Quantum Sparse Support Vector Machines

Tomasz Arodz* and Seyran Saeedi[†]

Department of Computer Science, Virginia Commonwealth University Richmond, VA, USA

Abstract

We present a quantum machine learning algorithm for training Sparse Support Vector Machine, a linear classifier that minimizes the hinge loss and the L_1 norm of the feature weights vector. Sparse SVM results in a classifier that uses only a small fraction of the input features in making decisions, and is especially suitable for cases where the total number of features is at the same order, or larger, than the number of training samples. The algorithm utilizes recently proposed quantum solvers for semidefinite programming and linear programming problems. We show that while for an arbitrary binary classification problem no quantum speedup is achieved by using quantum SDP/LP solvers during training, there are realistic scenarios in which using a sparse linear classifier makes sense in terms of the expected accuracy of predictions, and polynomial quantum speedup compared to classical methods can be achieved.

1 Introduction

Binary classification involves vector-scalar pairs $(x,y) \in \mathcal{X} \times \mathcal{Y}$, where $\mathcal{Y} = \{-1,1\}$ and $\mathcal{X} \subset \mathbb{R}^p$ is a compact subset of p-dimensional feature space. Each pair describes an object of study, for example a brain scan or a tissue sample of a medical patient. Individual components x^j of a vector x are called features. Each feature describes some numerical property of the object represented by x, for example signal intensity in a single voxel of a brain scan, or expression level of a single gene. The value of y tells us whether the object belongs to the positive or the negative class. In many scenarios the feature vectors are easy to obtain, but the class variable is not. For example, we can measure methylation status of each CpG basepair in patient's genome relatively easily, but deciding if the patient's prognosis is positive or negative is challenging.

In statistical learning [FHT01], we assume that samples (x,y) come from a fixed but unknown distribution D over $\mathcal{X} \times \mathcal{Y}$. For a given feature vector x, the probabilities of either class are given by conditional distribution $D_{y|x}$ over \mathcal{Y} , and for a given class y, the probability density of feature vectors in that class is given by conditional distribution $D_{x|y}$ over \mathcal{X} . While the underlying distributions D, $D_{y|x}$, and $D_{x|y}$ are unknown, we have access to a training set Z consisting of m samples $z_i = (y_i, x_i)$ drawn independently from D. In the binary classification problem the goal is to use the training set to learn how to predict classes y for feature vectors x, even if we did not see such a feature vector in the training set.

^{*}Corresponding author. e-mail: tarodz@vcu.edu

[†]e-mail: saeedis@vcu.edu

The training set can be used to construct a predictive model, in a form of a hypothesis function $h: \mathcal{X} \to \mathbb{R}$, where the sign of h(x) indicates the predicted class for input feature vector x. For a given sample (x,y), the prediction is considered correct if the signs of the predicted and the true class agree, that is, if yh(x) > 0. The predictive model should make as few errors as possible over samples z = (x,y) sampled from distribution D, that is, it should minimize $\int_{\mathcal{X} \times \mathcal{Y}} I[yh(x) \leq 0]D(z) \, \mathrm{d}z$, where I is an indicator function over Boolean domain returning 1 for true and 0 for false.

A simple but often effective class of hypotheses is the class of linear functions $h(x; \beta, b) = \beta^T x + b = \sum_{j=1}^p \beta_j x^j + b$. A linear predictive model is parameterized by a vector of feature weights $\beta \in \mathbb{R}^p$ and a bias term $b \in \mathbb{R}$. To simplify the notation, we often add one more dimension to \mathcal{X} with all samples having a value of one. The predictive model is then simply $h(x; \beta) = \beta^T x$, $\beta \in \mathbb{R}^{p+1}$, with β_{p+1} playing the role of bias.

Training of a linear model involves finding a suitable parameter vector β . For a single sample (x, y), the suitability of a model h with specific β will be captured by a loss function $\ell(y, h(x; \beta))$, which returns a nonnegative real number that we interpret as a measure of our dissatisfaction with the prediction $h(x; \beta)$. The natural $\theta/1$ loss, defined as $\ell(y, h(x; \beta)) = I[y\beta^T x \leq 0]$, is not a continuous function of the parameter vector β , and is flat almost everywhere, leading to problems with finding β that minimizes the loss. Instead of the $\theta/1$ loss, a convex function that upper-bounds it is often used in training classification models. For example, the least-square loss $\ell(y, h) = (y - h)^2$ is used in Fisher's Linear Discriminant and in Least-Squares Support Vector Machine (LS-SVM) classifier [SV99].

Once the loss function is chosen, the goal of training a model is to find the parameter vector β that minimizes the expected loss $\mathbb{E}_{z\sim D}\ell(y,h(x;\beta))$, referred to as risk of the model, $R(\beta) = \int_{\mathcal{X}\times\mathcal{Y}}\ell(y,h(x;\beta))D(z)\,\mathrm{d}z$. Since D is unknown, a surrogate goal is to search for β that leads to low loss on samples from the training set. For example, the *empirical risk minimization* strategy involves finding parameters β that minimize *empirical risk*, that is, the average loss on the training set, $\hat{R}(\beta) = \frac{1}{m} \sum_{i=1}^{m} \ell(y_i, h(x_i; \beta))$.

The model β that minimizes the empirical risk may have high generalization risk $R(\beta)$, that is, may fare poorly on samples outside of the training set, especially if the number of training samples m is smaller, or not much larger, than the number of features p, the features are not statistically independent, or the feature values are noisy. Often, the generalization error can be reduced if a penalty on the complexity of the model is introduced into the optimization problem. Typically, this penalty term, known as regularization term, is based on $\|\beta\|$, a norm of the vector of model parameters, leading to regularized empirical risk minimization strategy, which finds parameters that minimize $\hat{L}(\beta) = \hat{R}(\beta) + \lambda f(\|\beta\|)$. For example, most Support Vector Machines [CV95] use squared L_2 norm of β , $\|\beta\|_2^2$, as the regularizer.

With technical progress in many experimental disciplines, the ability to measure large number of features in large number of samples is rapidly increasing. There is ongoing interest in fast methods for solving large-scale classification problems. One of the approaches being explored is training the predictive model using a quantum algorithm that has access to the training set stored in quantum-accessible memory. In parallel to research on efficient architectures for quantum memory [Ble10], work on quantum machine learning algorithms and on quantum learning theory is under way (see for example Refs. [BWP+17, DB18, SP18] and [AW17] for review). A pioneering example of this approach is Quantum LS-SVM [RML14], which achieves exponential speedup compared to classical LS-SVM algorithm, although the trained model, that is, the feature weights vector β , is computed as a quantum state and is not directly available for inspection.

