

Bundle Methods for Machine Learning (Teo, Vishwanathan, Smola, Le) JMLR 2010, NIPS 2007

Presented by Kevin Duh
Bayes Reading Group 6/4/2010

1

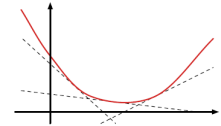
1-slide summary

- Many machine learning methods involve solving a minimum regularized risk objective

$$\min_w J(w) := \lambda \Omega(w) + R_{\text{emp}}(w),$$

$$\text{where } R_{\text{emp}}(w) := \frac{1}{m} \sum_{i=1}^m l(x_i, y_i, w)$$

- Cutting-plane algorithm & Bundle methods solve it iteratively by using a piece-wise lower bound



2

Why I chose this paper

- These optimization methods (invented in 1960s) are becoming popular in supervised learning
 - Very fast
 - Scale to large datasets
- Handles non-smooth convex optimization, so widely applicable
 - Can be used when LBFGS fails

3

Outline

- Background
- Cutting plane algorithm
- Bundle method
- Different loss functions

4

Warm-up

- Convex Set:
 - a set is a convex set if it contains the line segment joining any of its points

$$x, y \in S; a, b \geq 0; a + b = 1$$

$$\Rightarrow ax + by \in S$$
 - are these sets convex?



5

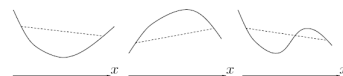
Background: Convex functions

- Convex function:
 - A function is convex if its domain is a convex set and the segment joining any two points on f do not have values lower than f

$$\forall x, y \in \text{dom}(f); a, b \geq 0, a + b = 1$$

$$af(x) + bf(y) \geq f(ax + by)$$

- What's convex, what's concave?



6

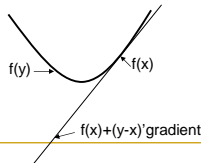
Convex & differentiable functions

- Gradient exists if f is differentiable

$$\nabla f(x) = \left(\frac{\partial f(x)}{\partial x_1}, \frac{\partial f(x)}{\partial x_2}, \dots, \frac{\partial f(x)}{\partial x_n} \right)$$

- 1st-order condition: a differentiable f is convex iff:

$$f(y) \geq f(x) + \langle y - x, \nabla f \rangle \quad \forall y$$



Gradient provides a global lower bound to f

7

Nonsmooth (non-differentiable) functions

- What if f is not differentiable, e.g.

- L1-regularizer or $|x|$
- Envelope function: $f(x) = \max_{s \in S} \phi_s(x)$

Derivative at $x > 1$ is 1
Derivative at $x < 1$ is -1
Derivative at $x = 0$?



- Subgradient: a vector s is a subgradient if

$$f(y) \geq f(x) + \langle y - x, s \rangle \quad \forall y$$

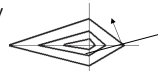
- There may exist many subgradients at a point
 - The set of subgradients is called the subdifferential
 - The methods we deal with only require one subgradient

8

Subgradient Method for optimizing non-smooth functions

- Similar as gradient descent, except:

- works on non-differentiable functions
- step-lengths not chosen by line-search, but fixed
- it is not a descent method
 - descent direction d can only be defined by $\langle d, s \rangle < 0$ for all s in sub-differential



- Pseudo-code:

- Repeat until convergence:
 - $s =$ subgradient at $f(x)$
 - $x = x - \text{stepsize} * s$
 - keep track of best x so far

9

Outline

- Background
- Cutting plane algorithm
- Bundle method
- Different loss functions

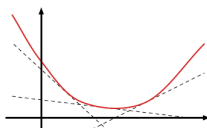
10

Cutting-plane algorithm for optimizing non-smooth functions

- Recall subgradient forms a lowerbound on f

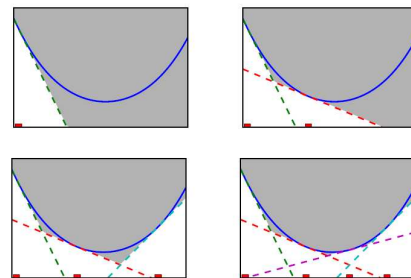
- Main Idea:

- If we have multiple subgradients at different points, we get a tighter lowerbound
- The lowerbound improves with each iteration, so minimizing the lowerbound eventually minimizes the desired objective



