CSCE 990 Lecture 8: Implementation* Stephen D. Scott	Introduction • We know that the convex quadratic program representing our SVM optimization problem has a unique global optimum
February 28, 2006	 How do we efficiently find it? Unlike "classical" optimization problems, we sometimes must deal with extremely large, dense matrices
*Most figures ©2002 MIT Press, Bernhard Schölkopf, and Alex Smola.	2
Outline	 Tricks of the Trade Stopping Criteria Recall that one property of the dual is that its optimum equals the primal's optimum
Outline • Tricks of the trade: stopping criteria, restart- ing, caching, shrinking the training set • Sparse greedy matrix approximation (SGMA)	 Tricks of the Trade Stopping Criteria Recall that one property of the dual is that its optimum equals the primal's optimum <u>KKT-Gap</u> is difference between primal and dual objective functions at a feasible solution
Outline • Tricks of the trade: stopping criteria, restart- ing, caching, shrinking the training set • Sparse greedy matrix approximation (SGMA) • Interior point algorithms	 Tricks of the Trade Stopping Criteria Recall that one property of the dual is that its optimum equals the primal's optimum <u>KKT-Gap</u> is difference between primal and dual objective functions at a feasible solution Theorem 6.27 can bound KKT-Gap for SVMs P10.1 Let <i>f</i> be feasible soln to problem of minimizing regularized risk functional <i>B</i>ma[f]. Then
Outline • Tricks of the trade: stopping criteria, restart- ing, caching, shrinking the training set • Sparse greedy matrix approximation (SGMA) • Interior point algorithms • Subset selection methods: Chunking, working set algorithms	Tricks of the Trade Stopping Criteria • Recall that one property of the dual is that its optimum equals the primal's optimum • <u>KKT-Gap</u> is difference between primal and dual objective functions at a feasible solution • Theorem 6.27 can bound KKT-Gap for SVMs P10.1 Let <i>f</i> be feasible soln to problem of minimizing regularized risk functional $R_{reg}[f]$. Then $R_{reg}[f] \ge R_{reg}[f^*] \ge R_{reg}[f] - Gap[f]/(Cm)$ where f^* is optimal feasible solution and \underline{m}
<section-header><list-item><list-item><list-item><list-item><list-item><list-item><list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></section-header>	Tricks of the Trade Stopping Criteria • Recall that one property of the dual is that its optimum equals the primal's optimum • <u>KKT-Gap</u> is difference between primal and dual objective functions at a feasible solution • Theorem 6.27 can bound KKT-Gap for SVMs P10.1 Let <i>f</i> be feasible soln to problem of minimizing regularized risk functional $R_{reg}[f]$. Then $R_{reg}[f] \ge R_{reg}[f^*] \ge R_{reg}[f] - Gap[f]/(Cm)$ where f^* is optimal feasible solution and $Gap[f] = \sum_{j=1}^{m} C \max\{0, 1-y_jf(x_j)\} + \alpha_j(y_jf(x_j)-1)$ for <i>C</i> -SVM and $Gap[f] = \sum_{j=1}^{m} \max\{0, \rho - y_jf(x_j)\} + \alpha_j(y_jf(x_j)-\rho)$ for ν -SVM