Quantum LS-SVM uses quadratic least-squares loss and squared- L_2 regularizer, and thus

translates to an unconstrained quadratic programming (QP) problem, which can be solved using the seminal HHL [HHL09] algorithm for solving quantum linear systems of equations. The least-squares loss, while often used in regression problems, leads to high-magnitude loss if h(x) has large magnitude, even if h(x) and y agree on the sign, that is, the prediction is correct. Most classification loss functions involve a nonincreasing $\mathbb{R} \to \mathbb{R}_+$ function of the product yh. If the sign of prediction h(x) and the target class y agree, then loss should not increase if |h| increases. One prominent example of a convex, monotonic loss is the $hinge\ loss$, defined as $\ell(y,h)=[1-yh]_+=\max(0,1-yh)$, which is used in the original variant of the Support Vector Machine (SVM) classifier [CV95]. Hinge loss leads to hinge risk $\hat{R}_{SVM}(\beta)=\frac{1}{m}\sum_{i=1}^{m}\max(0,1-y_i\beta^Tx_i)$. However, using the hinge loss in SVM leads to a quadratic problem with inequality constraints, and algorithms based on quantum manipulation of eigenvalues such as HHL and other recent methods [SCK16, SBJ18] are only applicable to unconstrained, or equality-constrained quadratic problems, since only these types of QP problems can be re-interpreted as linear systems of equations.

The L_2 regularizer used in LS-SVM penalizes large-magnitude feature weights, but is unlikely to set any feature weights to null. In many real-world scenarios involving classification problems with large number of features we expect that highly-accurate predictions can be made using just a few discriminative features. The remaining features either carry no information about the separation of classes, or the information is redundant. For example, classification problems involving gene expression measured using microarrays or RNA-seq may have tens of thousands of features, and brain scans can have million of voxels, but only a small number may be enough to separate subjects with one subtype of a disease from another subtype, an information that is useful in choosing treatment. In these scenarios, we expect that a well-performing model should be sparse; that there is a vector β composed mostly of zeros that achieves near-optimal risk $R(\beta)$. The key problem is to decide which feature weights should be non-zero.

To find sparse solutions to classification problems, a regularization term in the form of L_1 norm of β is often included in the objective function. L_1 regularization is especially useful when working with a training set with large number of features compared to the number of training samples, which is referred to as the p > m case. Optimization problems involving L_1 norm typically lead to inequality constraints that cannot be presently handled by quantum algorithms based on HHL.

In this paper, we focus on Sparse SVM (sSVM) [Ben99, KH00, ZRHT04], a linear classifier based on regularized empirical risk minimization involving hinge loss and L_1 regularizer, $\hat{L}_{sSVM}(\beta) = \hat{R}_{SVM}(\beta) + \lambda \|\beta\|_1$, where $\lambda > 0$ is a hyperparameter specifying the strength of regularization. Training of a Sparse SVM model can be transformed into an optimization problem with linear objective function and linear inequality constraints. We introduce Quantum Sparse SVM (QsSVM), which is based on recently proposed quantum algorithms for solving semidefinite programming (SDP) problems [BS17, AGGW17, BKLL+17, AG18], of which linear programs are a special case. We show that while for arbitrary binary classification problems no quantum speedup is achieved using quantum SDP/LP solvers, there are realistic families of cases in which using a sparse linear classifier makes sense in terms of the expected accuracy of predictions, and polynomial quantum speedup compared to classical methods can be guaranteed. Moreover, the quantum SDP/LP solvers underlying QsSVM return more information about the trained predictive model β than it is in the case of Quantum LS-SVM based on HHL method. This is especially important for sparse linear predictive models, which are often used not just to predict class variables from the feature vectors, but to gain insight into which features affect the class variable.

2 Quantum Sparse SVM

The training of Sparse SVM model using a training set with p features and m samples involves solving a minimization problem

$$\underset{\beta \in \mathbb{R}^p}{\operatorname{arg min}} \frac{1}{m} \sum_{i=1}^m \max(0, 1 - y_i \beta^T x_i) + \lambda \sum_{j=1}^p |\beta_j|. \tag{1}$$

Using standard techniques, this non-linear unconstrained optimization problem can be transformed to an equivalent constrained linear program with n=m+2p nonnegative variables and m linear inequality constraints

$$\min_{\xi,\beta^{+},\beta^{-}} \frac{1}{m} \sum_{i=1}^{m} \xi_{i} + \lambda \sum_{j=1}^{p} \beta_{j}^{+} + \lambda \sum_{j=1}^{p} \beta_{j}^{-}$$
s.t.
$$\sum_{j=1}^{p} y_{i} x_{i}^{j} \beta_{j}^{+} - \sum_{j=1}^{p} y_{i} x_{i}^{j} \beta_{j}^{-} \ge 1 - \xi_{i}, \quad i \in [m]$$

$$\xi_{i}, \beta_{j}^{+}, \beta_{j}^{-} \ge 0,$$
(2)

where $[m] = \{1, ..., m\}$. Under positive λ , we can read out the solution as $\beta_j = \beta_j^+ - \beta_j^-$. We also have $|\beta_j| = \beta_j^+ + \beta_j^-$. The value of the hinge loss of *i*-th training sample is equal to ξ_i .

Simplex-based and interior-point methods are currently the main approaches for solving practical LP problems, but alternative theoretical approaches are being studied, including those aimed at a broader problem of semidefinite programs. An SDP problem with known symmetric $n \times n$ matrices C, A_i for i = [m], and known scalar constants b_i , involves finding a positive semidefinite $n \times n$ matrix X

$$\begin{aligned} & \min_{X} \quad C \cdot X \\ & \text{s.t.} \quad A_i \cdot X \leq b_i, \quad i \in [m] \\ & \quad X \succeq 0 \end{aligned}$$

where \cdot represents element-wise multiplication and $X \succeq 0$ indicates that matrix X is positive semidefinite. A corresponding dual SDP problem involves finding a vector $\alpha \in \mathbb{R}^m$

$$\max_{\alpha} \quad b^{T} \alpha$$
 s.t.
$$C - \sum_{i=1}^{m} \alpha_{i} A_{i} \succeq 0.$$

LP involving nonnegative variables is a special case of SDP where all the matrices are diagonal.

