

Support Vector Machines

Regression and Classification

Introduction to Support Vector Machines (SVM)

The Core Idea: Optimal Hyperplanes

The main idea behind Support Vector Machines is to find an optimal hyperplane that best separates or fits the data.

For Classification (SVC):

- The hyperplane is a decision boundary that separates data points of different classes.
- "Optimal" means it has the **maximum margin** (distance) from the nearest data points of any class.

For Regression (SVR):

- The hyperplane is a function that best fits the data.
- "Optimal" means it has as many data points as possible within an ϵ -insensitive tube, balancing model complexity and prediction error.

The Kernel Trick

For non-linear data, SVMs use the **kernel trick** to map data into a high-dimensional feature space where a linear hyperplane can be found. All computations are done using a kernel function, avoiding explicit mapping.

Goal of Support Vector Regression (SVR)

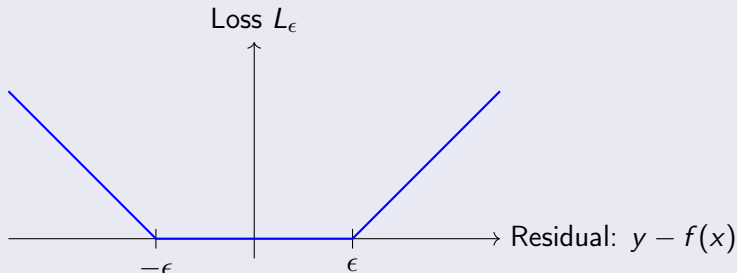
Fitting an ϵ -Insensitive Tube

Unlike traditional regression (e.g., Least Squares) which tries to minimize error for all points, SVR aims to find a function $f(x) = w^T x + b$ such that most data points (x_i, y_i) lie within an ϵ -tube.

The ϵ -Insensitive Loss Function

Errors are only penalized if a data point's residual, $|y - f(x)|$, is greater than ϵ .

$$L_\epsilon(y, f(x)) = \max(0, |y - f(x)| - \epsilon)$$



SVR: The Primal Formulation

Minimizing Complexity and Error

We want to minimize the norm of w to control model complexity. To handle points outside the ϵ -tube, we introduce non-negative slack variables ξ_i (for points above the tube) and ξ_i^* (for points below).

Primal Optimization Problem

$$\min_{w, b, \xi, \xi^*} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n (\xi_i + \xi_i^*)$$

subject to:

$$y_i - (w^T x_i + b) \leq \epsilon + \xi_i$$

$$(w^T x_i + b) - y_i \leq \epsilon + \xi_i^*$$

$$\xi_i, \xi_i^* \geq 0 \quad \text{for } i = 1, \dots, n$$

The constant $C > 0$ is a regularization parameter that controls the trade-off between the flatness of $f(x)$ and the tolerance for errors.

SVR: The Dual Formulation

Derivation via Lagrange Multipliers

We introduce Lagrange multipliers $\alpha_i, \alpha_i^*, \mu_i, \mu_i^* \geq 0$ and form the Lagrangian. Taking derivatives with respect to the primal variables (w, b, ξ_i, ξ_i^*) and setting to zero gives:

$$\frac{\partial \mathcal{L}}{\partial w} = 0 \quad \Rightarrow \quad w = \sum_{i=1}^n (\alpha_i - \alpha_i^*) x_i$$

$$\frac{\partial \mathcal{L}}{\partial b} = 0 \quad \Rightarrow \quad \sum_{i=1}^n (\alpha_i - \alpha_i^*) = 0$$

$$\frac{\partial \mathcal{L}}{\partial \xi_i} = 0 \quad \Rightarrow \quad C - \alpha_i - \mu_i = 0 \quad \Rightarrow \quad \alpha_i \leq C$$

$$\frac{\partial \mathcal{L}}{\partial \xi_i^*} = 0 \quad \Rightarrow \quad C - \alpha_i^* - \mu_i^* = 0 \quad \Rightarrow \quad \alpha_i^* \leq C$$

The Dual Optimization Problem

Substituting these back into the Lagrangian yields the dual problem to be

SVR, Representer Theorem, and Kernels I

Making Predictions

The dual formulation reveals two key insights:

1. Connection to Representer Theorem

The expression for the optimal weight vector, $w = \sum_{i=1}^n (\alpha_i - \alpha_i^*) x_i$, shows that w is a linear combination of the input data points. This is exactly what the **Representer Theorem** states: the solution to this type of regularized problem lies in the span of the training data.

2. The Kernel Trick in Action

The dual objective function depends only on the dot product $x_i^T x_j$. This allows us to apply the kernel trick by replacing it with a kernel function $K(x_i, x_j)$.

SVR, Representer Theorem, and Kernels II

Making Predictions

Prediction for a New Point z

The prediction function becomes:

$$f(z) = w^T z + b = \sum_{i=1}^n (\alpha_i - \alpha_i^*) x_i^T z + b = \sum_{i=1}^n (\alpha_i - \alpha_i^*) K(x_i, z) + b$$

The data points with non-zero coefficients $(\alpha_i - \alpha_i^*)$ are the **Support Vectors**.

