

Numerical Optimisation Project

Binary Classification Task on Spiral Data based on SVM

Bowen Ding(21118015)

April 1, 2022

Abstract

In this Project, we particularly explore a binary classification task on Spiral Data(a non-linear separable data set) based on Support Vector Machine(SVM). Under this case, we introduce kernel method and convert the primal SVM format into dual format to improve the computation efficiency. In the process of solving the optimal parameters, we utilise Sequential Minimal Optimisation(SMO) algorithm with two different working set selection tactics: SMO-Simplified and SVM-Light , afterwards we benchmark and compare the performance and convergence of them.

Contents

1	Introduction	1
2	Theory	2
2.1	Support Vector Machine	2
2.1.1	Hard-margin SVM	2
2.1.2	Soft-margin SVM	3
2.1.3	Kernel SVM	3
2.2	Numerical Optimisation Method for SVM	4
2.2.1	Sequential Minimal Optimisation Procedure	5
2.2.2	Working Set Selection	5
2.2.3	Updating Parameters	7
3	Experiment and Analysis	7
3.1	Analysis on Working Set Selection Strategy	8
3.2	Analysis on Local Convergence Property	8
3.3	Analysis on Parameters	10
3.4	Analysis on Complexity	11
3.4.1	Time Complexity	11
3.4.2	Space Complexity	11
3.5	Analysis on Classification Performance	11
4	Conclusion	12

1 Introduction

In this project, all of our following work is based on spiral data set. The spiral data set contains m two-dimensional samples and all of them are tagged with binary label. Label y_i of a specified sample $\mathbf{x}_i = [x_{i1}, x_{i2}]^T$ is tagged as '+1' if $\begin{cases} x_{i1} = r_i \sin 2r_i + \varepsilon \\ x_{i2} = r_i \cos 2r_i + \varepsilon \end{cases}$ while is '-1' if $\begin{cases} x_{i1} = r_i \sin (2r_i + \pi) + \varepsilon \\ x_{i2} = r_i \cos (2r_i + \pi) + \varepsilon \end{cases}$, where $r_i = 1 + (i - 1) \frac{2\pi}{m-1}$, $i = 1, \dots, m$ and ε is noise term submitted to $\mathcal{N}(0, 0.2)$.

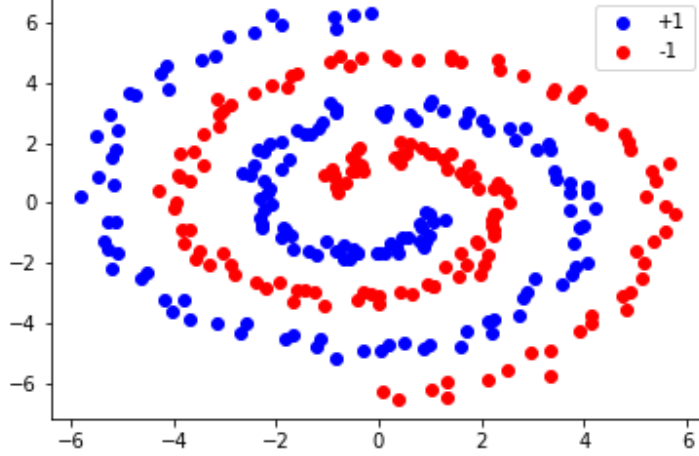


Figure 1: Visualisation of Spiral Data Set

In order to benchmark the binary classification performance of SVM, we separate the generated data into training set (300 samples) and test set (300 samples) uniformly. However, original primal SVM format implementation on this data set does not work due to some challenges [?]:

- It is obvious that this data set is not linear separable, which means we are unlikely to find a hyper-plane to classify '+1' and '-1' accurately. However, kernel trick helps solve this problem subtly by mapping the original sample points into a linear separable high dimensional space. So, we will implement kernel SVM on this problem.
- High or infinite dimension of mapped samples leads to the learnable parameter \mathbf{w} being high dimensional or infinite for primal format SVM, making the training process quite low efficiency or sometimes infeasible(for infinite dimension situation). Therefore, dual format SVM is utilised to guarantee the learnable parameter into finite dimensional(the size of train set).
- The dual kernel SVM essentially abstracts the original problem into a large dense quadratic programming problem with constraint qualification. In other words, the kernel matrix is an $m \times m$ fully dense matrix, where m is the size of training set. When m goes quite large, it will cost quite large computer memory to store the kernel matrix, let alone calculating Hessian matrix of the kernel matrix. Hence, Newton and Quasi Newton method cannot be applied in the optimisation stage of SVM straight-forwardly.

We organise the report as below. In section 2, we will give a brief introduction to SVM format used on this problem, as well as theoretical fundamentals of optimisation methods applied in SVM. Section 3 presents the comparison between two working set selection in terms of applicability, convergence, CPU time and memory efficiency based on numerical experiments. Section 4 gives a summary of our experiment result.

2 Theory

2.1 Support Vector Machine

2.1.1 Hard-margin SVM

Given training set $D = \{(\mathbf{x}_k, y_k)\}_{k=1}^m, y_k \in \{1, -1\}$. We want to have a hyper-plane $(\mathbf{w}, b) : \mathbf{w}^T \mathbf{x} + b = 0$ to separate all data accurately, where \mathbf{w} is normal vector of the hyper-plane and b is bias determining the intercept. Our aim is to find an optimal hyper-plane to separate all data correctly with the maximum margin. The margin refers to the minimal distance from samples in the data set to the hyper-plane, which is denoted as $r = \frac{|\mathbf{w}^T \mathbf{x} + b|}{\|\mathbf{w}\|} = \frac{1}{\|\mathbf{w}\|}$ and the data points lying on the margin are called support vectors. Therefore, the binary task can be described as a quadratic programming problem(primal hard-margin SVM format) [3]:

$$\begin{aligned} \min_{\mathbf{w}, b} \quad & \frac{\mathbf{w}^T \mathbf{w}}{2} \\ \text{s.t.} \quad & y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1, i = 1, \dots, m \end{aligned}$$

Note: $\max_{\mathbf{w}, b} \frac{2}{\|\mathbf{w}\|} \Leftrightarrow \min_{\mathbf{w}, b} \frac{\mathbf{w}^T \mathbf{w}}{2}$

2.1.2 Soft-margin SVM

Basically, primal hard-margin SVM aims to make sure the selected hyper-plane to classify all the data correctly, which required data space to be absolutely linearly separable. But for most data set, linearly separable condition is hard to satisfy. For such cases, we allow training errors occur. So we introduce slack variables $\xi_i \geq 0$ to weaken the constraints. Therefore, the soft-margin SVM can be shown as below:

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

The slack variable ξ_i measure in what extent the sample not observe the primal constraints. For example, \mathbf{x}_i will not be on the accurate side of the hyper-plane when $\xi_i \geq 1$. C refers to a large penalty parameter, making sure the optimal ξ_i tends to be 0 after several iterations of numerical optimisation.

