and the negative examples to be further than distance r (in the negative direction) from the hyperplane. Analogously to the combination of (12.6) and (12.7) into (12.8), we have

$$y_n(\langle \mathbf{w}, \mathbf{x}_n \rangle + b) \geq r. \quad (12.10)$$

In other words we can combine the requirements that examples are further than r from the hyperplane (in the positive and negative direction) into one single inequality.

Since we are interested only in the direction, we add an assumption to our model that the parameter vector \mathbf{w} is of unit length, that is, $\|\mathbf{w}\| = 1$ where we use the Euclidean norm $\|\mathbf{w}\| = \sqrt{\mathbf{w}^\top \mathbf{w}}$ (Section 3.1). Collecting the three requirements into one constrained optimization problem, we

A reader familiar with other presentations of the margin would notice that our definition of $\|\mathbf{w}\| = 1$ is different from the presentation in for example Schölkopf and Smola (2002). We will show that the two approaches are equivalent in Section 12.2.3.

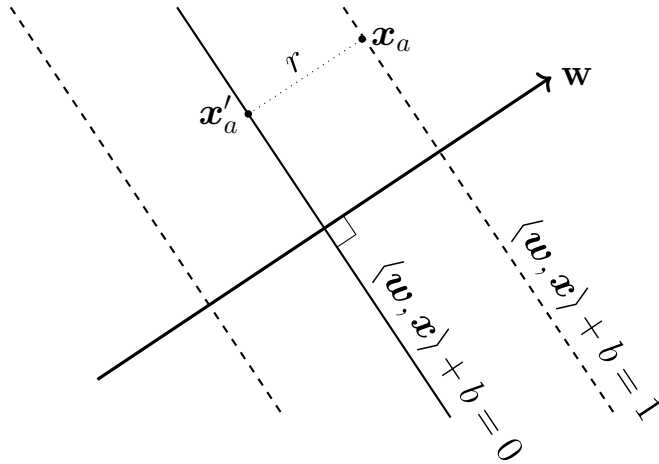


Figure 12.5
Derivation of the
margin: $r = \frac{1}{\|\mathbf{w}\|}$.

obtain the following

$$\max_{\mathbf{w}, b, r} \underbrace{r}_{\text{margin}} \quad \text{subject to} \quad \underbrace{y_n(\langle \mathbf{w}, \mathbf{x}_n \rangle + b) \geq r}_{\text{data fitting}}, \quad \underbrace{\|\mathbf{w}\| = 1}_{\text{normalization}}, \quad r > 0, \quad (12.11)$$

which says that we want to maximize the margin r , while ensuring that the data lies on the correct side of the hyperplane.

Remark. The idea of the margin turns out to be highly pervasive in machine learning. It was used by Vladimir Vapnik and Alexey Chervonenkis to show that when the margin is large, the “complexity” of the function class is low, and, hence, learning is possible (Vapnik, 2000). It turns out that the concept is useful for various different approaches for theoretically analyzing generalization error (Shalev-Shwartz and Ben-David, 2014; Steinwart and Christmann, 2008). \diamond

12.2.2 Traditional Derivation Of The Margin

In the previous section, we derived (12.11) by making the observation that we are only interested in the direction of \mathbf{w} and not its length, leading to the assumption that $\|\mathbf{w}\| = 1$. In this section, we derive the margin maximization problem by making a different assumption. Instead of choosing that the parameter vector is normalised, we choose a scale for the data. We choose this scale such that the value of the predictor $\langle \mathbf{w}, \mathbf{x} \rangle + b$ is 1 at the closest example. Let us also consider \mathbf{x}_a to be the example in the dataset that is closest to the hyperplane.

Figure 12.5 is the same as Figure 12.4, except that now we have rescaled the axes, such that we have the example \mathbf{x}_a exactly on the margin, i.e., $\langle \mathbf{w}, \mathbf{x}_a \rangle + b = 1$. Since \mathbf{x}'_a is the orthogonal projection of \mathbf{x}_a onto the

Recall that we currently consider linearly separable data.

hyperplane, it must by definition lie on the hyperplane, i.e.,

$$\langle \mathbf{w}, \mathbf{x}'_a \rangle + b = 0. \quad (12.12)$$

By substituting (12.9) into (12.12) we obtain

$$\left\langle \mathbf{w}, \mathbf{x}_a - r \frac{\mathbf{w}}{\|\mathbf{w}\|} \right\rangle + b = 0. \quad (12.13)$$

Multiplying out the inner product, we get

$$\langle \mathbf{w}, \mathbf{x}_a \rangle + b - r \frac{\langle \mathbf{w}, \mathbf{w} \rangle}{\|\mathbf{w}\|} = 0, \quad (12.14)$$

where we exploited the linearity of the inner product (see Section 3.2). Observe that the first term is unity by our assumption of scale, that is, $\langle \mathbf{w}, \mathbf{x}_a \rangle + b = 1$. From (3.18) in Section 3.1 we recall that $\langle \mathbf{w}, \mathbf{w} \rangle = \|\mathbf{w}\|^2$, and hence the second term reduces to $r\|\mathbf{w}\|$. Using these simplifications, we obtain

$$r = \frac{1}{\|\mathbf{w}\|}, \quad (12.15)$$

We can also think of the distance as the projection error that incurs when projecting \mathbf{x}_a onto the hyperplane. where we have derived the distance r in terms of the normal vector \mathbf{w} of the hyperplane. At first glance this equation is counterintuitive as we seem to have derived the distance from the hyperplane in terms of the length of the vector \mathbf{w} , but we do not yet know this vector. One way to think about it is to consider the distance r to be a temporary variable that we only use for this derivation. In fact, for the rest of this section we will refer to the distance to the hyperplane by $\frac{1}{\|\mathbf{w}\|}$. In Section 12.2.3 we will see that the choice that the margin equals 1 is equivalent to our previous assumption of $\|\mathbf{w}\| = 1$ in Section 12.2.1.

