

6 Support Vector Machines

The core idea of Support Vector Machines becomes clear when viewed geometrically. Figure 6.1 shows a slice of the iris dataset containing two linearly separable species. In the left panel two conventional linear classifiers are drawn; although both label every training point correctly, their separating lines sit very close to several observations, so even a small perturbation could lead to misclassifications on new data. The right panel plots the boundary produced by an SVM. The solid line still splits the classes cleanly, but it is positioned to maximise the distance to the nearest points, leaving the widest possible “street” (bounded by the dashed parallels). The samples that lie on these margins are the *support vectors*. This strategy, called large margin classification, usually delivers models that generalise better than those that merely separate the training set.

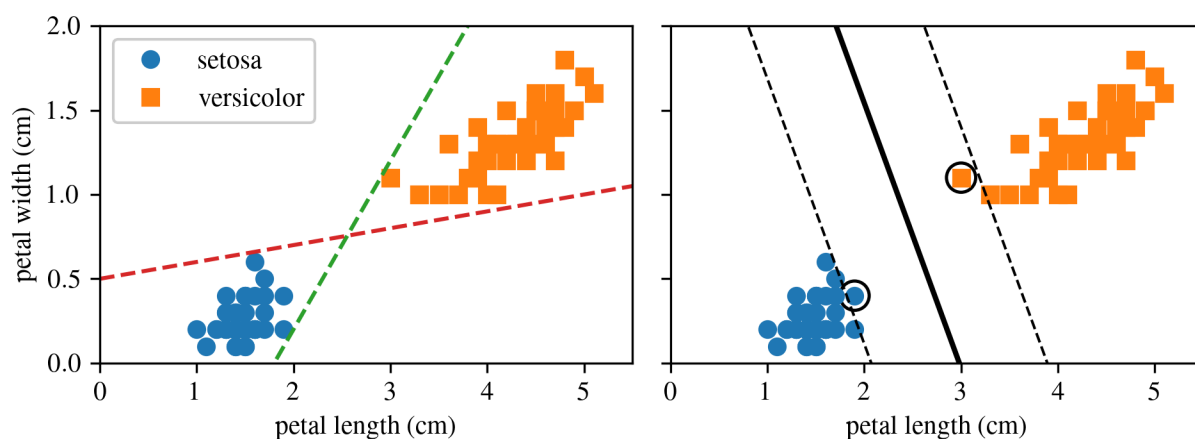
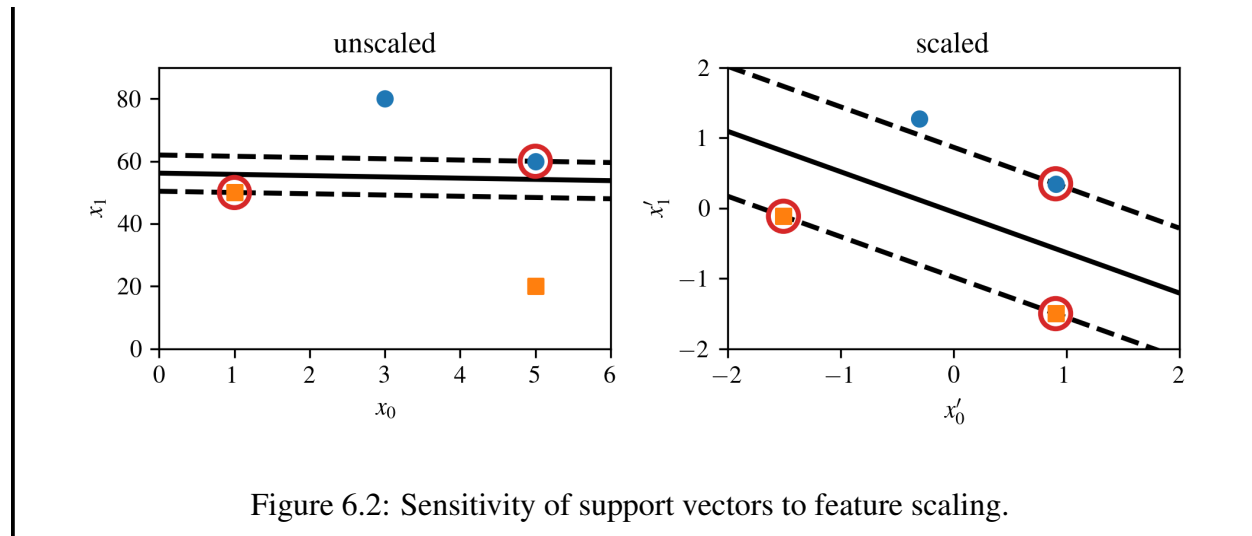


Figure 6.1: Large margin classification.

Observe that the addition of further training instances outside the “street” does not influence the decision boundary; it is entirely shaped by the instances situated on the boundary’s edge. These pivotal instances are termed support vectors and are highlighted with circles in Figure 6.1.

NOTE

The sensitivity of SVMs to feature scales is evident in Figure 6.2. In the left plot, the vertical dimension greatly outweighs the horizontal dimension, resulting in a nearly horizontal “street.” However, after applying feature scaling such as using Scikit-Learn’s `StandardScaler` the decision boundary becomes more appropriate, as illustrated in the right plot.

