

## 支持向量回归（SVR）原理介绍

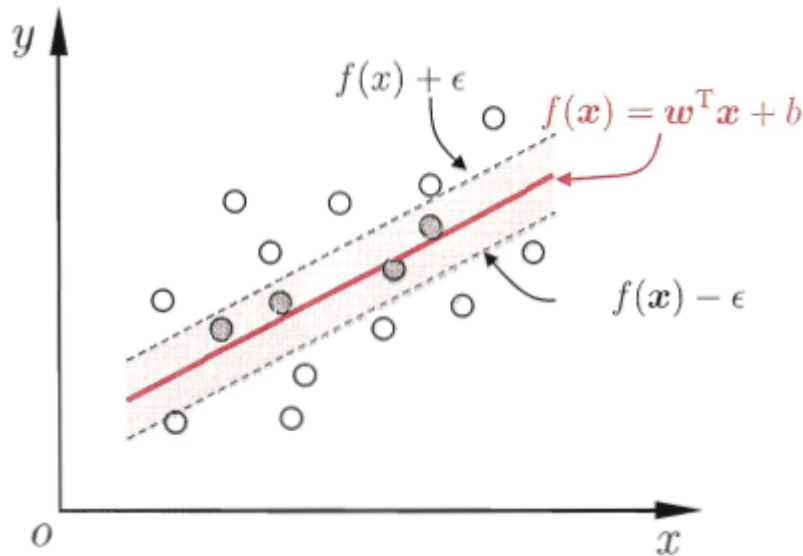
主要参考周志华老师的《机器学习》中对应内容。

考虑回归问题，给定训练样本  $D = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_m, y_m)\}$ ，我们希望得到如下回归模型：

$$f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b \quad (1)$$

使得  $f(\mathbf{x})$  和  $y$  尽可能接近。

传统的回归模型通常直接基于模型输出  $f(\mathbf{x})$  与真实输出  $y$  之间的差别来计算 loss，当且仅当  $f(\mathbf{x})$  和  $y$  完全相同时，loss 才为 0。与此不同，SVR 假设我们能容忍  $f(\mathbf{x})$  与  $y$  之间最多有  $\epsilon$  的偏差，即仅当  $f(\mathbf{x})$  与  $y$  之间的差别绝对值大于  $\epsilon$  时才计算 loss。如下图所示，这相当于以  $f(\mathbf{x})$  为中心，构建了一个宽度为  $2\epsilon$  的间隔带，若训练样本落入此间隔带，则认为是被预测正确的。



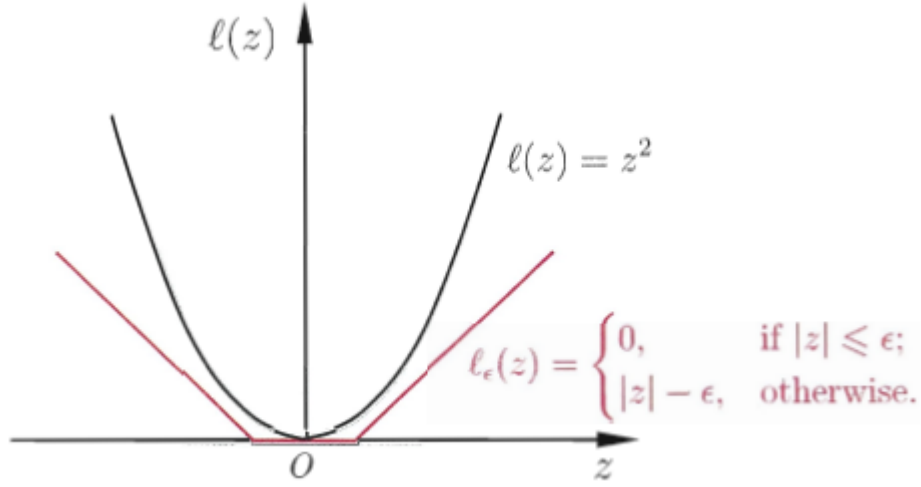
于是，SVR 问题形式化为

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m \ell_{\epsilon}(f(\mathbf{x}_i) - y_i) \quad (2)$$

其中  $C$  为正则化系数， $\ell_{\epsilon}$  是  $\epsilon$ -不敏感损失函数：

$$\ell_{\epsilon}(z) = \begin{cases} 0, & \text{if } |z| \leq \epsilon \\ |z| - \epsilon, & \text{otherwise} \end{cases} \quad (3)$$

其函数图像为：



样本点落在 $\epsilon$ -间隔带中的条件为：

$$f(x_i) - \epsilon \leq y_i \leq f(x_i) + \epsilon \quad (4)$$

我们引入松弛变量 $\xi_i$ 和 $\hat{\xi}_i$ （间隔带两侧的松弛程度不同），条件（4）变成：

$$f(x_i) - \epsilon - \hat{\xi}_i \leq y_i \leq f(x_i) + \epsilon + \xi_i \quad (5)$$

于是优化问题可以重写为：

$$\begin{aligned} \min_{w, b, \xi_i, \hat{\xi}_i} \quad & \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m (\xi_i + \hat{\xi}_i) \\ \text{s.t.} \quad & f(x_i) - y_i \leq \epsilon + \xi_i, \\ & y_i - f(x_i) \leq \epsilon + \hat{\xi}_i, \\ & \xi_i \geq 0, \hat{\xi}_i \geq 0, i = 1, 2, \dots, m. \end{aligned} \quad (6)$$

