

# 16 Lecture 16, Feb 25

## Announcements

- HW5 (Newton's method, handwritten digit recognition) due next Tue Mar 1 @ 11:59PM.
- Quiz 3 next Tue 3/1. In class, closed book.

## Last time

- EM algorithm.
- EM example: finite Gaussian mixture.
- MM algorithm.
- MM example: PET imaging.

## Today

- PET imaging: GPU computing.
- MM example: matrix completion.
- MM example: nonnegative matrix factorization (NNMF).
- Quasi-Newton method.

## GPU

GPUs are ubiquitous: servers, desktops, and laptops.

NVIDIA GPUs	Tesla M2090	GTX 580	GT 650M
			
Computers	servers, cluster	desktop	laptop
			
Main usage	scientific computing	gaming	gaming
Current version	K40	GTX 980	GTX 900M
Memory	6GB	1.5GB	1GB
Memory bandwidth	177GB/sec	192GB/sec	80GB/sec
Number of cores	512	512	384
Processor clock	1.3GHz	1.5GHz	0.9GHz
Peak DP performance	666Gflops		
Peak SP performance	1332Gflops	1581Gflops	691Gflops
Release price	\$2500	\$500	OEM

## MM example: Netflix and matrix completion

- Snapshot of the kind of data collected by Netflix. Only 100,480,507 ratings (1.2% entries of the 480K-by-18K matrix) are observed

ID	movie 1	movie 2	movie 3	movie 4	movie 5	movie 6	...	movie 17,770
user 1	5	3	4	3	3	NA	...	1
user 2	4	NA	NA	NA	NA	NA	...	NA
user 3	NA	NA	NA	NA	NA	NA	...	NA
user 4	4	NA	NA	NA	NA	2	...	4
user 5	NA	NA	NA	5	NA	NA	...	NA
user 6	3	NA	NA	5	1	NA	...	3
user 7	NA	NA	NA	NA	NA	NA	...	NA
user 8	5	NA	5	NA	NA	NA	...	NA
user 9	NA	NA	NA	NA	3	NA	...	NA
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
user 480,189	NA	5	NA	NA	NA	NA	...	NA

- Netflix challenge: impute the unobserved ratings for personalized recommendation. [http://en.wikipedia.org/wiki/Netflix\\_Prize](http://en.wikipedia.org/wiki/Netflix_Prize)



- *Matrix completion problem.* Observe a very sparse matrix  $\mathbf{Y} = (y_{ij})$ . Want to impute all the missing entries. It is possible only when the matrix is structured, e.g., of *low rank*.
- Let  $\Omega = \{(i, j) : \text{observed entries}\}$  index the observed entries and  $P_\Omega(\mathbf{M})$  denote the projection of matrix  $\mathbf{M}$  to  $\Omega$ . The problem

$$\min_{\text{rank}(\mathbf{X}) \leq r} \frac{1}{2} \|P_\Omega(\mathbf{Y}) - P_\Omega(\mathbf{X})\|_F^2 = \frac{1}{2} \sum_{(i,j) \in \Omega} (y_{ij} - x_{ij})^2$$

unfortunately is non-convex and difficult.

- *Convex relaxation* (Mazumder et al., 2010)

$$\min_{\mathbf{X}} f(\mathbf{X}) = \frac{1}{2} \|P_\Omega(\mathbf{Y}) - P_\Omega(\mathbf{X})\|_F^2 + \lambda \|\mathbf{X}\|_*,$$

where  $\|\mathbf{X}\|_* = \|\sigma(\mathbf{X})\|_1 = \sum_i \sigma_i(\mathbf{X})$  is the nuclear norm.

- Majorization step:

$$\begin{aligned}
f(\mathbf{X}) &= \frac{1}{2} \sum_{(i,j) \in \Omega} (y_{ij} - x_{ij})^2 + \frac{1}{2} \sum_{(i,j) \notin \Omega} 0 + \lambda \|\mathbf{X}\|_* \\
&\leq \frac{1}{2} \sum_{(i,j) \in \Omega} (y_{ij} - x_{ij})^2 + \frac{1}{2} \sum_{(i,j) \notin \Omega} (x_{ij}^{(t)} - x_{ij})^2 + \lambda \|\mathbf{X}\|_* \\
&= \frac{1}{2} \|\mathbf{X} - \mathbf{Z}^{(t)}\|_F^2 + \lambda \|\mathbf{X}\|_* \\
&= g(\mathbf{X} | \mathbf{X}^{(t)}),
\end{aligned}$$

where  $\mathbf{Z}^{(t)} = P_\Omega(\mathbf{Y}) + P_{\Omega^\perp}(\mathbf{X}^{(t)})$ . “fill in missing entries by current imputation”

- Minimization step:

Rewrite the surrogate function

$$g(\mathbf{X} | \mathbf{X}^{(t)}) = \frac{1}{2} \|\mathbf{X}\|_F^2 - \text{tr}(\mathbf{X}^T \mathbf{Z}^{(t)}) + \frac{1}{2} \|\mathbf{Z}^{(t)}\|_F^2 + \lambda \|\mathbf{X}\|_*$$