Similar to the argument to obtain (12.10), we want the positive examples to be further than 1 from the hyperplane, and the negative examples to be further than distance 1 (in the negative direction) from the hyperplane

$$y_n(\langle \mathbf{w}, \mathbf{x}_n \rangle + b) \geq 1. \quad (12.16)$$

Combining the margin maximization with the fact that examples need to be on the correct side of the hyperplane (based on their labels) gives us

$$\max_{\mathbf{w}, b} \frac{1}{\|\mathbf{w}\|} \quad (12.17)$$

$$\text{subject to } y_n(\langle \mathbf{w}, \mathbf{x}_n \rangle + b) \geq 1 \quad \text{for all } n = 1, \dots, N. \quad (12.18)$$

Instead of maximizing the reciprocal of the norm as in (12.17), we often minimize the squared norm. We also often include a constant $\frac{1}{2}$ that does not affect the optimal \mathbf{w}, b but yields a tidier form when we take the derivative. Then, our objective becomes

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 \quad (12.19)$$

The squared norm results in a convex quadratic programming problem for the SVM (Section 12.3.4).

$$\text{subject to } y_n(\langle \mathbf{w}, \mathbf{x}_n \rangle + b) \geq 1 \quad \text{for all } n = 1, \dots, N. \quad (12.20)$$

Equation (12.19) is known as the *hard margin SVM*. The reason for the expression “hard” is because the above formulation does not allow for any violations of the margin condition. We will see in Section 12.2.4 that this “hard” condition can be relaxed to accommodate violations.

hard margin SVM

12.2.3 Why We Can Set The Margin To 1

In Section 12.2.1 we argue that we would like to maximize some value r , which represents the distance of the closest example to the hyperplane. In Section 12.2.2 we scaled the data such that the closest example is of distance 1 to the hyperplane. Here we relate the two derivations, and show that they are actually equivalent.

Theorem 12.1. *Maximizing the margin r where we consider normalized weights as in (12.11),*

$$\max_{\mathbf{w}, b, r} \underbrace{r}_{\text{margin}} \quad \text{subject to} \quad \underbrace{y_n(\langle \mathbf{w}, \mathbf{x}_n \rangle + b) \geq r}_{\text{data fitting}}, \quad \underbrace{\|\mathbf{w}\| = 1}_{\text{normalization}}, \quad r > 0 \quad (12.21)$$

is equivalent to scaling the data such that the margin is unity

$$\min_{\mathbf{w}, b} \underbrace{\frac{1}{2} \|\mathbf{w}\|^2}_{\text{margin}} \quad \text{subject to} \quad \underbrace{y_n(\langle \mathbf{w}, \mathbf{x}_n \rangle + b) \geq 1}_{\text{data fitting}}. \quad (12.22)$$

Proof Consider (12.21), and note that because the square is a monotonic transformation for non-negative arguments, the maximum stays the same if we consider r^2 in the objective. Since $\|\mathbf{w}\| = 1$ we can reparameterize the equation with a new weight vector \mathbf{w}' that is not normalized by explicitly using $\frac{\mathbf{w}'}{\|\mathbf{w}'\|}$,

$$\max_{\mathbf{w}', b, r} r^2 \quad \text{subject to} \quad y_n \left(\left\langle \frac{\mathbf{w}'}{\|\mathbf{w}'\|}, \mathbf{x}_n \right\rangle + b \right) \geq r, \quad r > 0. \quad (12.23)$$

In (12.23) we have explicitly written that distances are non-negative. We can divide the first constraint by r ,

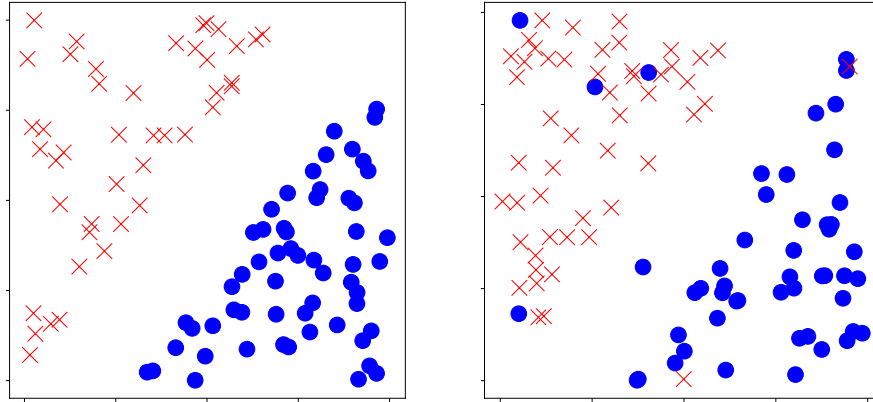
Note that $r > 0$ because we assumed linear separability, and hence there is no issue to divide by r .

$$\max_{\mathbf{w}', b, r} r^2 \quad \text{subject to} \quad y_n \left(\underbrace{\left\langle \frac{\mathbf{w}'}{\|\mathbf{w}'\|}, \mathbf{x}_n \right\rangle}_{\mathbf{w}''} + \underbrace{\frac{b}{r}}_{b''} \right) \geq 1, \quad r > 0 \quad (12.24)$$

renaming the parameters to \mathbf{w}'' and b'' . Since $\mathbf{w}'' = \frac{\mathbf{w}'}{\|\mathbf{w}'\| r}$, rearranging for r gives

$$\|\mathbf{w}''\| = \left\| \frac{\mathbf{w}'}{\|\mathbf{w}'\| r} \right\| = \left| \frac{1}{r} \right| \cdot \left\| \frac{\mathbf{w}'}{\|\mathbf{w}'\|} \right\| = \frac{1}{r}. \quad (12.25)$$

Figure 12.6 (left) linearly separable data, with a large margin. (right) non-separable data.



Substituting into (12.24), we obtain

$$\max_{\mathbf{w}'', b''} \frac{1}{\|\mathbf{w}''\|^2} \quad \text{subject to} \quad y_n (\langle \mathbf{w}'', \mathbf{x}_n \rangle + b'') \geq 1. \quad (12.26)$$

6364 The final step is to observe that maximizing $\frac{1}{\|\mathbf{w}''\|^2}$ yields the same solution
 6365 as minimizing $\frac{1}{2} \|\mathbf{w}''\|^2$. □

12.2.4 Soft Margin SVM: Geometric View

6366 We may wish to allow some examples to fall within the margin region,
 6367 or even to be on the wrong side of the hyperplane (as illustrated in Fig-
 6368 ure 12.6). This also naturally provides us with an approach that works
 6369 when we do not have linearly separable data.
 6370

soft margin SVM 6371 The resulting model is called the *soft margin SVM*. In this section, we
 6372 derive the resulting optimization problem using geometric arguments. In
 6373 Section 12.2.5, we will derive the same optimization problem using the
 6374 idea of a loss function. Using Lagrange multipliers (Section 7.2), we will
 6375 derive the dual optimization problem of the SVM in Section 12.3. This
 6376 dual optimization problem allows us to observe a third interpretation of
 6377 the SVM, as a hyperplane that bisects the line between convex hulls cor-
 6378 responding to the positive and negative data examples (Section 12.3.2).