11

Figure: lower bound improves after each iteration



12

The math

- Overall optimization goal:

$$\min_w J(w) := \lambda\Omega(w) + R_{\text{emp}}(w),$$

where $R_{\text{emp}}(w) := \frac{1}{m} \sum_{i=1}^m l(x_i, y_i, w)$

- Lower bound:

- Given sequence of iterates w and subgradients s , the (piecewise-linear) lower bound is:

$$J(w) \geq J_t^{\text{CP}}(w) := \max_{1 \leq t \leq t} \{J(w_{t-1}) + \langle w - w_{t-1}, s_t \rangle\}$$

- Because $J(w) \geq J(w') + \langle w - w', s' \rangle$
- At each iteration, compute $w_t := \underset{w}{\text{argmin}} J_t^{\text{CP}}(w)$.

Note: Slight change of notation starting now: function $f(x) \rightarrow J(w)$

13

Cutting-plane pseudocode

1. Compute $J(w_i)$ and its subgradient s_t
2. Compute error $\epsilon_t := \min_{0 \leq i \leq t} J(w_i) - J_t^{\text{CP}}(w_t)$
3. If error < threshold, stop
4. Update bound $J_t^{\text{CP}}(w) := \max_{1 \leq i \leq t} \{J(w_{i-1}) + \langle w - w_{i-1}, s_i \rangle\}$
5. Optimize it to get new iterate $w_t := \underset{w}{\text{argmin}} J_t^{\text{CP}}(w)$
6. Goto step 1

14

Why does this work?

Let w^* be optimal solution, then

$$J(w_i) \geq J(w^*) \Rightarrow \min_{0 \leq i \leq t} J(w_i) \geq J(w^*)$$

By construction,

$$J(w) \geq J_t^{\text{CP}}(w), \forall w \Rightarrow J(w^*) \geq J_t^{\text{CP}}(w_t)$$

So optimal point is sandwiched:

$$\min_{0 \leq i \leq t} J(w_i) \geq J(w^*) \geq J_t^{\text{CP}}(w_t)$$

This error is monotonically decreases.

When it reaches zero, we have w^*

$$\epsilon_t := \min_{0 \leq i \leq t} J(w_i) - J_t^{\text{CP}}(w_t)$$

15

Final word on cutting plane-algorithm

- It has nice stopping criteria
 - (better than subgradient method)
- Cost is solving linear programs per iteration:
 - Size of this subproblem grows with each iteration
 - But usually this can be solved quickly
- Speed depends critically on the set of cutting planes
 - Zig-zag behavior possible, slowing down convergence

16

Outline

1. Background
2. Cutting plane algorithm
3. Bundle method
4. Different loss functions

17

Standard Bundle Method

- Zig-zag in cutting-plane is caused by taking large steps and neglecting previous solutions
- Bundle methods extend cutting-plane by ensuring new iterate is not too far

$$w_t := \underset{w}{\text{argmin}} \left\{ \frac{\alpha}{2} \|w - \hat{w}_{t-1}\|^2 + J_t^{\text{CP}}(w) \right\}$$

18

(Standard) Bundle method pseudo-code

Algorithm 1 Proximal Bundle Method

```

1: input & initialization:  $\epsilon \geq 0, \rho \in (0, 1), w_1, t \leftarrow 0, \hat{w}_0 \leftarrow w_0$ 
2: loop
3:    $t \leftarrow t + 1$ 
4:   Compute  $J(w_{t-1})$  and  $s_t \in \partial_w J(w_{t-1})$ 
5:   Update model  $J_t^{\text{CP}}(w) := \max_{1 \leq i \leq t} \{J(w_{i-1}) + \langle w - w_{i-1}, s_i \rangle\}$ 
6:    $\hat{w}_t \leftarrow \operatorname{argmin}_w J_t^{\text{CP}}(w) + \frac{\epsilon}{2} \|w - \hat{w}_{t-1}\|^2$ 
7:    $\epsilon_t \leftarrow J(\hat{w}_{t-1}) - [J_t^{\text{CP}}(\hat{w}_t) + \frac{\epsilon}{2} \|\hat{w}_t - \hat{w}_{t-1}\|^2]$ 
8:   if  $\epsilon_t < \epsilon$  then return  $\hat{w}_t$ 
9:   Linesearch:  $\eta_t \leftarrow \operatorname{argmin}_{\eta \in \mathbb{R}} J(\hat{w}_{t-1} + \eta(\hat{w}_t - \hat{w}_{t-1}))$  (if expensive, set  $\eta_t = 1$ )
10:   $\hat{w}_t \leftarrow \hat{w}_{t-1} + \eta_t(\hat{w}_t - \hat{w}_{t-1})$ 
11:  if  $J(\hat{w}_{t-1}) - J(\hat{w}_t) \geq \rho \epsilon_t$  then
12:    SERIOUS STEP:  $\hat{w}_t \leftarrow \hat{w}_t$ 
13:  else
14:    NULL STEP:  $\hat{w}_t \leftarrow \hat{w}_{t-1}$ 
15:  end if
16: end loop

