

Homework #3

Due Mar 13 (Monday) at the beginning of class

Please show all work and intermediate steps. Late submission will lead to 0 credit.

Problem 1: Support Vector Machine

Support vector machine (SVM) is a popular model in machine learning used for classification. Mathematically, given a training dataset of m points

$$(x_1, y_1), \dots, (x_m, y_m)$$

where $x_i \in \mathbf{R}^n$ stands for the feature vector and $y_i \in \{1, -1\}$ stands for two classes. The goal is to find two parallel hyperplanes represented by (w, b) with maximal margin that separates the two classes of data, such that for class with $y_i = 1$, we have $w^T x_i + b \geq 1$ and for class with $y_i = -1$, we have $w^T x_i + b \leq -1$. Hence, we wish to satisfy $y_i(w^T x_i + b) \geq 1$ for $i = 1, \dots, m$.

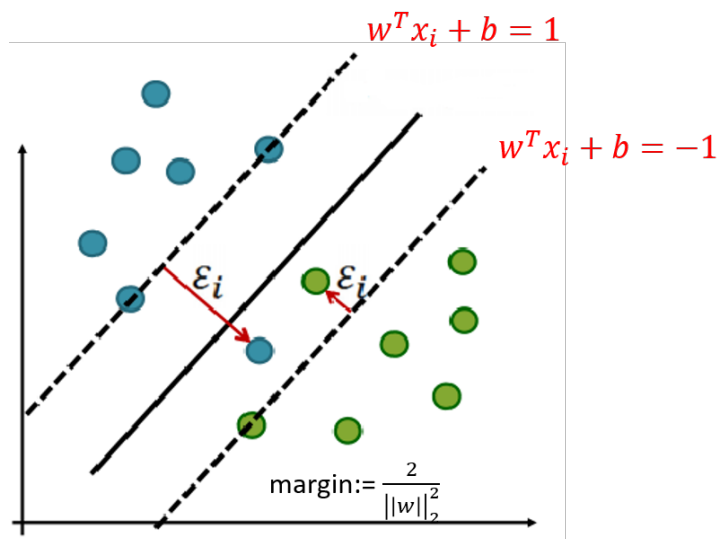


Figure 1: Illustration of SVM

If the data is not fully separable, we allow for small margin errors $\epsilon_i > 0, i = 1, \dots, m$, and we wish to also minimize these errors. This leads to solving the following optimization problem:

$$\begin{aligned} \min_{w, b, \epsilon} \quad & \frac{1}{2} \|w\|_2^2 + C \cdot \sum_{i=1}^m \epsilon_i \\ \text{s.t.} \quad & y_i(w^T x_i + b) \geq 1 - \epsilon_i, \quad i = 1, \dots, m \\ & \epsilon_i \geq 0, \quad i = 1, \dots, m \end{aligned} \tag{P}$$

where the parameter $C > 0$ plays a role of controlling the relative importance of minimizing the norm of w (i.e., maximizing the margin) and minimize the errors. Note that this problem is indeed a convex optimization problem.

Exercise 1.1 (Lagrange Duality) Let $\alpha \geq 0$ and $\beta \geq 0$ be the Lagrange multipliers associated with the two constraints. Show that the Lagrange dual problem of (P) is given by the quadratic program:

$$\begin{aligned} \max_{\alpha} \quad & \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j (x_i^T x_j) \\ & \sum_{i=1}^m \alpha_i y_i = 0 \\ & 0 \leq \alpha_i \leq C, \quad i = 1, \dots, m \end{aligned} \tag{D}$$

Moreover, show that the primal and dual optimal solutions satisfy that

$$\begin{aligned} \alpha_i = 0 & \Rightarrow y_i(w^T x_i + b) \geq 1 \\ \alpha_i = C & \Rightarrow y_i(w^T x_i + b) \leq 1 \\ 0 < \alpha_i < C & \Rightarrow y_i(w^T x_i + b) = 1 \end{aligned}$$

We call the data points with non-zero Lagrangian multipliers the *support vectors*.

Exercise 1.2 (Reformulation) Show that (P) can be equivalently written as an unconstrained convex problem

$$\min_{w,b} \frac{1}{m} \sum_{i=1}^m \max(1 - y_i(w^T x_i + b), 0) + \lambda \|w\|_2^2 \tag{P'}$$

where $\lambda > 0$ is some parameter.

Exercise 1.3 (Programming) Implement the Ellipsoid method to solve the problem (P') in Matlab or Python whichever you prefer. Your input should be the data matrix X , y and the parameter λ , and the maximum number of iterations T . Your output should be the best solution and objective function value obtained after running the algorithm within T iterations.

Note: When initializing the Ellipsoid method, you need to have some prior knowledge of the bound of the solution. You can either simply set a large enough ball or you may derive some simple bounds for w, b . For example, let $f(w, b)$ be the objective in (P') , then we know

$$\lambda \|w^*\|_2^2 \leq f(w^*, b^*) \leq f(0, 0) = 1$$

Hence, $\|w^*\|_2 \leq 1/\sqrt{\lambda}$.

Exercise 1.4 (Test on Real Dataset) Apply your algorithm with $T = 100$ iterations on the Wisconsin Diagnostic Breast Cancer (WDBC) dataset ($n = 30, m = 569$) provided (read [here](#) for detailed description of the dataset) with $\lambda = 1$.

- Plot the objective function values at current solution, i.e. $f(w_t)$ vs the number of iteration t ;
- On the same figure, plot the objective function values at best solution, i.e. $\min_{1 \leq \tau \leq t} f(w_\tau)$ vs the number of iteration t ;
- Compute the classification error: the ratio of misclassified points (i.e. $y_i(w^T x_i + b) < 1$).