CS838 Topics In Optimization Convex Geometry in High-Dimensional Data Analysis

Ben Recht Spring 2010

Logistics

- Class Tu-Th 1-2:15PM
- Office Hours CS4387, Tuesday 2:15-3:30
- Course webpage: <u>http://pages.cs.wisc.edu/~brecht/</u> <u>cs838.html</u>
- Readings will be posted here.

- **Scribing.** All students are required to scribe notes for at least one lecture. LaTeX template will be provided.
- **Project.** All students are required to prepare a 20-30 minute presentation on the themes of this course. This can be a literature review or an application of the course's techniques to your research.

Recommender Systems

More Top Picks for You







amazon.com

Because you enjoyed:

2001: A Space Odyssey Blue Velvet Bottle Rocket

We think you'll enjoy: Stalker

Add



S Not Interested







Netflix

• Rate some movies...





• Get some recommendations:







Netflix Prize

• One million big ones!



Data Sleuths in an Internet Age





For Today's Graduate, Just One Word: Statistics

By STEVE LOHR Published: August 5, 2009

MOUNTAIN VIEW, Calif. — At Harvard, Carrie Grimes majored in anthropology and archaeology and ventured to places like Honduras, where she studied Mayan settlement patterns by mapping where artifacts were found. But she was drawn to what she calls "all the computer and math stuff" that was part of the job.

Ø	SION IN TO RECOMMEND
13	TWITTER
Ψ	COMMENTS (50)
	E-MAIL
в	SEND TO PHONE

Netflix Prize

• One million big ones!



• Given 100 million ratings on a scale of 1 to 5, predict 3 million ratings to highest accuracy



- 17770 total movies x 480189 total users
- Over 8 billion total ratings
- How to fill in the blanks?

Abstract Setup: Matrix Completion



X_{ij} known for black cells X_{ij} unknown for white cells *Rows index movies Columns index users*

• How do you fill in the missing data?





- $k = Number of movies = 2 \times 10^4$
- $n = Number of users = 5 \times 10^5$
- $m = Number of Given Ratings = 10^8$
- kn ≈ 10¹⁰
- For r < 200, r(k+n) < 10⁸

Matrix Rank



• The rank of **X** is...

the dimension of the span of the rows the dimension of the span of the columns the smallest number r such that there exists an k x r matrix L and an n x r matrix R with X=LR*



or Singular Values

Affine Rank Minimization

 PROBLEM: Find the matrix of lowest rank that satisfies/approximates the underdetermined linear system

 $\mathcal{A}(\mathbf{X}) = \mathbf{b} \qquad \mathcal{A}: \mathbb{R}^{k \times n} \to \mathbb{R}^m$

 $\begin{array}{ll} \text{minimize} & \operatorname{rank}(\mathbf{X}) \\ \text{subject to} & \mathcal{A}(\mathbf{X}) = \mathbf{b} \end{array}$

• NP-HARD:

- Reduce to finding solutions to polynomial systems
- Hard to approximate
- Exact algorithms are awful

Heuristic: Gradient Descent

$$\mathcal{F}(\mathbf{L},\mathbf{R}) = \sum_{i=1}^{k} \sum_{k=1}^{r} L_{ik}^{2} + \sum_{j=1}^{n} \sum_{k=1}^{r} R_{jk}^{2} + \lambda \sum_{i,j} \left(\sum_{k} L_{ik} R_{jk} - X_{ij} \right)^{2}$$

- Just run gradient descent to minimize $\ensuremath{\mathcal{F}}$
- λ determines tradeoff between satisfying constraints and the size of the factors

Netflix Prize

Leaderboard

Mixture of hundreds of models, including gradient descent

2

53

Gradient descent							
on low-rank							
parameterization							

Team Name No Grand Prize candidates yet	Best Scor	provement	Last Submit Time
No Progress Prize candidates yet			
When Gravity and Dinosaurs Unite BellKor	0.8675 0.8682 0.8708	8.82 8.75 8.47	2008-03-01 07:03:35 2008-02-28 23:40:45 2008-02-06 14:12:44
18.2007 - RHSE = 0.8712			
KorBell acmehill Dan Tillberg basho Just a guy in a garage BigChaos Dinosaur Planet	0.8712 0.8720 0.8727 0.8729 0.8740 0.8748 0.8753	8.43 8.35 8.27 8.25 8.14 8.05 8.00	2007-10-01 23:25:23 2008-03-02 05:08:12 2008-03-02 08:42:29 2007-11-24 14:27:00 2008-02-06 12:16:40 2008-03-01 17:26:06 2007-10-04 04:56:45
amgl Remco molg JustWithSVD	0.8897 0.8899 0.8900 0.8900 0.8900 0.8900 0.8901	6.49 6.46 6.45 6.45 6.45 6.44	2007-12-23 18:44:03 2007-04-04 06:16:56 2007-12-23 18:54:46 2008-02-14 16:17:54 2008-02-28 09:56:20 2008-02-29 05:53:11
5020 The Clawn			