The dual variable α_i is null whenever $A_i \cdot X < b_i$, and for every non-zero α_i , the primal program constraint is satisfied with equality, $A_i \cdot X = b_i$. In the context of SVMs, where primal constraints and dual variables correspond to training samples, samples with $\alpha_i \neq 0$ are called support vectors. Those are all samples with non-zero loss, $\xi_i > 0$, as well as those with $\xi_i = 0$ for which the linear prediction is exactly equal to the class variable, $\beta^T x_i = y_i$; in both these cases, $y_i \beta^T x_i + \xi_i = 1$. On the other hand, samples with $y_i \beta^T x_i > 1$, which have null loss $\xi_i = 0$, are not support vectors, and $\alpha_i = 0$.

Arora and Kale [AK07, AK16] introduced a primal-dual algorithm for solving SDP/LP problems with input of size n and with m constraints with asymptotic computational complexity $\tilde{O}(mn \text{ poly}(R, r, 1/\varepsilon))$, where ε is the desired accuracy of the algorithm¹. The complexity depends not only on the size of the problem, n and m, but also on the size of the primal solution, R, and dual solution r; in the LP case these are captured by the L_1 norms of the primal and dual solution vectors.

Based on the Arora and Kale's approach, a quantum algorithm that uses quantum Gibbs sampling has been proposed recently by Brandão and Svore [BS17], with time complexity $\tilde{O}\left((\sqrt{mn}) \text{ poly}(R,r,1/\varepsilon)\right)$, which is a quadratic speed-up compared to the classical algorithm. Subsequent stream of improvements [AGGW17, BKLL⁺17] culminated up to this date with a quantum algorithm utilizing fast quantum OR lemma, provided by van Apeldoorn and Gilyén [AG18], with complexity $\tilde{O}\left((\sqrt{m}+\sqrt{n}) \text{ poly}(R,r,1/\varepsilon)\right)$.

The Quantum Sparse SVM (QsSVM) algorithm operates in the following way. The LP problem is seen as an SDP problem that involves diagonal matrices C and A_i , all of size n by n, where n = 2p + m. It also involves an m-dimensional vector b with all elements equal to -1. Access to the b vector is given by a unitary oracle

$$O_b |i\rangle |0\rangle = |i\rangle |b_i\rangle = |i\rangle |-1\rangle$$
,

where $i \in [m]$ and $|-1\rangle$ is a binary representation of -1 up to a chosen precision.

Access to matrix C and matrices A_i is given by unitary oracles

$$O_C |k, z\rangle = |k, z \oplus C[k, k]\rangle,$$

 $O_A |i, k, z\rangle = |i, k, z \oplus A_i[k, k]\rangle$

where z is a binary string with length depending on the chosen precision, \oplus represents bitwise XOR, and where $k \in [2p+m]$. For C, we have $C[k,k] = \frac{1}{m}$ for $k \in \{1,...,m\}$ and $C[k,k] = \lambda$ otherwise. For A_i , we have $A_i[k,k] = -1$ for $k \in \{1,...,m\}$, $A_i[k,k] = -y_ix_i^{(k-m)}$ for $k \in \{m+1,...,m+p\}$, and $A_i[k,k] = y_ix_i^{(k-m-p)}$ for k > m+p. The oracle O_A can be constructed from unitary oracles returning binary representations of $y_i \in \{-1,1\}$ and x_i^j and from efficient unitaries for bitstring multiplication and index addition and subtraction. The quantum oracles for y_i and x_i^j need to be implemented using quantum RAM. Quantum random access memory (qRAM) uses $\log N$ qubits to address any quantum superposition of N memory cell which may contains either quantum or classical information. For example, qRAM allows accessing classical data entries x_i^j in quantum superposition by a transformation

$$\frac{1}{\sqrt{mp}} \sum_{i=1}^{m} \sum_{j=1}^{p} |i,j\rangle |0...0\rangle \xrightarrow{\text{qRAM}} \frac{1}{\sqrt{mp}} \sum_{i=1}^{m} \sum_{j=1}^{p} |i,j\rangle |x_{i}^{j}\rangle,$$

where $|x_i^j\rangle$ is a binary representation up to a given precision. Discovering practical architectures for qRAM that allow query access in logarithmic time in terms of the number of items to be loaded is still an open challenge in quantum machine learning [BWP⁺17], with several approaches being considered [GLM08b, GLM08a, AGJO⁺15].

The quantum algorithm for training QsSVMs produces output in the form of samples from the normalized dual solution vector α , providing the identities of one support vector at a time, or sampling from a density operator proportional to solution X, which for X diagonal in computational basis provides the identities of the few non-zero-loss samples and the few non-zero feature weights β_i in the sparse solution.

 $[\]tilde{O}(f(n,m))$ hides factors that are logarithmic in n,m.

3 Complexity of Quantum Sparse SVM

Assuming oracle access to input, the computational complexity of quantum SDP solver utilized in Quantum Sparse SVM shows improved dependence on n and m, but polynomial dependence on R and r may erase any gains compared to classical LP solvers. Indeed, for any training set, the minimum of the objective function of the SparseSVM optimization problem (eq. 1) is bounded from above by one

$$1 \ge \min_{\beta} \frac{1}{m} \sum_{i=1}^{m} \max(0, 1 - y_i \beta^T x_i) + \lambda \sum_{i=1}^{p} |\beta_j|,$$

since an objective function value of one can be obtain by setting $\beta=0$, which leads to unit loss for each training sample, and thus unit average loss. For some training sets, one is the minimum of the objective function – for example if training samples come in pairs, (x,+1) and (x,-1). In this case, the norm of the primal solution is $R=\sum_i |\xi_i|+\sum_j |\beta_j|=m$, and the norm of the dual solution is $r=\sum_i |\alpha_i|=1$, since the dual objective function is just $\sum_i -\alpha_i$ and strong duality makes it equal to the value of the primal objective function.

In the complexity bound in quantum SDP/LP solvers, $\tilde{O}\left((\sqrt{m}+\sqrt{n})\operatorname{poly}(R,r,1/\varepsilon)\right)$, the polynomial term involving R and r has high exponent. The solver of van Apeldoorn and Gilyén [AG18] has complexity $\tilde{O}\left(\sqrt{m}(Rr/\varepsilon)^4+\sqrt{n}(Rr/\varepsilon)^5\right)$. Thus, $Rr=O\left(m\right)$ erases any speedups compared to classical solvers.