Let  $\sigma_i$  be the singular values of  $\mathbf{X}$  and  $\omega_i$  be the singular values of  $\mathbf{Z}^{(t)}$ . Observe

$$\begin{aligned}
\|\mathbf{X}\|_F^2 &= \|\sigma(\mathbf{X})\|_2^2 = \sum_i \sigma_i^2 \\
\|\mathbf{Z}^{(t)}\|_F^2 &= \|\sigma(\mathbf{Z}^{(t)})\|_2^2 = \sum_i \omega_i^2
\end{aligned}$$

and by the Fan-von Neuman’s inequality

$$\text{tr}(\mathbf{X}^T \mathbf{Z}^{(t)}) \leq \sum_i \sigma_i \omega_i$$

with equality achieved if and only if the left and right singular vectors of the two matrices coincide. Thus we can choose  $\mathbf{X}$  to have same singular vectors as  $\mathbf{Z}^{(t)}$  and

$$\begin{aligned}
g(\mathbf{X} | \mathbf{X}^{(t)}) &= \frac{1}{2} \sum_i \sigma_i^2 - \sum_i \sigma_i \omega_i + \frac{1}{2} \omega_i^2 + \lambda \sum_i \sigma_i \\
&= \frac{1}{2} \sum_i (\sigma_i - \omega_i)^2 + \lambda \sum_i \sigma_i,
\end{aligned}$$

with minimizer given by  $\sigma_i^{(t+1)} = (\omega_i - \lambda)_+$ . “Singular value thresholding”

- Algorithm for matrix completion:

Initialize  $\mathbf{X}^{(0)} \in \mathbb{R}^{m \times n}$

**repeat**

$$\mathbf{Z}^{(t+1)} \leftarrow P_{\Omega}(\mathbf{Y}) + P_{\Omega^\perp}(\mathbf{X}^{(t)})$$

Compute SVD:  $\mathbf{U}\text{diag}(\mathbf{w})\mathbf{V}^\top \leftarrow \mathbf{Z}^{(t+1)}$

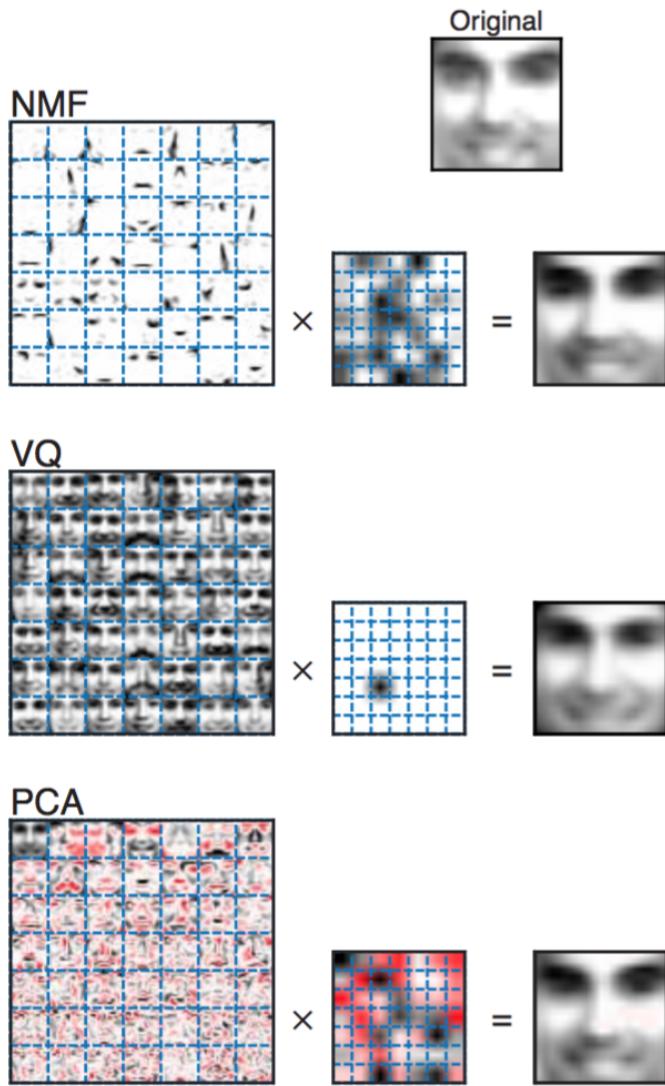
$$\mathbf{X}^{(t+1)} \leftarrow \mathbf{U}\text{diag}[(\mathbf{w} - \lambda)_+] \mathbf{V}^\top$$

**until** objective value converges

- “Golub-Kahan-Reinsch” algorithm takes  $4m^2n + 8mn^2 + 9n^3$  flops for a  $m \geq n$  matrix and is not going to work for 480K-by-18K Netflix matrix.

Notice only top singular values are needed since low rank solutions are achieved by large  $\lambda$ . Lanczos/Arnoldi algorithm is the way to go. Matrix-vector multiplication  $\mathbf{Z}^{(t)}\mathbf{v}$  is fast (why?)

## MM example: non-negative matrix factorization (NNMF)



**FIGURE 14.33.** Non-negative matrix factorization (NMF), vector quantization (VQ, equivalent to k-means clustering) and principal components analysis (PCA) applied to a database of facial images. Details are given in the text. Unlike VQ and PCA, NMF learns to represent faces with a set of basis images resembling parts of faces.

- Nonnegative matrix factorization (NNMF) was introduced by Lee and Seung (1999, 2001) as an analog of principal components and vector quantization with applications in data compression and clustering.
- In mathematical terms, one approximates a data matrix  $\mathbf{X} \in \mathbb{R}^{m \times n}$  with non-

negative entries  $x_{ij}$  by a product of two low-rank matrices  $\mathbf{V} \in \mathbb{R}^{m \times r}$  and  $\mathbf{W} \in \mathbb{R}^{r \times n}$  with nonnegative entries  $v_{ik}$  and  $w_{kj}$ . Consider minimization of the squared Frobenius norm

$$L(\mathbf{V}, \mathbf{W}) = \|\mathbf{X} - \mathbf{V}\mathbf{W}\|_{\text{F}}^2 = \sum_i \sum_j (x_{ij} - \sum_k v_{ik}w_{kj})^2, \quad v_{ik} \geq 0, w_{kj} \geq 0,$$

which should lead to a good factorization.

