Support Vector Machines: Algorithms and Applications

Tong Wen

LBNL

in collaboration with

Alan Edelman

MIT

Jan 17, 2003

- Support Vector Machines (SVMs) solve classification problems by learning from examples.
- Contents:
 - 1. Introduction to Support Vector Machines.
 - Fast SVM training algorithms.
 - 3. Financial applications of SVMs?
- Notations:
 - the *i*th training point $oldsymbol{x}_i$
 - the label of the *i*th training point y_i
 - the vector of Lagrange multipliers $\boldsymbol{\alpha}$
 - the *i*th Lagrange multiplier α_i
 - the vector giving the normal direction of a hyperplane W
 - the bias term b
 - the dimension of a training point n
 - the number of training points m
 - Η the Hessian matrix
 - Xthe matrix whose columns are the training points
 - Pthe matrix that spans the current search space
 - $Q \\ S_1$ the orthogonal complement of P
 - the set of SVs that are correctly separated
 - the set of SVs that are separated with errors S_2
 - $h(\cdot)$ a decision rule (function)

Introduction to Support Vector Machines

- An example.
- The primal SVM quadratic programming problem.
- The dual SVM quadratic programming problem.
- Properties of Support Vectors(SVs).

An Example

• The mechanism that generates each point *x* and its color *y*:



• The main theme of classification is to determine a decision rule:

 $h: \mathbb{X} \subseteq \mathbb{R}^n \to \mathbb{Y} \subseteq \mathbb{R}.$

For this example, $\mathbb{Y} = \{\mathbf{1}, -\mathbf{1}\}.$

• What is the optimal solution to our detection problem?

• By Bayes' Rule, the maximum likelihood decision rule is

$$h_{\mathsf{ML}}(x) = \underset{\hat{y} \in \{\pm 1\}}{\operatorname{argmax}} p(x \mid y = \hat{y}),$$

which minimize the risk

$$R[h] = E[(\frac{1}{2}|y - h(x)|)].$$



m_=(1,1), m_=(3,3), σ =1

- The decision rule $h_{ML}(\cdot)$ depends on P(x, y).
- In general, statistical information such as P(x, y) is not available. A distribution-free approach is needed.
- The SVM decision rule is determined based on the training examples:

$$S = \{(x_1, y_1), \ldots, (x_m, y_m)\}.$$

• The training error:

$$R_S[h] = \frac{1}{m} \sum_{i=1}^m \frac{1}{2} |y_i - h(x_i)|.$$

• Overfitting vs. the training error.

• In the previous example,

$$h_{\mathsf{ML}}(x) = \operatorname{sgn}(w^T x + b).$$

- Every decision rule geometrically corresponds to a separating boundary.
- If $f_{SB}(x) = 0$ defines the separating boundary, then the corresponding decision rule has a general form

$$h(x) = \operatorname{sgn}(f_{\mathsf{SB}}(x)).$$

• The SVM decision rule:

$$\begin{split} h_{\mathsf{SVM}}(x) &= \mathsf{sgn}(w^T x + b), \\ h_{\mathsf{SVM}}(x) &= \mathsf{sgn}(w^T \phi(x) + b) \\ &= \mathsf{sgn}(\sum_{i=1}^m \alpha_i k(x_i, x) + b). \end{split}$$

 SVMs use two parallel hyperplanes instead of one to determine the separating boundary. These two (oriented) hyperplanes are computed directly using the training set S.



The Primal SVM QP Problem: the Separable Case

• w and b are determined by the quadratic programming (QP) problem:

$$\underset{w,b}{\text{minimize}} \ \frac{1}{2} \|w\|_2^2$$

subject to

$$y_i(w^T x_i + b) \ge 1$$
, for $i = 1, ..., m$.

- We can always separate the training points by mapping them to a higher dimensional space $\mathbb{R}^{n'}$ (n' > n), such that in $\mathbb{R}^{n'}$ the two sets $\{\phi(x_i) \mid y_i = 1\}$ and $\{\phi(x_i) \mid y_i = -1\}$ can be separated by hyperplanes.
- In practice, people may not want to fully separate the training set because of overfitting.

