

Notes on SVM

Shasha Liao
Georgia Tech

November 5, 2020

1 Classification

- Decision boundary: $w^T x + b = 0$
- Positive class: (x, y) with $w^T x + b \geq c$ and $y > 1$.
- Negative class: (x, y) with $w^T x + b \leq -c$ and $y > 1$.
- Margin: $\gamma = \frac{2c}{\|w\|}$, the distance between two hyperplanes $w^T x + b = c$ and $w^T x + b = -c$.
- Correct classification condition: $(w^T x_i + b)y_i \geq c$ for all data points (x_i, y_i) , $i = 1, 2, \dots, m$.

2 Maximum margin classifier

Find decision boundary, i.e. w and b to maximize the margin.

$$\max_{w,b} \gamma = \frac{2c}{\|w\|} \text{ s.t. } (w^T x_i + b)y_i \geq c \text{ for all } i. \quad (1)$$

This becomes an optimization problem with constraints. And the solution is invariant under the scaling of c . So we set $c = 1$ and the problem (1) is equivalent to the following constrained convex quadratic programming problem:

$$\min_{w,b} \|w\| \text{ s.t. } (w^T x_i + b)y_i \geq 1 \text{ for all } i = 1, 2, \dots, m. \quad (2)$$

- Only a few of the constraints are relevant \rightarrow **support vectors**
- Kernel methods are introduced for nonlinear classification with nonlinear decision boundary.

To solve the optimization problem (2), we need to use Lagrangian multipliers.

2.1 Lagrangian Duality

The Lagrangian function:

$$L(w, b, \alpha) = \frac{1}{2}w^T w + \sum_{i=1}^m \alpha_i (1 - y^i (w^T x^i + b)). \quad (3)$$

The KKT conditions:

$$\frac{\partial L}{\partial w} = 0, \quad \frac{\partial L}{\partial b} = 0, \quad \alpha_i (1 - y^i (w^T x^i + b)) = 0, \alpha_i \geq 0, \forall i.$$

The last equation implies that all the points that are not on the margin will have the corresponding $\alpha_i = 0$. The first equation gives

$$w = \sum_{i=1}^m \alpha_i y^i x^i,$$

this tells us that w is a linear combination of the support vectors, i.e., the vectors on the margins.

From the KKT conditions we can solve for optimal w^* and b^* , plug in them into L we get the dual Lagrangian problem:

$$\max_{\alpha} L(w^*, b^*, \alpha) = \sum_i^m \alpha_i - \frac{1}{2} \sum_{i,j}^m \alpha_i \alpha_j y^i y^j ((x^i)^T x^j) \quad (4)$$

$$\text{s.t. } \alpha_i \geq 0, i = 1, 2, \dots, m \quad (5)$$

$$\sum_i^m \alpha_i y^i = 0. \quad (6)$$

This is a constrained quadratic programming. $L(w^*, b^*, \alpha)$ is nice and convex, and the global maximum can be found. For more details, see the slides. The optimal α_i can be found by solving the dual Lagrangian problem.

After that, we can find w, b in terms of α_i for i with x_i on the margins, i.e., x_i being support vector.

For a new test point z , compute

$$w^T z + b = \sum_{i \in \text{support vectors}} \alpha_i y^i ((x^i)^T z) + b.$$

If $w^T z + b > 0$, z is predicted to be in the positive class. Otherwise, z is predicted to be in the negative class.

3 Generalized SVM

- Kernelized SVM

Since the Lagrangian and the decision rule only depends on the pairwise inner products of the points, we can replace all the inner products with nonlinear kernel functions to obtain nonlinear decision boundary.

- Soft-margin SVM

When the data is not linearly separable, we allow a few points to violate the hard margin constraint and solve the following optimization problem:

$$\begin{aligned} \min_{w,b,\xi} \quad & \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m \xi^i \\ \text{s.t.} \quad & y^i(w^T x^i + b) \geq 1 - \xi^i, \xi^i \geq 0, \forall i. \end{aligned}$$

Now the examples are allowed to have functional margin less than 1, and if an example has a margin $1 - \xi^i$, we would pay a cost of the objective function being increased by $C\xi^i$. The parameter C controls the relative weighting between the twin goals of making the $\|w\|^2$ large (which makes the margin small) and of ensuring that most examples have functional margin at least 1.

- If C is large, the model will not allow examples to have functional margin less than 1. Only in this case, we achieve maximal margin classifier.
- If C is small, the model allows for incorrectly classified examples and has a wider margin. This is good for the situation when we have outliers and we don't want to affect the decision boundary.

4 Choosing SVM Parameters

Choosing C :

- Large $C \Rightarrow$ overfitting, higher variance, lower bias
- Small $C \Rightarrow$ underfitting, low variance, high bias

Choose σ^2 for Gaussian Kernel:

- Large $\sigma^2 \Rightarrow$ underfitting, low variance, high bias
- Small $\sigma^2 \Rightarrow$ overfitting, higher variance, lower bias

The above optimization problem can be solved using SMO algorithm.