A more realistic case in which we see R = O(m) is a regular XOR problem, for example involving two features and four training samples, [+1,+1] and [-1,-1] with y = +1 and [+1,-1], [-1,+1] with y = -1. For any $\beta \in \mathbb{R}^2$, if there is a sample with loss $1 - \delta$, there is another sample with loss $1 + \delta$. Thus, sum of ξ_i variables is one for any β , and again $\beta = 0$ is the minimizer of the regularized empirical risk. However, for XOR problems a linear classifier is known to be useless even if speed is not a concern.

Below, we explore two realistic families of classification problems in which using a sparse linear classifier makes sense in terms of accuracy of the model, and in which quantum speedup can be achieved under mild conditions.

3.1 Hard-margin Sparse SVM

A ν -margin linearly separable classification problem, for $\nu > 0$, is a problem defined by underlying distribution D, characterized by an underlying vector $\beta^* \in \mathbb{R}^p$ with $\|\beta\|_2 = 1$, in which the class y is a deterministic function of x, that is, the conditional distribution over classes is $D_{y|x} \in \{0,1\}$ for each vector x. Further,

- for each $x \in \mathcal{X}$, the conditional distribution over classes $y \in \mathcal{Y}$ is $D_{y|x} = [\operatorname{sign}(y\beta^{\star T}x)]_+ \in \{0,1\}$, that is, the separation between the positive and negative classes is linear,
- the distribution D(x,y) has null mass for $\left\{x: \left|\frac{1}{\nu}\beta^{\star T}x\right| \leq 1\right\}$, that is, a narrow band of width ν on both sides of the linear decision boundary is devoid of samples from either class.

A ν -margin linearly separable problem is called sparse if the number of non-zero components in the vector β^* defining D is small compared to the number of features, p.

For this special case, the Sparse SVM regularized empirical risk minimization is characterized by the following lemma.

Lemma 1. For $p \to \infty$, consider a family of p-dimensional ν -margin linearly separable sparse classification problems D_p over $\mathcal{X}_p \times \mathcal{Y}$, where $\mathcal{X}_p \in \mathbb{R}^p$, based on sparse underlying vectors β_p^* with at most p' = O(f(p)) non-zero components, for some slowly increasing function f. Then, for each p, there exist

- a primal linear program with optimal solution $\hat{\beta}_p$ such that each training sample is classified correctly by $h(x) = \hat{\beta}_p^T x$, and $R_p = \|\hat{\beta}_p\|_1 = \sum_j \left(\hat{\beta}_j^+ + \hat{\beta}_j^-\right) \leq \rho_p$,
- a corresponding dual linear program with a solution vector α_p such that $r_p = \|\alpha_p\|_1 \le \rho_p$, such that $\rho_p \le \sqrt{p'}/\nu$.

Proof. For any training set sampled from the D there is at least one solution β with $\hat{R}(\beta) = 0$, that is, with $1 - y_i \beta^T x_i \leq 0$, or with $\xi_i = 0$ for each $i \in [m]$. We can thus narrow the solution space to solutions with null empirical risk, and consider a constrained linear problem without the slack variables ξ_i

$$\min_{\beta^{+},\beta^{-}} \sum_{j=1}^{p} \beta_{j}^{+} + \sum_{j=1}^{p} \beta_{j}^{-}$$
s.t.
$$\sum_{j=1}^{F} y_{i} x_{i}^{j} \beta_{j}^{+} - \sum_{j=1}^{F} y_{i} x_{i}^{j} \beta_{j}^{-} \ge 1, \quad i \in [m]$$

$$\beta_{j}^{+}, \beta_{j}^{-} \ge 0.$$
(3)

For brevity, below we will refer to the vectors in the LP solution space as β , and to the minimizer as $\hat{\beta}$. These should be understood as reconstructed from β^+ and β^- . For example, $\hat{\beta} = \hat{\beta}^+ + \hat{\beta}^-$, and $\|\hat{\beta}\|_1 = \sum_i (\hat{\beta}_i^+ + \hat{\beta}_i^-)$.

The scaled vector β^*/ν is a feasible solution to problem 3, since it achieves null loss on any samples from D, including those in the training set. For any unit L_2 -norm vector β with p' non-zero entries, the highest L_1 norm is achieved if all the p' coordinates are equal to $1/\sqrt{p'}$. Thus, $\|\beta\|_1 \leq \sqrt{p'}$. The solution β^* has unit L_2 -norm, and thus $\|\frac{1}{\nu}\beta^*\|_1 \leq \sqrt{p'}/\nu$. On the given training set, the empirical minimizer $\hat{\beta}$ has lowest L_1 norm among all feasible parameter vectors. Thus, $\|\hat{\beta}\|_1 \leq \|\frac{1}{\nu}\beta_p^*\|_1 \leq \sqrt{p'}/\nu$. Dual program has objective $\sum_{i=1}^m -\alpha_i$, and from strong duality, $r_p = \sum_{i=1}^m |\alpha_i| = \|\hat{\beta}\|_1 = R_p \leq \sqrt{p'}/\nu$.

An immediate corollary follows. A family of ν -margin linearly separable sparse classification problems in which the number of non-zero coordinates in the solution vector solution grows as $p' = O(\log p)$ leads to the family of LP problems with primal and dual solution norms $R_p = r_p = \frac{1}{\nu} O(\sqrt{\log p})$. For these problems, quantum SDP solvers offer speedup compared to currently available classical solvers.

3.2 Soft-margin Sparse SVM

To move beyond linearly separable case, we will consider scenarios where the classes overlap, but the optimal decision boundary is not far from being linear, and the region of overlap is limited, so that the generalization error resulting from using linear classifier is not high. As a motivating example, consider a p-dimensional classification problem where samples x in each class are distributed as an isotropic multivariate normal distribution with diagonal covariance matrix with the same value σ on the diagonal, but with different means. Without loss of generality, we will assume that $\sigma = 1$ – it can be achieved by rescaling the feature values.

