Vectors $x_j$ for which $\xi_j = 0$, we have (7.31). Thus, the threshold can be obtained by averaging (7.32) over all Support Vectors $x_j$ (recall that they satisfy $\alpha_j > 0$) with $\alpha_j < C$.

In the above formulation, $C$ is a constant determining the trade-off between two conflicting goals: minimizing the training error, and maximizing the margin. Unfortunately, $C$ is a rather unintuitive parameter, and we have no a priori way to select it.[9] Therefore, a modification was proposed in [481], which replaces $C$ by

$\nu$-SVC

a parameter $\nu$; the latter will turn out to control the number of margin errors and Support Vectors.

As a primal problem for this approach, termed the $\nu$-SV classifier, we consider

$$\underset{\mathbf{w} \in \mathcal{H}, \boldsymbol{\xi} \in \mathbb{R}^m, \rho, b \in \mathbb{R}}{\text{minimize}} \quad \tau(\mathbf{w}, \boldsymbol{\xi}, \rho) = \frac{1}{2}\|\mathbf{w}\|^2 - \nu\rho + \frac{1}{m}\sum_{i=1}^{m}\xi_i \tag{7.40}$$

$$\text{subject to} \quad y_i(\langle \mathbf{x}_i, \mathbf{w} \rangle + b) \geq \rho - \xi_i \tag{7.41}$$

$$\text{and} \quad \xi_i \geq 0, \quad \rho \geq 0. \tag{7.42}$$

Note that no constant $C$ appears in this formulation; instead, there is a parameter $\nu$, and also an additional variable $\rho$ to be optimized. To understand the role of $\rho$, note that for $\boldsymbol{\xi} = 0$, the constraint (7.41) simply states that the two classes are separated by the *margin* $2\rho/\|\mathbf{w}\|$ (cf. Problem 7.4).

Margin Error

To explain the significance of $\nu$, let us first recall the term *margin error*: by this, we denote points with $\xi_i > 0$. These are points which are either errors, or lie within the margin. Formally, the fraction of margin errors is

$$R_{\text{emp}}^\rho[g] := \frac{1}{m}\left|\{i \mid y_i g(x_i) < \rho\}\right|. \tag{7.43}$$

Here, $g$ is used to denote the argument of the sgn in the decision function (7.25): $f = \text{sgn} \circ g$ (see footnote 5, p. 344). We are now in a position to state a result that

$\nu$-Property

explains the significance of $\nu$.

**Proposition 7.5 ([481])** *Suppose we run $\nu$-SVC with k on some data with the result that $\rho > 0$. Then*

*(i) $\nu$ is an upper bound on the fraction of margin errors.*

*(ii) $\nu$ is a lower bound on the fraction of SVs.*

*(iii) Suppose the data $(x_1, y_1), \ldots, (x_m, y_m)$ were generated iid from a distribution $\mathrm{P}(x, y) = \mathrm{P}(x)\mathrm{P}(y|x)$, such that neither $\mathrm{P}(x, y = 1)$ nor $\mathrm{P}(x, y = -1)$ contains any discrete component. Suppose, moreover, that the kernel used is analytic and non-constant. With probability 1, asymptotically, $\nu$ equals both the fraction of SVs and the fraction of errors.*

The proof can be found in Section A.2.

Before we get into the technical details of the dual derivation, let us take a look

---

9. As a default value, we use $C/m = 10$ unless stated otherwise.