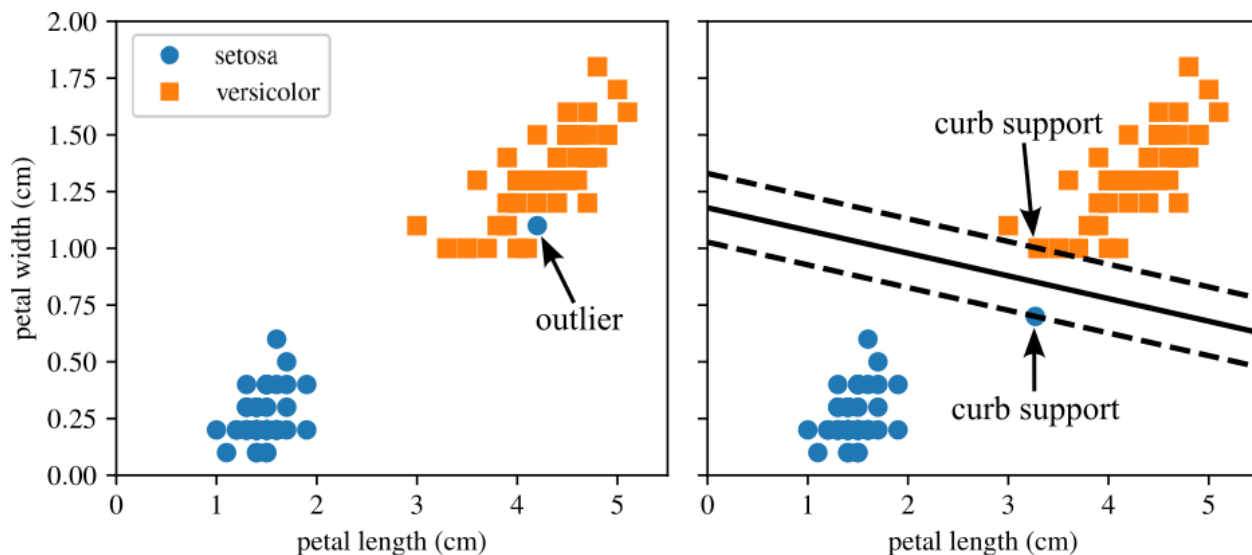


6.1 Linear SVM Classification

Hard margin classification demands that all instances be correctly classified without any margin violations. This strict approach faces two significant challenges:

- It is only feasible when the data is linearly separable.
- It is highly sensitive to outliers.

Figure 6.3 illustrates these challenges using the iris dataset with an added outlier. On the left, achieving a hard margin is impossible due to the outlier. On the right, although a decision boundary is found, it deviates substantially from the optimal boundary shown in Figure 6.1 and is less likely to perform well on new data.



To mitigate the limitations of hard margin classification, a more adaptable model, known as soft margin classification, is often employed. The goal here is to achieve an optimal balance between maximizing the margin width and minimizing margin violations, where instances might fall into the margin or on the incorrect side.

Scikit-Learn's SVM implementations facilitate this balance through the hyperparameter C . A smaller value of C results in a wider margin but allows more margin violations, which is beneficial for model flexibility. Conversely, a larger C value tightens the margin, reducing margin violations but at the risk of a less flexible model. Figure 6.4 demonstrates this trade-off: the left plot with a low C value

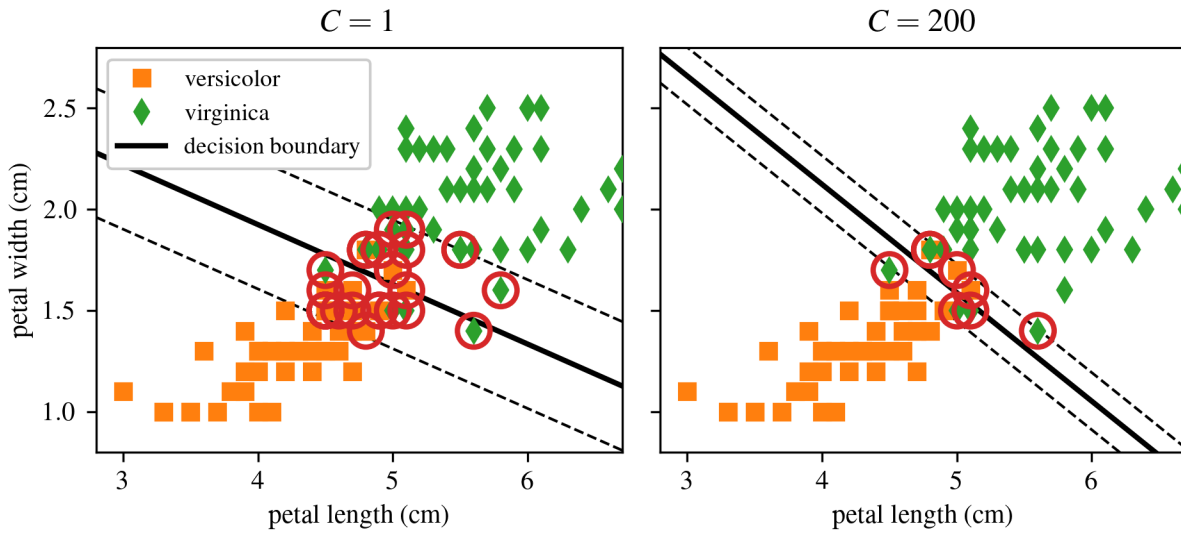


Figure 6.4: SVM margin sizes for different C values.