2.1.3 Kernel SVM

The introduction of soft margin can tolerate training error to some extent, but it still utilises a linear boundary to divide data points in original data space. But for the data with a distribution like the spiral data set where it is impossible to use a linear boundary to separate the data points into correct side. So, we have to introduce the kernel trick to map the data into a higher linearly dimensional data space where we can find a proper hyper-plane. Here, we define a mapping function $\Phi(\mathbf{x}) : R^n \rightarrow R^d$, $d > n$ and d can be taken as ∞ . Hence, the primal soft-margin SVM can be shown as:

$$\begin{aligned} \min_{\tilde{\mathbf{w}}, b, \xi} \quad & \frac{1}{2} \tilde{\mathbf{w}}^T \tilde{\mathbf{w}} + C \sum_{i=1}^m \xi_i \\ \text{s.t.} \quad & y_i (\tilde{\mathbf{w}}^T \Phi(\mathbf{x}_i) + b) \geq 1 - \xi_i \\ & \xi_i \geq 0 \\ & i = 1, \dots, m \end{aligned}$$

where $\tilde{\mathbf{w}} \in R^d$ and $\Phi(\mathbf{x}) = [\phi_1(\mathbf{x}), \dots, \phi_d(\mathbf{x})]^T$. However, when dimension of $\Phi(\mathbf{x})$ is very high even infinite, making quite hard to get optimal $\tilde{\mathbf{w}}$ because of dimensional curse. So we need another method to guarantee our learnable parameter's dimension is not quite high. In order to solve the problem above, we define the Lagrangian:

$$L(\tilde{\mathbf{w}}, b, \xi, \alpha, \lambda) = \frac{1}{2} \tilde{\mathbf{w}}^T \tilde{\mathbf{w}} + C \sum_{i=1}^m \xi_i + \sum_{i=1}^m \alpha_i (1 - y_i (\tilde{\mathbf{w}}^T \Phi(\mathbf{x}_i) + b) - \xi_i) + \sum_{i=1}^m \lambda_i (-\xi_i)$$

And we can get KKT conditions for primal soft-margin:

$$\left\{ \begin{array}{l} y_i (\tilde{\mathbf{w}}^T \Phi(\mathbf{x}_i) + b) \geq 1 - \xi_i \\ \xi_i \geq 0 \end{array} \right\} \text{primal feasible condition}$$

$$\left\{ \begin{array}{l} \alpha_i \geq 0 \\ \lambda_i \geq 0 \end{array} \right\} \text{dual feasible condition}$$

$$\left\{ \begin{array}{l} \alpha_i (1 - y_i (\tilde{\mathbf{w}}^T \Phi(\mathbf{x}_i) + b) - \xi_i) = 0 \\ \lambda_i \xi_i = 0 \end{array} \right\} \text{complementary slackness}$$

$$\left\{ \begin{array}{l} \frac{\partial L}{\partial \tilde{\mathbf{w}}} = \tilde{\mathbf{w}} - \sum_{i=1}^m \alpha_i y_i \Phi(\mathbf{x}_i) = 0 \\ \frac{\partial L}{\partial b} = \sum_{i=1}^m \alpha_i y_i = 0 \\ \frac{\partial L}{\partial \xi_i} = C - \alpha_i - \lambda_i = 0 \end{array} \right\} \text{dual - inner optimal condition}$$

$$i = 1, \dots, m$$

The KKT conditions are sufficient and necessary for this optimisation problem because primal soft-margin SVM is a quadratic programming problem essentially. Also, we know that the strong duality of the optimality holds when original problem is convex with feasible solutions, whose constraints are linear. Consequently, the dual format soft-margin SVM is deduced by strong duality as below

$$\min_{\tilde{\mathbf{w}}, b, \xi} \left(\max_{\alpha, \lambda \geq 0} L(\tilde{\mathbf{w}}, b, \xi, \alpha, \lambda) \right) = \max_{\alpha, \lambda \geq 0} \left(\min_{\tilde{\mathbf{w}}, b, \xi} L(\tilde{\mathbf{w}}, b, \xi, \alpha, \lambda) \right)$$

Apply dual-inner optimal condition into above format which is equivalent to primal soft-margin SVM, then we get

$$\max_{\alpha, \lambda \geq 0} \left(\min_{\tilde{\mathbf{w}}, b, \xi} L(\tilde{\mathbf{w}}, b, \xi, \alpha, \lambda) \right) = \max_{\alpha, \lambda \geq 0} \left(L(\tilde{\mathbf{w}}, b, \xi, \alpha, \lambda) |_{\tilde{\mathbf{w}} = \sum_{i=1}^m \alpha_i y_i \Phi(\mathbf{x}_i), \sum_{i=1}^m \alpha_i y_i = 0, C = \alpha_i + \lambda_i} \right)$$

which is equivalent to the optimality:

$$\begin{aligned} \min_{\alpha} \quad & \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i y_i \Phi(\mathbf{x}_i)^T \Phi(\mathbf{x}_j) \alpha_j y_j - \sum_{i=1}^m \alpha_i \\ \text{s.t.} \quad & \alpha_i \leq 0 \\ & \alpha_i - C \leq 0 \\ & \sum_{i=1}^m \alpha_i y_i = 0 \\ & i = 1, \dots, m \end{aligned}$$

Finally, we get the dual soft-margin format SVM. At the same time, we can calculate the kernel matrix $[K_{ij}]_{m \times m} = [\Phi(\mathbf{x}_i)^T \Phi(\mathbf{x}_j)]_{m \times m}$ previously to get rid of repeated computing on inner product of two high dimensional vectors. Hence, the kernel SVM has been deduced successfully, and can be shown as:

$$\begin{aligned} \min_{\alpha} f(\alpha) = \quad & \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i y_i K_{ij} \alpha_j y_j - \sum_{i=1}^m \alpha_i \\ \text{s.t.} \quad & \alpha_i \leq 0 \\ & \alpha_i - C \leq 0 \\ & \sum_{i=1}^m \alpha_i y_i = 0 \\ & i = 1, \dots, m \end{aligned}$$

Similarly, we can also construct Lagrangian for kernel format above:

$$L(\alpha, \lambda, \mu, b) = \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i y_i K_{ij} \alpha_j y_j - \sum_{i=1}^m \alpha_i - \sum_{i=1}^m \lambda_i \alpha_i + \sum_{i=1}^m \mu_i (\alpha_i - C) + b \sum_{i=1}^m y_i \alpha_i$$

Hence, we can get the corresponding KKT conditions:

$$\left\{ \begin{array}{l} 0 \leq \alpha_i \leq C, i = 1 \dots m \\ \sum_{i=1}^m y_i \alpha_i = 0 \\ \lambda_i \geq 0, i = 1 \dots m \\ \mu_i \geq 0, i = 1 \dots m \end{array} \right\} \text{dual feasible condition}$$

$$\left\{ \begin{array}{l} \lambda_i \alpha_i = 0, i = 1 \dots m \\ \mu_i (C - \alpha_i) = 0, i = 1 \dots m \end{array} \right\} \text{complementary slackness}$$

$$\frac{\partial L}{\partial \alpha_i} = \nabla f(\alpha)_i - \lambda_i - \mu_i + b y_i = 0, i = 1 \dots m$$

Actually, KKT conditions can serve as stopping condition when solving this problem because they are sufficient and necessary conditions for optimal point in this problem, which will be discussed more afterwards.