slack variable The key geometric idea is to introduce a *slack variable* ξ_n corresponding
 to each example (\mathbf{x}_n, y_n) that allows a particular example to be within the
 margin or even on the wrong side of the hyperplane (refer to Figure 12.7).
 We subtract the value of ξ_n from the margin, constraining ξ_n to be non-
 negative. To encourage correct classification of the samples, we add ξ_n to
 the objective

$$\min_{\mathbf{w}, b, \xi} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{n=1}^N \xi_n \quad (12.27)$$

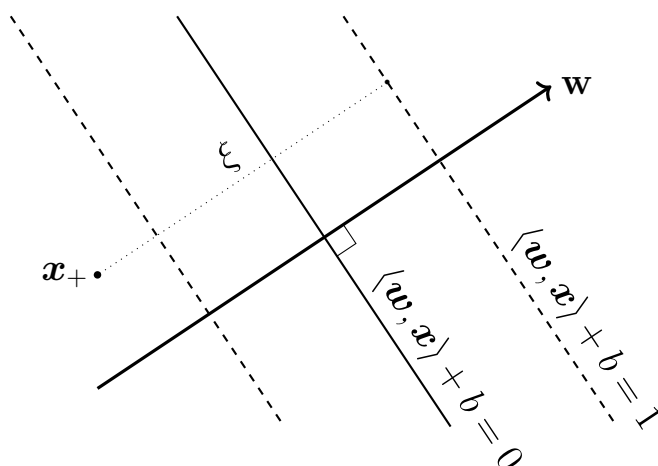


Figure 12.7 Soft Margin SVM allows examples to be within the margin or on the wrong side of the hyperplane. The slack variable ξ measures the distance of a positive example x_+ to the positive margin hyperplane $\langle w, x \rangle + b = 1$ when x_+ is on the wrong side.

$$\text{subject to } y_n(\langle w, x_n \rangle + b) \geq 1 - \xi_n \quad \text{for all } n = 1, \dots, N \quad (12.28)$$

$$\xi_n \geq 0 \quad \text{for all } n = 1, \dots, N. \quad (12.29)$$

In contrast to the optimization problem (12.19) from the previous section (the hard margin SVM), this one is called the *soft margin SVM*. The parameter $C > 0$ trades off the size of the margin and the total amount of slack that we have. This parameter is called the *regularization parameter* since, as we will see in the following section, the margin term in the objective function (12.27) is a regularization term. The margin term $\|w\|^2$ is called the *regularizer*, and in many books on numerical optimization, the regularization parameter multiplied with this term (Section 8.1.3). This is in contrast to our formulation in this section. Some care needs to be taken when interpreting the regularizer, as a large value of C implies low regularization, as we give the slack variables larger weight.

Remark. One detail to note is that in the formulation of the SVM (12.27) w is regularized but b is not regularized. We can see this by observing that the regularization term does not contain b . The unregularized term b complicates theoretical analysis (Steinwart and Christmann, 2008, Chapter 1) and decreases computational efficiency (Fan et al., 2008). \diamond

soft margin SVM

regularization
parameter

regularizer

There are alternative parametrizations of this regularization, which is why (12.27) is also often referred to as the C -SVM.

12.2.5 Soft Margin SVM: Loss Function View

Recall from Section 9.2.1 that when performing maximum likelihood estimation we usually consider the negative log likelihood. Furthermore since the likelihood term for linear regression with Gaussian noise is Gaussian, the negative log likelihood for each example is a squared error function (9.8). The squared error function is the term that is minimized when looking for the maximum likelihood solution. Let us consider the error function point of view, which is also known as the *loss function* point of view.

loss function

In contrast to Chapter 9 where we consider regression problems (the output of the predictor is a real number), in this chapter we consider binary classification problems (the output of the predictor is one of two labels $\{+1, -1\}$). Therefore the error function or the loss function for each single (example, label) pair needs to be appropriate for binary classification. For example, the squared loss that is used for regression (9.9b) is not suitable for binary classification.

Remark. The ideal loss function between binary labels is to count the number of mismatches between the prediction and the label. That is for a predictor f applied to an example \mathbf{x}_n , we compare the output $f(\mathbf{x}_n)$ with the label y_n . We define the loss to be zero if they match, and one if they do not match. This is denoted by $\mathbf{1}(f(\mathbf{x}_n) \neq y_n)$ and is called the zero-one loss. Unfortunately the zero-one loss results in a difficult optimization problem for finding the best parameters \mathbf{w}, b . \diamond

What is the loss function corresponding to the SVM? Consider the error between the output of a predictor $f(\mathbf{x}_n)$ and the label y_n . The loss should capture how much we care about the error that is made on the training data. An equivalent way to derive (12.27) is to use the *hinge loss*

hinge loss

$$\ell(t) = \max\{0, 1 - t\} \quad \text{where} \quad t = yf(\mathbf{x}) = y(\langle \mathbf{w}, \mathbf{x} \rangle + b). \quad (12.30)$$

If $f(\mathbf{x})$ is on the correct side (based on y) of the hyperplane, and further than distance 1, this means that $t \geq 1$ and the hinge loss returns a value of zero. If $f(\mathbf{x})$ is on the correct side but close to the hyperplane, that is, $0 < t < 1$, then the example \mathbf{x} is within the margin and the hinge loss returns a positive value. When the example is on the wrong side of the hyperplane ($t < 0$) the hinge loss returns an even larger value, which increases linearly. In other words, we pay a penalty once we are closer than the margin, even if the prediction is correct, and the penalty increases linearly. An alternative way to express the hinge loss is by considering it as two linear pieces

$$\ell(t) = \begin{cases} 0 & \text{if } t \geq 1 \\ 1 - t & \text{if } t < 1 \end{cases}, \quad (12.31)$$

as illustrated in Figure 12.8. The loss corresponding to the hard margin SVM 12.19 is defined as follows

$$\ell(t) = \begin{cases} 0 & \text{if } t \geq 1 \\ \infty & \text{if } t < 1 \end{cases}. \quad (12.32)$$

This loss can be interpreted as never allowing any examples inside the margin.