SVR in Practice: A Non-Linear Example I

We generate noisy sinusoidal data and fit an SVR model using a non-linear kernel (Radial Basis Function - RBF):

$$K(\mathbf{x}, \mathbf{z}) = \exp(-\gamma \|\mathbf{x} - \mathbf{z}\|_2^2)$$

The figure shows:

- The original noisy data points.
- The SVR prediction curve, which captures the underlying non-linear trend.
- The ϵ -tube around the prediction.
- The **Support Vectors** (points on or outside the tube) that define the model.

SVR in Practice: A Non-Linear Example II

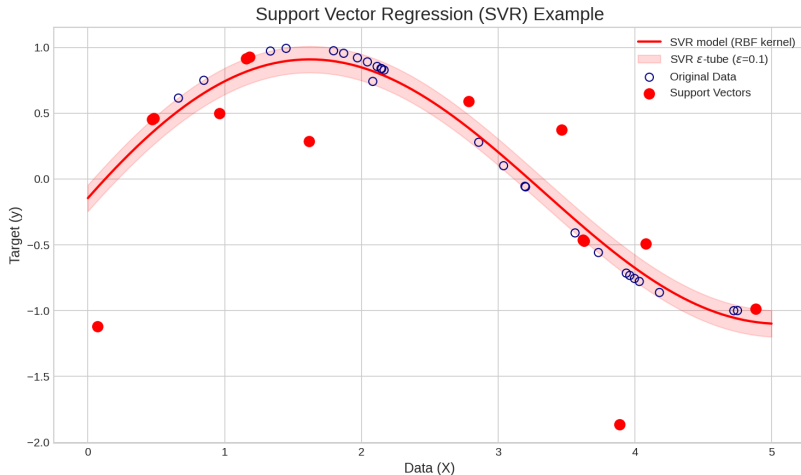


Figure: SVR with an RBF kernel fitting non-linear data.

Goal of Support Vector Classification (SVC)

Finding the Maximum-Margin Hyperplane

For a binary classification problem with labels $y_i \in \{-1, 1\}$, SVC aims to find a hyperplane (w, b) that separates the two classes with the largest possible margin.

The Margin

The margin is the distance between the two parallel hyperplanes $w^T x + b = 1$ and $w^T x + b = -1$. This distance can be shown to be $\frac{2}{\|w\|}$. Maximizing the margin is therefore equivalent to minimizing $\|w\|^2$.

The Hinge Loss Function I

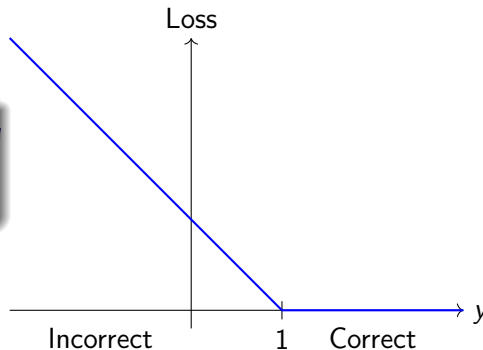
For SVC, the decision function is $f(x) = w^T x + b$. The product $y_i f(x_i)$ measures how correctly and confidently a point is classified.

- If $y_i f(x_i) \geq 1$: The point is correctly classified and outside the margin. No penalty.
- If $y_i f(x_i) < 1$: The point is inside the margin or misclassified. It incurs a linear penalty.

The Hinge Loss Function II

Hinge Loss Formula

$$L(y, f(x)) = \max(0, 1 - yf(x))$$



SVC: Primal and Dual Formulation I

Derivation with Soft Margins

To handle non-separable data, we introduce slack variables $\xi_i \geq 0$ (soft-margin SVC). The term $\sum \xi_i$ is an upper bound on the number of misclassifications and is equivalent to using the Hinge Loss.

Primal Optimization Problem

$$\min_{w, b, \xi} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i$$

subject to:

$$y_i(w^T x_i + b) \geq 1 - \xi_i \quad \text{and} \quad \xi_i \geq 0$$

SVC: Primal and Dual Formulation II

Derivation with Soft Margins

Dual Optimization Problem

Following a similar derivation with Lagrange multipliers, we arrive at the dual problem:

$$\max_{\alpha} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j K(x_i, x_j)$$

subject to:

$$\sum_{i=1}^n \alpha_i y_i = 0 \quad \text{and} \quad 0 \leq \alpha_i \leq C$$

The prediction for a new point z is $\text{sign}(\sum_{i=1}^n \alpha_i y_i K(x_i, z) + b)$.

SVR vs. SVC: Key Differences

Regression vs. Classification

While both are SVMs, their objectives and mechanics are fundamentally different.

Support Vector Regression (SVR)

- **Goal:** Fit a function to data.
- **Loss:** Uses ϵ -insensitive loss. Errors are ignored if they are within the ϵ -tube.
- **Constraints:** Two-sided constraints to keep data points inside the tube.
- **Support Vectors:** Points on the margin or outside the ϵ -tube.

Support Vector Classification (SVC)

- **Goal:** Find a separating boundary.
- **Loss:** Uses **Hinge Loss**, which penalizes points inside the margin or on the wrong side of the hyperplane.
- **Constraints:** One-sided constraint to ensure points are correctly classified with a margin.
- **Support Vectors:** Points on the margin or misclassified points.