2.2 Numerical Optimisation Method for SVM

Kernel SVM optimality problem can be represented as a quadratic optimization problem with linearly and bound constraints:

$$\begin{aligned} \min_{\alpha} f(\alpha) = \quad & \frac{1}{2} \alpha^T Q \alpha - \alpha^T e \\ \text{s.t.} \quad & \mathbf{y}^T \alpha = 0 \\ & 0 < \alpha_i < C, i = 1 \dots m \end{aligned}$$

where Hessian matrix $[Q_{ij}]_{m \times m} = [y_i K_{ij} y_j]_{m \times m}$, which is a positive semi-definite dense matrix; $\alpha = [\alpha_1, \dots, \alpha_m]^T$ and each component of it is a Lagrangian multiplier for a primal feasible condition with a range from 0 to C (penalty parameter for slack variables). We can find that the dimension of Hessian matrix is highly depended on the data number m of the data set. So some traditional optimisation algorithms like Newton or Quasi Newton methods are not suitable here attributed to heavy intensive computation based on Q and overload memory to store them [A two level]. Here we consider a decomposition problem—Sequential Minimal Optimisation (SMO) to solve this optimality problem.

2.2.1 Sequential Minimal Optimisation Procedure

SMO algorithm is a decomposition method. Instead of solving a quadratic programming problem depended on $m \times m$ fully dense matrix Q and optimise m dimension. In each iteration, SMO only focus on optimising two parameters α_i, α_j , where $i, j \in B, B \subseteq \{1, 2, \dots, m\}, |B| = 2$. B is called as working set and keeps updating with iteration until no suitable elements can be selected. The remaining of $\{1, 2, \dots, m\}$ is denoted as $N = \{1, 2, \dots, m\} \setminus B$. Hence, the original densely quadratic programming problem can be reduced as a sub-problem in each iteration:

$$\begin{aligned} \min_{\alpha_i, \alpha_j} W(\alpha_i, \alpha_j) &= \frac{1}{2}Q_{ii}\alpha_i^2 + \frac{1}{2}Q_{jj}\alpha_j^2 + Q_{ij}\alpha_i\alpha_j - (\alpha_i + \alpha_j) + (\boldsymbol{\alpha}_N^k)^T Q_{NB} \begin{bmatrix} \alpha_i \\ \alpha_j \end{bmatrix} \\ \text{s.t.} \quad &0 < \alpha_i < C \\ &0 < \alpha_j < C \\ &y_i\alpha_i + y_j\alpha_j = -\mathbf{y}_N^T \boldsymbol{\alpha}_N^k \end{aligned}$$

where α_i and α_j are selected based on the elements of working set B , Q_{ij} is the (i, j) element in Q matrix while Q_{NB} is the sub-matrix of Q based on the component in set N and set B . $\boldsymbol{\alpha}_N^k$ is the corresponding parameter vector of set N in the last iteration and it serves as a constant vector in current optimality problem. Consequently, the sub-problem is a constrained quadratic programming problem with only two parameter and it can be analytically solved [Platt, 1998], greatly saving the memory and time for computation in each iteration. The algorithm process of Generalised SMO algorithm is as below:

Generalised SMO algorithm

Input: training set $D = \{(\mathbf{x}_p, y_p)\}_{p=1}^m, \mathbf{x}_p \in R^n, y_p \in \{-1, 1\}$

Initialization: $\boldsymbol{\alpha}^0 = [\alpha_1^0, \alpha_2^0, \dots, \alpha_m^0]^T = [0, 0, \dots, 0]^T$

Repeat:

Select working set: Select working set $\{i^k, j^k\} \subseteq \{1, 2, \dots, m\}$, and the corresponding $\alpha_{i^k}^{k-1}, \alpha_{j^k}^{k-1}$ to be updated.

Update: By solving sub-problem analytically, we can get the optimal $\boldsymbol{\alpha}^k$, where

$$\begin{aligned} \alpha_{i^k}^k &= \text{update}(\alpha_{i^k}^{k-1}), \alpha_{j^k}^k = \text{update}(\alpha_{j^k}^{k-1}) \\ \alpha_p^k &= \alpha_p^{k-1} \text{ when } p \in \{1, 2, \dots, m\} \setminus \{i^k, j^k\} \end{aligned}$$

Until: Stopping conditions(KKT conditions) are satisfied with tolerance ε or when $\boldsymbol{\alpha}'s$ change is quite small.

In fact, selection of working set highly affects the convergence of the algorithm, so we mainly explore and compare two selecting strategies shown below.

2.2.2 Working Set Selection

1. SMO-Simple Algorithm

SMO-Simple algorithm is dedicated to employ a simpler heuristics strategy to choose element α_i and α_j by iterating over all α_i . At the same time, we select the corresponding α_j among the $m - 1$ elements in $\boldsymbol{\alpha}$ except α_i . We stop iteration until current α_i satisfies the KKT condition with a ε tolerance. However, this algorithm cannot guarantee the global optimal because we do not iterate over all the pair of α_i and α_j so that some more optimised working set is missed [2].