Complex Systems



© 60005 4 23% @NASDAQ:AAPL +1.36% @NASDAQ:MSFT +1.54





Predictions



Structure



Smoothness





Modeling Simplicity: Strategy

• Find a "natural" *convex* heuristic



 Use probabilistic analysis to prove the heuristic succeeds

 Provide efficient algorithms for solving the heuristic

Topics

- Sparsity
- Rank
- Smoothness

Themes

- Random Projections Preserve Geometry (encoding)
- Atomic Norms Recover Geometry (*decoding*)

Parsimonious Models



- Search for best linear combination of fewest atoms
- "rank" = fewest atoms needed to describe the model

• "natural" heuristic is the *atomic norm*:

$$||x||_{\mathcal{A}} \equiv \inf\left\{\sum_{k=1}^{r} |w_k| : x = \sum_{k=1}^{r} w_k \alpha_k\right\}$$

Mining for Biomarkers



Topic 1: Cardinality/Sparsity

 Vector x has cardinality s if it has at most s nonzeros. (x is *s-sparse*)

$$x = \sum_{k=1}^{s} w_k e_{i_k}$$

- Atoms are a discrete set of orthogonal points
- Typical Atoms:
 - standard basis
 - Fourier basis
 - Wavelet basis

Cardinality Minimization

• **PROBLEM:** Find the vector of lowest cardinality that satisfies/approximates the underdetermined linear system Ax = b $A: \mathbb{R}^n \to \mathbb{R}^m$

• NP-HARD:

- Reduce to EXACT-COVER [Natarajan 1995]
- Hard to approximate
- Known exact algorithms require enumeration

Proposed Heuristic

Cardinality Minimization:

 $\begin{array}{ll}\text{minimize} & \operatorname{card}(x)\\ \text{subject to} & Ax = b \end{array}$

Convex Relaxation:

minimize $||x||_1 = \sum_{i=1}^{D} |x_i|$ subject to Ax = b

- Long history (back to geophysics in the 70s)
- Flurry of recent work characterizing success of this heuristic: Candès, Donoho, Romberg, Tao, Tropp, etc., etc...
- "Compressed Sensing"

Compressed Sensing

- Model: most of the energy is at low frequencies
- Basis for JPG compression
- Use the fact that the image is sparse in DCT/wavelet basis to reduce number of measurements required for signal acquisition.
- decode using *I*₁ minimization

Integer Programming

 $\begin{array}{ll} \text{minimize} & \|x\|_{\infty} = \max_i |x_i| \\ \text{subject to} & Ax = b \end{array}$

Cardinality/Sparsity

- How many samples are required to reconstruct sparse vectors?
- Relationship to coding theory
- When can we guarantee the l1 heuristic works?
- What are efficient ways to compute minimum 11 norm solutions?

Topic 2: (Matrix) Rank

 Matrix X has rank r if it has at most r nonzero singular values.

$$X = \sum_{j=1}^r \sigma_j u_j v_j^* = \sum_{j=1}^r \sigma_j A_j$$

- Atoms are the set of all rank one matrices
- Not a discrete set

Singular Value Decomposition (SVD)

If X is a matrix of size k x n (k≤m) then there matrices
 U (k x k) and V (n x k) such that

$$\mathbf{X} = \mathbf{U}\Sigma\mathbf{V}^*$$
$$\mathbf{U}^*\mathbf{U} = I_m \qquad \mathbf{V}^*\mathbf{V} = I_m$$

- \sum a diagonal matrix, $\sigma_1 \ge \dots \ge \sigma_k \ge 0$
- σ_i^2 is an eigenvalue of **XX**^{*}. **U** are eigenvectors of **XX**^{*}.
- Fact: If X has rank r, then X has only r non-zero singular values.