The Primal SVM QP Problem: the Inseparable Case

 To accommodate the separation errors on S, nonnegative slack variables (errors) ε_i are introduced:

$$\underset{w,b,\varepsilon}{\text{minimize}} \quad \frac{1}{2} \|\boldsymbol{w}\|_2^2 + c \sum_i \varepsilon_i$$

subject to (for $i = 1, \ldots, m$)

$$y_i(w^T \phi(x_i) + b) \ge 1 - \varepsilon_i \text{ and } \varepsilon_i \ge 0.$$

- If the training point $\phi(x_i)$ is correctly separated by the soft-margin hyperplanes, then $\varepsilon_i = 0$.
- The coefficient c controls the trade-off between the size of the gap and how well S is separated.

• In practice, the dual problems are computed instead:

maximize
$$\alpha^T 1 - \frac{1}{2} \alpha^T H \alpha$$

subject to

$$y^T lpha = 0,$$

 $lpha \ge 0.$

maximize
$$\alpha^T \mathbf{1} - \frac{1}{2} \alpha^T H \alpha$$

subject to

- $y^T lpha = 0,$ $0 \le lpha \le c.$
- The Hessian matrix $H_{ij} = y_i(\phi(x_i)^T \phi(x_j)) y_j$ is a $m \times m$ symmetric positive semi-definite matrix.

• Since only inner products are involved in the dual problems, a positive definite kernel function $k(\cdot, \cdot)$ is used instead of $\phi(\cdot)$, where

```
k(\boldsymbol{x}_i, \boldsymbol{x}_j) = \phi(\boldsymbol{x}_i)^T \phi(\boldsymbol{x}_j) and H_{ij} = y_i k(\boldsymbol{x}_i, \boldsymbol{x}_j) y_j.
```

- Using kernel functions avoids the curse of dimensionality.
- Three positive definite kernels are used here:

linear: $k(s,t) = s^T t$ poly(d): $k(s,t) = (s^T t + 1)^d$ rbf(σ): $k(s,t) = e^{-\|s-t\|_2^2/(2\sigma^2)}$

• For simplicity, we use x_i in place of $\phi(x_i)$ in our later discussions.

• At optimality,

$$egin{aligned} & m{w} = \sum\limits_{i=1}^m y_i lpha_i x_i, \ & ext{and if } 0 < lpha_i < c \ (y_i (m{w}^T x_i + b) = 1) ext{ then} \ & b = y_i - m{w}^T x_i = y_i (1 - e_i^T H lpha). \end{aligned}$$

•
$$h_{\mathsf{SVM}}(x) = \mathsf{sgn}(w^T x + b).$$

- $S_1 \cup S_2$ contain the support vectors (SVs). $\begin{cases}
 0 < \alpha_i < c & \longleftrightarrow \mathbf{x}_i \in S_1, \ S_1 = \{\mathbf{x}_i | \ y_i(\mathbf{w}^T \mathbf{x}_i + b) = 1 \ (\varepsilon_i = 0)\}; \\
 \alpha_i = c & \longleftrightarrow \mathbf{x}_i \in S_2, \ S_2 = \{x_i | \ y_i(\mathbf{w}^T \mathbf{x}_i + b) = 1 - \varepsilon_i \text{ and } \varepsilon_i > 0\}; \\
 \alpha_i = 0 & \longleftrightarrow \mathbf{x}_i \in S_3, \ S_3 = \{\mathbf{x}_i | \ y_i(\mathbf{w}^T \mathbf{x}_i + b) > 1 \ (\varepsilon_i = 0)\}.
 \end{cases}$
- SVs are the training points that determine the two separating hyperplanes through linear equations.
- The solution to the dual problem is sparse.

Solving the Dual SVM QP Problem Efficiently

- Solving the dual problem is computationally challenging.
- The decomposition strategy.
- A projected conjugate gradient algorithm.
- Speed considerations.
- Numerical experiments.
- Conclusion and future work.

The Decomposition Strategy

- The sparsity property.
- An observation from the primal problem point of view:



- How to apply this strategy to solve the dual problem efficiently?
 - 1. Constructing a subproblem.
 - 2. Solving the subproblems.
 - 3. The optimality conditions.
- A projected conjugate gradient algorithm.