2. SVM-Light Algorithm Actually, SVM-Light [1] is a popular decomposition methods which can select working set more systematically. When the size of working set equals to two, it is equivalent to a SMO modification format proposed by Keerthi et al. [4]. Working set of this algorithm is selected by KKT violation. The KKT conditions for Kernel format SVM can be rewritten as:

$$\begin{cases} \nabla f(\boldsymbol{\alpha})_i + by_i \geq 0, \text{ if } \alpha_i < C \\ \nabla f(\boldsymbol{\alpha})_i + by_i \leq 0, \text{ if } \alpha_i > 0 \end{cases} \quad (KKT \ 1)$$

According to $y_i = \pm 1$, we can rewrite the (KKT 1) above as:

$$\begin{cases} \nabla f(\boldsymbol{\alpha})_i + b \geq 0, \text{ if } \alpha_i < C, y_i = 1 \\ \nabla f(\boldsymbol{\alpha})_i - b \geq 0, \text{ if } \alpha_i < C, y_i = -1 \\ \nabla f(\boldsymbol{\alpha})_i + b \leq 0, \text{ if } \alpha_i > 0, y_i = 1 \\ \nabla f(\boldsymbol{\alpha})_i - b \leq 0, \text{ if } \alpha_i > 0, y_i = -1 \end{cases} \quad (KKT \ 2)$$

From (KKT 2), we can easily deduce the range for b:

$$\begin{cases} \max_{i \in I_{up}(\boldsymbol{\alpha})} -y_i \nabla f(\boldsymbol{\alpha})_i \leq b \leq \min_{j \in I_{low}(\boldsymbol{\alpha})} -y_j \nabla f(\boldsymbol{\alpha})_j \\ I_{up}(\boldsymbol{\alpha}) \equiv \{i | \alpha_i < C, y_i = 1 \text{ or } \alpha_i > 0, y_i = -1\} \\ I_{low}(\boldsymbol{\alpha}) \equiv \{i | \alpha_i < C, y_i = -1 \text{ or } \alpha_i > 0, y_i = 1\} \end{cases}$$

Therefore, we conclude that α is optimal if and only if

$$\max_{i \in I_{up}(\alpha)} -y_i \nabla f(\alpha)_i \leq \min_{j \in I_{low}(\alpha)} -y_j \nabla f(\alpha)_j$$

Consequently, α is not optimal when $\exists i \in I_{up}(\alpha), j \in I_{low}(\alpha), -y_i \nabla f(\alpha)_i > -y_j \nabla f(\alpha)_j$. Such (i, j) is called as a "violating pair" [4]. Hence, α is optimal when no "violating pairs" are found. To guarantee the asymptotic convergence, we select the working set containing the maximal violating pair in each iteration. i.e the indices of two elements in the working set should be $i \in \arg \max_{t \in I_{up}(\alpha)} -y_t \nabla f(\alpha)_t$ and $j \in \arg \min_{t \in I_{low}(\alpha)} -y_t \nabla f(\alpha)_t$. And here we will give a brief interpretation about why the maximal violating pair can ensure the global asymptotic convergence.

Theorem (Asymptotic Convergence of SVM-Light Algorithm) If $\{\alpha_k\}$ is the sequence generated by the decomposition method (like SVM-Light based SMO algorithm), any convergent subsequence converges to the optimal point α^* of kernel SVM format optimality [5].

Proof [5] Solving the sub-problem quadratic programming problem for kernel SVM based on SVM-Light essentially means to find the minimisation in a linear equality constraint $\alpha_i y_i + \alpha_j y_j = -\mathbf{y}_N^T \alpha_N^k := \zeta$ within a rectangle area due to $0 \leq \alpha_i, \alpha_j \leq C$. We use $\alpha(t) = [\alpha_1(t), \alpha_2(t), \dots, \alpha_m(t)]^T$ to denote the change of α in the linear equality constraint where

$$\alpha_h(t) = \begin{cases} \alpha_h^k + \frac{t}{y_h}, & h = i \\ \alpha_h^k - \frac{t}{y_h}, & h = j \\ \alpha_h^k, & \text{Otherwise} \end{cases}$$

Hence, the previous sub-problem with parameter α_i and α_j can be converted to a optimality problem with only parameter t :

$$\begin{aligned} \min_t \quad & \varphi(t) := f(\alpha(t)) \\ \text{s.t.} \quad & 0 \leq \alpha_i(t) \leq C \\ & 0 \leq \alpha_j(t) \leq C \end{aligned}$$

We assume \bar{t} is the optimal point of the problem and $\alpha^{k+1} = \alpha(\bar{t})$, $\bar{t} = \frac{\|\alpha^{k+1} - \alpha^k\|}{\sqrt{2}}$. We use the second-order Taylor expansion to approximate the $\varphi(t)$: $\varphi(t) = \varphi(0) + \varphi'(0)t + \varphi''(0)\frac{t^2}{2}$. Because

$$\begin{aligned} \varphi'(t) &= \sum_{h=1}^m \nabla f(\alpha(t))_h \alpha_h'(t) \\ &= y_i \nabla f(\alpha(t))_i - y_j \nabla f(\alpha(t))_j \\ &= y_i \left(\sum_{h=1}^m Q_{ih} \alpha_h(t) - 1 \right) - y_j \left(\sum_{h=1}^m Q_{jh} \alpha_h(t) - 1 \right) \\ \varphi''(t) &= Q_{ii} + Q_{jj} - 2y_i y_j Q_{ij} \end{aligned}$$

Therefore, when $t = 0$,

$$\begin{aligned} \varphi'(0) &= y_i \nabla f(\alpha^k)_i - y_j \nabla f(\alpha^k)_j \\ \varphi''(0) &= K_{ii} + K_{jj} - 2K_{ij} \\ &= \Phi(\mathbf{x}_i)^T \Phi(\mathbf{x}_i) + \Phi(\mathbf{x}_j)^T \Phi(\mathbf{x}_j) - 2\Phi(\mathbf{x}_i)^T \Phi(\mathbf{x}_j) \\ &= \|\Phi(\mathbf{x}_i) - \Phi(\mathbf{x}_j)\|^2 \end{aligned}$$

So, considering Q is S.P.D and $\varphi''(0) \geq 0$, we talk about the following cases.

Case 1 $\varphi''(0) = 0$, i.e $\Phi(\mathbf{x}_i) = \Phi(\mathbf{x}_j)$, hence

$$\begin{aligned} \varphi'(0) &= y_i \sum_{h=1}^m Q_{ih} \alpha_h^k - y_j \sum_{h=1}^m Q_{jh} \alpha_h^k \\ &= y_i \left(\sum_{h=1}^m y_i y_h \Phi(\mathbf{x}_i)^T \Phi(\mathbf{x}_h) \alpha_h^k - 1 \right) - y_j \left(\sum_{h=1}^m y_j y_h \Phi(\mathbf{x}_j)^T \Phi(\mathbf{x}_h) \alpha_h^k - 1 \right) \\ &= y_j - y_i \end{aligned}$$

Potentially, the descent is guaranteed when

$$\varphi'(0) = y_i \nabla f(\alpha^k)_i - y_j \nabla f(\alpha^k)_j \neq 0$$

In this circumstance, $y_i \neq y_j$ and $|\varphi'(0)| = 2$. According to $\varphi''(0) = \varphi''(t)|_{t=0} = 0$, it is shown that $\varphi'(t)$ is a linear function, so $\varphi(\bar{t}) \leq \varphi(0)$ with $|\bar{t}| \leq C$. Then we get

$$\varphi(\bar{t}) - \varphi(0) = -|\varphi'(0)\bar{t}| \leq -\frac{2}{C}\bar{t}^2 = -\frac{\|\alpha^{k+1} - \alpha^k\|^2}{C} (*1)$$