In the two isotropic Gaussians classes case, the optimal solution is known to be a hyperplane, and the projection on the line orthogonal to the hyperplane, $u = \beta^{\star T} x$, results in the two classes forming univariate Gaussians $\mathcal{N}_{\mu}(u)$ and $\mathcal{N}_{-\mu}(u)$. The multivariate scenario, and the corresponding univariate case are depicted schematically in Figure 1a and 1b. We can further simplify the setup by considering a new variable $v = yu = y\beta^{\star T} x$; then both classes are distributed according to $D_{v|+} = D_{v|-} = \mathcal{N}_{\mu}(v)$. The hinge loss $\ell(y, h(x)) = [1 - yh(x)]_+$ in this setting is simply $L(v) = \max(0, 1 - y\beta^T x) = \max(0, 1 - v) = [1 - v]_+$, and the generalization risk associated with hinge loss is $\int_{-\infty}^{\infty} L(v) \mathcal{N}_{\mu}(v) dv$.

As before, we will consider a scenario in which as the number of features p grows, the number of features needed for accurate predictions, p', grows much more slowly. These discriminative features will have means +c and -c in the positive and the negative class, respectively – though the situation does not change if the signs of the means are swapped for some of the discriminative features. With increasing number p' of discriminative features, each with means differing by 2c, the distance between the means of the two multivariate isotropic Gaussians increases at the rate of at least $2c\sqrt{p'}$, and after the projection into single dimension to form $D_{v|+} = D_{v|-} = \mathcal{N}_{\mu}(v)$ as described above, the value of μ increases as $c\sqrt{p'}$.

To move beyond this idealized Gaussian scenario, we will consider problems governed by distributions D(x,y) that give rise to univariate class conditional distributions $D_{v|+}$ and $D_{v|-}$ that have tails in the region of non-zeros loss, $v \leq 1$, bounded from above by Gaussian tails, with the Gaussian mean μ diverging at a rate $c\sqrt{p'}$ as the number of discriminative features p' increases, and the tails are truncated beyond some constant $-\Delta$, also increasing with p'. Figure 1c shows this generalized scenario, which is formalized by the definition below.

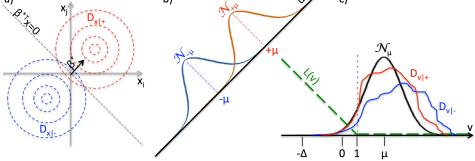
A (Δ, μ) -truncated subgaussian classification problem, for $\mu > 1$, $\Delta > 0$, is defined by distribution D such that there is an underlying vector $\beta^* \in \mathbb{R}^p$ with $\|\beta\|_2 = 1$, for which

- the conditional distributions $D_{x|+}$ and $D_{x|-}$ of the samples from the positive and negative class, respectively, give rise to univariate distributions $D_{v|+}$ and $D_{v|-}$ on a line resulting from the projection $v = y\beta^{*T}x$,
- the tails of $D_{v|+}$ and $D_{v|-}$ are bounded from above, in the region $v \in (-\infty, 1]$, by the probability density function $\mathcal{N}_{\mu}(v)$,
- the tails of $D_{v|+}$ and $D_{v|-}$ have zero mass for $v < -\Delta$.

A p-dimensional (Δ, μ) -truncated subgaussian problem is called sparse if the number of non-zero components in the vector β^* is small compared to the number of features, p.

Figure 1: Multivariate two-class problem and its univariate views.

b) c)



For the hinge loss, the generalization risk $R(\beta^*)$ associated with model $h(x) = \beta^{*T}x$ on the (Δ, μ) -truncated subgaussian problem D is bounded through the following lemma.

Lemma 2. Let D be a (Δ, μ) -truncated subgaussian classification problem with underlying vector β^* leading to univariate distributions $D_{v|+}$ and $D_{v|-}$ as described above. Let $L = \max(0, 1 - v)$ be a univariate random variable capturing hinge loss of the model $h(x) = \beta^{*T} x$ for samples from D. Then, the expectation and standard deviation of L are bounded by

$$R(\beta^*) = \mathbb{E}[L] \le \frac{1}{\sqrt{2\pi}} e^{-\frac{(1-\mu)^2}{2}} = \mathcal{N}_{\mu}(1),$$
 (4)

$$Var[L] \le \left[(1 - \mu)^2 + 1 \right] \left[1 + \operatorname{erf} \left(\frac{1 - \mu}{\sqrt{2}} \right) \right]. \tag{5}$$

Also, values of L are in the range $[0, \Delta + 1]$.

Proof. The proof relies on properties of integrals of $x^k \mathcal{N}_0(x)$. The technical details are given in the Appendix.

The result above gives the bound on the expected value of the hinge loss for the model $h(x) = \beta^{\star T} x$ on the distribution D, that is, it bounds from above the the generalization risk of that model, $R(\beta^{\star}) = \mathbb{E}[L]$. However, it does not give an upper bound on the empirical risk for the model $h(x) = \beta^{\star T} x$ on a specific training set with m samples and p features, sampled from p. This bound is given be the following lemma.

Lemma 3. Let D be a (Δ, μ) -truncated subgaussian problem based on β^* . Let $\hat{R}(\beta^*)$ be the empirical risk associated with model β^* over a m-sample training set sampled i.i.d. from D. Then, with probability at least $1 - \delta$

$$\hat{R}(\beta^*) \leq \frac{1}{\sqrt{2\pi}} e^{-\frac{(1-\mu)^2}{2}}$$

$$+ 4 \frac{\sqrt{\log(2/\delta)}}{\sqrt{m}} \left[(1-\mu)^2 + 1 \right] \left[1 + \operatorname{erf} \left(\frac{1-\mu}{\sqrt{2}} \right) \right]$$

$$+ 4 \frac{(\Delta+1)\log(2/\delta)}{m}$$

$$(6)$$

Proof. Consider m values $l_1, ..., l_m$ drawn from a univariate random variable L taking values in range in $[a, b] = [0, \Delta + 1]$, and with finite variance $s = \mathbb{V}\operatorname{ar}[L]$ and finite mean $R = \mathbb{E}[L]$. Let $\hat{R} = \frac{1}{m} \sum_{i=1}^{m} l_i$ be the empirical mean. Bernstein's inequality states that

$$\mathbb{P}(|\hat{R} - R| \ge t) \le 2 \exp\left(\frac{mt^2}{2(s^2 + (b - a)t)}\right).$$

That is, with probability at least $1 - \delta$,

$$\hat{R} \le R + 4s\sqrt{\frac{\log(2/\delta)}{m}} + \frac{4(b-a)\log(2/\delta)}{m}.$$

We thus have

$$\hat{R}(\beta^*) \le \mathbb{E}[L] + 4\mathbb{V}\operatorname{ar}[L]\sqrt{\frac{\log(2/\delta)}{m}} + \frac{4(\Delta+1)\log(2/\delta)}{m}$$

The bound follows from plugging in the bounds on expected value (eq. 4) and variance (eq. 5) of the loss.