NOTE

Overfitting in an SVM model can often be addressed by reducing the C value, which increases regularization.

In earlier chapters we placed every model parameter in a single vector θ : the first entry θ_0 acted as the bias, while $\theta_1, \dots, \theta_n$ were the feature weights, and each input was augmented with a constant bias feature $x_0 = 1$. In this chapter we adopt the notation most common for SVMs. The bias is written as b , the weight vector as \mathbf{w} , and no extra bias feature is appended to the input vectors.

6.1.1 Decision Function and Predictions

A linear SVM classifier determines the class of a new instance x by calculating the decision function

$$\mathbf{w}^T \mathbf{x} + b = w_1 x_1 + \dots + w_n x_n + b. \quad (6.1)$$

If the outcome is positive, the predicted class \hat{y} is the positive class (1); otherwise, it is the negative class (0). This is written as

$$\hat{y} = \begin{cases} 0 & \text{if } \mathbf{w}^T \mathbf{x} + b < 0, \\ 1 & \text{if } \mathbf{w}^T \mathbf{x} + b \geq 0. \end{cases} \quad (6.2)$$

Figure 6.5 plots the decision function for a model with two features, so the surface is a plane in \mathbb{R}^2 . The thick solid line marks the decision boundary where the function equals zero. The dashed lines show the loci where the function equals 1 and -1 ; they run parallel to the boundary and sit at equal distance from it, outlining the margin. Training a linear SVM adjusts \mathbf{w} and b to make this margin as wide as possible while either prohibiting margin violations (hard margin) or keeping them small through a penalty term (soft margin).

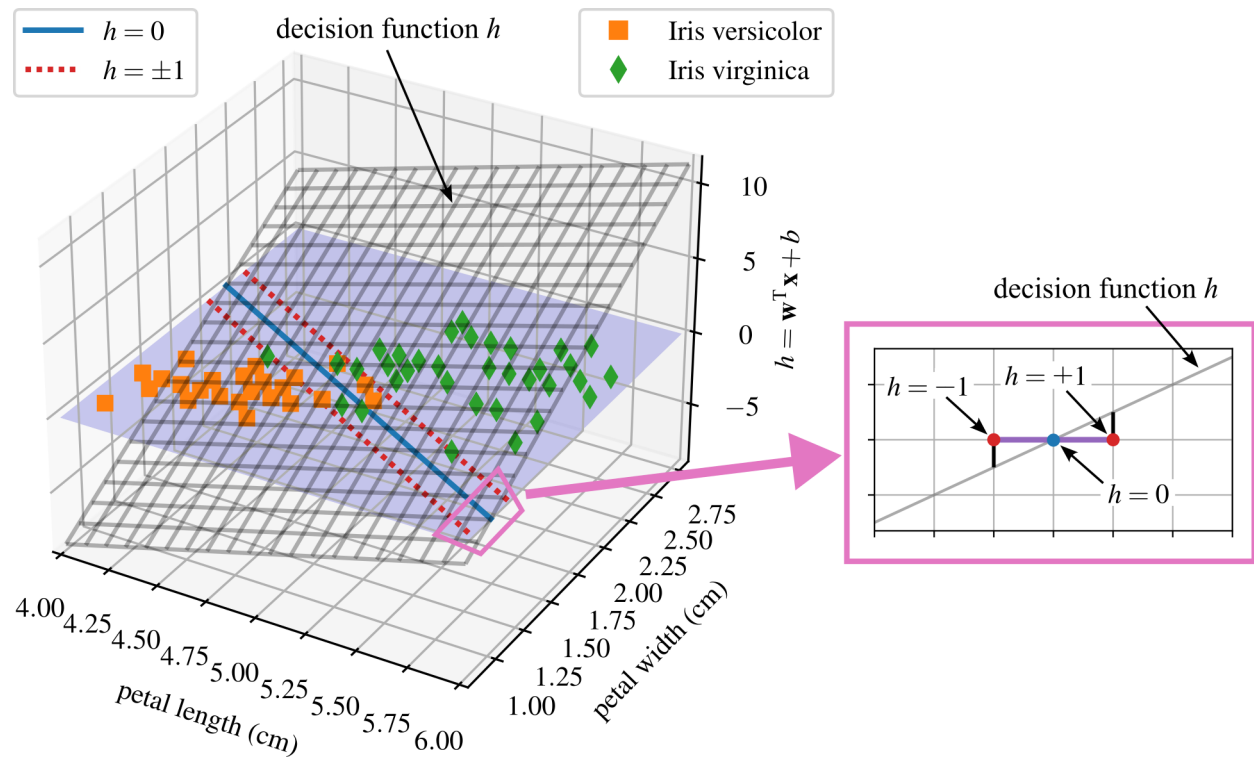


Figure 6.5: Decision function for the Iris Dataset showing how the decision function h cuts through the feature space.

6.1.2 Training Objective

The slope of the decision function corresponds to the norm of the weight vector, $\|\mathbf{w}\|$. Halving this slope causes the decision boundary margins, where the decision function equals ± 1 , to double in distance from the decision boundary. Effectively, reducing the norm of \mathbf{w} by half doubles the margin. This geometric interpretation is perhaps simpler to visualize in two dimensions, as shown in Figure 6.6. Thus, minimizing $\|\mathbf{w}\|$ maximizes the margin.