Tricks of the Trade Restarting with Different Parameters Stopping Criteria (cont'd) Recall the C-SV classifier: • Can halt if gap is relatively smaller than ϵ : $\underset{\mathbf{w}\in\mathcal{H},b\in\mathbb{R},\boldsymbol{\xi}\in\mathbb{R}^m}{\text{minimize}} \quad \tau(\mathbf{w},\boldsymbol{\xi}) = \frac{1}{2} \|\mathbf{w}\|^2 + \frac{C}{m} \sum_{i=1}^m \xi_i$ $\operatorname{Gap}[f] \le \epsilon \left(\frac{|R_{\operatorname{reg}}| + |R_{\operatorname{reg}}[f] - \operatorname{Gap}[f]|}{2} \right)$ $y_i(\langle \mathbf{x}_i, \mathbf{w} \rangle + b) \ge 1 - \xi_i, i = 1, \dots, m$ $\xi_i \ge 0, i = 1, \dots, m$ i.e. the gap is small compared to its "midpoint" • This can be thought of as minimizing the regularized risk functional • Alternatively, can see that if f_i is solution at $R_{\mathsf{reg}}[f,C] := C \sum_{i=1}^{m} c(x_i, y_i, f(x_i)) + \Omega[f]$ iteration *i*, $\min_{i} \{R_{\mathsf{reg}}[f_i]\} \ge R_{\mathsf{reg}}[f^*] \ge \max_{i} \{R_{\mathsf{reg}}[f_i] - \mathsf{Gap}[f_i]\}$ where $c(\cdot)$ measures the average margin error of the training set which is useful since Gap can increase when f_i improves • How do we set C? 5 6 **Restarting with Different Parameters Restarting with Different Parameters** (cont'd) (cont'd) • Finally, • If f_C minimzes $R_{reg}[f, C]$, then for all C' > C

 $R_{\mathsf{reg}}[f_C, C'] \ge R_{\mathsf{reg}}[f_{C'}, C'] \ge R_{\mathsf{reg}}[f_{C'}, C] \ge R_{\mathsf{reg}}[f_C, C]$

• Thus

$$R_{\text{reg}}[f_{C'}, C'] \leq R_{\text{reg}}[f_C, C']$$

$$= C' \sum_{i=1}^m c(x_i, y_i, f_C(x_i)) + \Omega[f_C]$$

$$= \frac{C'}{C} \left(C \sum_{i=1}^m c(x_i, y_i, f_C(x_i)) + \Omega[f_C] \right)$$

$$= \left(\frac{C'}{C} \right) R_{\text{reg}}[f_C, C]$$

7

- $\left(\frac{C}{C'}\right) R_{\mathsf{reg}}[f_{C'}, C'] \le R_{\mathsf{reg}}[f_C, C] \le R_{\mathsf{reg}}[f_{C'}, C']$
- I.e. changes in $R_{\text{reg}}[f_C, C]$ are bounded by changes in C, so f_C is reasonable starting point for search for $f_{C'}$
- Thus can start with large C (i.e. focus on minimizing margin errors) and steadily decrease C to increase regularization
- By scaling parameters appropriately, can dramatically speed up training

Tricks of the Trade Caching

- If kernel matrix *K* is too large to store in memory, may need to store most on disk, caching a relatively small amount in memory
 - 1. Row cache: cache m_c rows, each with m entries, and replace with LRU. Works well with e.g. SMO
 - 2. <u>Element cache</u>: store individual elements of K. Works well if most $\alpha_i = 0$, but significant overhead involved
 - 3. <u>Function cache</u>: Cheap way to update $f(x_j)$ (prediction of classifier on x_j). If e.g. first n alphas are changed in current update:

$$f^{\text{new}}(x_j) = \sum_{i=1}^m \alpha_i^{\text{new}} k(x_i, x_j) + b$$

= $f^{\text{old}}(x_j) + \left[\sum_{i=1}^n \left(\alpha_i^{\text{new}} - \alpha_i^{\text{old}}\right) k(x_i, x_j)\right]$

Sparse Greedy Matrix Approximation

- Cost of computing/storing/using for optimization entire $m \times m$ Gram matrix K is $\Theta(m^2)$
- \bullet Problematic when m (size of training set) is large, e.g. 10^5
- SGMA builds a sparse approximation \tilde{K} of K
 - \tilde{K} still $m\times m,$ but represented by matrix $\alpha\in\mathbb{R}^{m\times n}$ for $n\ll m$
- Recall that for $x_i \in \mathcal{X}$, we can think of $k(x_i, \cdot) = k_i(\cdot)$ as a function that computes the dot product in feature space of $\Phi(x_i)$ and $\Phi(\cdot)$

- So
$$K_{ij} = k_i(x_j)$$

• SGMA approximates $k_i(\cdot)$ with

$$\tilde{k}_i(\cdot) := \sum_{j=1}^n \alpha_{ij} k(x_j, \cdot)$$

(w.l.o.g. assume that we use the first n training patterns x_1, \ldots, x_n in the approximation)