Case 2 $\varphi''(0) > 0$. Assume $\bar{t} = \gamma t^*$ where $0 < \gamma < 1$ and $t^* = -\frac{\varphi'(0)}{\varphi''(0)}$ is the unconstrained minimum of φ . So we can get that

$$\varphi(\bar{t}) - \varphi(0) = \gamma \frac{-\varphi'(0)^2}{\varphi''(0)} + \frac{\gamma^2 \varphi'(0)^2}{2 \varphi''(0)} \leq -\frac{\gamma^2 \varphi'(0)^2}{2 \varphi''(0)} = -\frac{\varphi''(0)}{2} \bar{t}^2 = -\frac{\varphi''(0)}{4} \|\alpha^{k+1} - \alpha^k\|^2 C \quad (*)$$

In summary, note that $\varphi(0) = f(\alpha^k)$, $\varphi(\bar{t}) = f(\alpha^{k+1})$ and with use of $\bar{t} = \frac{\|\alpha^{k+1} - \alpha^k\|}{\sqrt{2}}$, $\varphi''(t) = Q_{ii} + Q_{jj} - 2y_i y_j Q_{ij}$ and $(*)$, we can get $\sigma \equiv \min \left\{ \frac{2}{C}, \min_{i,j} \left\{ \frac{Q_{ii} + Q_{jj} - 2y_i y_j Q_{ij}}{2} : Q_{ii} + Q_{jj} - 2y_i y_j Q_{ij} > 0 \right\} \right\}$ to make sure $f(\alpha^{k+1}) - f(\alpha^k) \leq -\frac{\sigma}{2} \|\alpha^{k+1} - \alpha^k\|^2$. Hence, the asymptotic convergence of SVM-Light is proved.

Q.E.D

2.2.3 Updating Parameters

After selecting the proper working set $B = \{i, j\}$, we solve the sub-problem mentioned in section 1.2.1 to update α in each iteration, which means the aim is to find the required minimisation taking place in the rectangle area $S = [0, C] \times [0, C]$ along a line $\alpha_i y_i + \alpha_j y_j = -\mathbf{y}_N \alpha_N^k := \zeta$, where ζ is a constant in current iteration $k+1$ and $y_i, y_j = \pm 1$. Hence, α_i can be represented as $\alpha_i = (\zeta - \alpha_j y_j) y_i$ and the sub-problem can be re-written as:

$$\begin{aligned} \min_{\alpha_i} \psi(\alpha_i) &= \frac{1}{2} Q_{ii} \alpha_i^2 + \frac{1}{2} Q_{jj} (\zeta - \alpha_i y_i)^2 + y_j Q_{ij} \alpha_i (\zeta - \alpha_i y_i) - (\zeta - \alpha_i y_i) y_j - \alpha_i + \left(\alpha_N^k \right)^T Q_{NB} \begin{bmatrix} \alpha_i \\ (\zeta - \alpha_i y_i) y_j \end{bmatrix} \\ \text{s.t.} \quad & L \leq \alpha_i \leq H \\ & \begin{cases} L = \max(0, \alpha_i^k - \alpha_j^k); H = \min(C, C + \alpha_i^k - \alpha_j^k), y_i \neq y_j \\ L = \max(0, \alpha_i^k + \alpha_j^k - C); H = \min(C, \alpha_i^k + \alpha_j^k), y_i = y_j \end{cases} \end{aligned}$$

By setting $\frac{\partial \psi}{\partial \alpha_i} = 0$, we can deduce that

$$\alpha_i^{k+1, unc} = \alpha_i^k + \frac{y_i (E_i - E_j)}{\eta}$$

where $E_t = \left(\sum_{j=1}^N \alpha_j y_j K(\mathbf{x}_j, \mathbf{x}_t) + b \right) - y_t = \left(\frac{[Q\alpha]_t}{y_t} + b \right) - y_t$, $t = 1, 2$ implying the distance between prediction and truth, and $\eta = K_{ii} + K_{jj} - 2K_{ij} = \frac{Q_{ii}}{y_i^2} + \frac{Q_{jj}}{y_j^2} - 2\frac{Q_{ij}}{y_i y_j}$. Then, considering the constraint $L \leq \alpha_i \leq H$, we clip the result $\alpha_i^{k+1, unc}$:

$$\alpha_i^{k+1} = \begin{cases} H, & \alpha_i^{k+1, unc} > H \\ \alpha_i^{k+1, unc}, & L \leq \alpha_i^{k+1, unc} \leq H \\ L, & \alpha_i^{k+1, unc} < L \end{cases}$$

Then, $\alpha_j^{k+1} = \alpha_j^k + y_i y_j (\alpha_i^k - \alpha_i^{k+1})$ can be updated. At last, we upgrade b :

$$b^{k+1} = \begin{cases} b_i^{k+1} & , 0 < \alpha_i, \alpha_j < C \\ \frac{b_i^{k+1} + b_j^{k+1}}{2} & , otherwise \end{cases}$$

where

$$\begin{aligned} b_i^{k+1} &= -E_i - y_i Q_{ii} (\alpha_i^{k+1} - \alpha_i^k) - y_j \frac{Q_{ji}}{y_j y_i} (\alpha_j^{k+1} - \alpha_j^k) + b^k \\ b_j^{k+1} &= -E_j - y_i \frac{Q_{ij}}{y_i y_j} (\alpha_i^{k+1} - \alpha_i^k) - y_j Q_{jj} (\alpha_j^{k+1} - \alpha_j^k) + b^k \end{aligned}$$

Therefore, the updating process is finished. The algorithm keeps iterating until KKT conditions are satisfied with tolerance ε .

3 Experiment and Analysis

In this section, we implement proper SVM format based on SMO algorithm with two working set selection tactics and then test on spiral data set. We select one of the strategy with a better performance in convergence and then train the algorithm with different regularisation parameter C and kernel parameter μ to benchmark convergence performance influenced by parameters. Also, we talk about the performance of the selected method in terms of complexity, CPU time and memory used. At last, we use accuracy as the metric to verify the validity of the implemented SVM.

In this part, we generate 600 sample points in total and split half of the spiral data set as training set (size 300) and set left samples as test set (size 300). It is obvious that the spiral data set is a non-linearly separable data set and we decide to use the kernel format SVM with RBF kernel for this task. RBF kernel $K_\mu(\mathbf{p}, \mathbf{q}) = e^{-\mu \|\mathbf{p} - \mathbf{q}\|^2}$ possesses the great advantage for mapping the data samples into the infinity dimensional space, showing that RBF kernel is with a strong and flexible mapping ability than other kernels.