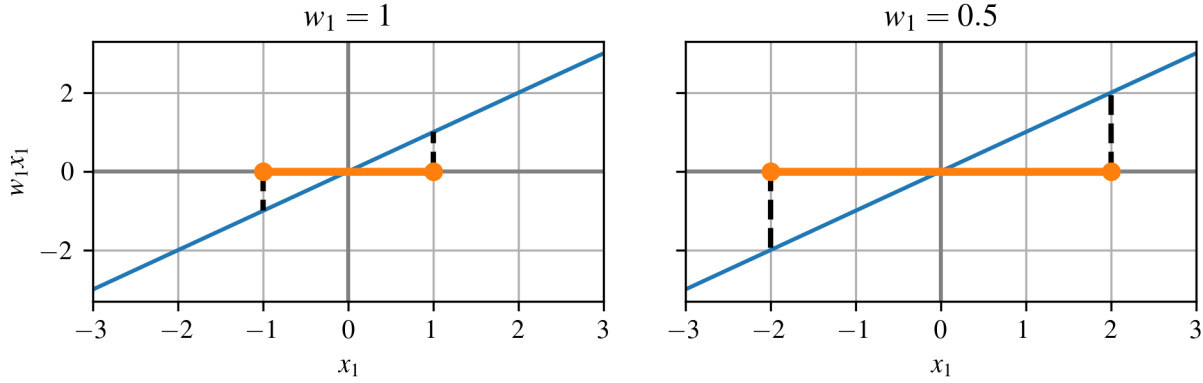


Figure 6.6: The margin is dependent on the value of the weight vector where a smaller weight vector results in a larger margin and vice versa.

To achieve a large margin while enforcing that no data points fall within the margin (hard margin), we ensure the decision function exceeds $+1$ for all positive training instances and is less than -1 for all negative instances. Let $t^{(i)}$ equal -1 for negative instances ($y^{(i)} = 0$) and $+1$ for positive ones ($y^{(i)} = 1$). The constraints then require

$$t^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + b) \geq 1 \quad (6.3)$$

for all training instances. This forms the basis of the hard margin linear SVM classifier optimization problem:

$$\begin{aligned} & \underset{\mathbf{w}, b}{\text{minimize}} && \frac{1}{2} \mathbf{w}^T \mathbf{w} \\ & \text{subject to} && t^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + b) \geq 1 \text{ for } i = 1, 2, \dots, m \end{aligned} \quad (6.4)$$

NOTE

The objective function minimized is $\frac{1}{2} \mathbf{w}^T \mathbf{w}$, equivalent to $\frac{1}{2} \|\mathbf{w}\|^2$. This formulation is chosen over minimizing $\|\mathbf{w}\|$ directly because $\frac{1}{2} \|\mathbf{w}\|^2$ offers a straightforward derivative, simply \mathbf{w} , facilitating gradient calculations. In contrast, $\|\mathbf{w}\|$ lacks differentiability at $\mathbf{w} = 0$, posing challenges for optimization algorithms, which typically require smooth, differentiable functions to ensure effective optimization.

To formulate the soft margin objective, it is necessary to introduce a slack variable $\zeta^{(i)} \geq 0$ for each instance. This variable, $\zeta^{(i)}$, quantifies the permissible margin violation for the i^{th} instance. Consequently, we face dual objectives: minimizing the slack variables to reduce margin violations and minimizing $\frac{1}{2} \mathbf{w}^T \mathbf{w}$ to maximize the margin. The hyperparameter C plays a crucial role here, enabling a balance between these competing objectives. The introduction of C transforms our task into a constrained optimization problem.

$$\begin{aligned}
& \underset{\mathbf{w}, b, \zeta}{\text{minimize}} && \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^m \zeta^{(i)} \\
& \text{subject to} && t^{(i)} (\mathbf{w}^T \mathbf{x}^{(i)} + b) \geq 1 - \zeta^{(i)} \text{ and } \zeta^{(i)} \geq 0 \text{ for } i = 1, 2, \dots, m
\end{aligned} \tag{6.5}$$

6.1.3 Quadratic Programming

Both hard margin and soft margin problems are examples of convex quadratic optimization problems with linear constraints, commonly referred to as Quadratic Programming (QP) problems. QP involves solving optimization problems where the objective is a quadratic function and the constraints are linear. This form of programming, established in the 1940s, predates and is distinct from “computer programming,” and is sometimes more descriptively termed “quadratic optimization” to avoid confusion.

A variety of techniques available through off-the-shelf solvers can address these QP problems, though they extend beyond the scope of this text. The general formulation of a QP problem is as follows:

$$\begin{aligned}
& \underset{\mathbf{p}}{\text{minimize}} && \frac{1}{2} \mathbf{p}^T \mathbf{H} \mathbf{p} + \mathbf{f}^T \mathbf{p} \\
& \text{subject to} && \mathbf{A} \mathbf{p} \leq \mathbf{b}
\end{aligned} \tag{6.6}$$

Here, \mathbf{p} is an n_p -dimensional vector (where n_p is the number of parameters), \mathbf{H} is an $n_p \times n_p$ matrix, \mathbf{f} is an n_p -dimensional vector, \mathbf{A} is an $n_c \times n_p$ matrix (with n_c being the number of constraints), and \mathbf{b} is an n_c -dimensional vector.