For a given training set $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$ we would like to minimize the total loss, while regularizing the objective with ℓ_2 regularization (see Section 8.1.3). Using the hinge loss (12.30) gives us the uncon-

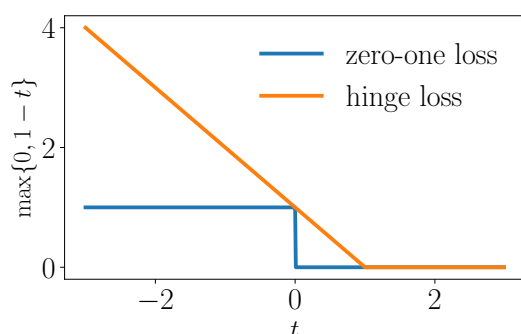


Figure 12.8 Hinge Loss is a convex envelope of zero-one loss.

strained optimization problem

$$\min_{\mathbf{w}, b} \underbrace{\frac{1}{2} \|\mathbf{w}\|^2}_{\text{regularizer}} + C \underbrace{\sum_{n=1}^N \max\{0, 1 - y_n(\langle \mathbf{w}, \mathbf{x}_n \rangle + b)\}}_{\text{error term}}. \quad (12.33)$$

The first term in (12.33) is called the regularization term or the *regularizer* (see Section 9.2.3), and the second term is called the *loss term* or the *error term*. Recall from Section 12.2.4 that the term $\frac{1}{2} \|\mathbf{w}\|^2$ is actually the term arising from the margin. In other words, margin maximization can be interpreted as a regularizer.

In principle, the unconstrained optimization problem in (12.33) can be directly solved with (sub-)gradient descent methods as described in Section 7.1. To see that (12.33) and (12.27) are equivalent, observe that the hinge loss (12.30) essentially consists of two linear parts, as expressed in (12.31). Therefore, we can equivalently replace minimization of the hinge loss with two constraints, i.e.,

$$\min_t \max\{0, 1 - t\} \quad (12.34)$$

is equivalent to

$$\begin{aligned} \min_{\xi, t} \quad & \xi \\ \text{subject to} \quad & \xi \geq 0 \\ & \xi \geq 1 - t. \end{aligned} \quad (12.35)$$

By substituting this into (12.33) and rearranging one of the constraints, we obtain exactly the soft margin SVM (12.27).

12.3 Dual Support Vector Machine

The description of the SVM in the previous sections, in terms of the variables \mathbf{w} and b , is known as the primal SVM. Recall that we are considering

input vectors \mathbf{x} , which have dimension D , i.e., we are looking at input examples with D features. Since \mathbf{w} is of the same dimension as \mathbf{x} , this means that the number of parameters (the dimension of \mathbf{w}) of the optimization problem grows linearly with the number of features.

In the following, we consider an equivalent optimization problem (the so-called dual view) which is independent of the number of features. We see a similar idea appear in Chapter 10 where we express the learning problem in a way that does not scale with the number of features. This is useful for problems where we have more features than the number of examples. Instead the number of parameters increases with the number of examples in the training set. The dual SVM also has the additional advantage that it easily allows kernels to be applied, as we shall see at the end of this chapter. The word “dual” appears often in mathematical literature, and in this particular case it refers to convex duality. The following subsections are essentially an application of convex duality as discussed in Section 7.2.

12.3.1 Convex Duality Via Lagrange Multipliers

In Chapter 7 we used λ as Lagrange multipliers. In this section we follow the notation commonly chosen in SVM literature, and use α and γ .

Recall the primal soft margin SVM (12.27). We call the variables \mathbf{w} , b and ξ corresponding to the primal SVM the primal variables. We use $\alpha_n \geq 0$ as the Lagrange multiplier corresponding to the constraint (12.28) that the examples are classified correctly and $\gamma_n \geq 0$ as the Lagrange multiplier corresponding to the non-negativity constraint of the slack variable, see (12.29). The Lagrangian is then given by

$$\begin{aligned} \mathcal{L}(\mathbf{w}, b, \xi, \alpha, \gamma) = & \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{n=1}^N \xi_n \\ & - \underbrace{\sum_{n=1}^N \alpha_n (y_n (\langle \mathbf{w}, \mathbf{x}_n \rangle + b) - 1 + \xi_n)}_{\text{constraint (12.28)}} - \underbrace{\sum_{n=1}^N \gamma_n \xi_n}_{\text{constraint (12.29)}} . \end{aligned} \quad (12.36)$$

Differentiating the Lagrangian (12.36) with respect to the three primal variables \mathbf{w} , b and ξ respectively, we obtain

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}} = \mathbf{w} - \sum_{n=1}^N \alpha_n y_n \mathbf{x}_n , \quad (12.37)$$

$$\frac{\partial \mathcal{L}}{\partial b} = \sum_{n=1}^N \alpha_n y_n , \quad (12.38)$$

$$\frac{\partial \mathcal{L}}{\partial \xi_n} = C - \alpha_n - \gamma_n . \quad (12.39)$$

We now find the maximum of the Lagrangian by setting each of these partial derivatives to zero. By setting (12.37) to zero we find

$$\mathbf{w} = \sum_{n=1}^N \alpha_n y_n \mathbf{x}_n, \quad (12.40)$$

which is a particular instance of the *representer theorem* (Kimeldorf and Wahba, 1970). Equation (12.40) says that the optimal weight vector in the primal is a linear combination of the examples. Recall from Section 2.6.1 that this means that the solution of the optimization problem lies in the span of training data. Additionally the constraint obtained by setting 12.38 to zero implies that the optimal weight vector is an affine combination of the examples. The representer theorem turns out to hold for very general settings of regularized empirical risk minimization (Hofmann et al., 2008; Argyriou and Dinuzzo, 2014). The theorem has more general versions (Schölkopf et al., 2001), and necessary and sufficient conditions on its existence can be found in Yu et al. (2013).

representer theorem

Remark. The representer theorem (12.40) also provides an explanation of the name Support Vector Machine. The examples \mathbf{x}_n whose corresponding parameters $\alpha_n = 0$ do not contribute to the solution \mathbf{w} at all. The other examples, where $\alpha_n > 0$, are called *support vectors* since they “support” the hyperplane. \diamond

The representer theorem is actually a collection of theorems saying that the solution of minimizing empirical risk lies in the subspace (Section 2.4.3) defined by the examples. support vectors

By substituting the expression for \mathbf{w} into the Lagrangian (12.36), we obtain the dual

$$\begin{aligned} \mathfrak{D}(\xi, \alpha, \gamma) = & \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N y_i y_j \alpha_i \alpha_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle - \sum_{i=1}^N y_i \alpha_i \left\langle \sum_{j=1}^N y_j \alpha_j \mathbf{x}_j, \mathbf{x}_i \right\rangle \\ & + C \sum_{i=1}^N \xi_i - b \sum_{i=1}^N y_i \alpha_i + \sum_{i=1}^N \alpha_i - \sum_{i=1}^N \alpha_i \xi_i - \sum_{i=1}^N \gamma_i \xi_i. \end{aligned} \quad (12.41)$$