Tricks of the Trade

Shrinking the Training Set

- Recall that only x_i for which $\alpha_i > 0$ affect the solution
- Thus can speed up training by dropping non-SVs from training set
- Don't want to do this too early, but at various points during optimization, can discard parts of training set
- Will cover subset selection schemes later

10

Sparse Greedy Matrix Approximation (cont'd)

• Goodness of the approximation will be the squared norm between the functions in feature space:

$$\left|k_{i}(\cdot)-\tilde{k}_{i}(\cdot)\right|_{\mathcal{H}}^{2}=\left\langle k_{i}(\cdot)-\tilde{k}_{i}(\cdot),k_{i}(\cdot)-\tilde{k}_{i}(\cdot)\right\rangle_{\mathcal{H}}$$

• Let's hark back to Section 2.2.2, where if $f(\cdot) = \sum_{j=1}^{m} \alpha_j k(x_j, \cdot)$, then

$$\langle f, f \rangle_{\mathcal{H}} = \sum_{j,\ell=1}^{m} \alpha_j \alpha_\ell k(x_j, x_\ell)$$

• If we let $\alpha_0'=$ 1, $\alpha_j'=-\alpha_{ij},\; x_j'=x_j$ and $x_0'=x_i,$ then

$$f(\cdot) = k_i(\cdot) - \tilde{k}_i(\cdot) = k(x_i, \cdot) - \sum_{j=1}^n \alpha_{ij} k(x_j, \cdot)$$
$$= \sum_{j=0}^n \alpha'_j k(x'_j, \cdot)$$

and

$$\langle f, f \rangle_{\mathcal{H}} = K_{ii} - 2\sum_{j=1}^{n} \alpha_{ij} K_{ij} + \sum_{j,\ell=1}^{n} \alpha_{ij} \alpha_{i\ell} K_{j\ell}$$

Sparse Greedy Matrix Approximation (cont'd) • Given basis functions $k_i(\cdot)$, goal is to find $\alpha \in \mathbb{R}^{m \times n}$ to minimize $\operatorname{Err}(\alpha) := \sum_{i=1}^{m} \left\ k_i(\cdot) - \tilde{k}_i(\cdot) \right\ _{\mathcal{H}}^2$ $= \sum_{i=1}^{m} \left(K_{ii} - 2 \sum_{j=1}^{n} \alpha_{ij} K_{ij} + \sum_{j,\ell=1}^{n} \alpha_{ij} \alpha_{i\ell} K_{j\ell} \right)$ • Differentiating wrt α and setting to 0 yields $\alpha_{\text{opt}} = K^{mn} (K^{nn})^{-1}$ where K^{mn} is first n columns of K and K^{nn} is upper left submatrix of K T10.2 The approximation \tilde{K} of K is PD, $\operatorname{Err}(\alpha_{\text{opt}}) = \operatorname{tr}(K) - \operatorname{tr}(\tilde{K})$, and $\tilde{K} = \alpha K^{nn} \alpha^{T}$	Sparse Greedy Matrix Approximation (cont'd) In general, it is intractable to choose as a ba- sis the proper subset of the <i>m</i> functions that minimizes the objective function Alternative approach: given current set of <i>n</i> basis functions $k_1(\cdot), \ldots, k_n(\cdot)$, consider adding one more function - Random - Greedy: best out of all $m - n$ remaining - "Semi-greedy": best out of a subset of size N (e.g. $N = 59$) Given that $\alpha^{m,n}$ is the optimal matrix for the n chosen basis functions, and consider adding $k_{n+1}(\cdot)$. Then $\operatorname{Err}(\alpha^{m,n+1}) = \operatorname{Err}(\alpha^{m,n}) - \eta^{-1} K^{mn}\mathbf{v} - \bar{\mathbf{k}} ^2$ where $\bar{\mathbf{k}} = (K_{1,n+1}, \dots, K_{m,n+1}), \eta = (K_{n+1,n+1} - \mathbf{k}^{\top}(K^{nn})^{-1}\mathbf{k}), \mathbf{v} = ((K^{nn})^{-1}\mathbf{k}), \text{ and}$ $\mathbf{k} = (K_{n+1,1}, \dots, K_{n+1,n})$
 Thus given K and basis functions, can easily compute α_{opt}, Err(α_{opt}), and K How do we choose basis functions? 	 Use the above to compute Err for each of the N new candidates and take the best; repeat until Err sufficiently small Can update α in O(n²) time
Sparse Greedy Matrix Approximation Experiments • Gaussian kernel, USPS data • First ten basis functions correspond to 9 of the 10 digits 90176 83024 01650	Predictor-Corrector Methods • Considers lower-order approximations when do- ing optimization of complex functions • Predict update from lower-order approx, then correct it with higher-order version • E.g. solving $f(x) = f_0 + ax + bx^2/2 = 0$: 1. Start with $f_0 + ax = 0$, giving update of $x^{\text{pred}} = -f_0/a$ 2. Substitute x^{pred} into original equation and solve: $f_0 + ax^{\text{corr}} + \frac{b}{2} \left(\frac{f_0}{a}\right)^2 = 0$ so $x^{\text{corr}} = -\frac{f_0}{a} \left(1 + \frac{\tilde{b}f_0}{2a^2}\right)$