Equation 6.6 defines a standard quadratic program with constraints of the form $\mathbf{A} \mathbf{p} \leq \mathbf{b}$. For training a hard margin linear SVM, this setup can be used by choosing parameters that encode the SVM objective and constraints. The solution vector \mathbf{p} contains both the bias term and the feature weights. Using this knowledge, a standard QP solver can be applied directly to find the optimal SVM classifier.

Example 6.1 Support Vector Machine Classification

This example applies a linear Support Vector Machine (SVM) classifier to distinguish Iris-Virginica flowers based on petal length and width. The decision boundary is derived after scaling the data, and margins are visualized. Misclassified instances within the margin are identified and marked. The model’s performance is assessed using a confusion matrix and F1 score.

6.2 Nonlinear SVM Classification

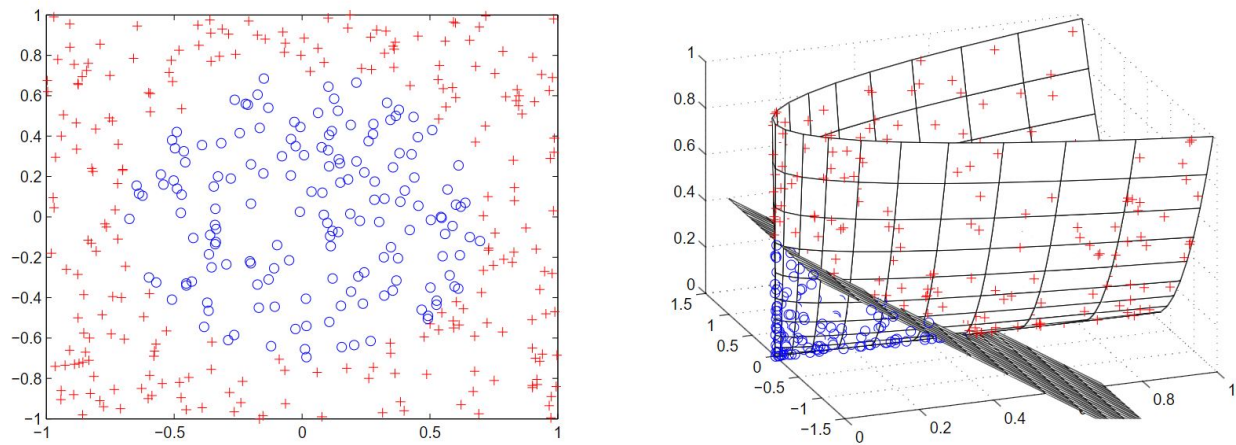


Figure 6.7: Nonlinear SVM example illustration ^a

While linear SVM classifiers are quite effective and perform exceptionally well in various scenarios, many datasets are far from being linearly separable. One strategy to address non-linear datasets is to introduce additional features, such as polynomial feature. Adding features can sometimes transform the dataset into one that is linearly separable. A representation of this technique is shown in 6.7.

A simple example of converting non-linearly separable variables is shown in figure 6.8 where the left plot displays a simple dataset with a single feature x_1 . Clearly, this dataset is not linearly separable. However, by adding another feature $x_2 = (x_1)^2$, the dataset becomes perfectly linearly separable in two dimensions.

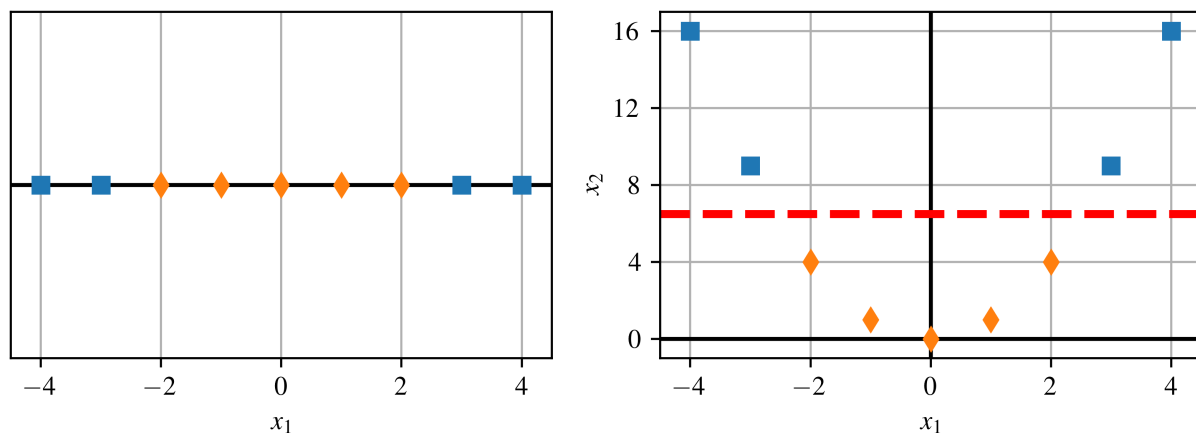


Figure 6.8: Illustration of SVM in higher dimensions

^aMachine Learner, CC BY-SA 4.0 <<https://creativecommons.org/licenses/by-sa/4.0/>>, via Wikimedia Commons

You can implement this idea in Scikit-Learn by creating a pipeline that applies a Polynomial Features transformer (introduced in the Regression Chapter), followed by a `StandardScaler` and a `LinearSVC`. The approach works nicely on the moons dataset, a toy binary-classification problem in which the samples trace two interleaving half-circles, as illustrated in Figure 6.9. You can generate this dataset with the function `make_moons()`.