Applying the Decomposition Strategy

- Projecting the dual QP problem onto a subspace.
 - Let us restrict α in the affine space determined by α_0 and *P*. That is, $\alpha = \alpha_0 + P\hat{\alpha}$.
 - The projected dual problem is

maximize
$$\hat{F}(\hat{\alpha}) \equiv \hat{\alpha}^T P^T r_0 - \frac{1}{2} \hat{\alpha}^T P^T H P \hat{\alpha}$$

subject to

$$y^T P \hat{\alpha} = 0,$$

 $0 \le \alpha_0 + P \hat{\alpha} \le c,$

where r_0 is the gradient of $F(\alpha)$ at α_0 .

- *P* is a general matrix so far. It can be specified directly or by its orthogonal complement *Q*, where $P^T Q = 0$.
- How to determine *P* or *Q*, and how to solve each subproblem?

•
$$\alpha = \alpha_0 + P\hat{\alpha}$$
.

$$P = E_p \text{ and } Q = E_q \ (m = p + q)$$
maximize $\hat{F}(\hat{\alpha}) \equiv \hat{\alpha}^T \hat{r}_0 - \frac{1}{2} \hat{\alpha}^T \hat{H} \hat{\alpha}$
subject to
 $\hat{y}^T \hat{\alpha} = 0,$
 $0 \le \hat{\alpha}_0 + \hat{\alpha} \le c,$
where $\hat{r}_0 = P^T r_0, \hat{y} = P^T y, \hat{H} = P^T H P$, etc.
$$P = I - \frac{\hat{y}\hat{y}^T}{p}, Q = \hat{y} \text{ and } \hat{\alpha} = P \hat{\alpha}$$
maximize $\hat{F}(\hat{\alpha}) \equiv \hat{\alpha}^T P^T \hat{r}_0 - \frac{1}{2} \hat{\alpha}^T P^T \hat{H} P \hat{\alpha}$
subject to
 $0 \le \hat{\alpha}_0 + P \hat{\alpha} \le c.$

•
$$\hat{H}$$
 is a $p \times p$ matrix ($p \ll m$).

Solving the Subproblems by the CG Method

• If $0 \leq \hat{lpha_0} + P \hat{lpha} \leq c$ never becomes active, we have

$$P^T \hat{H} P \hat{\alpha} = P^T \hat{r}_0,$$

or equivalently

 $P\hat{H}P\hat{\alpha} = P\hat{r}_0$ (*P* is symmetric).

• How about $0 \le (\hat{\alpha}_0)_i + (P\hat{\alpha})_i \le c$ becomes active during the iteration?

• An equivalent problem:

$$Q = [E_q, y]$$
 and $P = I - Q(Q^T Q)^{-1}Q^T$.

• Define

$$d = Pr$$

and

$$\boldsymbol{\beta} = -[I_{q \times q}, \boldsymbol{0}](Q^T Q)^{-1} Q^T \boldsymbol{r},$$

where $r = 1 - H\alpha$ and *P* is a $m \times m$ matrix.

• α is optimal iff

$$d = 0$$
 and $\beta \ge 0$.

• Since $(Q^T Q)^{-1} = \begin{bmatrix} I + \frac{E_q^T y y^T E_q}{p} & \frac{-E_q^T y}{p} \\ \frac{-y^T E_q}{p} & \frac{1}{p} \end{bmatrix}$, it is easy to compute d and β .

$$\beta = E_q^T r + \frac{E_q^T y((E_q^T y) E_q^T r - y^T r)}{p}.$$

• Each time when active constraints are to be relaxed, we start with the one that has the most negative β_i .

• A summary of our algorithm:

```
\alpha = 0;
initialize E_p;
while d \neq 0 or \beta < 0
at most k steps of the CG iterations for the subproblem
update \alpha and r:
\alpha = \alpha + E_p(I - \frac{yy^T}{2})\hat{\alpha}.
```

$$egin{aligned} &lpha &= lpha + E_p(I - \frac{p}{p}) lpha, \ &r &= r - HE_p(I - \frac{yy^T}{p}) \hat{lpha}; \end{aligned}$$

relax at most l active constraints with the most negative β_i ; update E_p ;

compute the columns of H corresponding to the relaxed constraints; update HE_p and $\hat{H};$

end

• Memory requirement: (*H* is not precomputed!)