Note that there are no longer any terms involving the primal variable \mathbf{w} . By setting (12.38) to zero, we obtain $\sum_{n=1}^N y_n \alpha_n = 0$. Therefore, the term involving b also vanishes. Recall that inner products are symmetric and linear (see Section 3.2). Therefore, the first two terms in (12.41) are over the same objects. These terms (coloured blue) can be simplified, and we obtain the Lagrangian

$$\mathfrak{D}(\xi, \alpha, \gamma) = -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N y_i y_j \alpha_i \alpha_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle + \sum_{i=1}^N \alpha_i + \sum_{i=1}^N (C - \alpha_i - \gamma_i) \xi_i. \quad (12.42)$$

The last term in this equation is a collection of all terms that contain slack variables ξ_i . By setting (12.39) to zero, we see that the last term in (12.41) is also zero. Furthermore, by using the same equation and recalling that the Lagrange multipliers γ_i are non-negative, we conclude that $\alpha_i \leq C$.

We now obtain the dual optimization problem of the SVM, which is expressed exclusively in terms of the Lagrange multipliers α_i . Recall from Lagrangian duality (Theorem 7.1) that we maximize the dual problem. This is equivalent to minimizing the negative dual problem, such that we end up with the *dual SVM*

dual SVM

$$\begin{aligned} \min_{\alpha} \quad & \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N y_i y_j \alpha_i \alpha_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle - \sum_{i=1}^N \alpha_i \\ \text{subject to} \quad & \sum_{i=1}^N y_i \alpha_i = 0 \\ & 0 \leq \alpha_i \leq C \quad \text{for all } i = 1, \dots, N. \end{aligned} \quad (12.43)$$

The equality constraint in (12.43) is obtained from setting (12.38) to zero. The inequality constraint $\alpha_i \geq 0$ is the condition imposed on Lagrange multipliers of inequality constraints (Section 7.2). The inequality constraint $\alpha_i \leq C$ is discussed in the previous paragraph.

The set of inequality constraints in the SVM are called “box constraints” because they limit the vector $\alpha = [\alpha_1, \dots, \alpha_N]^\top \in \mathbb{R}^N$ of Lagrange multipliers to be inside the box defined by 0 and C on each axis. These axis-aligned boxes are particularly efficient to implement in numerical solvers (Dostál, 2009, Chapter 5).

Once we obtain the dual parameters α we can recover the primal parameters \mathbf{w} by using the representer theorem (12.40). Let us call the optimal primal parameter \mathbf{w}^* . However there remains the question on how to obtain the parameter b^* . Consider an example (\mathbf{x}_n) that lies exactly on the margin’s boundary, that is, $\langle \mathbf{w}^*, \mathbf{x}_n \rangle + b = y_n$. Recall that y_n is either +1 or −1, and therefore the only unknown is b which can be computed by

$$b^* = y_n - \langle \mathbf{w}^*, \mathbf{x}_n \rangle. \quad (12.44)$$

It turns out examples that lie exactly on the margin are examples whose dual parameters lie strictly inside the box constraints, $0 < \alpha_i < C$. This is derived using the Karush Kuhn Tucker conditions, for example in Schölkopf and Smola (2002).

Remark. In principle there may be no examples that lie exactly on the margin. In this case we should compute $|y_n - \langle \mathbf{w}^*, \mathbf{x}_n \rangle|$ for all support vectors and take the median value of this absolute value difference to be the value of b^* . A derivation of this fact can be found in <http://fouryears.eu/2012/06/07/the-svm-bias-term-conspiracy/>. \diamond

12.3.2 Soft Margin SVM: Convex Hull View

Another approach to obtain the SVM is to consider an alternative geometric argument. Consider the set of examples \mathbf{x}_n with the same label. We would like to build a convex boundary around this set of examples that is the smallest possible. This is called the convex hull and is illustrated in Figure 12.9.

Let us first build some intuition about a convex combination of points.

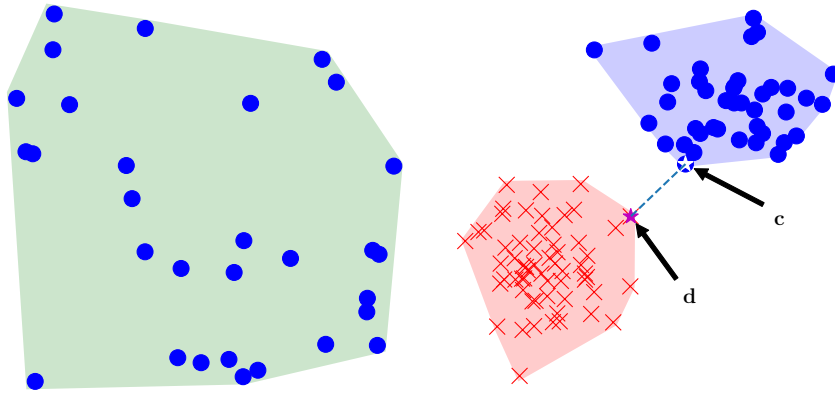


Figure 12.9 (left) Convex hull of points (right) A convex hull around positive and negative examples.

Consider two points x_1 and x_2 and corresponding non-negative weights $\alpha_1, \alpha_2 \geq 0$ such that $\alpha_1 + \alpha_2 = 1$. The equation $\alpha_1 x_1 + \alpha_2 x_2$ describes each point on a line between x_1 and x_2 . Consider what happens when we add a third point x_3 along with a weight $\alpha_3 \geq 0$ such that $\sum_{n=1}^3 \alpha_n = 1$. The convex combination of these three points x_1, x_2, x_3 span a two dimensional area. The convex hull of this area is the triangle formed by the edges corresponding to each pair of points. As we add more points, and the number of points become greater than the number of dimensions, some of the points will be inside the convex hull, as we can see in the left of Figure 12.9.