引入拉格朗日乘子 $\mu_i \geq 0, \hat{\mu}_i \geq 0, \alpha_i \geq 0, \hat{\alpha}_i \geq 0$ ，构造拉格朗日函数为：

$$\begin{aligned} L(w, b, \alpha, \hat{\alpha}, \xi, \hat{\xi}, \mu, \hat{\mu}) \\ = \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m (\xi_i + \hat{\xi}_i) - \sum_{i=1}^m \mu_i \xi_i - \sum_{i=1}^m \hat{\mu}_i \hat{\xi}_i \\ + \sum_{i=1}^m \alpha_i (f(x_i) - y_i - \epsilon - \xi_i) + \sum_{i=1}^m \hat{\alpha}_i (y_i - f(x_i) - \epsilon - \hat{\xi}_i) \end{aligned} \quad (7)$$

令 $L(w, b, \alpha, \hat{\alpha}, \xi, \hat{\xi}, \mu, \hat{\mu})$ 对 $w, b, \xi_i$ 和 $\hat{\xi}_i$ 的偏导为零，可得：

$$w = \sum_{i=1}^m (\hat{\alpha}_i - \alpha_i) x_i \quad (8)$$

$$0 = \sum_{i=1}^m (\hat{\alpha}_i - \alpha_i) \quad (9)$$

$$C = \alpha_i + \mu_i \quad (10)$$

$$C = \hat{\alpha}_i + \hat{\mu}_i \quad (11)$$

将式（8）～（11）代入式（7），即可得到SVR的对偶问题

$$\begin{aligned}
& \max_{\alpha, \hat{\alpha}} \sum_{i=1}^m y_i (\hat{\alpha}_i - \alpha_i) - \epsilon (\hat{\alpha}_i + \alpha_i) \\
& \quad - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m (\hat{\alpha}_i - \alpha_i) (\hat{\alpha}_j - \alpha_j) \mathbf{x}_i^T \mathbf{x}_j \\
& \text{s.t.} \quad \sum_{i=1}^m (\hat{\alpha}_i - \alpha_i) = 0 \\
& \quad 0 \leq \alpha_i, \hat{\alpha}_i \leq C
\end{aligned} \tag{12}$$

上述过程满足KKT条件，有

$$\begin{cases}
\alpha_i (f(\mathbf{x}_i) - y_i - \epsilon - \xi_i) = 0 \\
\hat{\alpha}_i (y_i - f(\mathbf{x}_i) - \epsilon - \hat{\xi}_i) = 0 \\
\alpha_i \hat{\alpha}_i = 0, \xi_i \hat{\xi}_i = 0 \\
(C - \alpha_i) \xi_i = 0, (C - \hat{\alpha}_i) \hat{\xi}_i = 0
\end{cases} \tag{13}$$

(这里有一个问题：式 (13) 中的第三行是如何通过KKT条件得到的?)

可以看出，当且仅当  $f(\mathbf{x}_i) - y_i - \epsilon - \xi_i = 0$  时  $\alpha_i$  能取非零值，当且仅当  $y_i - f(\mathbf{x}_i) - \epsilon - \hat{\xi}_i = 0$  时  $\hat{\alpha}_i$  能取非零值。换言之，仅当样本  $(\mathbf{x}_i, y_i)$  不落入  $\epsilon$ -间隔带中，相应的  $\alpha_i$  和  $\hat{\alpha}_i$  才能取非零值。此外，约束  $f(\mathbf{x}_i) - y_i - \epsilon - \xi_i = 0$  和  $y_i - f(\mathbf{x}_i) - \epsilon - \hat{\xi}_i = 0$  不能同时成立，因此  $\alpha_i$  和  $\hat{\alpha}_i$  中至少有一个为零。

将式 (8) 代入式 (1)，可得SVR的解：

$$f(\mathbf{x}) = \sum_{i=1}^m (\hat{\alpha}_i - \alpha_i) \mathbf{x}_i^T \mathbf{x} + b \tag{14}$$

能使式 (14) 中的  $(\hat{\alpha}_i - \alpha_i) \neq 0$  的样本即为SVR的支持向量，它们必落在  $\epsilon$ -间隔带之外。显然，SVR的支持向量仅是训练样本的一部分，即其解仍具有稀疏性。

由KKT条件可以看出，对每个样本  $(\mathbf{x}_i, y_i)$  都有  $(C - \alpha_i) \xi_i = 0$  且  $\alpha_i (f(\mathbf{x}_i) - y_i - \epsilon - \xi_i) = 0$ 。于是，在得到  $\alpha_i$  后，若  $0 < \alpha_i < C$ ，则必有  $\xi_i = 0$ ，进而有

$$b = y_i + \epsilon - \sum_{i=1}^m (\hat{\alpha}_i - \alpha_i) \mathbf{x}_i^T \mathbf{x} \tag{15}$$

因此，在求解式 (12) 得到  $\alpha_i$  后，理论上来说，可任意选取满足  $0 < \alpha_i < C$  的样本通过式 (15) 求得  $b$ 。实践中常采用一种更鲁棒的方法：选取多个（或所有）满足条件  $0 < \alpha_i < C$  的样本求解  $b$  后取平均值。

同样，我们可以在SVR中引入核函数，此时式 (8) 变为：

$$\mathbf{w} = \sum_{i=1}^m (\hat{\alpha}_i - \alpha_i) \phi(\mathbf{x}_i) \tag{16}$$

而SVR可表示为：

$$f(\mathbf{x}) = \sum_{i=1}^m (\hat{\alpha}_i - \alpha_i) \kappa(\mathbf{x}, \mathbf{x}_i) + b \tag{17}$$

其中  $\kappa(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$  为核函数。