- $L(\mathbf{V}, \mathbf{W})$  is not convex, but bi-convex. The strategy is to alternately update  $\mathbf{V}$  and  $\mathbf{W}$ .
- The key is the majorization, via convexity of the function  $(x_{ij} - x)^2$ ,

$$\left( x_{ij} - \sum_k v_{ik}w_{kj} \right)^2 \leq \sum_k \frac{a_{ikj}^{(t)}}{b_{ij}^{(t)}} \left( x_{ij} - \frac{b_{ij}^{(t)}}{a_{ikj}^{(t)}} v_{ik}w_{kj} \right)^2,$$

where

$$a_{ikj}^{(t)} = v_{ik}^{(t)}w_{kj}^{(t)}, \quad b_{ij}^{(t)} = \sum_k v_{ik}^{(t)}w_{kj}^{(t)}.$$

- This suggests the alternating multiplicative updates

$$\begin{aligned} v_{ik}^{(t+1)} &= v_{ik}^{(t)} \frac{\sum_j x_{ij}w_{kj}^{(t)}}{\sum_j b_{ij}^{(t)}w_{kj}^{(t)}} \\ b_{ij}^{(t+.5)} &= \sum_k v_{ik}^{(t+1)}w_{kj}^{(t)} \\ w_{kj}^{(t+1)} &= w_{kj}^{(t)} \frac{\sum_i x_{ij}v_{ik}^{(t+1)}}{\sum_i b_{ij}^{(t+.5)}v_{ik}^{(t+1)}}. \end{aligned}$$

- In MATLAB or JULIA notation, the update is extremely simple

$$\begin{aligned} \mathbf{B}^{(t)} &= \mathbf{V}^{(t)}\mathbf{W}^{(t)} \\ \mathbf{V}^{(t+1)} &= \mathbf{V}^{(t)} \odot (\mathbf{X}\mathbf{W}^{(t)T}) \oslash (\mathbf{B}^{(t)}\mathbf{W}^{(t)T}) \\ \mathbf{B}^{(t+.5)} &= \mathbf{V}^{(t+1)}\mathbf{W}^{(t)} \\ \mathbf{W}^{(t+1)} &= \mathbf{W}^{(t)} \odot (\mathbf{X}^T\mathbf{V}^{(t+1)}) \oslash (\mathbf{B}^{(t+.5)T}\mathbf{V}^{(t)}), \end{aligned}$$

where  $\odot$  denotes elementwise multiplication and  $\oslash$  denotes elementwise division. If we start with  $v_{ik}, w_{kj} > 0$ , parameter iterates stay positive.

- Database #1 from the MIT Center for Biological and Computational Learning (CBCL) (<http://cbcl.mit.edu>) reduces to a matrix  $X$  containing  $m = 2,429$  gray-scale face images with  $n = 19 \times 19 = 361$  pixels per face. Each image (row) is scaled to have mean and standard deviation 0.25.
- CPU: Intel Core 2 @ 3.2GHZ (4 cores), implemented using BLAS in the GNU Scientific Library (GSL)
- GPU: NVIDIA GeForce GTX 280 (240 cores), implemented using cuBLAS.

Rank $r$	Iters	CPU		GPU		Speedup
		Time	Function	Time	Function	
10	25459	1203	106.2653503	55	106.2653504	22
20	87801	7564	89.56601262	163	89.56601287	46
30	55783	7013	78.42143486	103	78.42143507	68
40	47775	7880	70.05415929	119	70.05415950	66
50	53523	11108	63.51429261	121	63.51429219	92
60	77321	19407	58.24854375	174	58.24854336	112

## Quasi-Newton methods (KL 14.9)

- Quasi-Newton is probably the most successful “black-box” optimizers in use. E.g., implemented in the general purpose optimization routine `optim()` in R.
- Consider the general Newton scheme for minimizing  $f(\mathbf{x})$ ,  $\mathbf{x} \in \mathcal{X} \subset \mathbb{R}^p$ :

$$\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - s[\mathbf{A}^{(t)}]^{-1} \nabla f(\mathbf{x}^{(t)}),$$

where  $\mathbf{A}^{(t)}$  a pd approximation of the Hessian  $d^2 f(\mathbf{x}^{(t)})$ .

- Pros: fast (quadratic) convergence
- Cons: compute and store Hessian at each iteration (usually  $O(np^2)$  cost in statistical problems), solving a linear system ( $O(p^3)$  cost in general), human efforts (derive and implement gradient and Hessian, pd approximation, ...)
- Any pd  $\mathbf{A}$  gives a descent algorithm – tradeoff between convergence rate and cost per iteration.

- Setting  $\mathbf{A} = \mathbf{I}$  leads to the *steepest descent* algorithm. Picture: slow convergence (zig-zagging) of steepest descent (with exact line search) for minimizing a convex quadratic function (ellipse).