We are now ready to analyze the behavior of empirical risk of models β_p^* on problems D_p as the number of all features p and the number of discriminative features p' grow.

Lemma 4. For $p \to \infty$, consider a family of p-dimensional (Δ_p, μ_p) -truncated subgaussian problems D_p with underlying vectors β_p^* . Assume that the vector β_p^* is sparse, it only has $p' = 1 + 2 \log p$ non-zero coefficients. Further, assume that the mean μ_p diverges with the number of discriminative features p' as $\mu_p > c\sqrt{p'}$ for some c > 1. As p grows, we allow scattering of the samples farther into the region dominated by the other class – specifically, we allow $\Delta_p \leq 2 \log p$. Then, with probability at least $1 - \delta$, we have

$$\hat{R}(\beta_p^*) \le \frac{1}{\sqrt{2\pi}p} + 4\frac{(2\log p + 1)\log(2/\delta)}{m}.\tag{7}$$

Proof. Under the assumption that μ_p grows at least as $c\sqrt{p'} = c\sqrt{1+2\log p}$, we have $\mu_p \ge 1+\sqrt{2\log p}$, which leads to the bound on the first term of eq. (6), and to the second term approaching null limit. The technical details of the proof are given in the Appendix.

Sparse SVM involves regularized empirical risk, that is, minimization of a weighted sum of the empirical risk and the L_1 norm of the model β . Under the scenario of slowly increasing number of discriminative features, the Sparse SVM regularized empirical risk minimization is characterized by the following lemma.

Lemma 5. For $p \to \infty$, consider a family of p-dimensional classification problems D_p as described in Lemma 4. For each D_p , consider the SparseSVM regularized empirical minimization problem (eq. 1)

$$\underset{\beta}{\text{arg min}} \quad \frac{1}{m} \sum_{i=1}^{m} \max(0, 1 - y_i \beta^T x_i) + \lambda \|\beta\|_1,$$

involving m-sample training set sampled from D_p . Then, for each p, with probability $1 - \delta$, there exist an empirical minimizer $\hat{\beta}_p$ of the problem above that can be found using a linear program (eq. 2), with L_1 norms of the primal and dual solutions, R_p and r_p , respectively, bounded as

$$R_p \le \frac{1}{\sqrt{2\pi}} \frac{m}{p} + 4(1 + 2\log p)\log(2/\delta)] + \lambda\sqrt{1 + 2\log p},$$
 (8)

$$r_p \le \frac{1}{\sqrt{2\pi}} \frac{1}{p} + \frac{4(1+2\log p)\log(2/\delta)}{m} + \lambda\sqrt{1+2\log p}.$$
 (9)

Proof. As shown in the hard-margin case in the previous section, we have $\|\beta_p^{\star}\|_1 \leq \sqrt{p'}$ for β_p^{\star} with unit L_2 -norm and p' non-zero coefficients. On the training set, the minimizer $\hat{\beta}_p$ has lowest objective function of all possible β , including β_p^{\star} . Thus, we have

$$\hat{R}(\hat{\beta}_p) + \lambda \|\hat{\beta}_p\|_1 \le \hat{R}(\beta_p^*) + \lambda \|\beta_p^*\|_1 \le \frac{1}{\sqrt{2\pi}} \frac{1}{p} + 4 \frac{(1 + 2\log p)\log(2/\delta)}{m} + \lambda \sqrt{1 + 2\log p}.$$

From strong duality, we have $r_p = \hat{R}(\hat{\beta}_p) + \lambda ||\hat{\beta}_p||_1$. The norm R_p of the primal solution does not involve averaging the losses $\max(0, 1 - y_i \beta^T x_i)$. Instead, the losses are added up, that is, R_p includes the term $m\hat{R}(\hat{\beta}_p)$ instead of the empirical risk.

The immediate corollary is that if we are dealing with p > m scenario, in which L_1 regularization is especially useful, that is, when m/p = O(1), neither R_p nor r_p grow with m, and both grow with p as $O(\log p)$. Utilizing the quantum SDP/LP solver proposed of van Apeldoorn and Gilyén [AG18], training QsSVMs, which translates to solving an LP problem (eq. 2) with m constraints and n = 2p + m variables, has computational complexity of

$$\tilde{O}\left(\sqrt{m}\left(\log^2 p/\varepsilon\right)^4 + \sqrt{m+2p}\left(\log^2 p/\varepsilon\right)^5\right) = \tilde{O}\left(\sqrt{m+2p}\operatorname{poly}\left(\log p, 1/\varepsilon\right)\right),\,$$

assuming efficient oracle access to data. Thus, using quantum SDP solvers offers speedup compared to classical solvers. More generally, the computational complexity is $\tilde{O}\left(\sqrt{m+2p}\operatorname{poly}\left(\frac{m\log p}{p},\log p,1/\varepsilon\right)\right)$, leading to speedup even in some cases beyond the p>m scenario, such as $m=O\left(p\log p\right)$. The scenarios in which number of features is larger or at least comparable to the number of samples is of great practical importance – it is common in biomedical data analysis, for example in classification of molecular profiles such as gene expression or methylation, or classification of 3D brain scans.

Acknowledgments

TA is supported by NSF grant IIS-1453658. Early results that are part of this work were presented as poster at the Conference on Quantum Machine Learning Plus in Innsbruck, Austria, 2018. We are grateful to Ronald de Wolf for inspiring comments on our work.

References

- [AG18] Joran van Apeldoorn and András Gilyén. Improvements in quantum SDP-solving with applications. arXiv preprint arXiv:1804.05058, 2018.
- [AGGW17] Joran van Apeldoorn, András Gilyén, Sander Gribling, and Ronald de Wolf. Quantum SDP-solvers: Better upper and lower bounds. In *Proc. IEEE 58th Annual Symposium on Foundations of Computer Science (FOCS'17)*, pages 403–414. IEEE, 2017.
- [AGJO⁺15] Srinivasan Arunachalam, Vlad Gheorghiu, Tomas Jochym-O'Connor, Michele Mosca, and Priyaa Varshinee Srinivasan. On the robustness of bucket brigade quantum RAM. New Journal of Physics, 17(12):123010, 2015.
- [AK07] Sanjeev Arora and Satyen Kale. A combinatorial, primal-dual approach to semidefinite programs. In *Proc. 39th Annual ACM Symposium on Theory of Computing (STOC'07)*, pages 227–236. ACM, 2007.
- [AK16] Sanjeev Arora and Satyen Kale. A combinatorial, primal-dual approach to semidefinite programs. *Journal of the ACM*, 63(2):12, 2016.
- [AW17] Srinivasan Arunachalam and Ronald de Wolf. A survey of quantum learning theory. ACM SIGACT News, 48(2):41–67, 2017.
- [Ben99] Kristin Bennett. Combining support vector and mathematical programming methods for induction. Advances in Kernel Methods: Support Vector Learning, pages 307–326, 1999.