3.1 Analysis on Working Set Selection Strategy

We compare the convergence situation when we utilise different working set selection strategies: SMO-Simple algorithm and SVM-Light algorithm. In this experiment, we set the regularisation parameter C as 0.5 , kernel parameter μ as 10, tolerance ε as 0.001 and maximum iteration number as 1000. Then we plot the convergence plot for α based on the posterior α^* .

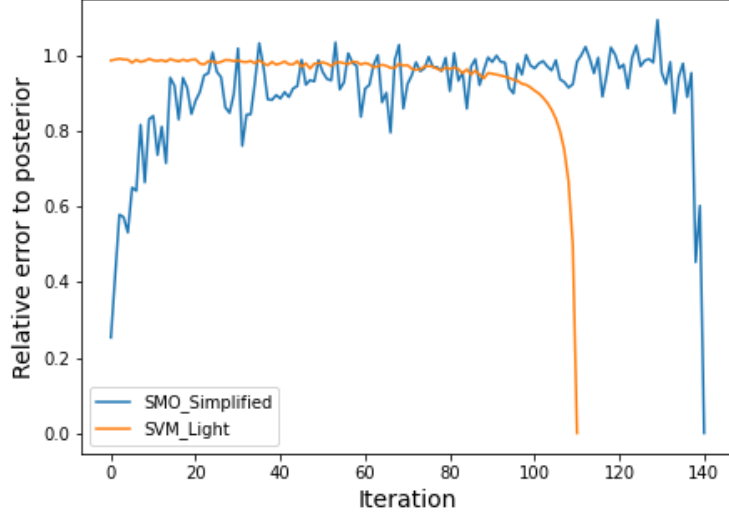


Figure 2: Convergence plot based on different working set selection strategies

The relative error to α^* is measured based on Q matrix, which is $\frac{(\alpha^{k+1} - \alpha^*)^T Q (\alpha^{k+1} - \alpha^*)}{(\alpha^k - \alpha^*)^T Q (\alpha^k - \alpha^*)}$. From the plot above, we can easily find that SVM-Light algorithm possesses a faster convergence, which converges after 112 iterations compared with converging after 143 iterations of SMO-Simplified algorithm. Also, the converging process of SMO-Simplified highly fluctuates because of the randomness with choosing j element in working set in every iteration, which is avoided in SVM-Light algorithm. Therefore, SVM-Light's working set selection strategy is more robust so the experiments afterwards are mainly based on SMO-Light algorithm.

3.2 Analysis on Local Convergence Property

We have known that for SVM-Light algorithm, the asymptotic convergence can be assured. Now we talk about the local convergence rate of the algorithm. In this part, we still use the parameters set in section 3.1.

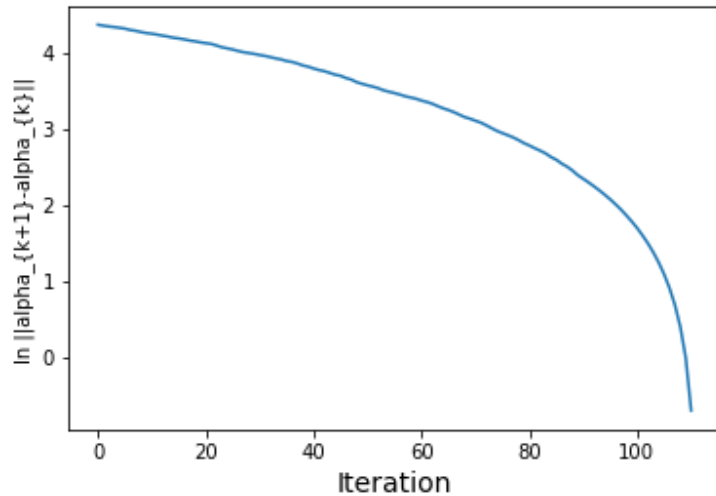


Figure 3: Exponential Convergence plot for SMO-Light

From Figure 2, we find that $\frac{(\alpha^{k+1}-\alpha^*)^T Q(\alpha^{k+1}-\alpha^*)}{(\alpha^k-\alpha^*)^T Q(\alpha^k-\alpha^*)} < 1$ as $k \rightarrow \infty$ and from Figure 3, we find that when $k \rightarrow \infty$, the exponential convergence plot approximates to a straight line. Hence, we empirically estimate this algorithm possesses linear convergence. Theoretically, the linear convergence can be guaranteed based on two assumptions [4]:

Assumption 1: The Kernel Matrix K is positive definite.

This assumption ensure that Q matrix is positive definite as well and the optimality problem becomes a strictly convex programming problem with a unique global optimum α^* .

Assumption 2: Nondegeneracy For the optimal solution α^* , $\nabla f(\alpha^*)_i + b^* y_i \neq 0$ holds when $\alpha_i^* = 0, C$.

b^* denotes the value of $\max_{i \in I_{up}(\alpha^*)} -y_i \nabla f(\alpha^*)_i \leq b \leq \min_{j \in I_{low}(\alpha^*)} -y_j \nabla f(\alpha^*)_j$, if the algorithm takes infinite iterations. In fact, only elements whose $-y_i \nabla f(\alpha^*)_i$ are b^* can still form violating pairs. This assumption guarantees that when k is larger enough, all bounded variables are constants, so at the optimal point $\alpha^* = \arg \min_{\alpha} f(\alpha) = \frac{1}{2} \alpha^T Q \alpha - \alpha^T e$ with constraint $y^T \alpha^* = \Delta$ where Δ is a constant. By using Lagrange multiplier b^* to solve the optimal point, we can get a linear system like:

$$\begin{bmatrix} Q & y \\ y^T & 0 \end{bmatrix} \begin{bmatrix} \alpha^* \\ b^* \end{bmatrix} = \begin{bmatrix} e \\ \Delta \end{bmatrix} \quad (LS1)$$

In order to guarantee the decent direction, we solve the following problem in each iteration

$$\min_d \frac{1}{2} d^T Q_{BB} d + \nabla f(\alpha^k)_B^T d \text{ s.t. } y_B^T d = 0$$

Similarity, we can get

$$\begin{bmatrix} Q_{BB} & y_B \\ y_B^T & 0 \end{bmatrix} \begin{bmatrix} d^* \\ b^* \end{bmatrix} = \begin{bmatrix} -\nabla f(\alpha^k)_B \\ 0 \end{bmatrix} \quad (LS2)$$

With (LS1), we get $Q(\alpha^k - \alpha^*) = Q\alpha^k - e + b^* y = \nabla f(\alpha^k) + b^* y$, as k is large enough.