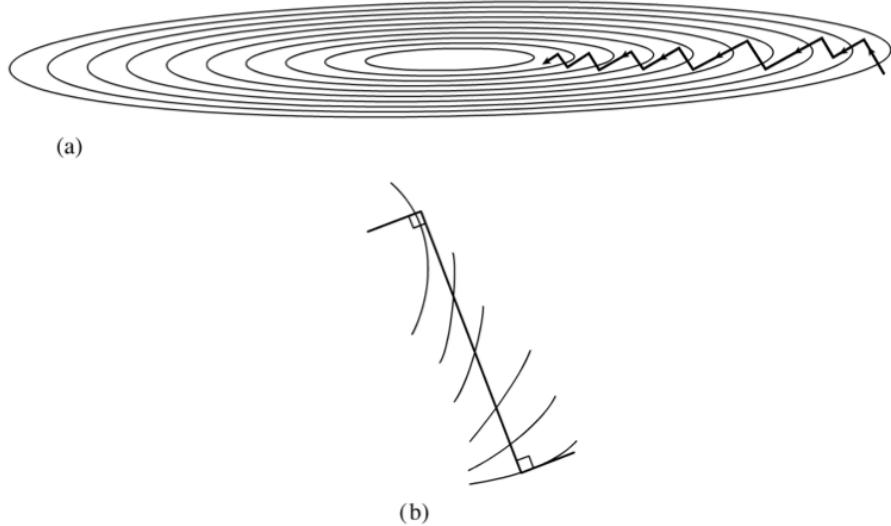


Figure 10.6.1. (a) Steepest descent method in a long, narrow “valley.” While more efficient than the strategy of Figure 10.5.1, steepest descent is nonetheless an inefficient strategy, taking many steps to reach the valley floor. (b) Magnified view of one step: A step starts off in the local gradient direction, perpendicular to the contour lines, and traverses a straight line until a local minimum is reached, where the traverse is parallel to the local contour lines.

How many steps does the Newton’s method take for a convex quadratic  $f$ ?

- Idea of Quasi-Newton: No analytical Hessian is required (still need gradient). Update approximate Hessian  $\mathbf{A}$  according to the *secant condition*

$$\nabla f(\mathbf{x}^{(t-1)}) - \nabla f(\mathbf{x}^{(t)}) \approx [d^2 f(\mathbf{x}^{(t)})](\mathbf{x}^{(t-1)} - \mathbf{x}^{(t)}).$$

Instead of computing  $\mathbf{A}$  from scratch at each iteration, we update an approximation  $\mathbf{A}$  to  $d^2 f(\mathbf{x})$  which satisfies (1) p.d., (2) the secant condition, and (3) closest to the previous approximation.

- Super-linear convergence, compared to the quadratic convergence of Newton’s method. But each iteration only takes  $O(p^2)$ !
- History: Davidon was a physicist at Argonne National Lab in 50s and proposed the very first idea of quasi-Newton method in 1959. Later studied, implemented, and popularized by Fletcher and Powell. Davidon’s original paper was not

accepted for publication  $\ominus$ ; 30 years later it appeared as the first article in the first issue of *SIAM Journal of Optimization* (Davidon, 1991).

## William C. Davidon

From Wikipedia, the free encyclopedia

**William Cooper Davidon** (1927–November 8, 2013) was an American professor of physics and mathematics, and peace activist. He was the mastermind of the March 8, 1971 F.B.I. office break-in, in Media, Pennsylvania, and the informal leader of the Citizens' Commission to Investigate the FBI.

### Contents

- 1 Life
- 2 Activism
- 3 Family
- 4 References
- 5 External links

### Life

Davidon was born in Fort Lauderdale, Florida in 1927. He attended Purdue University, and graduated from the University of Chicago with a Ph.D. in 1957.<sup>[1]</sup>

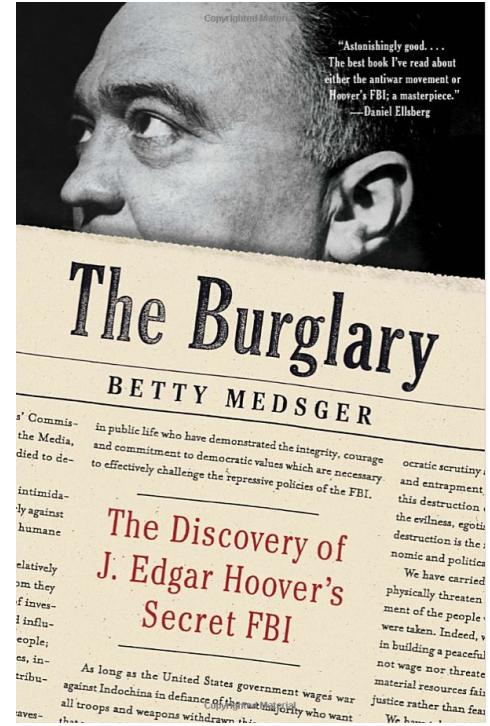
From 1954 to 1956, he was a research associate at the Enrico Fermi Institute. From 1956 to 1961, he was an associate physicist at the Argonne National Laboratory. He was professor of physics at Haverford College, beginning in 1961, and then Professor of Mathematics. He retired in 1991. He was a 1966 Fulbright Scholar.<sup>[2]</sup>

Davidon moved to Highlands Ranch, Colorado, in 2010. He died November 8, 2013, of Parkinson's disease.

William C. Davidon



<b>Born</b>	1927 Fort Lauderdale
<b>Died</b>	November 8, 2013 Highlands Ranch, Colorado
<b>Nationality</b>	American
<b>Occupation</b>	physics professor
<b>Known for</b>	Citizens' Commission to Investigate the FBI



- Davidon-Fletcher-Powell (DFP) rank-2 update. Solve

$$\begin{aligned} & \text{minimize} && \|\mathbf{A} - \mathbf{A}^{(t)}\|_F \\ & \text{subject to} && \mathbf{A} = \mathbf{A}^T \\ & && \mathbf{A}(\mathbf{x}^{(t)} - \mathbf{x}^{(t-1)}) = \nabla f(\mathbf{x}^{(t)}) - \nabla f(\mathbf{x}^{(t-1)}) \end{aligned}$$

for the next approximation  $\mathbf{A}^{(t+1)}$ . The solution is a low rank (rank 1 or rank 2) update of  $\mathbf{A}^{(t)}$ . The inverse is too thanks to Sherman-Morrison-Woodbury!  $O(p^2)$  operations. Need to store a  $p$ -by- $p$  dense matrix. Positive definiteness is guaranteed by the same trick you are using in HW5!