- [BKLL⁺17] Fernando GSL Brandão, Amir Kalev, Tongyang Li, Cedric Yen-Yu Lin, Krysta M Svore, and Xiaodi Wu. Quantum SDP solvers: Large speed-ups, optimality, and applications to quantum learning. arXiv preprint arXiv:1710.02581, 2018.
- [Ble10] Miles Blencowe. Quantum computing: Quantum RAM. Nature, 468(7320):44, 2010.
- [BS17] Fernando GSL Brandão and Krysta M Svore. Quantum speed-ups for solving semidefinite programs. In *Proc. IEEE 58th Annual Symposium on Foundations of Computer Science (FOCS'17)*, pages 415–426. IEEE, 2017.
- [BWP⁺17] Jacob Biamonte, Peter Wittek, Nicola Pancotti, Patrick Rebentrost, Nathan Wiebe, and Seth Lloyd. Quantum machine learning. *Nature*, 549(7671):195, 2017.
- [CV95] Corinna Cortes and Vladimir Vapnik. Support-vector networks. *Machine Learning*, 20(3):273–297, 1995.
- [DB18] Vedran Dunjko and Hans J Briegel. Machine learning & artificial intelligence in the quantum domain: a review of recent progress. Reports on Progress in Physics, 81(7):074001, 2018.
- [FHT01] Jerome Friedman, Trevor Hastie, and Robert Tibshirani. The elements of statistical learning. Springer, 2001.
- [GLM08a] Vittorio Giovannetti, Seth Lloyd, and Lorenzo Maccone. Architectures for a quantum random access memory. *Physical Review A*, 78(5):052310, 2008.
- [GLM08b] Vittorio Giovannetti, Seth Lloyd, and Lorenzo Maccone. Quantum random access memory. *Physical Review Letters*, 100(16):160501, 2008.
- [HHL09] Aram W Harrow, Avinatan Hassidim, and Seth Lloyd. Quantum algorithm for linear systems of equations. *Physical Review Letters*, 103(15):150502, 2009.
- [KH00] Vojislav Kecman and Ivana Hadzic. Support vectors selection by linear programming. In *Proc. International Joint Conference on Neural Networks (IJCNN'00)*, volume 5, pages 193–198. IEEE, 2000.
- [RML14] Patrick Rebentrost, Masoud Mohseni, and Seth Lloyd. Quantum support vector machine for big data classification. *Physical Review Letters*, 113(13):130503, 2014.
- [SBJ18] Sathyawageeswar Subramanian, Steve Brierley, and Richard Jozsa. Implementing smooth functions of a Hermitian matrix on a quantum computer. arXiv preprint arXiv:1806.06885, 2018.
- [SCK16] Rolando Somma, Andrew Childs, and Robin Kothari. Quantum linear systems algorithm with exponentially improved dependence on precision. In *APS Meeting Abstracts*, 2016.
- [SP18] M. Schuld and F. Petruccione. Supervised Learning with Quantum Computers. Springer Nature, 2018.
- [SV99] Johan AK Suykens and Joos Vandewalle. Least squares support vector machine classifiers. *Neural Processing Letters*, 9(3):293–300, 1999.

[ZRHT04] Ji Zhu, Saharon Rosset, Trevor Hastie, and Robert Tibshirani. 1-norm support vector machines. In *Advances in Neural Information Processing Systems* (NIPS'04), pages 49–56. MIT Press, 2004.

A Complete Proofs of Lemmas 2 and 4

Lemma 2. Let D be a (Δ, μ) -truncated subgaussian classification problem with underlying vector β^* leading to univariate distributions $D_{v|+}$ and $D_{v|-}$ as described above. Let $L = \max(0, 1-v)$ be a univariate random variable capturing hinge loss of the model $h(x) = \beta^{*T}x$ for samples from D. Then, the expectation and standard deviation of L are bounded by

$$R(\beta^*) = \mathbb{E}[L] \le \frac{1}{\sqrt{2\pi}} e^{-\frac{(1-\mu)^2}{2}} = \mathcal{N}_{\mu}(1),$$
$$\mathbb{V}\text{ar}[L] \le \left[(1-\mu)^2 + 1 \right] \left[1 + \text{erf} \left(\frac{1-\mu}{\sqrt{2}} \right) \right].$$

Also, values of L are in the range $[0, \Delta + 1]$.

Proof. Let

$$G(x) = \mathcal{N}_0(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}},$$

$$G_k(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x t^k e^{-\frac{t^2}{2}} dt.$$

Then, we have

$$G_0(x) = \frac{1}{2} \left[1 + \operatorname{erf} \left(\frac{x}{\sqrt{2}} \right) \right],$$

$$G_1(x) = -G(x) = -\frac{1}{\sqrt{2\pi}} e^{-\frac{t^2}{2}},$$

$$G_2(x) = G_0(x) - xG(x).$$

Let p_+ and p_- by the probabilities, under D, of the positive and the negative class, respectively. For expected value, we have

$$\mathbb{E}[L] = \int_{-\infty}^{\infty} \max(0, 1 - v) [p_{+}D_{v|+}(v) + p_{-}D_{v|-}(v)] dv$$

$$\leq \int_{-\infty}^{1} (1 - v) [p_{+}\mathcal{N}_{\mu}(v) + p_{-}\mathcal{N}_{\mu}(v)] dv$$

$$= \int_{-\infty}^{1} (1 - v) \frac{1}{\sqrt{2\pi}} e^{-\frac{(v - \mu)^{2}}{2}} dv$$

$$= \int_{-\infty}^{1 - \mu} (1 - (t + \mu)) \frac{1}{\sqrt{2\pi}} e^{-\frac{t^{2}}{2}} dt$$

$$= (1 - \mu)G_{0}(1 - \mu) - G_{1}(1 - \mu) \leq \frac{1}{\sqrt{2\pi}} e^{-\frac{(1 - \mu)^{2}}{2}} = \mathcal{N}_{\mu}(1),$$

where we used substitution $t = v - \mu$, and the last inequality comes from $1 - \mu < 0$.