In general, building a convex boundary of points (called the *convex hull*) can be done by introducing non-negative weights $\alpha_n \geq 0$ corresponding to each example x_n . Then the convex hull can be described as the set

$$\text{conv}(\mathbf{X}) = \left\{ \sum_{n=1}^N \alpha_n x_n \right\} \quad \text{with} \quad \sum_{n=1}^N \alpha_n = 1 \quad \text{and} \quad \alpha_n \geq 0, \quad (12.45)$$

for all $n = 1, \dots, N$. If the two clouds of points corresponding to the positive and negative classes are well separated, then the convex hulls do not overlap. Given the training data $(x_1, y_1), \dots, (x_N, y_N)$ we form two convex hulls, corresponding to the positive and negative classes respectively. We pick a point c , which is in the convex hull of the set of positive examples, and is closest to the negative class distribution. Similarly we pick a point d in the convex hull of the set of negative examples, and is closest to the positive class distribution. Refer to the right of Figure 12.9. We draw a vector from d to c

$$w = c - d. \quad (12.46)$$

Picking the points c and d as above, and requiring them to be closest to each other is the same as saying that we want to minimize the length/norm of w , such that we end up with the corresponding optimization

problem

$$\operatorname{argmin}_w \|w\| = \operatorname{argmin}_w \frac{1}{2} \|w\|^2. \quad (12.47)$$

Since c must be in the positive convex hull, it can be expressed as a convex combination of the positive examples, i.e., for non-negative coefficients α_n^+

$$c = \sum_{y_n=+1} \alpha_n^+ x_n. \quad (12.48)$$

Similarly, for the examples with negative labels we obtain

$$d = \sum_{y_n=-1} \alpha_n^- x_n. \quad (12.49)$$

By substituting (12.46), (12.48) and (12.49) into (12.47), we obtain the following objective function

$$\min_{\alpha} \frac{1}{2} \left\| \sum_{y_n=+1} \alpha_n^+ x_n - \sum_{y_n=-1} \alpha_n^- x_n \right\|^2. \quad (12.50)$$

Let α be the set of all coefficients, i.e., the concatenation of α^+ and α^- . Recall that we require that for each convex hull that their coefficients sum to one,

$$\sum_{y_n=+1} \alpha_n^+ = 1 \quad \text{and} \quad \sum_{y_n=-1} \alpha_n^- = 1. \quad (12.51)$$

This implies the constraint

$$\sum_{n=1}^N y_n \alpha_n = 0. \quad (12.52)$$

This result can be seen by multiplying out the individual classes

$$\begin{aligned} \sum_{n=1}^N y_n \alpha_n &= \sum_{y_n=+1} (+1) \alpha_n^+ + \sum_{y_n=-1} (-1) \alpha_n^- \\ &= \sum_{y_n=+1} \alpha_n^+ - \sum_{y_n=-1} \alpha_n^- = 1 - 1 = 0. \end{aligned} \quad (12.53)$$

The objective function (12.50) and the constraint (12.52), along with the assumption that $\alpha \geq 0$, give us a constrained (convex) optimization problem. This optimization problem can be shown to be the same as that of the dual hard margin SVM (Bennett and Bredensteiner, 2000a).

Remark. To obtain the soft margin dual, we consider the reduced hull. The *reduced hull* is similar to the convex hull but has an upper bound to the size of the coefficients α . The maximum possible value of the elements of α restricts the size that the convex hull can take. In other words, the bound on α shrinks the convex hull to a smaller volume (Bennett and Bredensteiner, 2000b). \diamond

12.3.3 Kernels

Consider the formulation of the dual SVM (12.43). Notice that the inner product in the objective occurs only between examples \mathbf{x}_i and \mathbf{x}_j . There are no inner products between the examples and the parameters. Therefore if we consider a set of features $\phi(\mathbf{x}_i)$ to represent \mathbf{x}_i , the only change in the dual SVM will be to replace the inner product. This modularity, where the choice of the classification method (the SVM) and the choice of the feature representation $\phi(\mathbf{x})$ can be considered separately, provides flexibility for us to explore the two problems independently.

Since $\phi(\mathbf{x})$ could be a non-linear function, we can use the SVM (which assumes a linear classifier) to construct nonlinear classifiers. This provides a second avenue, in addition to the soft margin, for users to deal with a dataset that is not linearly separable. It turns out that there are many algorithms and statistical methods, which have this property that we observed in the dual SVM: the only inner products are those that occur between examples. Instead of *explicitly* defining a non-linear feature map $\phi(\cdot)$ and computing the resulting inner product between examples \mathbf{x}_i and \mathbf{x}_j , we define a similarity function $k(\mathbf{x}_i, \mathbf{x}_j)$ between \mathbf{x}_i and \mathbf{x}_j . For a certain class of similarity functions called *kernels*, the definition of the similarity function *implicitly* defines a non-linear feature map $\phi(\cdot)$. Kernels are by definition functions $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ for which there exists a Hilbert space \mathcal{H} and $\phi : \mathcal{X} \rightarrow \mathcal{H}$ a feature map such that

$$k(\mathbf{x}_i, \mathbf{x}_j) = \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle_{\mathcal{H}}. \quad (12.54)$$

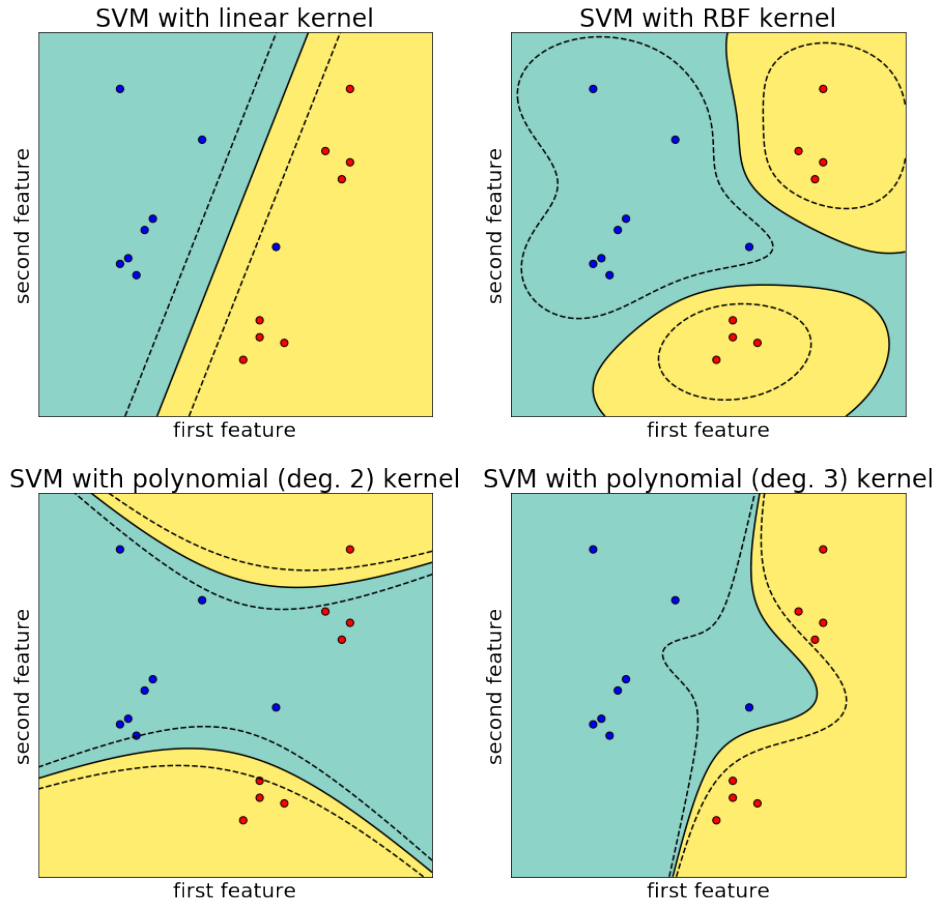
There is a unique reproducing kernel Hilbert space associated with every kernel k (Aronszajn, 1950; Berlinet and Thomas-Agnan, 2004). In this unique association $\phi(\mathbf{x}) = k(\cdot, \mathbf{x})$ is called the canonical feature map. This is known as the *kernel trick* (Schölkopf and Smola, 2002; Shawe-Taylor and Cristianini, 2004), as it hides away the explicit non-linear feature map. The matrix $\mathbf{K} \in \mathbb{R}^{N \times N}$, resulting from the inner products or the application of $k(\cdot, \cdot)$ to a dataset, is called the *Gram matrix*, and is often just referred to as the *kernel matrix*. Kernels must be symmetric and positive semi-definite, i.e., every kernel matrix \mathbf{K} must be symmetric and positive semi-definite (Section 3.2.3):