```

This paper argues that some parameters are hard to tune:
 ζ_t, ρ

19

Proposed Bundle Method (BMRM)

- Regularized risk minimization objective already has a regularization term:

$$\min_w J(w) := \lambda \Omega(w) + R_{\text{emp}}(w),$$
- So optimize this subproblem: $J_t(w) := \lambda \Omega(w) + R_t^{\text{CP}}(w)$
 $w_t := \min_w J_t(w)$
- No need for serious/null step

20

Algorithm 2 BMRM

```

1: input & initialization:  $\epsilon \geq 0, w_0, t \leftarrow 0$ 
2: repeat
3:    $t \leftarrow t + 1$ 
4:   Compute  $a_t \in \partial_w R_{\text{emp}}(w_{t-1})$  and  $b_t \leftarrow R_{\text{emp}}(w_{t-1}) - \langle w_{t-1}, a_t \rangle$ 
5:   Update model:  $R_t^{\text{CP}}(w) := \max_{1 \leq i \leq t} \{w, a_i\} + b_i$ 
6:    $w_t \leftarrow \operatorname{argmin}_w J_t(w) := \lambda \Omega(w) + R_t^{\text{CP}}(w)$ 
7:    $\epsilon_t \leftarrow \min_{0 \leq i \leq t} J(w_i) - J_t(w_t)$ 
8:   until  $\epsilon_t \leq \epsilon$ 
9: return  $w_t$ 

```

(subgradient of R_{emp}) $a_t \in \partial_w R_{\text{emp}}(w_{t-1})$
 (offset) $b_t := R_{\text{emp}}(w_{t-1}) - \langle w_{t-1}, a_t \rangle$
 (piecewise linear lower bound of R_{emp}) $R_t^{\text{CP}}(w) := \max_{1 \leq i \leq t} \{w, a_i\} + b_i$
 (piecewise convex lower bound of J) $J_t(w) := \lambda \Omega(w) + R_t^{\text{CP}}(w)$
 (iterate) $w_t := \min_w J_t(w)$
 (approximation gap) $\epsilon_t := \min_{0 \leq i \leq t} J(w_i) - J_t(w_t)$

21

In more detail: how to solve subproblem in step 6

- Reformulate as constrained optimization:

$$w_t = \operatorname{argmin}_w J_t(w) := \lambda \Omega(w) + \max_{1 \leq i \leq t} \{w, a_i\} + b_i$$

$$\min_{w, \xi} \lambda \Omega(w) + \xi$$

subject to $\xi \geq \langle w, a_i \rangle + b_i$ for $i = 1, \dots, t$

- Then call linear/quadratic program depending on regularizer

□ # constraints = #iterations, unrelated to #samples!

□ Dual program for L2 regularizer:

$$\alpha_t = \operatorname{argmax}_{\alpha \in \mathbb{R}^t} \{-\frac{1}{2\lambda} \alpha^\top A^\top A \alpha + \alpha^\top b \mid \alpha \geq 0, \|\alpha\|_1 = 1\}$$

22

Convergence Analysis

Theorem 4 Assume that $\max_{w \in \partial_w R_{\text{emp}}(w)} \|w\| \leq G$ for all $w \in \operatorname{dom} J$. Also assume that Ω^* has bounded curvature, i.e., $\|\partial_w^2 \Omega^*(\mu)\| \leq H^*$ for all $\mu \in \{-\lambda^{-1} \sum_{i=1}^{t+1} \alpha_i a_i \mid \alpha_i \geq 0, \forall i \text{ and } \sum_{i=1}^{t+1} \alpha_i = 1\}$. In this case we have