For variance, we have

$$\begin{aligned} \operatorname{Var}[L] &= \mathbb{E}[L^2] - \mathbb{E}[L]^2 \leq \mathbb{E}[L^2] = \int_{-\infty}^{\infty} \max(0, 1 - v)^2 [p_+ D_{v|+}(v) + p_- D_{v|-}(v)] \, \mathrm{d}v \\ &\leq \int_{-\infty}^{1} (1 - v)^2 [p_+ \mathcal{N}_{\mu}(v) + p_- \mathcal{N}_{\mu}(v)] \, \mathrm{d}v \\ &\leq \int_{-\infty}^{1} (1 - v)^2 \frac{1}{\sqrt{2\pi}} e^{-\frac{(v - \mu)^2}{2}} \, \mathrm{d}v \\ &= \int_{-\infty}^{1 - \mu} (1 - \mu - t)^2 \frac{1}{\sqrt{2\pi}} e^{-\frac{t^2}{2}} \, \mathrm{d}t \\ &= \int_{-\infty}^{1 - \mu} ((1 - \mu)^2 - 2(1 - \mu)t + t^2) \frac{1}{\sqrt{2\pi}} e^{-\frac{t^2}{2}} \, \mathrm{d}t \\ &= (1 - \mu)^2 G_0(1 - \mu) - 2(1 - \mu)G_1(1 - \mu) + G_2(1 - \mu) \\ &= (1 - \mu)^2 G_0(1 - \mu) + 2(1 - \mu)G(1 - \mu) + G_0(1 - \mu) - (1 - \mu)G(1 - \mu) \\ &= [(1 - \mu)^2 + 1]G_0(1 - \mu) + (1 - \mu)G(1 - \mu) \\ &\leq [(1 - \mu)^2 + 1]G_0(1 - \mu) = \left[(1 - \mu)^2 + 1 \right] \left[1 + \operatorname{erf} \left(\frac{1 - \mu}{\sqrt{2}} \right) \right]. \end{aligned}$$

The range of L follows immediately from the null mass of $D_{v|y}$ for $v \leq -\Delta$, and from null loss for any $v \geq 1$.

Lemma 4. For $p \to \infty$, consider a family of p-dimensional (Δ_p, μ_p) -truncated subgaussian problems D_p with underlying vectors β_p^{\star} . Assume that the vector β_p^{\star} is sparse, it only has $p' = 1 + 2\log p$ non-zero coefficients. Further, assume that the mean μ_p diverges with the number of discriminative features p' as $\mu_p > c\sqrt{p'}$ for some c > 1. As p grows, we allow scattering of the samples farther into the region dominated by the other class – specifically, we allow $\Delta_p \leq 2\log p$. Then, with probability at least $1 - \delta$, we have

$$\hat{R}(\beta_p^\star) \leq \frac{1}{\sqrt{2\pi}p} + 4\frac{(2\log p + 1)\log(2/\delta)}{m}.$$

Proof. We will start with the eq. (6)

$$\hat{R}(\beta^*) \le \frac{1}{\sqrt{2\pi}} e^{-\frac{(1-\mu_p)^2}{2}} + 4\frac{\sqrt{\log(2/\delta)}}{\sqrt{m}} \left[(1-\mu_p)^2 + 1 \right] \left[1 + \operatorname{erf} \left(\frac{1-\mu_p}{\sqrt{2}} \right) \right] + 4\frac{(\Delta_p + 1)\log(2/\delta)}{m}$$

and provide bounds on the first and second term.

We have the following limit $\lim_{x\to\infty} c\sqrt{1+kx}/[1+\sqrt{kx}]=c$. Thus, under the assumption that μ_p grows at least as $c\sqrt{p'}=c\sqrt{1+2\log p}$, for c>1, for sufficiently large p, we have $\mu_p\geq c\sqrt{1+2\log p}\geq 1+\sqrt{2\log p}$. That is, $1-\mu_p\leq -\sqrt{2\log p}$.

For the first term of the bound in eq. (6), e^x is an increasing function of x, we thus have the following upper bound

$$\frac{1}{\sqrt{2\pi}}e^{-\frac{(1-\mu_p)^2}{2}} \le \frac{1}{\sqrt{2\pi}}e^{-\log p} = \frac{1}{\sqrt{2\pi}}\frac{1}{p}.$$

For the second term in eq. (6), we can show that $1 + \operatorname{erf}\left(\frac{1-\mu_p}{\sqrt{2}}\right)$ approaches null with the rate faster than $\frac{1}{p^2}$. Since $1 - \mu_p \le -\sqrt{2\log p}$ and $\operatorname{erf}(x)$ is an increasing function of x, we have $1 + \operatorname{erf}\left(\frac{1-\mu_p}{\sqrt{2}}\right) \le 1 + \operatorname{erf}\left(-\sqrt{2\log p}\right)$. We also have

$$\frac{\mathrm{d}\left[1 + \mathrm{erf}\left(-\sqrt{2\log p}\right)\right]}{\mathrm{d}p} = -\frac{\sqrt{\frac{2}{\pi}}}{p^3\sqrt{\log p}}.$$

Thus, from the L'Hôpital's rule,

$$\lim_{p \to \infty} \frac{1 + \operatorname{erf}\left(-\sqrt{2\log p}\right)}{p^{-2}}$$

$$= \lim_{p \to \infty} \frac{\frac{d}{dp} \left[1 + \operatorname{erf}\left(-\sqrt{2\log p}\right)\right]}{\frac{d}{dp} p^{-2}}$$

$$= \lim_{p \to \infty} \frac{-\frac{\sqrt{\frac{2}{\pi}}}{p^3 \sqrt{\log p}}}{-2\frac{1}{p^3}}$$

$$= \lim_{p \to \infty} \frac{1}{2\sqrt{\pi}} \frac{1}{\sqrt{\log p}} = 0.$$

Thus, $\left[(1-\mu_p)^2+1\right]\left[1+\operatorname{erf}\left(\frac{1-\mu_p}{\sqrt{2}}\right)\right]=O\left(\left[(-\sqrt{2\log p})^2+1\right]/p^2\right)=O\left(\log p/p^2\right)$. The second term in eq. (6) quickly approaches null as p grows.