- Broyden-Fletcher-Goldfarb-Shanno (BFGS) rank 2 update is considered by many the most effective among all quasi-Newton updates. BFGS imposes secant con-

dition on the inverse of Hessian  $\mathbf{H} = \mathbf{A}^{-1}$ .

$$\begin{aligned} & \text{minimize} && \|\mathbf{H} - \mathbf{H}^{(t)}\|_{\text{F}} \\ & \text{subject to} && \mathbf{H} = \mathbf{H}^{\top} \\ & && \mathbf{H}[\nabla f(\mathbf{x}^{(t)}) - \nabla f(\mathbf{x}^{(t-1)})] = \mathbf{x}^{(t)} - \mathbf{x}^{(t-1)}. \end{aligned}$$

Again  $\mathbf{H}^{(t+1)}$  is a rank two update of  $\mathbf{H}^{(t)}$ .  $O(p^2)$  operations. Still need to store a dense  $p$ -by- $p$  matrix.

- Limited-memory BFGS (L-BFGS). Only store the secant pairs. Particularly useful for large scale optimization.
- L-BFGS-B: with box-constraints. Implemented in the general purpose optimization routine `optim()` in R.
- How to set the initial approximation  $\mathbf{A}^{(0)}$ ? Identity or Hessian (if pd) or Fisher information matrix at starting point.

## (Linear) Conjugate gradient method

- (Linear) Conjugate gradient is the top-notch iterative method for solving large, structured linear systems  $\mathbf{Ax} = \mathbf{b}$ . Earlier we talked about Jacobi and Gauss-Seidel as the classical iterative solvers.

**Table 1. Kershaw's results for a fusion problem.**

Method	Number of iterations
Gauss Seidel	208,000
Block successive overrelaxation methods	765
Incomplete Cholesky conjugate gradients	25

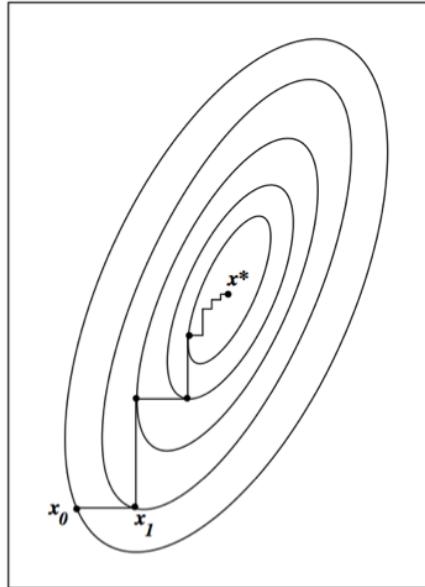
- *Linear* conjugate gradient method: for solving large linear systems of equations.  
History: Hestenes and Stiefel in 50s.

- Solve linear equation  $\mathbf{A}\mathbf{x} = \mathbf{b}$ , where  $\mathbf{A} \in \mathbb{R}^{n \times n}$  is pd, is equivalent to

$$\text{minimize } f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^\top \mathbf{A} \mathbf{x} - \mathbf{b}^\top \mathbf{x}.$$

Note  $\nabla f(\mathbf{x}) = \mathbf{A}\mathbf{x} - \mathbf{b} =: r(\mathbf{x})$ .

- Consider a simple idea: coordinate descent, that is to update  $x_j$  alternatingly. Same as the Gauss-Seidel iteration.



**Figure 9.1**

Coordinate search method makes slow progress on this function of two variables.

- A set of vectors  $\{\mathbf{p}^{(0)}, \dots, \mathbf{p}^{(l)}\}$  is said to be conjugate wrt  $\mathbf{A}$  if

$$\mathbf{p}_i^\top \mathbf{A} \mathbf{p}_j = 0, \quad \text{for all } i \neq j.$$

- *Conjugate direction* method: Given a set of conjugate vectors  $\{\mathbf{p}^{(0)}, \dots, \mathbf{p}^{(l)}\}$ , at iteration  $t$ , we search along the conjugate direction  $\mathbf{p}^{(t)}$

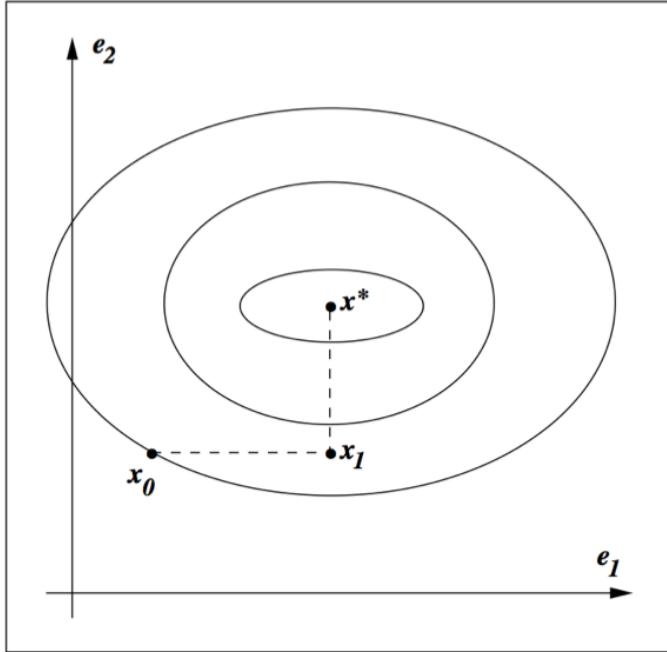
$$\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} + \alpha^{(t)} \mathbf{p}^{(t)},$$

where

$$\alpha^{(t)} = -\frac{\mathbf{r}^{(t) \top} \mathbf{p}^{(t)}}{\mathbf{p}^{(t) \top} \mathbf{A} \mathbf{p}^{(t)}}.$$

- Theorem:  $\mathbf{x}^{(t)}$  converges to the solution in at most  $n$  steps.