$$\forall \mathbf{z} \in \mathbb{R}^N \quad \mathbf{z}^\top \mathbf{K} \mathbf{z} \geq 0. \quad (12.55)$$

Some popular examples of kernels for multivariate real-valued data $\mathbf{x}_i \in \mathbb{R}^D$ are the polynomial kernel, the Gaussian radial basis function kernel, and the rational quadratic kernel. Figure 12.10 illustrates the effect of different kernels on separating hyperplanes on an example dataset.

Remark. Unfortunately for the fledgling machine learner, there are multiple meanings of the word kernel. In this chapter, the word kernel comes from the idea of the Reproducing Kernel Hilbert Space (RKHS) (Aronszajn, 1950; Saitoh, 1988). We have discussed the idea of the kernel in

Figure 12.10
Support Vector Machine with different kernels. Note that while the decision boundary is nonlinear, the underlying problem being solved is for a linear separating hyperplane (albeit with a nonlinear kernel).



linear algebra (Section 2.7.3), where the kernel is another word for the nullspace. The third common use of the word kernel in machine learning is the smoothing kernel in kernel density estimation (Section 11.5). \diamond

Since the explicit representation $\phi(\mathbf{x})$ is mathematically equivalent to the kernel representation $k(\mathbf{x}_i, \mathbf{x}_j)$ a practitioner will often design the kernel function, such that it can be computed more efficiently than the inner product between explicit feature maps. For example, consider the polynomial kernel, where the number of terms in the explicit expansion grows very quickly (even for polynomials of low degree) when the input dimension is large. The kernel function only requires one multiplication per input dimension, which can provide significant computational savings.

Another useful aspect of the kernel trick is that there is no need for the original data to be already represented as multivariate real-valued data. Note that the inner product is defined on the output of the function $\phi(\cdot)$, but does not restrict the input to real numbers. Hence, the function $\phi(\cdot)$ and the kernel function $k(\cdot, \cdot)$ can be defined on any object, e.g., sets,

sequences, strings and graphs (Ben-Hur et al., 2008; Gärtner, 2008; Shi et al., 2009; Vishwanathan et al., 2010).

12.3.4 Numerical Solution

We conclude our discussion of SVMs by looking at how to express the problems derived in this chapter in terms of the concepts presented in Chapter 7. We consider two different approaches for finding the optimal solution for the SVM. First we consider the loss view of SVM 8.1.2 and express this as an unconstrained optimization problem. Then we express the constrained versions of the primal and dual SVMs as quadratic programs in standard form 7.3.2.

Consider the loss function view of the SVM (12.33). This is a convex unconstrained optimization problem, but the hinge loss (12.30) is not differentiable. Therefore, we apply a subgradient approach for solving it. However, the hinge loss is differentiable almost everywhere, except for one single point at the hinge $t = 1$. At this point, the gradient is a set of possible values that lie between 0 and -1 . Therefore, the subgradient g of the hinge loss is given by

$$g(t) = \begin{cases} -1 & t < 1 \\ [-1, 0] & t = 1 \\ 0 & t > 1 \end{cases} . \quad (12.56)$$

Using this subgradient above, we can apply the optimization methods presented in Section 7.1.

Both the primal and the dual SVM result in a convex quadratic programming problem (constrained optimization). Note that the primal SVM in (12.27) has optimization variables that have the size of the dimension D of the input examples. The dual SVM in (12.43) has optimization variables that have the size of the number N of examples.

To express the primal SVM in the standard form (7.35) for quadratic programming, let us assume that we use the dot product (3.6) as the inner product. We rearrange the equation for the primal SVM (12.27), such that the optimization variables are all on the right and the inequality of the constraint matches the standard form. This yields the optimization

$$\min_{\mathbf{w}, b, \boldsymbol{\xi}} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{n=1}^N \xi_n \quad (12.57)$$

$$\text{subject to} \quad \begin{aligned} -y_n \mathbf{x}_n^\top \mathbf{w} - y_n b - \xi_n &\leq -1 \\ -\xi_n &\leq 0 \end{aligned} \quad (12.58)$$

for all $n = 1, \dots, N$. By concatenating the variables $\mathbf{w}, b, \mathbf{x}_n$ into one single vector, and carefully collecting the terms, we obtain the following matrix form of the soft margin SVM. In the following optimization problem, the minimization is over $[\mathbf{w}^\top, b, \boldsymbol{\xi}^\top]^\top \in \mathbb{R}^{D+1+N}$, and we have used

Recall from Section 3.2 that in this book, we use the phrase dot product to mean the inner product on Euclidean vector space.

