Contents lists available at ScienceDirect

Information Processing Letters

www.elsevier.com/locate/ipl

A note on sparse least-squares regression

Christos Boutsidis^{a,*}, Malik Magdon-Ismail^b

^a Mathematical Sciences Department, IBM T.J. Watson Research Center, United States ^b Computer Science Department, Rensselaer Polytechnic Institute, United States

ARTICLE INFO

Article history: Received 29 May 2013 Received in revised form 2 September 2013 Accepted 15 November 2013 Available online 25 December 2013 Communicated by X. Wu

Keywords: Algorithms Least squares Regression Sparse approximation Sparsification Regularization Truncated SVD

1. Introduction

Fix inputs $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $\mathbf{b} \in \mathbb{R}^m$. We study leastsquares regression: $\min_{\mathbf{x} \in \mathbb{R}^n} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2$. It is well known that the minimum norm solution vector can be found using the pseudo-inverse of \mathbf{A} : $\mathbf{x}^* = \mathbf{A}^{\dagger}\mathbf{b} = (\mathbf{A}^T\mathbf{A})^{-1}\mathbf{A}^T\mathbf{b}$. When \mathbf{A} is ill-conditioned, \mathbf{A}^{\dagger} becomes unstable to perturbations and overfitting can become a serious problem. For example, when the smallest non-zero singular value of \mathbf{A} is close to zero, the largest singular value of \mathbf{A}^{\dagger} can be extremely large and the solution vector $\mathbf{x}^* = \mathbf{A}^{\dagger}\mathbf{b}$ obtained via a numerical algorithm is not the optimal, due to numerical instability issues. Practitioners deal with such situations using *regularization*.

Popular regularization techniques are the Lasso [8], the Tikhonov regularization [4], and the truncated SVD [6]. The lasso minimizes $\|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2 + \lambda \|\mathbf{x}\|_1$, and Tikhonov regularization minimizes $\|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2 + \lambda \|\mathbf{x}\|_2^2$ (in both cases $\lambda > 0$ is the regularization parameter). The truncated SVD

* Corresponding author. E-mail addresses: cboutsi@us.ibm.com (C. Boutsidis), magdon@cs.rpi.edu (M. Magdon-Ismail).

ABSTRACT

We compute a *sparse* solution to the classical least-squares problem $\min_{\mathbf{x}} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2$, where **A** is an arbitrary matrix. We describe a novel algorithm for this sparse least-squares problem. The algorithm operates as follows: first, it selects columns from **A**, and then solves a least-squares problem only with the selected columns. The column selection algorithm that we use is known to perform well for the well studied column subset selection problem. The contribution of this article is to show that it gives favorable results for sparse least-squares as well. Specifically, we prove that the solution vector obtained by our algorithm is close to the solution vector obtained via what is known as the "SVD-truncated regularization approach".

© 2013 Elsevier B.V. All rights reserved.

1.1. Preliminaries The compact (or thin) Singular Value Decomposition (SVD) of a matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ of rank ρ is

minimizes $\|\mathbf{A}_k \mathbf{x} - \mathbf{b}\|_2$, where $k < \operatorname{rank}(\mathbf{A})$ is a rank parameter and $\mathbf{A}_k \in \mathbb{R}^{m \times n}$ is the best rank-*k* approximation

to A obtained via the SVD. So, the truncated SVD solu-

tion is $\mathbf{x}_{\nu}^{*} = \mathbf{A}_{\nu}^{\dagger} \mathbf{b}$. Notice that these regularization methods

impose parsimony on x in different ways. A combinato-

rial approach to regularization is to explicitly impose the

sparsity constraint on x, requiring it to have few non-zero

elements. We give a new deterministic algorithm which, for r = O(k), computes an $\hat{\mathbf{x}}_r \in \mathbb{R}^n$ with at most r non-zero

