Sparse Convex Optimization Methods for Machine Learning

PhD Defense Talk 2011 / 10 / 04 Martin Jaggi

Examiner: Emo Welzl Co-Examiners: Bernd Gärtner, Elad Hazan, Joachim Giesen, Joachim Buhmann



 $D \subset \mathbb{R}^n$









The Linearized Problem

$$\min_{\boldsymbol{y}\in D} f(\boldsymbol{x}) + \langle \boldsymbol{y} - \boldsymbol{x}, d_{\boldsymbol{x}} \rangle$$

Algorithm 1 Greedy on a Compact Convex Set

Pick an arbitrary starting point
$$x^{(0)} \in D$$

for $k = 0 \dots \infty$ do
Let $d_x \in \partial f(x^{(k)})$ be a subgradient to f at $x^{(k)}$
Compute $s := approx \arg\min \langle y, d_x \rangle$
Let $\alpha := \frac{2}{k+2}$
Update $x^{(k+1)} := x^{(k)} + \alpha(s - x^{(k)})$
end for

Theorem: Algorithm obtains **accuracy** $O(\frac{1}{k})$ after k steps.

 ${}^{i}\mathcal{X}$

 $D \subset \mathbb{R}^n$

f(

 $|\mathcal{X}|$



History & Related Work

	Domain	Known Stepsize	Approx. Subproblem	Primal-Dual Guarantee
Frank & Wolfe 1956	linear inequality constraints	×	×	×
Dunn 1978, 1980	general bounded convex domain	×	\checkmark	×
Zhang 2003	convex hulls	×	\checkmark	×
Clarkson 2008, 2010	unit simplex	\checkmark	×	\checkmark
Hazan 2008	semidefinite matrices of bounded trace	\checkmark	\checkmark	\checkmark
J. PhD Thesis	general bounded convex domain	\checkmark	\checkmark	\checkmark

Sparse Approximation





for $k = 0...\infty$ do Let $d_x \in \partial f(x^{(k)})$ be a subgradient to f at $x^{(k)}$ Compute $i := \arg\min_i (d_x)_i$ Let $\alpha := \frac{2}{k+2}$ Update $x^{(k+1)} := x^{(k)} + \alpha(\mathbf{e}_i - x^{(k)})$ end for

Corollary:

Algorithm gives an \mathcal{E} -approximate solution of **sparsity** $O\left(\frac{1}{\varepsilon}\right)$.

[Clarkson SODA '08]



Sparsity as a function of the approximation quality

lower bound: $\Omega(\frac{1}{2})$



Smallest enclosing ball Linear Classifiers _INEAR Classifiers (such as Support Vector Machines, l_2 -loss) 61603453 • Model Predictive Control 2 403/20212 Б 232086 Mean Variance Portfolio Optimization 81511251 $x^T C x - t \cdot b^T x$ $x \in \Delta_n$

Sparse Approximation

$$D := conv(\{\pm \mathbf{e}_i \mid i \in [n]\})$$

for $k = 0...\infty$ do Let $d_x \in \partial f(x^{(k)})$ be a subgradient to f at $x^{(k)}$ Compute $i := \arg \max_i |(d_x)_i|$, and let $s := \mathbf{e}_i \cdot \operatorname{sign} ((-d_x)_i)$ Let $\alpha := \frac{2}{k+2}$ Update $x^{(k+1)} := x^{(k)} + \alpha(s - x^{(k)})$ end for

Corollary: Algorithm gives an \mathcal{E} -approximate solution of **sparsity** $O\left(\frac{1}{\varepsilon}\right)$.



Sparsity as a function of the approximation quality

lower bound: $\Omega(\frac{1}{2})$

 $\min_{\|x\|_1 \le 1} f(x)$

 ℓ_1 -ball



• ℓ_1 -regularized regression

 $\min_{\|x\|_1 \le t} \|Ax - b\|_2^2$

Sparse Recovery

Low Rank Approximation

 $D := conv(\{vv^T \mid v \in \mathbb{R}^n, \\ \|v\|_2 = 1\})$

 $\min_{x \in D} f(x)$

 $X \in Sym^{n \times n}$ $X \succeq 0$ Tr(X) = 1

spectahedron

for $k = 0...\infty$ do Let $D_X \in \partial f(X^{(k)})$ be a subgradient to f at $X^{(k)}$ Let $\alpha := \frac{2}{k+2}$ Compute $v := v^{(k)} = \text{ApproxEV}(D_X, \alpha C_f)$ Update $X^{(k+1)} := X^{(k)} + \alpha(vv^T - X^{(k)})$ end for

Corollary:

Algorithm gives an \mathcal{E} -approximate solution of **rank** $O\left(\frac{1}{2}\right)$.

[Hazan LATIN '08]

lower bound:

 $\Omega(\frac{1}{\varepsilon})$



Trace norm regularized problems

 $\min_{\|X\|_* \le t} f(X)$

Low-Rank Matrix Recovery

Max norm regularized problems

Matrix Factorizations for recommender systems



The Netflix challenge: 17k Movies 500k Customers 100M Observed Entries $(\approx 1\%)$

=: X

s eduivalent to: $\min_{X \succeq 0} f(X)$ s.t. Tr(X) = t

 $\min_{U,V} \sum_{(i,j)\in\Omega} (Y_{ij} - (UV^T)_{ij})^2$ s.t. $||U||_{Fro}^2 + ||V||_{Fro}^2 = t$

A Simple Alternative Optimization Duality

The Problem

 $\min_{x \in D} f(x)$

The Dual

 $\omega(x) := \min_{\substack{\boldsymbol{y} \in D}} f(x) + \langle \boldsymbol{y} - x, d_x \rangle$

Weak Duality

 $\omega(x) \le f(x^*) \le f(x')$



Pathwise Optimization

The Parameterized Problem

 $\min_{x \in D} f_t(x)$

"Better than necessary"

 $g_t(x) \le \frac{\varepsilon}{2}$

The difference

$$g_{t'}(x) - g_t(x) \le \varepsilon \left(1 - \frac{1}{2}\right)$$

"Still good enough"

$$g_{t'}(x) \le \varepsilon$$

 $f_{t'}(x)$ 200 $\omega_{t'}$ \boldsymbol{x} 100 0 t+' $\Leftrightarrow |t' - t| \le \varepsilon \cdot P_f$

"Continuity in the parameter"

Theorem: There are $O\left(\frac{1}{\varepsilon}\right)$ many intervals of piecewise constant ε -approx. solutions.

[Giesen, J, Laue ESA 2010]

Applications

 Smallest enclosing ball of moving points

- SVMs, MKL (with 2 base kernels) $\min_{x \in \Delta_n} x^T (K+t\mathbf{1}) x$
- Model Predictive Control
- Mean Variance Portfolio Optimization

 $\min_{x \in \Delta_n} x^T C x - t \cdot b^T x$

est accuracy 4032021 2 [4] 23210860 81011251 robust PCA

Recommender Systems

Thanks

co-authors:

Bernd Gärtner Joachim Giesen Soeren Laue Marek Sulovský

3D visualization:

Robert Carnecky