Intuition: Look at graph.



**Figure 5.1** Successive minimizations along the coordinate directions find the minimizer of a quadratic with a diagonal Hessian in  $n$  iterations.

- *Conjugate gradient* method. Idea: generate  $\mathbf{p}^{(t)}$  using only  $\mathbf{p}^{(t-1)}$

$$\mathbf{p}^{(t)} = -\mathbf{r}^{(t)} + \beta^{(t)} \mathbf{p}^{(t-1)},$$

where  $\beta^{(t)}$  is determined by the conjugacy condition  $\mathbf{p}^{(t-1)^\top} \mathbf{A} \mathbf{p}^{(t)} = 0$

$$\beta^{(t)} = \frac{\mathbf{r}^{(t)^\top} \mathbf{A} \mathbf{p}^{(t-1)}}{\mathbf{p}^{(t-1)^\top} \mathbf{A} \mathbf{p}^{(t-1)}}.$$

- CG algorithm (preliminary version):

Given  $\mathbf{x}^{(0)}$

Initialize:  $\mathbf{r}^{(0)} \leftarrow \mathbf{A}\mathbf{x}^{(0)} - \mathbf{b}$ ,  $\mathbf{p}^{(0)} \leftarrow -\mathbf{r}^{(0)}$ ,  $t = 0$

**while**  $\mathbf{r}^{(0)} \neq \mathbf{0}$  **do**

$$\alpha^{(t)} \leftarrow -\frac{\mathbf{r}^{(t)^\top} \mathbf{p}^{(t)}}{\mathbf{p}^{(t)^\top} \mathbf{A} \mathbf{p}^{(t)}}$$

$$\mathbf{x}^{(t+1)} \leftarrow \mathbf{x}^{(t)} + \alpha^{(t)} \mathbf{p}^{(t)}$$

$$\mathbf{r}^{(t+1)} \leftarrow \mathbf{A}\mathbf{x}^{(t+1)} - \mathbf{b}$$

$$\beta^{(t+1)} \leftarrow \frac{\mathbf{r}^{(t+1)^\top} \mathbf{A} \mathbf{p}^{(t)}}{\mathbf{p}^{(t)^\top} \mathbf{A} \mathbf{p}^{(t)}}$$

$$\mathbf{p}^{(t+1)} \leftarrow -\mathbf{r}^{(t+1)} + \beta^{(t+1)} \mathbf{p}^{(t)}$$

```

 $t \leftarrow t + 1$ 
end while

```

- Theorem: With CG algorithm

1.  $\mathbf{r}^{(t)}$  are mutually orthogonal.
2.  $\{\mathbf{r}^{(0)}, \dots, \mathbf{r}^{(t)}\}$  is contained in the *Krylov* subspace of degree  $t$  for  $\mathbf{r}^{(0)}$ , denoted by
$$\mathcal{K}(\mathbf{r}^{(0)}; t) = \text{span}\{\mathbf{r}^{(0)}, \mathbf{A}\mathbf{r}^{(0)}, \mathbf{A}^2\mathbf{r}^{(0)}, \dots, \mathbf{A}^t\mathbf{r}^{(0)}\}.$$
3.  $\{\mathbf{p}^{(0)}, \dots, \mathbf{p}^{(t)}\}$  is contained in  $\mathcal{K}(\mathbf{r}^{(0)}; t)$ .
4.  $\mathbf{p}^{(0)}, \dots, \mathbf{p}^{(t)}$  are conjugate wrt  $\mathbf{A}$ .

The iterates  $\mathbf{x}^{(t)}$  converge to the solution in at most  $n$  steps.

- CG algorithm (economical version): saves one matrix-vector multiplication.

Given  $\mathbf{x}^{(0)}$

Initialize:  $\mathbf{r}^{(0)} \leftarrow \mathbf{A}\mathbf{x}^{(0)} - \mathbf{b}$ ,  $\mathbf{p}^{(0)} \leftarrow -\mathbf{r}^{(0)}$ ,  $t = 0$

**while**  $\mathbf{r}^{(0)} \neq \mathbf{0}$  **do**

$$\alpha^{(t)} \leftarrow \frac{\mathbf{r}^{(t)\top} \mathbf{r}^{(t)}}{\mathbf{p}^{(t)\top} \mathbf{A} \mathbf{p}^{(t)}}$$

$$\mathbf{x}^{(t+1)} \leftarrow \mathbf{x}^{(t)} + \alpha^{(t)} \mathbf{p}^{(t)}$$

$$\mathbf{r}^{(t+1)} \leftarrow \mathbf{r}^{(t)} + \alpha^{(t)} \mathbf{A} \mathbf{p}^{(t)}$$

$$\beta^{(t+1)} \leftarrow \frac{\mathbf{r}^{(t+1)\top} \mathbf{r}^{(t+1)}}{\mathbf{r}^{(t)\top} \mathbf{r}^{(t)}}$$

$$\mathbf{p}^{(t+1)} \leftarrow -\mathbf{r}^{(t+1)} + \beta^{(t+1)} \mathbf{p}^{(t)}$$

$$t \leftarrow t + 1$$

**end while**

- Computation cost per iteration is one matrix vector multiplication:  $\mathbf{A}\mathbf{p}^{(t)}$ .

Consider PageRank problem,  $\mathbf{A}$  has dimension  $n \approx 10^{10}$  but is highly structured (sparse + low rank). Each matrix vector multiplication takes  $O(n)$ .

- Theorem: If  $\mathbf{A}$  has  $r$  distinct eigenvalues,  $\mathbf{x}^{(t)}$  converges to solution  $\mathbf{x}^*$  in at most  $r$  steps.