$$\mathbf{A} = \underbrace{(\mathbf{U}_{k}, \mathbf{U}_{\rho-k})}_{\mathbf{U}_{\mathbf{A}} \in \mathbb{R}^{m \times \rho}} \underbrace{\begin{pmatrix} \boldsymbol{\Sigma}_{k} & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\Sigma}_{\rho-k} \end{pmatrix}}_{\boldsymbol{\Sigma}_{\mathbf{A}} \in \mathbb{R}^{\rho \times \rho}} \underbrace{\begin{pmatrix} \mathbf{V}_{k}^{\mathrm{T}} \\ \mathbf{V}_{\rho-k}^{\mathrm{T}} \\ \mathbf{v}_{\mathbf{A}}^{\mathrm{T}} \in \mathbb{R}^{\rho \times n} \end{pmatrix}}_{\mathbf{V}_{\mathbf{A}}^{\mathrm{T}} \in \mathbb{R}^{\rho \times n}}.$$

entries such that $\|\mathbf{A}\hat{\mathbf{x}}_r - \mathbf{b}\|_2 \approx \|\mathbf{A}\mathbf{x}_k^* - \mathbf{b}\|_2$.

Here, $\mathbf{U}_k \in \mathbb{R}^{m \times k}$ and $\mathbf{U}_{\rho-k} \in \mathbb{R}^{m \times (\rho-k)}$ contain the left singular vectors of **A**. Similarly, $\mathbf{V}_k \in \mathbb{R}^{n \times k}$ and $\mathbf{V}_{\rho-k} \in$







^{0020-0190/\$ –} see front matter @ 2013 Elsevier B.V. All rights reserved. http://dx.doi.org/10.1016/j.ipl.2013.11.011

Algorithm 1: Deterministic s	sparse	rearession
------------------------------	--------	------------

- 1: **Input:** $\mathbf{A} \in \mathbb{R}^{m \times n}$, $\mathbf{b} \in \mathbb{R}^{m}$, target rank $k < \operatorname{rank}(\mathbf{A})$, and parameter $0 < \varepsilon < 1/2$.
- 2: Obtain $\mathbf{V}_k \in \mathbb{R}^{n \times k}$ from the SVD of **A** and compute $\mathbf{E} = \mathbf{A} \mathbf{A} \mathbf{V}_k \mathbf{V}_k^{\mathrm{T}} \in \mathbb{R}^{m \times n}$.
- 3: Set $\mathbf{C} = \mathbf{A}\mathbf{\Omega}\mathbf{S} \in \mathbb{R}^{m \times r}$, with $r = \lceil \frac{9k}{\varepsilon^2} \rceil$ and
- $[\boldsymbol{\Omega}, \mathbf{S}] = \mathsf{DeterministicSampling}(\mathbf{V}_k^{\mathrm{T}}, \mathbf{E}, r),$
- 4: Set $\mathbf{x}_r = \mathbf{C}^{\dagger}\mathbf{b} \in \mathbb{R}^r$, and $\hat{\mathbf{x}}_r = \mathbf{\Omega}^{\mathbf{S}}\mathbf{x}_r \in \mathbb{R}^n$ ($\hat{\mathbf{x}}_r$ has at most r nonzeros at the indices of the selected columns in \mathbf{C}).
- 5: **Return** $\hat{\mathbf{x}}_r \in \mathbb{R}^n$.