- Results (Fig. 10.3, p. 294) comparable to PCA, which requires computation of all *m* basis functions
- 3. Use x^{corr} as update
- Algorithm 6.5, p. 165

Interior Point Methods

- An <u>interior point</u> is a (x, α) pair satisfying both primal and dual constraints
- In general, start with a quadratic optimization problem:

 $\begin{array}{l} \underset{x}{\text{minimize}} \quad \frac{1}{2} x^{\top} K x + c^{\top} x \\ \text{s.t.} \quad A x + d + \xi = 0, \xi \geq 0 \end{array} (1) \\ \text{where } K \text{ is } m \times m \text{ PD matrix, } x, c \in \mathbb{R}^m, \ A \in \mathbb{R}^{n \times m}, \text{ and } d, \xi \in \mathbb{R}^n \end{array}$

• Apply Theorem 6.26 to get KKT conditions:

 $Kx + A^{\top}\alpha + c = 0$ $Ax + d + \xi = 0$ $\alpha^{\top}\xi = 0$ $\alpha, \xi \ge 0$

- First two constraints are linear, but third is not, so we'll use predictor-corrector method
- First change third constraint to $\alpha^{\top}\xi = \mu > 0$ and decrease μ over time

17

Interior Point Methods (cont'd)

- Need to ensure that updates are not too large and negative to make the variables negative
- Shrink length of $(\Delta x, \Delta \alpha, \Delta \xi)$ by λ such that

$$\min \left\{ \frac{\alpha_1 + \lambda \Delta \alpha_1}{\alpha_1}, \cdots, \frac{\alpha_n + \lambda \Delta \alpha_n}{\alpha_n}, \\ \frac{\xi_1 + \lambda \Delta \xi_1}{\xi_1}, \cdots, \frac{\xi_n + \lambda \Delta \xi_n}{\xi_n} \right\} \ge \epsilon$$

for e.g. $\epsilon=0.05$

• To update μ , after getting new values of α and ξ , set

$$\mu = \frac{\alpha^{\top} \xi}{n} \left(\frac{1 - \lambda}{10 + \lambda} \right)^2$$

Interior Point Methods

Linearization

- Start with initial values of x, $\alpha,$ $\xi,$ and μ
- We'll compute updates Δx, etc. by expanding e.g. x to x + Δx:

$$\begin{split} \mathbf{K} \Delta x + \mathbf{A}^{\top} \Delta \alpha &= -\mathbf{K} x - \mathbf{A}^{\top} \alpha - c &=: \ \rho_p \\ \mathbf{A} \Delta x + \Delta \xi &= -\mathbf{A} x - d - \xi &=: \ \rho_d \\ \alpha_i^{-1} \xi_i \Delta \alpha_i + \Delta \xi_i &= \mu \alpha_i^{-1} - \xi_i - \alpha_i^{-1} \Delta \alpha_i \Delta \xi_i &=: \ \rho_{\mathsf{K}\mathsf{K}\mathsf{T}_i} \end{split}$$