## Pre-conditioned conjugate gradient (PCG)

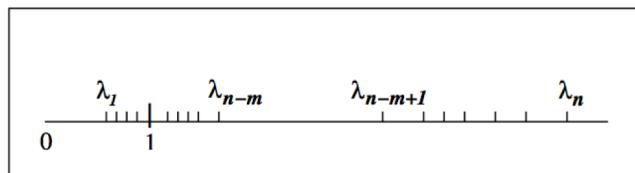
- Summary of conjugate gradient method for solving  $\mathbf{A}\mathbf{x} = \mathbf{b}$  or equivalently minimizing  $\frac{1}{2}\mathbf{x}^T\mathbf{A}\mathbf{x} - \mathbf{b}^T\mathbf{x}$ :
  - Each iteration needs one matrix vector multiplication:  $\mathbf{A}\mathbf{p}^{(t+1)}$ . For structured  $\mathbf{A}$ , often  $O(n)$  cost per iteration.
  - Guaranteed to converge in  $n$  steps.

- Two important bounds for conjugate gradient algorithm:

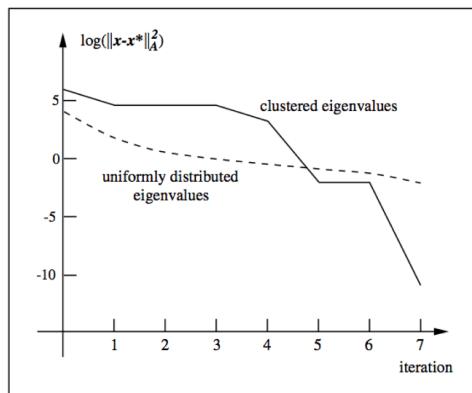
Let  $\lambda_1 \leq \dots \leq \lambda_n$  be the ordered eigenvalues of a pd  $\mathbf{A}$ .

$$\begin{aligned}\|\mathbf{x}^{(t+1)} - \mathbf{x}^*\|_{\mathbf{A}}^2 &\leq \left(\frac{\lambda_{n-t} - \lambda_1}{\lambda_{n-t} + \lambda_1}\right)^2 \|\mathbf{x}^{(0)} - \mathbf{x}^*\|_{\mathbf{A}}^2 \\ \|\mathbf{x}^{(t+1)} - \mathbf{x}^*\|_{\mathbf{A}}^2 &\leq 2 \left(\frac{\sqrt{\kappa(\mathbf{A})} - 1}{\sqrt{\kappa(\mathbf{A})} + 1}\right)^t \|\mathbf{x}^{(0)} - \mathbf{x}^*\|_{\mathbf{A}}^2,\end{aligned}$$

where  $\kappa(\mathbf{A}) = \lambda_n/\lambda_1$  is the condition number of  $\mathbf{A}$ .



**Figure 5.3** Two clusters of eigenvalues.



**Figure 5.4** Performance of the conjugate gradient method on (a) a problem in which five of the eigenvalues are large and the remainder are clustered near 1, and (b) a matrix with uniformly distributed eigenvalues.

- Messages:
  - Roughly speaking, if the eigenvalues of  $\mathbf{A}$  occur in  $r$  distinct clusters, the CG iterates will *approximately* solve the problem after  $O(r)$  steps.
  - $\mathbf{A}$  with a small condition number ( $\lambda_1 \approx \lambda_n$ ) converges fast.
- *Pre-conditioning*: Change of variables  $\hat{\mathbf{x}} = \mathbf{C}\mathbf{x}$  via a nonsingular  $\mathbf{C}$  and solve

$$(\mathbf{C}^{-T} \mathbf{A} \mathbf{C}^{-1}) \hat{\mathbf{x}} = \mathbf{C}^{-T} \mathbf{b}.$$

Choose  $\mathbf{C}$  such that

- $\mathbf{C}^{-T} \mathbf{A} \mathbf{C}^{-1}$  has small condition number, or
- $\mathbf{C}^{-T} \mathbf{A} \mathbf{C}^{-1}$  has clustered eigenvalues
- Inexpensive solution of  $\mathbf{C}^T \mathbf{C} \mathbf{y} = \mathbf{r}$
- Preconditioned CG does not make use of  $\mathbf{C}$  explicitly, but rather the matrix  $\mathbf{M} = \mathbf{C}^T \mathbf{C}$ .
- Preconditioned CG (PCG) algorithm:

```

Given  $\mathbf{x}^{(0)}$ , pre-conditioner  $\mathbf{M}$ 
 $\mathbf{r}^{(0)} \leftarrow \mathbf{A}\mathbf{x}^{(0)} - \mathbf{b}$ 
solve  $\mathbf{M}\mathbf{y}^{(0)} = \mathbf{r}^{(0)}$  for  $\mathbf{y}^{(0)}$ 
 $\mathbf{p}^{(0)} \leftarrow -\mathbf{r}^{(0)}$ ,  $t = 0$ 
while  $\mathbf{r}^{(0)} \neq \mathbf{0}$  do
   $\alpha^{(t)} \leftarrow \frac{\mathbf{r}^{(t)\top} \mathbf{y}^{(t)}}{\mathbf{p}^{(t)\top} \mathbf{A} \mathbf{p}^{(t)}}$ 
   $\mathbf{x}^{(t+1)} \leftarrow \mathbf{x}^{(t)} + \alpha^{(t)} \mathbf{p}^{(t)}$ 
   $\mathbf{r}^{(t+1)} \leftarrow \mathbf{r}^{(t)} + \alpha^{(t)} \mathbf{A} \mathbf{p}^{(t)}$ 
  Solve  $\mathbf{M}\mathbf{y}^{(t+1)} \leftarrow \mathbf{r}^{(t+1)}$  for  $\mathbf{y}^{(t+1)}$ 
   $\beta^{(t+1)} \leftarrow \frac{\mathbf{r}^{(t+1)\top} \mathbf{y}^{(t+1)}}{\mathbf{r}^{(t)\top} \mathbf{r}^{(t)}}$ 
   $\mathbf{p}^{(t+1)} \leftarrow -\mathbf{y}^{(t+1)} + \beta^{(t+1)} \mathbf{p}^{(t)}$ 
   $t \leftarrow t + 1$ 
end while

```

Remark: Only extra cost in the pre-conditioned CG algorithm is the need to solve the linear system  $\mathbf{M}\mathbf{y} = \mathbf{r}$ .