SVD = Filter Bank

• Multiply a vector ${f z}$ by ${f X}={f U}\Sigma{f V}^*$

Collaborative Filterings

- Z is a linear combination of eigenusers, $v_1, ..., v_k$.
- $u_1, ..., u_k$ are the eigenratings

Which Algorithm?

Affine Rank Minimization:

minimize $\operatorname{rank}(\mathbf{X})$ subject to $\mathcal{A}(\mathbf{X}) = \mathbf{b}$

Convex Relaxation:

minimize $\|\mathbf{X}\|_* = \sum_{i=1}^k \sigma_i(\mathbf{X})$ subject to $\mathcal{A}(\mathbf{X}) = \mathbf{b}$

- Proposed by Fazel (2002).
- Nuclear norm is the "numerical rank" in numerical analysis
- The "trace heuristic" from controls if **X** is p.s.d.

- 2x2 matrices
- plotted in 3d
- $\left\| \left[\begin{array}{cc} x & 0 \\ 0 & z \end{array} \right] \right\|_* \le 1$
- Projection onto x-z plane is l₁ ball

- 2x2 matrices
- plotted in 3d
- $\left\| \left[\begin{array}{cc} x & y \\ y & z \end{array} \right] \right\|_* \le 1$
- Not polyhedral...

So how do we compute it? And when does it work?

Computationally: Gradient Descent!

$$\mathcal{F}(\mathbf{L}, \mathbf{R}) = \sum_{i=1}^{k} \sum_{j=1}^{r} L_{ij}^{2} + \sum_{i=1}^{n} \sum_{j=1}^{r} R_{ij}^{2} + \lambda \|\mathcal{A}(\mathbf{LR}^{*}) - \mathbf{b}\|^{2}$$

- "Method of multipliers"
- Schedule for $\boldsymbol{\lambda}$ controls the noise in the data
- Same global minimum as nuclear norm

Topic 2: Rank

- How many samples are required to reconstruct low-rank matrices?
- Fast algorithms for SVD as compressed sensing
- When can we guarantee the nuclear norm heuristic works?
- What are efficient ways to compute minimum nuclear norm solutions?

Topic 3: Approximation

 Try to write a function as a sum of (non-orthogonal) bases:

$$f(x) \approx \sum_{k=1}^{n} c_k \phi_k(\mathbf{x}; \theta_k)$$

- Atoms are sets of basis functions
- Not a discrete set, infinite dimensional space.

• Solution: Approximate $f(\mathbf{x})$ by $f_n(\mathbf{x}) = \sum_{k=1}^{k} c_k \phi_k(\mathbf{x}; \theta_k)$ optimize sample

• Approximate
$$f(\mathbf{x})$$
 by $f_n(\mathbf{x}) = \sum_{k=1}^n c_k \phi_k(\mathbf{x}; \theta_k)$

For large class of f, sampling θ_k i.i.d. and optimizing c_k yields

$$\|f - f_n\| = O\left(\frac{1}{\sqrt{n}}\right)$$

Analysis via convex hull norm where the atoms are $\phi(\mathbf{x}; heta)$

 $\phi(\mathbf{x}; \,\omega, b) = \cos(\omega' \mathbf{x} + b)$ $\omega \sim \mathcal{N}(0, 1)$ $b \sim \text{unif}[-\pi, \pi]$

$$k(\mathbf{x}, \mathbf{y}) = \exp(-\gamma \|\mathbf{x} - \mathbf{y}\|^2)$$

Radial Basis Functions

```
% Approximates Gaussian Process regression
% with Gaussian kernel of variance gamma
% lambda: regularization parameter
% dataset: X is dxN, y is 1xN
% test: xtest is dx1
% D: dimensionality of random feature
% training
  w = randn(D, size(X,1));
  b = 2*pi*rand(D,1);
  Z = cos(sqrt(gamma)*w*X + repmat(b,1,size(X,2)));
  alpha = (lambda*eye(size(X,2)+Z*Z') \setminus (Z*y);
```

% testing

```
ztest = alpha(:)'*cos( sqrt(gamma)*w*xtest(:) + ...
+ repmat(b,1,size(X,2)) );
```

Topic 3: Approximation

- How many bases are required to approximate complicated behavior?
- What are efficient ways to fit functions in infinite dimensional function spaces?
- What are fast ways to fit functions when we are overwhelmed by data?