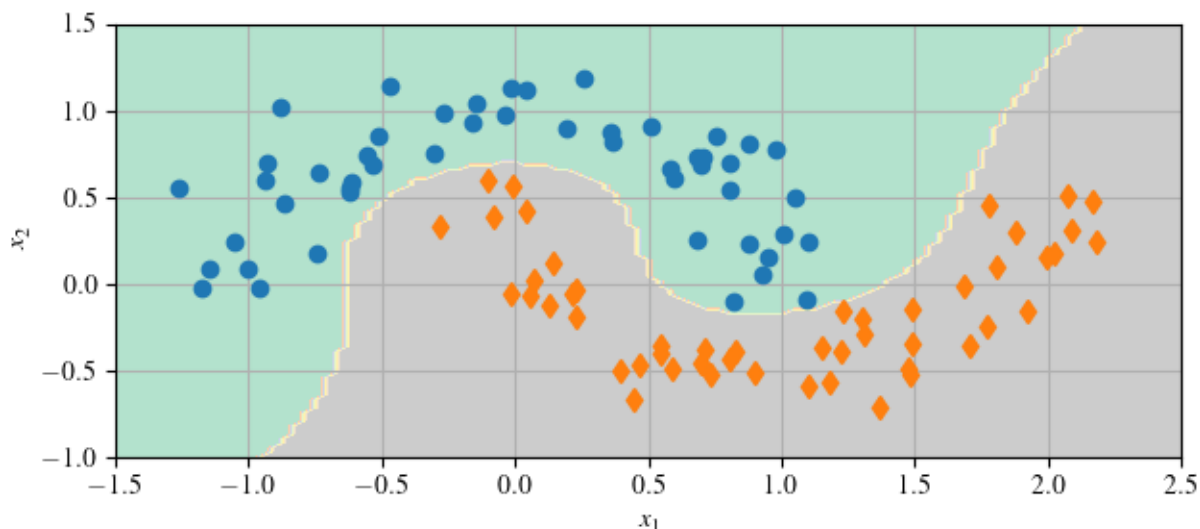


Figure 6.9: Demonstration of a SVM classifier with polynomial features.

Example 6.2 Nonlinear Classification

This example demonstrates nonlinear classification leveraging `LinearSVC` in a pipeline with `preprocessing.PolynomialFeatures` using the `make_moons` dataset. A polynomial feature transformation is combined with a linear SVM to classify the data, and the resulting decision boundaries are visualized.

6.2.1 Kernel trick

While adding polynomial features is straightforward and enhances the performance of various Machine Learning algorithms (not limited to SVMs), it presents limitations. Specifically, lower-degree polynomials may not adequately handle complex datasets, and higher-degree polynomials significantly increase the feature count, slowing down the model.

SVMs offer a unique solution through a remarkable mathematical technique known as the kernel trick, which allows for the benefits of high-degree polynomial features without actually expanding the feature space, thereby avoiding a rapid increase in computation. This kernel trick is incorporated within the `SVC` class.

In the dual form of an SVM the explicit dot product of two input vectors

$$\mathbf{x}^T \mathbf{z} \quad (6.7)$$

is replaced by a kernel function

$$k(\mathbf{x}, \mathbf{z}), \quad (6.8)$$

where \mathbf{x} and \mathbf{z} are n -dimensional feature vectors representing any two data points in the training set. This substitution lets the algorithm construct a linear separator in a (possibly infinite-dimensional) feature space while all calculations still occur in the original input coordinates.

Figure 6.10 shows configured SVM classifiers using a 3rd and 3th degree polynomial kernel; on the left and right respectively. Adjusting the polynomial degree can help manage model fit: reducing the degree may prevent overfitting, whereas increasing it may be necessary for underfitting scenarios. The ‘coef0’ hyperparameter is crucial as it determines the influence of high versus low-degree polynomials in the model.

NOTE

A typical method for determining optimal hyperparameter settings involves utilizing grid search techniques. Starting with a broad, coarse grid search to quickly narrow down potential candidates, followed by a more detailed, finer grid search centered on these promising values often yields faster results. Additionally, understanding the function and influence of each hyperparameter aids in efficiently targeting the most relevant areas of the hyperparameter space.