Hence,

$$-YQ(\alpha^k - \alpha^*) = -Y\nabla f(\alpha^k) - b^* e$$

where $Y \equiv \text{diag}(y)$. Therefore, a maximal violating pair is obtained by the following formulation, since b^* does not affect the order of sequence:

$$\arg \max_i \left(-y_i Q(\alpha^k - \alpha^*)_i \right) = \arg \max_i \left(-y_i \nabla f(\alpha^k)_i \right) \arg \min_i \left(-y_i Q(\alpha^k - \alpha^*)_i \right) = \arg \min_i \left(-y_i \nabla f(\alpha^k)_i \right)$$

In order to deduce the local convergence rate, we also have to introduce two lemmas, which has been proved by Lin [4].

Lemma 1

If $v_1 \geq \dots \geq v_l$, then $\sum_{i=1}^l v_i^2 \geq \frac{(v_1 - v_l)^2}{2}$

Lemma 2 If Q is invertible, then for any x such that (1) $e^T x = 0$; (2) $v \equiv Qx$, $v \equiv Qx \max_i ((Qx)_i) = v^1 >$

$v^l = \min_i ((Qx)_i)$, and $v^1 v^l \geq 0$, we have

$$\min(|v^1|, |v^l|) \leq \left(1 - \frac{e^T Q^{-1} e}{\sum_{i,j} |Q_{ij}^{-1}|} \right) \max(|v^1|, |v^l|)$$

On the basis of the assumptions and lemmas above, we talk about the difference between α^{k+1} and α^k :

$$\begin{aligned} & (\alpha^{k+1} - \alpha^*)^T Q(\alpha^{k+1} - \alpha^*) - (\alpha^k - \alpha^*)^T Q(\alpha^k - \alpha^*) \quad (\text{diff}) \\ &= 2(d^k)^T (Q(\alpha^k - \alpha^*))_B + (d^k)^T Q_{BB} d^k \\ &= (d^k)^T (2(Q(\alpha^k - \alpha^*))_B - \nabla f(\alpha^k)_B - b^k y_B) \\ &= (d^k)^T ((Q(\alpha^k - \alpha^*))_B + (b^* - b^k) y_B) \\ &= (d^k)^T ((Q(\alpha^k - \alpha^*))_B + (b^k - b^*) y_B) \\ &= - \left[(Q(\alpha^k - \alpha^*))_B + (b^* - b^k) y_B \right]^T Q_{BB}^{-1} \left[- (Q(\alpha^k - \alpha^*))_B + (b^* - b^k) y_B \right] \end{aligned}$$

we define that

$$\hat{Q} \equiv Y_B Q_{BB}^{-1} Y_B \quad \text{and} \quad v \equiv -Y(Q(\alpha^k - \alpha^*))$$

where $Y_B \equiv \text{diag}(y_B)$, then $v_B \equiv -Y_B(Q(\alpha^k - \alpha^*))_B$

So, the (diff) formulation can be rewritten as

$$- \left[v_B + (b^* - b^k) e_B \right]^T \hat{Q} \left[v_B + (b^* - b^k) e_B \right]$$

We define

$$v^1 \equiv \max_i (v_i) = \max_{i \in B} (v_i), v^l \equiv \min_i (v_i) = \min_{i \in B} (v_i)$$

and denote $\min(\text{eig}(\cdot))$ and $\max(\text{eig}(\cdot))$ to be the minimal and maximal eigenvalues of a matrix respectively. Then,

$$\begin{aligned} & - \left[v_B + (b^* - b^k) e_B \right]^T \hat{Q} \left[v_B + (b^* - b^k) e_B \right] \\ & \geq \min \left(\text{eig}(\hat{Q}) \right) \left[v_B + (b^* - b^k) e_B \right]^T \left[v_B + (b^* - b^k) e_B \right] \\ & \stackrel{\text{lemma1}}{\geq} \min \left(\text{eig}(\hat{Q}) \right) \frac{(v_1 - v_l)^2}{2} \\ & \stackrel{\text{lemma2}}{\geq} \frac{\min \left(\text{eig}(\hat{Q}) \right)}{2} \left(\frac{y^T Q^{-1} y}{\sum_{i,j} |Q_{ij}^{-1}|} \right)^2 \max(|v^1|, |v^l|)^2 \quad (IEQ1) \\ & \geq \frac{\min \left(\text{eig}(\hat{Q}) \right)}{2l} \left(\frac{y^T Q^{-1} y}{\sum_{i,j} |Q_{ij}^{-1}|} \right)^2 (Q(\alpha^k - \alpha^*))^T Q(\alpha^k - \alpha^*) \\ & \geq \frac{\min \left(\text{eig}(\hat{Q}) \right)}{2l \max(\text{eig}(Q^{-1}))} \left(\frac{y^T Q^{-1} y}{\sum_{i,j} |Q_{ij}^{-1}|} \right)^2 (Q(\alpha^k - \alpha^*))^T Q^{-1} Q(\alpha^k - \alpha^*) \\ & \geq \frac{\min \left(\text{eig}(\hat{Q}) \right)}{2l \max(\text{eig}(Q^{-1}))} \left(\frac{y^T Q^{-1} y}{\sum_{i,j} |Q_{ij}^{-1}|} \right)^2 (\alpha^k - \alpha^*)^T Q(\alpha^k - \alpha^*) \end{aligned}$$

Therefore, a constant $c \equiv 1 - \min_B \left(\frac{\min(\text{eig}(\hat{Q}))}{2l \max(\text{eig}(Q^{-1}))} \left(\frac{y^T Q^{-1} y}{\sum_{i,j} |Q_{ij}^{-1}|} \right)^2 \right)$, $0 < c < 1$ exists to make sure when k is large,

$$(\alpha^{k+1} - \alpha^*)^T Q(\alpha^{k+1} - \alpha^*) \leq c (\alpha^k - \alpha^*)^T Q(\alpha^k - \alpha^*)$$

holds, which show that the algorithm is with linear convergence.

3.3 Analysis on Parameters

In this part, we perform experiments on different regularisation parameter $C = 0.001, 0.1, 0.5, 5, 10$ and kernel parameter $\mu = 1, 5, 10, 15$ and then we compare the convergence and performance on them.

