

Support vector machines *

In a p -dimensional space, a hyperplane is a flat affine subspace of dimension $p-1$.
 $B_0 + B_1 x_1 + \dots + B_p x_p = 0$ defines a p -dimensional hyperplane, in the sense that if a point $x = (x_1, \dots, x_p)^T$ in p -dimensional space satisfies, then x lies on the hyperplane. if > 0 above the plane, < 0 below it.
 • we compute the perpendicular distance from each point to the plane, the smallest is the margin, and we want to maximize it. we can classify test on which side of the maximal margin hyperplane it lies.
 The hyperplane thus only depends on the support vectors, but not the others.

To get the maximal margin classifier maximize M
 B_0, B_1, \dots, B_p, M subject to $\sum_{j=1}^p B_j^2 = 1$

The support vector classifier (soft margin classifier), rather than seeking the largest possible margin, we allow some observations to be on the incorrect side of the margin, as even the incorrect side of the hyperplane maximize M

$$B_0, B_1, \dots, B_p, \epsilon_1, \dots, \epsilon_m, M \quad \text{subject to } \sum_{j=1}^p B_j^2 = 1$$

$\epsilon_1, \dots, \epsilon_m$ are slack variables that

allow observations to be on the wrong side of the margin or hyperplane.

if $\epsilon_i = 0$ then i th observation is on the right side, $\epsilon_i > 0$ wrong side of the margin, $\epsilon_i > 1$ wrong side of hyperplane
 • for $c > 0$ no more than c observations can be on the wrong side of the hyperplane. when c is small, we seek narrow margins that are rarely violated (low bias but high variance). when c is larger, the margin is wider and we allow more violations (more biased, less variance).

$$y_i (B_0 + B_1 x_{i1} + B_2 x_{i2} + \dots + B_p x_{ip}) \geq M(1 - \epsilon_i)$$

$$\epsilon_i \geq 0, \sum_{i=1}^m \epsilon_i \leq c \quad (c \text{ chosen by CV})$$

The support vectors with soft margin are the ones that lie directly on the margin, on the wrong side of the margin, hence when c is large \rightarrow margin is larger, \rightarrow many observations are support vector \rightarrow reduce variance.

The support vector machines (for non linear class)

we can change the feature space $(x_1, x_2, \dots, x_p \rightarrow x_1, x_1^2, x_2, x_2^2, \dots)$. The solution is linear in the new space but not in the former one where it is non-linear.

To enlarge the space in a computational friendly way we can use kernel.

The inner product of two observations x_i, x_j is given by $\langle x_i, x_j \rangle = \sum_{j=1}^p x_{ij} x_{ij}$

we can represent the linear machine classifier as $f(x) = B_0 + \sum_{i=1}^m \alpha_i \langle x, x_i \rangle$ where there are m parameters α_i , $i = 1, \dots, m$ one per training observation.
 to estimate $\alpha_1, \dots, \alpha_m$ and B_0 , we need the $\binom{m}{2}$ inner product between all pairs of training. But α_i is non zero only for the support vector.

Now suppose that every time the inner product appears in the representation, as in the calculation, we replace with a generalization of the inner product \rightarrow a kernel

Linear kernel $K(x_i, x_j) = \sum_{j=1}^p x_{ij} x_{ij}$ (classifier quantifies the similarities of a pair of observations using Pearson's standard correlation)

Polynomial kernel $K(x_i, x_j) = (1 + \sum_{j=1}^p x_{ij} x_{ij})^d$. when the support vector classifier is combined with a non linear kernel, the result is a support vector machine, of the form $f(x) = B_0 + \sum_{i \in S} \alpha_i K(x, x_i)$

Radial kernel $K(x_i, x_j) = \exp(-\gamma \sum_{j=1}^p (x_{ij} - x_{ij})^2)$. It has local behavior, only nearby training observations have an effect on the class label.

SVM with more than two classes

- **one versus one**: if there are $k > 2$ classes, we construct $\binom{k}{2}$ SVM, each of which compares a pair of classes. The final class is the most assigned.
- **one versus all**, we fit k SVM, each time comparing one of the k classes to the remaining $k-1$ classes.

Logistic regression

$f(x) = B_0 + B_1 x_1 + \dots + B_p x_p$ as minimize $\sum_{i=1}^n \max[0, 1 - y_i f(x_i)] + \lambda \sum_{j=1}^p B_j^2$ when λ is large, B_0, \dots, B_p are small \rightarrow max violation \rightarrow low variance high bias
 \Rightarrow generalization minimize $\sum_{i=1}^n L(x_i, y_i, \beta) + \lambda P(\beta)$

Support vector regression

Seeks coefficients that minimize a different type of loss, where only residuals larger in absolute value than some positive constant contribute to the loss function. This is an extension of the margin used in support vector classifier to the regression setting.