• Thus we get $\Delta \xi_i = \rho_{\mathsf{KKT}_i} - \xi_i \Delta \alpha_i / \alpha_i$, and $A\Delta x - \xi \Delta \alpha / \alpha = \rho_d - \rho_{\mathsf{KKT}}$:

 $\begin{bmatrix} K & A^{\mathsf{T}} \\ A & -D \end{bmatrix} \begin{bmatrix} \Delta x \\ \Delta \alpha \end{bmatrix} = \begin{bmatrix} \rho_p \\ \rho_d - \rho_{\mathsf{KKT}} \end{bmatrix}$ where $D := \operatorname{diag}(\alpha_1^{-1}\xi_1, \dots, \alpha_n^{-1}\xi_n)$

- We've eliminated $\Delta \xi$, so we can solve for $\Delta x^{\rm pred}$ and $\Delta \alpha^{\rm pred}$
- Now update ρ_{KKT} with $\Delta \alpha^{\text{pred}}$ (ρ_p, ρ_d are unchanged) and solve again to get Δx^{corr} and $\Delta \alpha^{\text{corr}}$, then solve for $\Delta \xi$

18

Interior Point Methods Application to SVMs

 Recall the dual optimization problem for C-SVMs:

$$\begin{array}{ll} \underset{\boldsymbol{\alpha} \in \mathbb{R}^m}{\text{maximize}} & W(\boldsymbol{\alpha}) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j k(x_i, x_j) \\ \text{s.t.} & 0 \leq \alpha_i \leq C, \ i = 1, \dots, m \\ & \sum_{i=1}^m \alpha_i y_i = 0 \end{array}$$

• Can put into the form of (1):

$$\begin{array}{l} \underset{\alpha,t\in\mathbb{R}^m}{\text{minimize}} \quad \frac{1}{2}\alpha^{\top}Q\alpha + c^{\top}\alpha\\ \text{s.t.} \qquad A\alpha = 0 \qquad (2)\\ \alpha + t = u\\ \alpha, t \ge 0 \end{array}$$

where $Q_{ij} = y_i y_j k(x_i, x_j), c = (-1, ..., -1) \in \mathbb{R}^m$, $u = (C, ..., C) \in \mathbb{R}^m$, $A = (y_1, ..., y_m)$

- Then dualize, linearize, and derive updates (Section 10.3.1)
- For large m, can use SGMA in optimization algorithm rather than all of Q

Subset Selection Methods

Working Set Algorithms

- A way of dealing with large data sets
- Focus on only a subset of the training patterns at any time, freezing the α variables for the other patterns
- Let $S_w \subset [m] = \{1, \dots, m\}$ be working set and $S_f = [m] \setminus S_w$ be fixed set

• Now split the problem (2) into
$$Q = \begin{bmatrix} Q_{ww} & Q_{fw} \\ Q_{wf} & Q_{ff} \end{bmatrix}$$
,
 $c = (c_w, c_f), A = (A_w, A_f), \text{ and } u = (u_w, u_f)$:
minimize $\frac{1}{2} \alpha_w^\top Q_{ww} \alpha_w + [c_w + Q_{wf} \alpha_f]^\top \alpha_w$
 $+ [\frac{1}{2} \alpha_f^\top Q_{ff} \alpha_f + c_f^\top \alpha_f]$
s.t. $A_w \alpha_w = -A_f \alpha_f$
 $\alpha_w + t_w = u_w$
 $\alpha_w, t_w \ge 0$

• Minimizing this also decreases (2)

21

Sequential Minimal Optimization

- Extreme case of subset selection, with working set of size 2
- With only two active variables α_i and α_j , can analytically solve optimization problem

