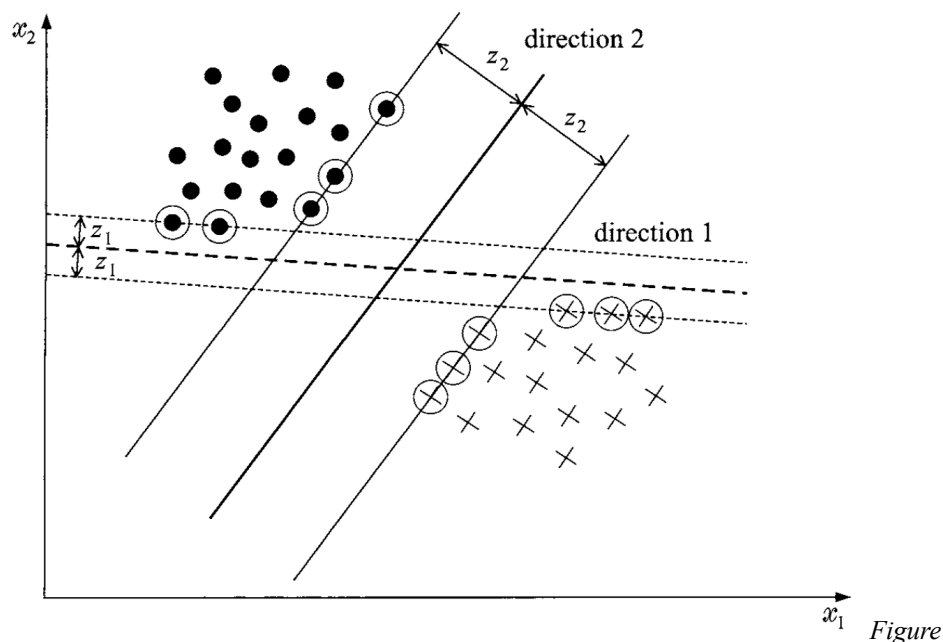


A16 – Support Vector Machines

Introduction

The Perceptron algorithm calculates the decision line between two linearly separable classes clustered in feature space. We note that the solution found by the Perceptron is not unique. Running the algorithm with a different set of initial weights will give you a different but equally valid solution.

What if we want to find the best decision line between the two classes? By “best” we mean that the decision line should be the one that offers the best generalization when tested with novel data. Visually, the best decision line should have the largest margin between the two classes as shown in Figure 1. For linearly separable classes the solution to the best decision line is unique.



1. Decision line along direction 2 has a larger margin (expressed as perpendicular distance z) than that of direction 1. The feature vectors that are closest to the decision lines are known as the support vectors. Figure from Chapter 3, Theodoridis 'Pattern Recognition'.

If the best decision line is give by

$$g(\mathbf{x}) = w_o + \mathbf{x}^T \mathbf{w} = 0 \quad (1)$$

the perpendicular distance z is given by

$$z = \frac{|g(\mathbf{x})|}{\|\mathbf{w}\|} \quad (2)$$

We can scale \mathbf{w} and w_o such that the value of $g(\mathbf{x})$ is 1 for ω_1 and -1 for ω_2 . This is equivalent to imposing that

1. The margin given by $\frac{1}{\|\mathbf{w}\|} + \frac{1}{\|\mathbf{w}\|} = \frac{2}{\|\mathbf{w}\|}$ is maximized;
2. And that
$$\begin{aligned} w_o + \mathbf{x}^T \mathbf{w} &\geq 1, \forall \mathbf{x} \in \omega_1 \\ w_o + \mathbf{x}^T \mathbf{w} &\leq -1, \forall \mathbf{x} \in \omega_2 \end{aligned}$$

Recall that for each data point \mathbf{x}_i we indicate the class by y_i which is +1 for ω_1 and -1 for ω_2 . Our two conditions above can therefore be efficiently rewritten as follows:

minimize :

$$J(\mathbf{w}) = \frac{1}{2} \|\mathbf{w}\|^2 \quad (3)$$

subject to:

$$y_i (w_o + \mathbf{x}_i^T \mathbf{w}) \geq 1, i = 1, 2, \dots, N \quad (4)$$

Using Lagrange multipliers λ_i we can convert the margin maximization problem given by:

$$\begin{aligned} \min L_p &= \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^N \lambda_i y_i (w_o + \mathbf{x}_i^T \mathbf{w}) + \sum_{i=1}^N \lambda_i \\ \text{subject to} \quad & \mathbf{w} = \sum_{i=1}^N \lambda_i y_i \mathbf{x}_i \\ & \sum_{i=1}^N \lambda_i y_i = 0 \\ & \lambda_i \geq 0 \end{aligned} \quad (5)$$

where N is the number of training points, to a problem involving only dot products given by

$$\begin{aligned}
 \max L_D(\lambda_i) &= \sum_{i=1}^N \lambda_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \lambda_i \lambda_j y_i y_j (\mathbf{x}_i \cdot \mathbf{x}_j) \\
 \text{subject to} \quad & \sum_{i=1}^l \lambda_i y_i = 0 \\
 & \lambda_i \geq 0.
 \end{aligned} \tag{6}$$

Finding the decision line with the widest margin from the nearest feature vectors using Equation 6 is an optimization problem. Using quadratic problem solver packages we can solve for λ_i and use the expression for \mathbf{w} in Equation 5 to solve for the weights \mathbf{w} and bias w_0 .

Procedure

1. Similar to the Perceptron, SVM is applied to pairs of classes. Read through Reference 2 slides and for each pair of fruits set up the matrices \mathbf{H} , \mathbf{B} and \mathbf{A} , and vectors \mathbf{f} , \mathbf{a} and \mathbf{b} found in slides 19 to 22.
2. Use quadratic programming packages available from Python (numpy and cvxopt) or Matlab (quadprog in Optimization toolbox) using the vectors and matrices setup in Step 1 to solve for the Lagrange multipliers λ_i .
3. Solve for the weights \mathbf{w} and bias w_0 and plot the decision lines in feature space with your fruit data.

Reference

1. Theodoridis, Chapter 3, Pattern Recognition
2. Veksler, O., CS 434a/541a: Pattern Recognition Lecture 11 slides.