 $\mathbb{R}^{n \times (\rho - k)}$ contain the right singular vectors. The singular values of **A**, which we denote as $\sigma_1(\mathbf{A}) \ge \sigma_2(\mathbf{A}) \ge$ $\cdots \ge \sigma_{\rho}(\mathbf{A}) > 0$ are contained in $\boldsymbol{\Sigma}_k \in \mathbb{R}^{k \times k}$ and $\boldsymbol{\Sigma}_{\rho-k} \in$ $\mathbb{R}^{(\rho-k)\times(\rho-k)}$. We use $\mathbf{A}^{\dagger} = \mathbf{V}_{\mathbf{A}}\boldsymbol{\Sigma}_{\mathbf{A}}^{-1}\mathbf{U}_{\mathbf{A}}^{\mathrm{T}} \in \mathbb{R}^{n \times m}$ to denote the Moore–Penrose pseudo-inverse of **A** with $\Sigma_{\mathbf{A}}^{-1}$ denoting the inverse of $\boldsymbol{\Sigma}_{\mathbf{A}}$. Let $\mathbf{A}_{k} = \mathbf{U}_{k}\boldsymbol{\Sigma}_{k}\mathbf{V}_{k}^{\mathrm{T}} \in \mathbb{R}^{m \times n}$ and $\mathbf{A}_{\rho-k} = \mathbf{A} - \mathbf{A}_{k} = \mathbf{U}_{\rho-k}\boldsymbol{\Sigma}_{\rho-k}\mathbf{V}_{\rho-k}^{\mathrm{T}} \in \mathbb{R}^{m \times n}$. For $k < \operatorname{rank}(\mathbf{A})$, the SVD gives the best rank \dot{k} approximation to **A** in both the spectral and the Frobenius norm: for $\tilde{\mathbf{A}} \in \mathbb{R}^{m \times n}$, let rank($\tilde{\mathbf{A}}$) $\leq k$; then, for $\xi = 2$, F, $\|\mathbf{A} - \mathbf{A}_k\|_{\xi} \leq \|\mathbf{A} - \tilde{\mathbf{A}}\|_{\xi}$. Also, $\|\mathbf{A} - \mathbf{A}_k\|_2 = \|\boldsymbol{\Sigma}_{\rho-k}\|_2 = \sigma_{k+1}(\mathbf{A})$, and $\|\mathbf{A} - \mathbf{A}_k\|_F^2 =$ $\|\boldsymbol{\Sigma}_{\rho-k}\|_{F}^{2} = \sum_{i=k+1}^{\rho} \sigma_{i}^{2}(\mathbf{A}).$ The Frobenius and the spectral norm of **A** are defined as: $\|\mathbf{A}\|_{F}^{2} = \sum_{i,j} \mathbf{A}_{ij}^{2} = \sum_{i=1}^{\rho} \sigma_{i}^{2}(\mathbf{A});$ and $\|\mathbf{A}\|_2 = \sigma_1(\mathbf{A})$. Let **X** and **Y** be matrices of appropriate dimensions; then, $\|\mathbf{X}\mathbf{Y}\|_{F} \leq \min\{\|\mathbf{X}\|_{F}\|\mathbf{Y}\|_{2}, \|\mathbf{X}\|_{2}\|\mathbf{Y}\|_{F}\}$. This is a stronger version of the standard submultiplicativity property $\|\mathbf{X}\mathbf{Y}\|_{F} \leq \|\mathbf{X}\|_{F} \|\mathbf{Y}\|_{F}$, which we will refer to as "spectral submultiplicativity".

Given $k < \rho = \text{rank}(\mathbf{A})$, the truncated rank-k SVD regularized weights are

$$\mathbf{x}_k^* = \mathbf{A}_k^{\dagger} \mathbf{b} = \mathbf{V}_k \boldsymbol{\Sigma}_k^{-1} \mathbf{U}_k^{\mathrm{T}} \mathbf{b} \in \mathbb{R}^n,$$

and note that $\|\mathbf{b} - \mathbf{A}_k \mathbf{A}_k^{\dagger} \mathbf{b}\|_2 = \|\mathbf{b} - \mathbf{U}_k \mathbf{U}_k^{\mathrm{T}} \mathbf{b}\|_2$.

Finally, for r < n, let $\Omega = [\mathbf{z}_{i_1}, ..., \mathbf{z}_{i_r}] \in \mathbb{R}^{n \times r}$ where $\mathbf{z}_i \in \mathbb{R}^m$ are standard basis vectors; Ω is a *sampling matrix* because $\mathbf{A}\Omega \in \mathbb{R}^{m \times r}$ is a matrix whose columns are sampled (with possible repetition) from the columns of **A**. Let $\mathbf{S} \in \mathbb{R}^{r \times r}$ be a diagonal *rescaling matrix* with positive entries; then, we define the sampled and rescaled columns from **A** by $\mathbf{C} = \mathbf{A}\Omega\mathbf{S}$: Ω samples some columns from **A** and then **S** rescales them.