 $\begin{array}{ll} \underset{\alpha_i,\alpha_j}{\text{minimize}} & \frac{1}{2} \left[\alpha_i^2 Q_{ii} + \alpha_j^2 Q_{jj} + 2\alpha_i \alpha_j Q_{ij} \right] + c_i \alpha_i + c_j \alpha_j \\ \text{s.t.} & s \alpha_i + \alpha_j = \gamma \\ & 0 \le \alpha_i \le C_i, 0 \le \alpha_j \le C_j \end{array}$

where $Q_{ij} = y_i y_j K_{ij}$, $s = y_i y_j$, $\gamma = y_i y_j \alpha_i^{\text{old}} + \alpha_j^{\text{old}}$, $c_i = y_i (f(x_i) - b - y_i) - \alpha_i K_{ii} - \alpha_j s K_{ij}$, C_i, C_j parameters weighting x_i, x_j

• Above values come from working set version of (2)

Subset Selection Methods (cont'd)

- When choosing the working set, want to base choice on what will speed up convergence
- Pick patterns whose Lagrange multipliers violate KKT conditions
- Want generally small working set (< 100), and balanced number of +1 and -1 labels
- In addition, can choose:
 - those with largest contribution to KKT gap (P10.1)
 - 2. those with largest negative gradient of objective function at current solution

22

Sequential Minimal Optimization (cont'd)

- Let $\xi = sc_j c_i + \gamma sQ_{jj} \gamma Q_{ij}$ and $\chi = Q_{ii} + Q_{jj} 2sQ_{ij}$
- If $y_i = y_j$, let $L = \max\{0, \gamma C_j\}$ and $H = \min\{C_i, \gamma\}$
- If $y_i \neq y_j$, let $L = \max\{0, \gamma\}$ and $H = \min\{C_i, C_j \gamma\}$
 - 1. If $\chi = 0$, set $\alpha_i = L$ if $\xi > 0$ and $\alpha_i = H$ otherwise
 - 2. If $\chi > 0$, set $\alpha_i = \min\{\max\{L, \xi/\chi\}, H\}$
 - 3. Set $\alpha_j = \gamma s\alpha_i$

Sequential Minimal Optimization Updating b (cont'd)

Sequential Minimal Optimization

Updating b

- \bullet When choosing patterns to bring in, need to know our prediction on them, i.e. we need a current value of b
- When all α s optimal, can apply KKT conditions to sove for b; choose some $\alpha_i \in (0, C_i)$ and solve: $y_i(\langle \mathbf{w}, \Phi(x) \rangle + b) = 1$
 - But the α s aren't yet optimal!
- Thus we will estimate *b* by choosing the midpoint of a range of possible values

• Partition set of training patterns X into
$$\begin{split} I_0 &= \{i \mid \alpha_i \in (0, C_i)\} \\ I_{+,0} &= \{i \mid \alpha_i = 0, y_i = +1\} \\ I_{+,C} &= \{i \mid \alpha_i = C_i, y_i = +1\} \end{split}$$

 $I_{-,0} = \{i \mid \alpha_i = 0, y_i = -1\}$ $I_{-,C} = \{i \mid \alpha_i = C_i, y_i = -1\}$

• Now define

$$e_{\mathsf{hi}} := \min_{i \in I_0 \cup I_{+,0} \cup I_{-,C}} \{f(x_i) - y_i\}$$

$$e_{\mathsf{lo}} := \min_{i \in I_0 \cup I_{-,0} \cup I_{+,C}} \{f(x_i) - y_i\}$$

(f is based on setting b with a SV)

- Using KKT conditions (see Keerthi et al.), can show that optimality occurs iff $e_{hi} \ge 0 \ge e_{lo}$
- Further, using as a bias term $b_{\rm hi} = b e_{\rm hi}$ for the "hi" sets and $b_{\rm lo} = b e_{\rm lo}$ for the "lo" sets will yield optimality
- Thus can update b as $(b_{hi} + b_{lo})/2$
- Complete SMO pseudocode on p. 313

26

Topic summary due in 1 week!

25