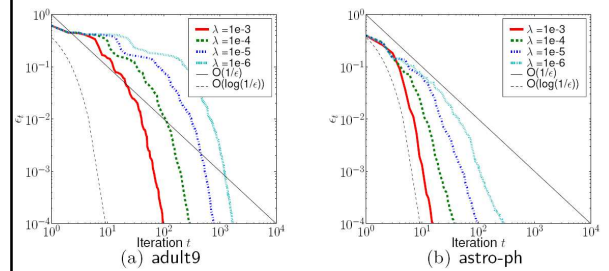
$$\epsilon_t - \epsilon_{t+1} \geq \frac{\epsilon_t}{2} \min(1, \lambda \epsilon_t / 4G^2 H^*). \quad (27)$$

Furthermore, if $\|\partial_w^2 J(w)\| \leq H$, then we have **Every iteration the error is halved!**

$$\epsilon_t - \epsilon_{t+1} \geq \begin{cases} \epsilon_t / 2 & \text{if } \epsilon_t \geq 4G^2 H^* / \lambda \\ \lambda / 8H^* & \text{if } 4G^2 H^* / \lambda \geq \epsilon_t \geq H/2 \\ \lambda \epsilon_t / 4HH^* & \text{otherwise} \end{cases} \quad (28)$$

23

Experiments in actual speed



24

Outline

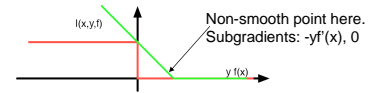
1. Background
2. Cutting plane algorithm
3. Bundle method
4. Different loss functions

25

Binary classification

- Accuracy-based loss: $\Delta(y, y') = \begin{cases} 0 & \text{if } y = y' \\ 1 & \text{otherwise} \end{cases}$

- Convex upper-bounds:
 - Soft margin loss: $l(x, y, f) = \max(0, 1 - yf(x))$



- Logistic: $\log(1 + \exp(-yf(x)))$
- MCE (Katagiri et. al.) – sigmoid, with adjustable parameter
- Gaussian process classifier, MAP solution:
 - Minimize $\frac{1}{2} \mathbf{f}^T \mathbf{K}^{-1} \mathbf{f} - \sum_{i=1}^n \log p(y_i | f_i)$.

26

Structured Prediction

- Similar to previous slide, convex upper bounds for structured loss $\Delta(y, y')$

- CRF: $l(x, y, w) = \log \sum_{y' \in \mathcal{Y}} \exp(\langle w, \phi(x, y') \rangle) - \langle w, \phi(x, y) \rangle$

$$\partial_w l(x, y, w) = \mathbf{E}_{y' \sim p(y'|x)} [\phi(x, y')] - \phi(x, y).$$

- Structured SVM:

$$l(x, y, w) = \max_{y' \in \mathcal{Y}} \Gamma(y, y') \langle w, \phi(x, y') - \phi(x, y) \rangle + \Delta(y, y')$$

$$\partial_w l(x, y, w) = \Gamma(y, \tilde{y}(x)) [\phi(x, \tilde{y}(x)) - \phi(x, y)]$$

$$\tilde{y}(x) := \operatorname{argmax}_{y'} \Gamma(y, y') \langle w, \phi(x, y') - \phi(x, y) \rangle + \Delta(y, y')$$

For bundle methods, just collect the vectors $\partial_w l(x, y, w)$ and give to the LP/QP

27

ROC score

- AUC is not continuous in w :

$$AUC(x, y, w) = \frac{1}{m_+ m_-} \sum_{y_i < y_j} \mathbf{I}(\langle w, x_i \rangle < \langle w, x_j \rangle),$$

- But this nonsmooth convex bound is:

$$R_{\text{emp}}(w) = \frac{1}{m_+ m_-} \sum_{y_i < y_j} \max(0, 1 + \langle w, x_i - x_j \rangle)$$

- We can directly calculate subgradients in closed-form, but we can also obtain from an algorithm if it's more efficient

- See algorithm 7 in JMLR paper

28

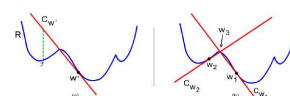
Do you see the pattern?

1. Give me any problem, with any evaluation metric (may be difficult to optimize)
2. Think of a convex upper bound for the metric
 - This bound does not need to be smooth
 - Just need to get subgradients from it
3. Solve with:
 - Gradient descent, LFBGS; or
 - Subgradient method, Bundle method, etc.
4. Done: submit NIPS paper

29

Discussions

- Fear not non-smooth convex functions
- What about non-convex optimization?
 - EM-style training where M is solved by bundle
 - Yu & Joachims, Learning Structural SVMs w/ Latent Variables (ICML09)
 - Modified bundle method:
 - Do & Artieres, Large margin training of HMMs w/ partially-observed states (ICML09)



30