Approximate SDP solvers, Matrix Factorizations, the Netflix Prize, and PageRank

Mittagseminar Martin Jaggi, Oct 6 2009

Sparse Approximation

f() convex

The Problem

 $\min f(x)$ $x \in \mathbb{R}^n$ $x \ge 0$ $\mathbf{1}^T x = 1$

vectors living in the simplex

 $\min f(X)$ $X \in \mathbb{S}^{n \times n}$ $X \succeq 0$ Tr(X) = 1

symetric matrices living in the spectahedron

The Problem	$\min f(x)$	$\min f(X)$	
	$x \in \mathbb{R}^n$	$X \in \mathbb{S}^{n \times n}$	
	$x \ge 0$	$X \succeq 0$	
	$1^T x = 1$	Tr(X) = 1	

The Algorithm

$$x^{(k+1)} := (1-\lambda)x^{(k)} + \lambda e_i$$

 $i := \arg \max - \nabla f(\mathbf{x}^{(k)})_i$

"Coordinate Descent"

 $\mathbf{X}^{(k+1)} := (1-\lambda) \mathbf{X}^{(k)} + \lambda \mathbf{v} \mathbf{v}^T$

 $\boldsymbol{v} := \underset{\|v\|=1}{\operatorname{arg\,max}} v^T \left(-\nabla f(\boldsymbol{X^{(k)}}) \right) v$

 $v^{(1)}$

 $v^{(k)}$

,largest' Eigenvector

 $X^{(k)}$

Sparsity = k Rank = k

 $\lambda = 1/k$

 $= UU^T$

$$\begin{array}{c|c} x^{(k+1)} := (1-\lambda) x^{(k)} + \lambda e_i \\ \text{The Algorithm} \\ i := \arg\max_i -\nabla f(x^{(k)})_i \\ \hline v := \arg\max_i v^T \left(-\nabla f(x^{(k)}) v \right) \\ \|v\| = 1 \end{array}$$

The Convergence

After $O\left(\frac{1}{\epsilon}\right)$ steps the primal-dual error is $\leq \epsilon$.

[Clarkson SODA '08]

After $O\left(\frac{1}{\epsilon}\right)$ steps the primal-dual error is $\leq \epsilon$.

[Hazan LATIN '08]

Approximate Eigenvector computation

Instead of

$$v := rgmax_{\|v\|=1} v^T M v$$

 $M := -\nabla f(\mathbf{X}^{(k)})$

it is enough to work with

$$v: v^T M v \ge \lambda_{\max} - \epsilon^2 \qquad \|v\| =$$

Such a v can be found by doing $O\left(\frac{1}{\epsilon}\right)$ Lanzcos steps. Alternative: Power method



How to solve general Semidefinite Programs?

Optimization Version:

 $\min \ Tr(CX)$ $Tr(A_iX) \le b_i$ $_{i \in [m]}$ $X \in \mathbb{S}^{n \times n}$ $X \succeq 0$

Feasibility Version:

Find X s.t. $Tr(A_iX) \leq b_i$ $i \in [m+1]$ $X \in \mathbb{S}^{n \times n}$ $X \succeq 0$ Tr(X) = 1

 $\min f(X)$

 $\overline{X} \in \mathbb{S}^{n \times n}$

Tr(X) = 1

4/10

 $X \succeq 0$

$$f(X) := \frac{1}{M} \log \left(\sum_{i=1}^{m+1} e^{M(Tr(A_i X) - b_i X)} \right)$$

Soft Max"

By this trick, Hazan's algorithm is able to satisfy all constraints up to an error $\leq\epsilon$.

Matrix Factorizations

and machine learning



Matrix Factorizations for recommender systems



[Short IEEE article, Wikipedia: Netflix Prize]



Figure 3. The first two vectors from a matrix decomposition of the Netflix Prize data. Selected movies are placed at the appropriate spot based on their factor vectors in two dimensions. The plot reveals distinct genres, including clusters of movies with strong female leads, fraternity humor, and quirky independent films.

Matrix Factorizations

and machine learning

Applications:

Customer i ↔ Product j

(Amazon, Netflix, Migros Cumulus etc...)

(Symmetry?, k=?)

Word i ↔ Document

(Search engines, Latent Semantic Analysis)

many other applications

 (e.g. dimensionality reduction, clustering)



Accuracy vs Model complexity





[Srebro NIPS '05]

$$f(X) := \sum_{ij \in S} \left((X - Y)_{ij} \right)^2$$

Low Norm Matrix Factorization

$$\begin{split} \min_{U,V} & f(UV^T) \\ \text{s.t.} & \|U\|_{\text{Fro}}^2 + \|V\|_{\text{Fro}}^2 \\ & = Tr(UU^T) + Tr(VV^T) \\ & = Tr(Z) \end{split}$$

is equivalent to

 $\min_{Z} f(Z)$

 $Z \in \mathbb{S}^{(n+m) \times (n+m)}$ $Z \succeq 0$ Tr(Z) = t

Perfectly fits for Hazan's Algorithm.

 V^T)

 $f(X) := \sum (X - Y)_{ij}^2$

m

 $ij \in S$

n

 $= \begin{pmatrix} U \\ V \end{pmatrix}$

Z

Low Norm Matrix Factorization

Have to be careful, principal EV is not always the largest EV! Add a constant to the diagonal in that case.

We need the largest Eigenvector of $M := -\nabla f(\mathbf{Z}^{(k)})$

Largest Eigenvector of the bipartite weighted graph with adjacency matrix M.





- Hazan's new approximate SDP solver applies to Low Norm Matrix Factorization
- Easy to parallelize (Power method)
- Algorithm maintains sparsity structure of the given matrix, needs no additional memory
- Speed is comparable to existing methods, and much better than generic SDP solvers

Thanks