 $X = [x_1, \ldots, x_m]$ and HE_p .

- Flops: Solving a sequence of smaller problems is cheaper when the solution is sparse.
- Two parameters: k and l.
- How to initialize E_p ?

Speed Considerations: Being Adaptive to the Memory Hierarchy

• A model of the modern computer memory hierarchy:



- Linear algebra operations: level 1 3.
- ATLAS BLAS.
- Thanks to MATLAB's inclusion of ATLAS BLAS.

Speed Considerations: Setting k and l Adaptively

- The spectrum of H.
 - The size of $S_1 \cup S_2$.
 - the convergence of the CG method.
- k is used to control the steps of the CG iterations.
 - There is no need to compute every subproblem exactly.
 - The rank deficiency of H.
- *l* is used to control the number of active constraints relaxed at each time.
 - When *l* active constraints are relaxed, $[x_{i_1}, \ldots, x_{i_l}]^T X$ is computed to generate the corresponding columns of *H*.
 - If *l* is set large, then there is more chance of some previously relaxed constraints becoming active again.
- k and l is set to be \hat{k} , the estimated number of the leading eigenvalues of H.
 - In our algorithm, \hat{k} is set to be the number of iterations for the CG method to reduce the residual norm to a scale of 0.01 for the initial subproblem.

• The two training sets:

Training sets	Digit17	Face
Size of the training set	13007	31022
The positive training points	6742	2901
The negative training points	6265	28121
The dimension of the training points	784	361
The format of the training points	sparse	dense
Separability	easier	harder

- Comparing against SVM^{light} (C) and SvmFu (C++).
- The choices of k(s, t) and c:

Training sets: Digit17-6000 and Digit17						
I(1)	I(0.1)	p2(0.1)	p2(0.001)			
r10(10)	r10(1)	r3(1)	r3(0.1)			
Training sets: Face-6000, Face-13901 and Face						
I(4)	I(1)	I(0.1)	p2(1)			
p2(0.001)	r10(10)	r3(10)	r3(1)			

Using Default Settings **Using Estimated Optimal Settings** Digit17-6000 Digit17-6000 50 30 FMSvm SVM^{light} SvmFu FMSvm SVM^{light} SvmFu * * 45 25 40 35 20 15 10 10 15 10 5 0 I(1) 0 I(1) l(0.1) p2(0.1) p2(0.001) r10(10) r10(1) r3(1) r3(0.1) kernel (c) l(0.1) p2(0.1) p2(0.001) r10(10) r10(1) r3(1) r3(0.1) kernel (c) Face-6000 Face-6000 700 2000 FMSvm SVM^{light} FMSvm SVM^{light} * 1800 600 SvmFu SvmFu 1600 500 1400 sp1200 1000 800 800 seconds 300 CPU seconds 600 200 400 100 200 0∟ I(4) 0 I(4) l(0.1) p2(0.1) p2(0.001) r10(10) r3(10) kernel (c) l(0.1) p2(0.1) p2(0.001) r10(10) r3(10) kernel (c) l(1) r3(1) l(1) r3(1)



200

0 (4)

l(1)

l(0.1)

10 0 I(1)

l(0.1) p2(0.1) p2(0.001) r10(10) r10(1) r3(1) r3(0.1) kernel (c)

Comparison of FMSvm with ${\rm SVM}^{\mathit{light}}$ and SvmFu with the estimated optimal parameter settings



r3(1)

Conclusion and Future Work

- Conclusion:
 - 1. efficient;
 - 2. easy to use;
 - 3. highly portable.
- Future work:
 - 1. Try FMSvm on more training problems and platforms.
 - 2. Improve FMSvm's performance and understand better its convergence behavior.

Applying SVMs to Financial Markets

- Identifying problems:
 - A mathematical model is available, but only an approximation.
 - Decisions are made based on experience.
- Constructing training sets:
 - Domain knowledges.
 - Data are from different categories.
 - Scaling the training data.
- Using SVMs to predict the movement of Dow Jones Industrials Average Index (DJIA):
 - The geometric Brownian motion solution.
 - The SVM solution.
 - The comparison between the two solutions.
- Conclusion and future work.