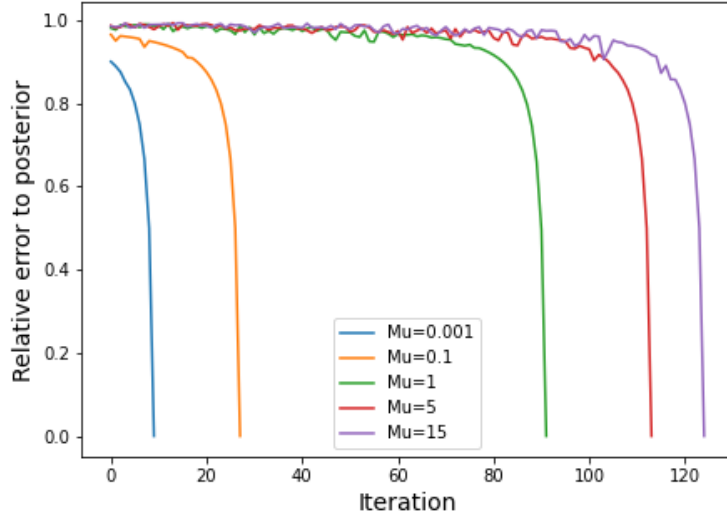


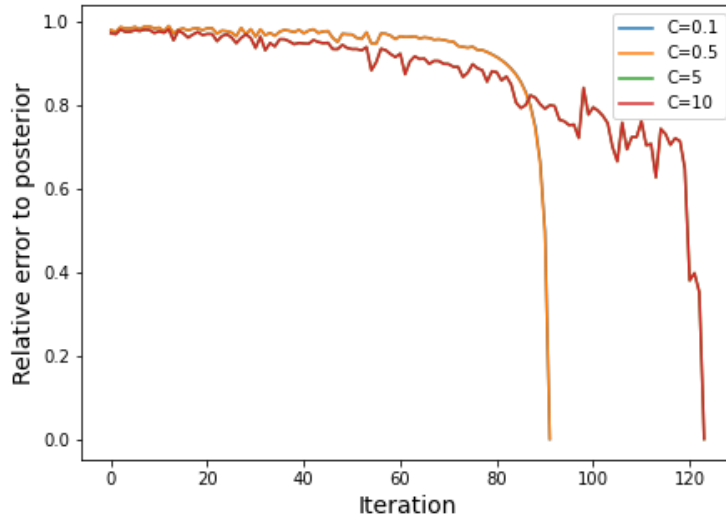
Figure 4: Convergence plot for multiple μ

From Figure 4, we can find that iteration speed to stationary point is slower with μ increasing, this is because the larger μ brings much more computation complexity when we create a kernel matrix with large parameter. However, from table 1 we find that when $\mu = 0.001$, the training and test error are around 0.5, which means the

Table 1: Training and test error for multiple μ

μ	Train error(%)	Test error(%)
0.001	0.5100	0.5267
0.1	0.0000	0.0000
1	0.0000	0.0000
5	0.0000	0.0000
15	0.0000	0.0067

SVM classification does not work here. This is attributes to the Kernel Matrix K is almost not S.P.D under this case, not ensuring the optimality problem has the unique global minimiser. When $\mu = 15$, test error begins to increase, hence the overfitting occurs because of too complicated mapping brought by RBF kernel with a high parameter.

Figure 5: Convergence plot for multiple C

Form the graph above, we find that different parameters has the same iteration situation: when $C=0.1$ or 0.5 , iteration stops at 92th iteration; when $C=5$ or 10 , iteration stops at 134th iteration. This is because parameter C just serves as a range bound when we select maximal violating pairs. So the situation where the feasible α are same might occur when we choose different C , leading to the same iteration situation.

3.4 Analysis on Complexity

3.4.1 Time Complexity

When we implement the SMO-Light algorithm, we mainly use two for-loops to find the maximal violating pair. In the worst case, we have to travel all m points of data set in each loop. Hence the time complexity for SVM-Light is $O(m^2)$ [1].

3.4.2 Space Complexity

In SMO-Light algorithm, the working set's size q is 2. So the highest memory requirements are depends on how many rows of Q are stored for each iteration. In this case, $O(qm)$ floating numbers are needed to stored for updating parameter, computing gradients, solving sub-problems phases. Also, Q_{BB} needs $O(q^2)$ space and α^k needs $O(m)$ space. So the total space complexity per iteration should be $O(qm + q^2 + m)$ [1].

3.5 Analysis on Classification Performance

In order to verify SVM-Light's validity, we select $\mu = 1$ and $C = 0.5$ as the parameters, which are the parameter pair with the best convergence performance. And then we measure the training and test error based on this and visualise the classification result.

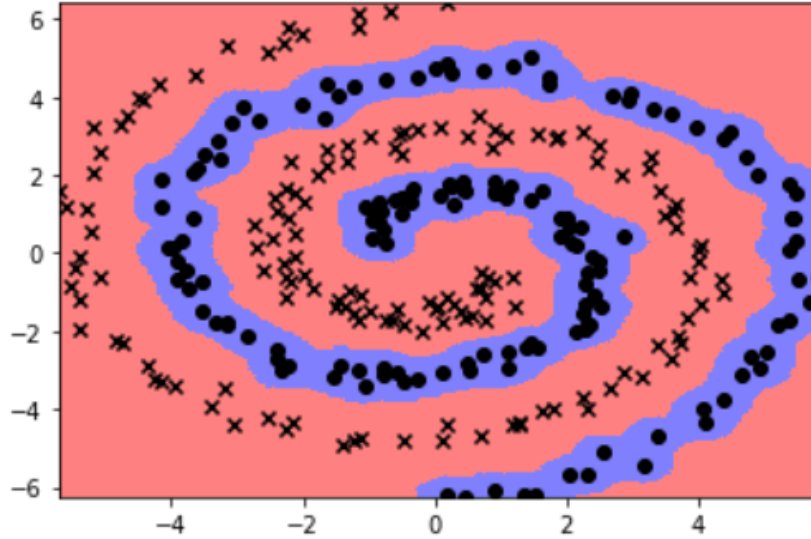


Figure 6: Classification result

The training error and test error are measured as $\frac{\text{number of correct prediction}}{\text{data set size}}$, which are both 0 in this case. This means the SVM-Light Algorithm works very well and classify all test data points correctly.

4 Conclusion

Based on experiments and theories, we can conclude that the SVM-Light is with a global asymptotic convergence and possess a faster convergence than SMO-Simplified. The local convergence could be linear for SVM-Light and we had better to choose μ between 0.1 and 5 to ensure a faster convergence and avoid non-S.P.D kernel matrix or overfitting happen. Finally, SVM-Light performs very well on spiral data set, dividing all test data correctly.

References

- [1]
- [2] A.Ng. Cs229 lecture notes. In *CS229 lecture notes*, page 1(1):1–3. 2020.
- [3] Hang Li. Svm. In *Statistical learning methods*. Tsinghua University Press, 2012.
- [4] Chih-Jen Lin. Linear convergence of a decomposition method for support vector machines. 2001.
- [5] Chih-Jen Lin. Asymptotic convergence of an smo algorithm without any assumptions. *IEEE Transactions on Neural Networks*, 13(1):248–250, 2002.