- Pre-conditioning is more like an art than science. Some choices include
  - Incomplete Cholesky.  $\mathbf{A} \approx \tilde{\mathbf{L}}\tilde{\mathbf{L}}^T$ , where  $\tilde{\mathbf{L}}$  is a sparse approximate Cholesky factor. Then  $\tilde{\mathbf{L}}^{-1}\mathbf{A}\tilde{\mathbf{L}}^{-T} \approx \mathbf{I}$  (perfectly conditioned) and  $\mathbf{M}\mathbf{y} = \tilde{\mathbf{L}}\tilde{\mathbf{L}}^T\mathbf{y} = \mathbf{r}$  is easy to solve.
  - Banded pre-conditioners.
  - Choose  $\mathbf{M}$  as a coarsened version of  $\mathbf{A}$ .
  - *Subject knowledge.* Knowledge about the structure and origin of a problem is often the key to devising efficient pre-conditioner. For example, see recent work of Stein et al. (2012) for pre-conditioning large covariance matrices. <http://pubs.siam.org/doi/abs/10.1137/110834469>

## More buzzwords and softwares

Here are a few variants of CG that we should at least know the names and what they are for (so we can Google later in need).

- MINRES (minimum residual method): symmetric indefinite  $\mathbf{A}$ .
- Bi-CG (bi-conjugate gradient): unsymmetric  $\mathbf{A}$ .
- Bi-CGSTAB (Bi-CG stabilized): improved version of Bi-CG.
- GMRES (generalized minimum residual method): current *de facto* method for unsymmetric  $\mathbf{A}$ . E.g., PageRank problem.
- Lanczos method: top eigen-pairs of a large symmetric matrix.
- Arnoldi method: top eigen-pairs of a large unsymmetric matrix.
- Lanczos bidiagonalization algorithm: top singular triples of large matrix.

Remark: For Lanczos/Arnoldi methods, the critical computation is still matrix vector multiplication  $\mathbf{A}\mathbf{v}$ .

Softwares:

- MATLAB:

- Iterative methods for solving linear equations:  
`pcg`, `bicg`, `bicgstab`, `gmres`, ...
- Iterative methods for top eigen-pairs and singular pairs:  
`eigs`, `svds`, ...
- Pre-conditioner:  
`cholinc`, `luinc`, ...

Get familiar with the reverse communication interface (RCI) for utilizing iterative solvers:

```
x = gmres(A, b)
x = gmres(@Afun, b)
eigs(A)
eigs(@Afun)
```

- Consider the PageRank problem. We want to find the top left eigenvector of the transition matrix

$$\mathbf{P} = p\mathbf{R}^+ \mathbf{A} + z\mathbf{1}_n^\top,$$

where  $\mathbf{R} = \text{diag}(r_1, \dots, r_n)$  and  $z_j = (1 - p)/n$  if  $r_i > 0$  and  $1/n$  if  $r_i = 0$ . Size of  $\mathbf{P}$  can be huge:  $n \approx 40$  billion web pages. How to call the `gmres` or `eigs` function?

- R: Try Google and good luck ...

## (Nonlinear) Conjugate gradient method

- *Linear* conjugate gradient method is for solving linear system  $\mathbf{A}\mathbf{x} = \mathbf{b}$ , or equivalently, minimizing  $\frac{1}{2}\mathbf{x}^T \mathbf{A} \mathbf{x} - \mathbf{b}^T \mathbf{x}$ .
- *Nonlinear* conjugate gradient is for nonlinear optimization

$$\text{minimize } f(\mathbf{x}).$$

- History: Fletcher and Reeves in 60s.
- Fletcher-Reeves CG for nonlinear minimization:

Given  $\mathbf{x}^{(0)}$

Evaluate  $\nabla f^{(0)} = \nabla f(\mathbf{x}^{(0)})$

Set  $\mathbf{p}^{(0)} \leftarrow -\nabla f^{(0)}$ ,  $t \leftarrow 0$

**while**  $\nabla f^{(t)} \neq \mathbf{0}$  **do**

    Compute  $\alpha^{(t)}$  and set  $\mathbf{x}^{(t+1)} \leftarrow \mathbf{x}^{(t)} + \alpha^{(t)} \mathbf{p}^{(t)}$

    Evaluate  $\nabla f^{(t+1)} = \nabla f(\mathbf{x}^{(t+1)})$

$$\beta^{(t+1)} \leftarrow \frac{df^{(t+1)} \nabla f^{(t+1)}}{df^{(t)} \nabla f^{(t)}}$$

$$\mathbf{p}^{(t+1)} \leftarrow -\nabla f^{(t+1)} + \beta^{(t+1)} \mathbf{p}^{(t)}$$

$$t \leftarrow t + 1$$

**end while**

- Most cost is evaluation of objective function and its gradient. No matrix operations are needed. Appealing for large nonlinear optimization problems.
- Line search (choose  $\alpha^{(t)}$ ) is necessary to get a descending algorithm.