2. Results

Our sparse solver to minimize $\|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2$ takes as input the sparsity parameter r (i.e., the solution vector \mathbf{x} is allowed at most r non-zero entries), and selects r rescaled columns from \mathbf{A} (denoted by \mathbf{C}). We then solve the leastsquares problem to minimize $\|\mathbf{C}\mathbf{x} - \mathbf{b}\|_2$. The result is a dense vector $\mathbf{C}^{\dagger}\mathbf{b}$ with r dimensions. The sparse solution $\hat{\mathbf{x}}_r$ will be zero at indices corresponding to columns not selected in \mathbf{C} , and we use $\mathbf{C}^{\dagger}\mathbf{b}$ to compute the other entries of $\hat{\mathbf{x}}_r$.

Theorem 1. Let $\mathbf{A} \in \mathbb{R}^{m \times n}$, $\mathbf{b} \in \mathbb{R}^m$, rank $k < \operatorname{rank}(\mathbf{A})$, and $0 < \varepsilon < 1/2$. Algorithm 1 runs in time $O(\operatorname{mn}\min\{m, n\} + nk^3/\varepsilon^2)$ and returns $\hat{\mathbf{x}}_r \in \mathbb{R}^n$ with at most $r = \lceil 9k/\varepsilon^2 \rceil$ non-zero entries such that:

Algorithm 2: Deterministic Sampling (110111

- 1: Input: $\mathbf{V}^{\mathrm{T}} = [\mathbf{v}_1, ..., \mathbf{v}_n] \in \mathbb{R}^{k \times n}$; $\mathbf{E} = [\mathbf{e}_1, ..., \mathbf{e}_n] \in \mathbb{R}^{m \times n}$; and
- *r* > *k*. 2: **Output:** Sampling and rescaling matrices $\boldsymbol{\Omega} \in \mathbb{R}^{n \times r}$, $\mathbf{S} \in \mathbb{R}^{r \times r}$.
- 3: Initialize $\mathbf{B}_0 = \mathbf{0}_{k \times k}$, $\boldsymbol{\Omega} = \mathbf{0}_{n \times r}$, $\mathbf{S} = \mathbf{0}_{r \times r}$.
- 4: for $\tau = 0$ to r 1 do
- 5: Set $L_{\tau} = \tau \sqrt{rk}$.
- 6: Pick index $i \in \{1, 2, ..., n\}$ and t such that
- $U(\mathbf{e}_i) \leqslant \frac{1}{t} \leqslant L(\mathbf{v}_i, \mathbf{B}_{\tau}, \mathbf{L}_{\tau}).$
- 7: Update $\mathbf{B}_{\tau+1} = \mathbf{B}_{\tau} + t\mathbf{v}_i\mathbf{v}_i^{\mathrm{T}}$. Set $\boldsymbol{\Omega}_{i,\tau+1} = 1$ and $\mathbf{S}_{\tau+1,\tau+1} = 1/\sqrt{t}$.
- 8: end for
- 9: **Return:** $\boldsymbol{\Omega} \in \mathbb{R}^{n \times r}$, $\mathbf{S} \in \mathbb{R}^{r \times r}$.

$$\|\mathbf{A}\hat{\mathbf{x}}_{r} - \mathbf{b}\|_{2} \leq \|\mathbf{A}\mathbf{x}_{k}^{*} - \mathbf{b}\|_{2} + (1 + \varepsilon) \cdot \|\mathbf{b}\|_{2} \cdot \frac{\|\mathbf{A} - \mathbf{A}_{k}\|_{F}}{\sigma_{k}(\mathbf{A})}$$

This upper bound is "small" when **A** is "effectively" low-rank, i.e., $\|\mathbf{A} - \mathbf{A}_k\|_F / \sigma_k(\mathbf{A}) \ll 1$. Also, a trivial bound is $\|\mathbf{A}\hat{\mathbf{x}}_r - \mathbf{b}\|_2 \leq \|\mathbf{b}\|_2$ (error when $\hat{\mathbf{x}}_r$ is the all-zeros vector), because $\|\mathbf{CC}^{\mathsf{T}}\mathbf{b} - \mathbf{b}\|_2 \leq \|\mathbf{C0}_{r \times 1} - \mathbf{b}\|_2 = \|\mathbf{b}\|_2$.