- The weak form of the market-efficiency hypothesis.
- The equations:

$$dS = \mu S dt + \sigma S dz.$$

By Itô's lemma, we have

$$d\ln S = (\mu - \frac{\sigma^2}{2})dt + \sigma dz.$$

That is,

$$\ln \frac{S_T}{S_0} \sim N((\mu - \frac{\sigma^2}{2})T, \sigma^2 T).$$

 We want to predict on Thursday (after the market is closed) whether or not the closing index of DJIA on Friday is higher than its opening index on Monday.

The Dow Jones Industrial Average Index



• The geometric Brownian motion model:

$$\frac{dS}{S} = \mu dt + \sigma dz.$$

• Its discrete-time version:

$$\frac{\Delta S}{S} = \mu \Delta t + \sigma \epsilon \sqrt{\Delta t}.$$
$$\frac{\Delta S}{S} \sim N(\mu \Delta t, \sigma^2 \Delta t).$$

• Stock price on Friday: $(R = \frac{S_{\rm Fr} - S_{\rm Th}}{S_{\rm Th}} \sim N(\mu, \sigma^2))$

$$S_{\mathsf{Fr}} = (1+R)S_{\mathsf{Th}}$$
$$= (1+\mu+\sigma\epsilon)S_{\mathsf{Th}}$$

• We want to know $S_{\mathsf{Fr}} \geq S_{\mathsf{Mo}}$.

$$\epsilon \stackrel{>}{<} \frac{S_{\rm Mo}/S_{\rm Th} - 1 - \mu}{\sigma} = \xi.$$

• Since $\begin{cases} \Pr(\epsilon > \xi) > \Pr(\epsilon < \xi) \text{ if } \xi < 0\\ \Pr(\epsilon > \xi) < \Pr(\epsilon < \xi) \text{ if } \xi > 0, \end{cases}$ $h_{\text{GBM}}(\xi) = \text{sgn}(-\xi).$

• Estimating μ and σ : only μ matters.



- The test set: December 7, 1990 to December 28, 2001.
- The performance: 96 errors out 511 trials. The accuracy of h_{GBM} on this test set is 81.21%.

The SVM Solution

- $h_{SVM}(x) = sgn(w^T x + b)$, where "experience" is incorporated in w and b.
- Our first training set: July 5, 1954 to November 30,1990.
 x = [341.3 343.6 339.4 340.9 23300 ... 340.4 343.0 339.2 341.1 30000]^T
- Scaling the raw training vectors:

 $m{x}_s = [1\ 1.007\ 0.994\ 0.999\ 1\ \dots\ 0.997\ 1.005\ 0.994\ 0.999\ 1.288]^T$

• Results: (There are totally 1704 training vectors, and the test set contains 511 vectors.)

poly(4)							
С	0.2	0.5	0.8	1	1.2	1.4	1.6
Errors	106	99	96	96	96	98	97
rbf(1)							
С	20	40	60	80	100	120	140
Errors	109	101	101	98	96	96	99

The SVM Solution

- Our intuition is that if there is any available information correlated to the movement of DJIA, then adding it to the training set may improve the performance.
- We add the daily federal funds rate to the training set. We choose this rate simply because it is a daily rate.
- Again, we use the relative values with respect to the Monday rate in the training vectors, instead of real rates.
- Results:

poly(3)								
С	0.05	1	1.5	2.5	10	20	30	35
Errors	103	102	96	96	94	93	89	89

• When c = 30, it takes around 2.6 hours (on newton.mit.edu) to train the SVM comparing with the scale of seconds when using the old training set.

- The Geometric Brownian Motion model is only an approximation.
- The semi-strong form of the market-efficiency hypothesis.



Conclusions and Future Work

- Conclusions:
 - 1. SVMs can learn well from time series data.
 - 2. There are potentials to achieve better performance using SVMs.
- Future work:
 - 1. We want to find good applications of SVMs (need domain knowledge).
 - 2. How to construct a training set based on data from different categories is also an interesting problem.