the notation: \mathbf{I}_m to represent the identity matrix of size $m \times m$, $\mathbf{0}_{m,n}$ to represent the matrix of zeros of size $m \times n$, and $\mathbf{1}_{m,n}$ to represent the matrix of ones of size $m \times n$. The soft margin SVM can be written in the following vector form:

$$\min_{\mathbf{w}, b, \xi} \frac{1}{2} \begin{bmatrix} \mathbf{w} \\ b \\ \xi \end{bmatrix}^\top \begin{bmatrix} \mathbf{I}_D & \mathbf{0}_{D, N+1} \\ \mathbf{0}_{N+1, D} & \mathbf{0}_{N+1, N+1} \end{bmatrix} \begin{bmatrix} \mathbf{w} \\ b \\ \xi \end{bmatrix} + [\mathbf{0}_{D+1, 1} \quad C\mathbf{1}_{N, 1}]^\top \begin{bmatrix} \mathbf{w} \\ b \\ \xi \end{bmatrix} \quad (12.59)$$

$$\text{subject to } \begin{bmatrix} -\mathbf{Y}\mathbf{X} & -\mathbf{y} & -\mathbf{I}_N \\ \mathbf{0}_{N, D+1} & & -\mathbf{I}_N \end{bmatrix} \begin{bmatrix} \mathbf{w} \\ b \\ \xi \end{bmatrix} \leq \begin{bmatrix} -\mathbf{1}_{N, 1} \\ \mathbf{0}_{N, 1} \end{bmatrix}, \quad (12.60)$$

where \mathbf{y} is the vector of labels $[y_1, \dots, y_N]^\top$, $\mathbf{Y} = \text{diag}(\mathbf{y})$ is an N by N matrix where the elements of the diagonal are from \mathbf{y} , and $\mathbf{X} \in \mathbb{R}^{N \times D}$ is the matrix obtained by concatenating all the examples.

We can similarly perform a collection of terms for the dual version of the SVM (12.43). To express the dual SVM in standard form, we first have to express the kernel matrix \mathbf{K} such that each entry is $K_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$. Or if we are using an explicit feature representation $K_{ij} = \langle \mathbf{x}_i, \mathbf{x}_j \rangle$. For convenience of notation we introduce a matrix with zeros everywhere except on the diagonal, where we store the labels, that is, $\mathbf{Y} = \text{diag}(\mathbf{y})$. The dual SVM can be written as

$$\min_{\alpha} \frac{1}{2} \alpha^\top \mathbf{Y} \mathbf{K} \mathbf{Y} \alpha - \mathbf{1}_{N, 1}^\top \alpha \quad (12.61)$$

$$\text{subject to } \begin{bmatrix} \mathbf{y}^\top \\ -\mathbf{y}^\top \\ -\mathbf{I}_N \\ \mathbf{I}_N \end{bmatrix} \alpha \leq \begin{bmatrix} \mathbf{0}_{N+2, 1} \\ C\mathbf{1}_{N, 1} \end{bmatrix}. \quad (12.62)$$

Remark. In Section 7.3.1 and 7.3.2 we introduced the standard forms of the constraints to be inequality constraints. We will express the dual SVM's equality constraint as two inequality constraints, i.e.,

$$\mathbf{A}\mathbf{x} = \mathbf{b} \quad \text{is replaced by} \quad \mathbf{A}\mathbf{x} \leq \mathbf{b} \quad \text{and} \quad \mathbf{A}\mathbf{x} \geq \mathbf{b} \quad (12.63)$$

Particular software implementations of convex optimization methods may provide the ability to express equality constraints. \diamond

Since there are many different possible views of the SVM, there are many approaches for solving the resulting optimization problem. The approach presented here, expressing the SVM problem in standard convex optimization form, is not often used in practice. The two main implementations of SVM solvers are (Chang and Lin, 2011) (which is open source) and (Joachims, 1999). Since SVMs have a clear and well defined optimization problem, many approaches based on numerical optimization

techniques (Nocedal and Wright, 2006) can be applied (Shawe-Taylor and Sun, 2011).

12.4 Further Reading

The SVM is one of many approaches for studying binary classification. Other approaches include the perceptron, logistic regression, Fisher discriminant, nearest neighbor, naive Bayes, and random forest (Bishop, 2006; Murphy, 2012). A short tutorial on SVMs and kernels on discrete sequences can be found in Ben-Hur et al. (2008). The development of SVMs is closely linked to empirical risk minimization 8.1, and hence the SVM has strong theoretical properties (Vapnik, 2000; Steinwart and Christmann, 2008). The book about kernel methods (Schölkopf and Smola, 2002) includes many details of support vector machines and how to optimize them. A broader book about kernel methods (Shawe-Taylor and Cristianini, 2004) also includes many linear algebra approaches for different machine learning problems.

An alternative derivation of the dual SVM can be obtained using the idea of the Legendre-Fenchel transform (Section 7.3.3). The derivation considers each term of the unconstrained formulation of the SVM (12.33) separately and calculates their convex conjugates (Rifkin and Lippert, 2007). Readers interested in the functional analysis view (also the regularization methods view) of SVMs are referred to the work by Wahba (1990). Theoretical exposition of kernels (Manton and Amblard, 2015; Aronszajn, 1950; Schwartz, 1964; Saitoh, 1988) require a basic grounding of linear operators (Akhiezer and Glazman, 1993). The idea of kernels have been generalized to Banach spaces (Zhang et al., 2009) and Krein spaces (Ong et al., 2004; Loosli et al., 2016).

Observe that the hinge loss has three equivalent representations, as shown by (12.30) and (12.31), as well as the constrained optimization problem in (12.35). The formulation (12.30) is often used when comparing the SVM loss function with other loss functions (Steinwart, 2007). The two piece formulation (12.31) is convenient for computing subgradients, as each piece is linear. The third formulation (12.35), as seen in Section 12.3.4, enables the use of convex quadratic programming (Section 7.3.2) tools.

Since binary classification is a well studied task in machine learning, other words are also sometimes used, such as discrimination, separation or decision. To further add to the confusion, there are three quantities that can be the output of a binary classifier. First is the output of the linear function itself. This output can be used for ranking the examples, and binary classification can be thought of as picking a threshold on the ranked examples (Shawe-Taylor and Cristianini, 2004). The second quantity that is often considered the output of a binary classifier is after the output is passed through a non-linear function to constrain its value to a bounded range.

6608 A common non-linear function is the sigmoid function (Bishop, 2006).
6609 When the non-linearity results in well calibrated probabilities (Gneiting
6610 and Raftery, 2007; Reid and Williamson, 2011), this is called class proba-
6611 bility estimation. The third output of a binary classifier is the final binary
6612 decision, which is the one most commonly assumed to be the output of
6613 the classifier.