In the heart of Algorithm 1 lies a method for selecting columns from **A** (Algorithm 2), which was originally developed in [1] for column subset selection, where one selects columns **C** from **A** to minimize $||\mathbf{A} - \mathbf{CC}^{\dagger}\mathbf{A}||_{F}$. Here, we adopt the same algorithm for least-squares.

The main tool used to prove Theorem 1 is a new "structural" result that may be of independent interest.

Lemma 2. Fix $\mathbf{A} \in \mathbb{R}^{m \times n}$, $\mathbf{b} \in \mathbb{R}^n$, rank $k < \operatorname{rank}(\mathbf{A})$, and sparsity r > k. Let $\mathbf{x}_k^* = \mathbf{A}_k^{\dagger} \mathbf{b} \in \mathbb{R}^n$, where $\mathbf{A}_k \in \mathbb{R}^{m \times n}$ is the rank-k SVD approximation to \mathbf{A} . Let $\boldsymbol{\Omega} \in \mathbb{R}^{n \times r}$ and $\mathbf{S} \in \mathbb{R}^{r \times r}$ be any sampling and rescaling matrices with $\operatorname{rank}(\mathbf{V}_k^T \boldsymbol{\Omega} \mathbf{S}) = k$. Let $\mathbf{C} = \mathbf{A} \boldsymbol{\Omega} \mathbf{S} \in \mathbb{R}^{m \times r}$ be a matrix of sampled rescaled columns of \mathbf{A} and let $\hat{\mathbf{x}}_r = \boldsymbol{\Omega} \mathbf{S} \mathbf{C}^{\dagger} \mathbf{b} \in \mathbb{R}^n$ (having at most r non-zeros). Then,

$$\begin{split} \|\mathbf{A}\hat{\mathbf{x}}_{r} - \mathbf{b}\|_{2} &\leq \left\|\mathbf{A}\mathbf{x}_{k}^{*} - \mathbf{b}\right\|_{2} \\ &+ \left\|(\mathbf{A} - \mathbf{A}_{k})\boldsymbol{\varOmega}\mathbf{S}\left(\mathbf{V}_{k}^{\mathrm{T}}\boldsymbol{\varOmega}\mathbf{S}\right)^{\dagger}\boldsymbol{\varSigma}_{k}\mathbf{U}_{k}^{\mathrm{T}}\mathbf{b}\right\|_{2}. \end{split}$$

The lemma says that if the sampling matrix satisfies a simple rank condition, then solving the regression on the sampled columns gives a sparse solution to the original problem with a performance guarantee.

2.1. Algorithm description

Algorithm 1 selects *r* columns from **A** to form **C** and the corresponding sparse vector $\hat{\mathbf{x}}_r$. The core of Algorithm 1 is the subroutine DeterministicSampling, which is a method to simultaneously sample the columns of two matrices, while controlling their spectral and Frobenius norms. DeterministicSampling takes inputs $\mathbf{V}^T \in \mathbb{R}^{k \times n}$ and $\mathbf{E} \in \mathbb{R}^{m \times n}$; the matrix **V** is orthonormal, $\mathbf{V}^T \mathbf{V} = \mathbf{I}_k$. (In our application, $\mathbf{V}^T = \mathbf{V}_k^T$ and $\mathbf{E} = \mathbf{A} - \mathbf{A}_k$.) We view \mathbf{V}^T and $\mathbf{E} = [\mathbf{e}_1, \dots, \mathbf{e}_n]$.

Given *k* and *r* and the iterator $\tau = 0, 1, 2, ..., r - 1$, define $L_{\tau} = \tau - \sqrt{rk}$. For a symmetric matrix $\mathbf{B} \in \mathbb{R}^{k \times k}$ with eigenvalues $\lambda_1, ..., \lambda_k$ and $L \in \mathbb{R}$, define functions

$$\phi(\mathbf{L}, \mathbf{B}) = \sum_{i=1}^{k} \frac{1}{\lambda_i - \mathbf{L}},$$

and

Download English Version:

https://daneshyari.com/en/article/428542

Download Persian Version:

https://daneshyari.com/article/428542

Daneshyari.com