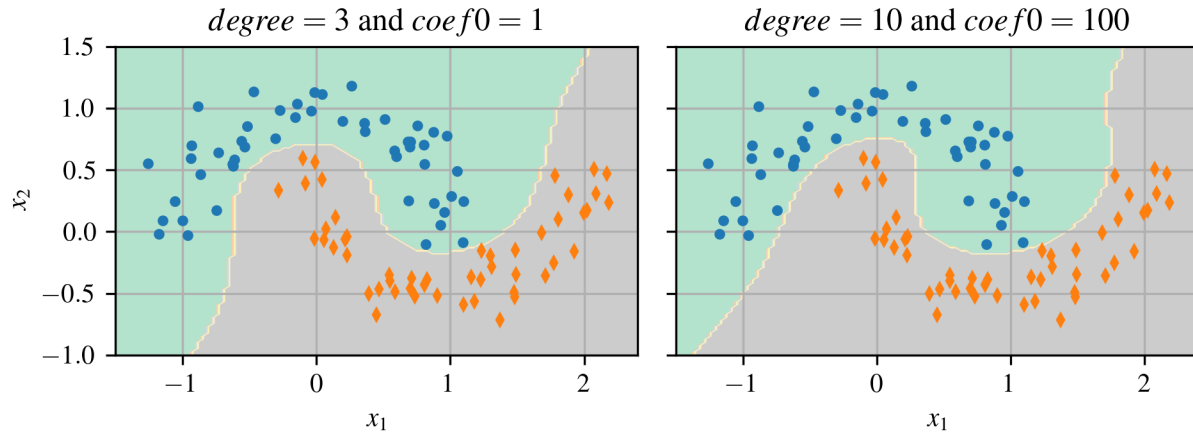


Figure 6.10: SVM polynomial kernel

6.2.2 Standard Kernels

Three kernels cover the vast majority of practical cases.

- The polynomial kernel captures interactions between input features up to a chosen degree d :

$$k_{\text{poly}}(\mathbf{x}, \mathbf{z}) = (\mathbf{x}^T \mathbf{z} + c)^d \quad (6.9)$$

where $c \geq 0$ is a constant offset. It is effective when domain knowledge suggests that a low-order combination of variables explains the target. Keep d modest (typically $d \leq 5$) and standardise inputs to avoid exploding feature dimensions and overfitting.

- The radial basis function (RBF) kernel builds smooth, highly flexible decision boundaries:

$$k_{\text{rbf}}(\mathbf{x}, \mathbf{z}) = \exp(-\gamma \|\mathbf{x} - \mathbf{z}\|^2) \quad (6.10)$$

with width parameter $\gamma > 0$. A large γ makes the surface too flat (underfitting), while a small γ lets every point carve its own pocket (overfitting). Tune C and γ jointly, typically on a logarithmic grid, after x-score standardising the features.

- The sigmoid kernel adds an S-shaped non-linearity to the inner product:

$$k_{\text{sig}}(\mathbf{x}, \mathbf{z}) = \tanh(\kappa \mathbf{x}^T \mathbf{z} + \theta) \quad (6.11)$$

where κ controls the slope and θ the offset. Although useful for some sparse or text data, this kernel is not always positive-semidefinite, so ensure your software handles the resulting optimisation safely.

Example 6.3 Polynomial Kernel Trick

This example uses the kernel trick to enable an SVM to classify non-linearly separable data using SVC (not LinearSVC). A third-degree polynomial kernel is applied to the moons dataset using Scikit-Learn's Pipeline and SVC tools, producing a curved decision boundary that cleanly separates the classes.

NOTE

Begin with the RBF kernel as a strong default, explore polynomial kernels when you expect specific interaction orders, and treat the sigmoid option as experimental unless you have evidence it helps.

6.3 Computational Complexity

Scikit-Learn provides two main SVM classifiers that trade accuracy for speed in different ways. LinearSVC is optimized for large, high-dimensional data when a linear decision boundary is adequate; SVC sacrifices runtime for the expressive power of kernels and thus handles complex, nonlinear patterns. Table 1 highlights how these priorities translate into computational costs and practical constraints and compares them to SGDClassifier for reference.

LinearSVC builds on the `liblinear` solver and is limited to *linear* decision boundaries. Because it does not apply the kernel trick, its training cost grows almost linearly with both the number of samples m and features n , i.e. $O(mn)$. Convergence is controlled by the tolerance parameter `tol` (denoted ϵ in the literature); the default value is usually sufficient, but smaller tolerances can be specified when higher accuracy is critical.

SVC, in contrast, relies on `libsvm` and *does* support kernel functions. Its computational burden is markedly heavier, between $O(m^2n)$ and $O(m^3n)$ in practice, so training becomes prohibitive once the dataset reaches the hundreds-of-thousands range. Nevertheless, SVC excels on smaller or medium-sized problems that demand nonlinear decision surfaces. Runtime also scales with the

average count of non-zero features per instance, meaning sparse high-dimensional inputs remain tractable.

Table 1: Key characteristics of three Scikit-Learn SVM classifier implementations.

True class	time complexity	out-of-core support	scaling required	kernel trick
LinearSVC	$O(m \times n)$	no	yes	no
SVC	$O(m^2 \times n)$ to $O(m^3 \times n)$	no	yes	yes
SGDClassifier	$O(m \times n)$	yes	yes	no

6.4 SVM Regression

Support Vector Regression (SVR) adapts the SVM idea to regression by surrounding the prediction curve with a tube of width ε ; points outside this tube incur a penalty and the optimizer keeps the tube as flat as possible while allowing a limited number of violations governed by the regularization constant C . Conceptually this reverses the classification goal: instead of widening a margin to separate two classes, SVR encourages the data to lie inside the margin, making ε the principal knob that controls how loose the fit may be for both linear and kernel-based models.

6.4.1 Linear SVR

The presence of additional training instances within the margin does not influence the predictions of the model, rendering it ε -insensitive. For linear SVM Regression tasks, the LinearSVR class from Scikit-Learn can be utilized. For instance, Figure 6.11 demonstrates two linear SVM Regression models trained on random linear data; one features a wide margin ($\varepsilon = 1.5$), and the other a narrower margin ($\varepsilon = 0.5$).

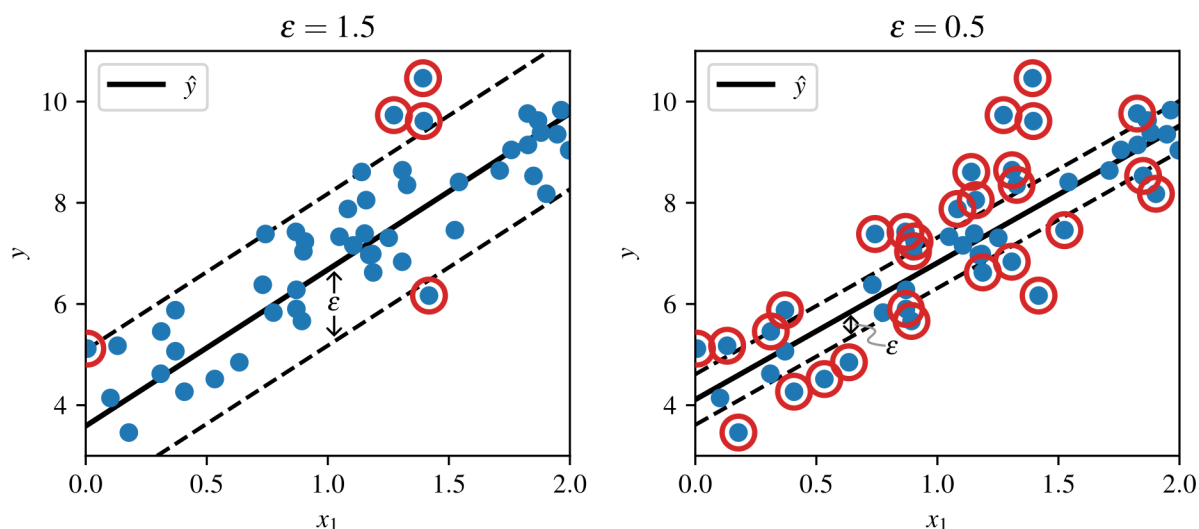


Figure 6.11: SVM Regression models with different ε values.

For data that follow an approximately linear trend, the LinearSVR class in Scikit-Learn solves

the primal problem directly. It scales in $O(mn)$ time with m samples and n features, making it a practical choice for large data sets where a simple linear fit is adequate.

6.4.2 Kernel SVR

When the relationship is non-linear, the SVR class employs the kernel trick:

$$f(x) = \sum_{i=1}^m (\alpha_i - \alpha_i^*) k(x_i, x) + b, \quad (6.12)$$

where $k(\cdot, \cdot)$ is typically RBF or polynomial. Figure 6.12 shows a 2nd-degree polynomial kernel capturing quadratic structure under various regularisation levels.

The Scikit-learn SVR class, supporting the kernel trick and acting as the regression counterpart to the SVC class, performs well with small to medium-sized datasets but slows considerably as dataset size increases. In contrast, the LinearSVR class, akin to the LinearSVC class, scales linearly with the size of the training set.

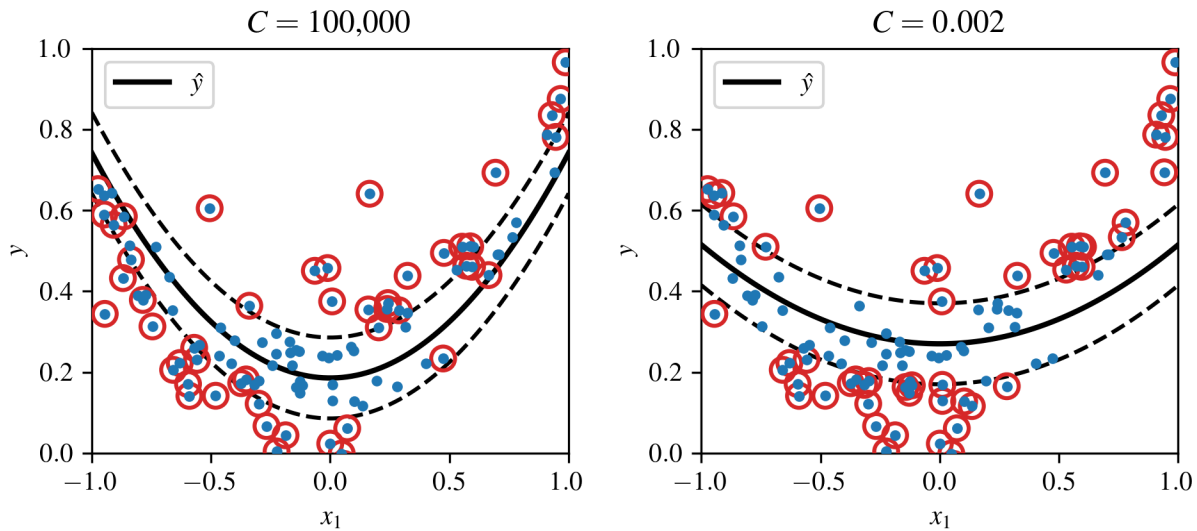


Figure 6.12: SVM regression with a 2nd-degree polynomial kernel, showcasing different regularization levels.

Example 6.4 SVM Polynomial Regression

This example fits a non-linear regression curve to noisy quadratic data using Support Vector Regression (SVR) with a polynomial kernel. It highlights how SVR models can handle non-linear relationships by introducing a margin of tolerance.

Examples**Example 6.1****Example 6.2****Example 6.3****Example 6.4**