



• Alternating convex optimization

given nonconvex problem with variables $(x_1, x_2, \dots, x_n) \in \mathbb{R}^n$

$I_1, \dots, I_k \subset \{1, \dots, n\}$ are index subset with $\bigcup_i I_j = \{1, \dots, n\}$

Suppose problem is convex in subset of variables $x_i, i \in I_j$

alternating convex optimization method: cycle through j , in each step optimizing over variable $x_i, i \in I_j$

eg- $f(x, A, b) = \|Ax - b\|$

f is convex in A when x, b fixed

x when A, b fixed

b when A, x fixed

f is jointly convex in X and B when A fixed

$$I_1 = \{A, B\}$$

A and B when X fixed

$$I_2 = \{B, X\}$$

f is not convex in A and X when b fixed

min. $\|Ax - b\|$

s.t. $A \in \mathbb{A}$

$b \in \mathbb{B}$

$x \in X$

Repeat {

min. $\|Ax - b\|$

s.t. $x \in X$

min. $\|Ax - b\|$

s.t. $A \in \mathbb{A}, b \in \mathbb{B}$

}

special case: bi-convex problem

$X = (u, v)$ problem is convex in u/v with v/u fixed

alternate optimizing over u and v

eg. Non-negative matrix factorization

min. $\|Ax - Y\|_F$ $X \in \mathbb{R}^{m \times k}$ $Y \in \mathbb{R}^{k \times n}$ data $A \in \mathbb{R}^{m \times n}$

s.t. $X_{ij} \geq 0, Y_{ij} \geq 0$

when k is small: express A as product of low-rank matrices

if there's no constraint, this problem is solved by truncating the SVD of A

when $k=1$:

Take the positive part of A , and take the 1st singular vectors of SVD of A
Perron-frobenius theorem:

for a pos matrix, the singular vectors associated with 1st singular value

can be chosen to be positive

alternating CONVEX optimization :

Solve QPs to optimize over X , then Y , then X ...

There are a handful of non-convex problems that can be solved exactly.

For example: can minimize any quadratic function, convex or not,

Subject to a single quadratic inequality, convex not

$$\text{eg: } \begin{array}{ll} \min & x^T P x \\ \text{st.} & \|x\|_2 \leq 1 \end{array} \Rightarrow \text{mini eigenvalue}$$

• Conjugate gradient method

• Three classes of methods for linear equations

Methods to solve linear system $Ax = b$ $A \in \mathbb{R}^{m \times n}$

1. dense direct (factor-solve method)

runtime depends only on size, (almost) independent of data

work well for n up to a few thousands

2. Sparse direct (factor solve method)

runtime depends on size, sparsity pattern, (almost) independent of data

work well for n up to 10^4 or 10^5 or more

require good heuristic for ordering

3. indirect (iterative methods)

runtime depends on data, size, sparsity, required accuracy

require tuning, preconditioning

good choices in many cases, only choice for $n=10^6$ or larger

• Symmetric positive definite linear system

SPD system of equations: $Ax = b$ $A \in \mathbb{R}^{n \times n}$ $A = A^T > 0$

eg Newton step $\nabla f(x) dx = -\nabla^2 f(x)$

(regularized) least square $(A^T A + \lambda I) x = A^T b$

minimize convex quadratic function $\frac{1}{2} x^T A x - b^T x$

• Conjugate gradient overview

the SPD system $Ax=b$ ($A \in \mathbb{R}^{n \times n}$, $b \in \mathbb{R}^n$) $\left\{ \begin{array}{l} \text{in theory in } n \text{ iterations} \\ \text{each iteration requires a few inner products in } \mathbb{R}^n, \\ \text{and one matrix-vector multiplication } z \rightarrow Az \end{array} \right.$

With roundoff error, CG can work poorly

for some A, b can get good approximation in $\ll n$ iterations

• Solution and Error

$x^* = A^{-1}b$ is the solution

x^* minimizes $f(x) = \frac{1}{2} x^T Ax - b^T x$

$\nabla f(x) = Ax - b$ (the residual)

$$\begin{aligned} f(x) - f^* &= \frac{1}{2} x^T Ax - \frac{1}{2} x^* T Ax^* + b^T x^* \\ &= \frac{1}{2} x^T Ax - b^T x - \frac{1}{2} b^T A^T A x^* + b^T A^T b \\ &= \frac{1}{2} x^T Ax + \frac{1}{2} b A^T b - b^T x \\ &= \frac{1}{2} x^T Ax + \frac{1}{2} x^{*T} A x^* - x^{*T} Ax \\ &= \frac{1}{2} (x - x^*)^T A (x - x^*) = \frac{1}{2} \|x - x^*\|_A^2 \end{aligned}$$

a relative measure (comparing to 0)

$$\tau = \frac{f(x) - f^*}{f(b) - f^*} = -\frac{\|x - x^*\|_A}{\|b - x^*\|_A}$$

• Residual

$r = b - Ax$ is called the residual at x

$$\nabla f = Ax - b = A(x - x^*) = -r$$

$$f(x) - f^* = \frac{1}{2} (x - x^*)^T A (x - x^*) = \frac{1}{2} r^T A^T r = \frac{1}{2} \|r\|_A^2$$

a commonly used measure of relative accuracy $y = \|r\|/\|b\|$

b is the residual norm if $x^* = 0$

if $\|r\|/\|b\| \geq 1$, the solution is overperformed by $x \neq 0$

• Krylov Subspace (controllability subspace)

$$\begin{aligned}K_k &= \text{span}\{b, Ab, \dots, A^{k-1}b\} \\&= \{p(A)b \mid p \text{ polynomial, } \deg(p) \leq k\}\end{aligned}$$

define Krylov sequence $x^{(0)}, x^{(1)}, \dots$ as

$$x^{(k)} = \underset{x \in K_k}{\operatorname{argmin}} f(x) = \underset{x \in K_k}{\operatorname{argmin}} \|x - x^*\|_A$$

the CG algorithm (among others) generates the Krylov sequence

Properties of Krylov Sequence

$f(x^{(k+1)}) \leq f(x^{(k)})$ (optimizing over a bigger and bigger subspace), but $\|r\|$ can increase
 $x^{(k)} = x^*$ even if $K_n \neq \mathbb{R}$

Cayley-Hamilton theorem:

$$\chi(S) = \det(SI - A) = S^n + a_1 S^{n-1} + \dots + a_n$$

by Cayley-Hamilton theorem:

$$\chi(A) = A^n + a_1 A^{n-1} + \dots + a_n I = 0$$

$$\therefore A^{-1} = -\frac{a_1}{a_n} A^{n-1} - \frac{a_2}{a_n} A^{n-2} - \dots - \frac{a_{n-1}}{a_n} I$$

the solution of $Ax=b$ is in the span of Krylov subspace $K_n = \{b, Ab, \dots, A^{n-1}b\}$

i.e. $b \in$

b is an eigenvector of A

$x^{(k)} = p_k(A) b$, where p_k is a polynomial with $\deg(p_k) \leq k$

$x^{(k+1)} = x^{(k)} + \alpha_k r^{(k)} + \beta_k (x^{(k)} - x^{(k-1)})$ for some α_k, β_k (basis of CG algorithm)

Spectral analysis of Krylov Sequence

$$A = Q \Lambda Q^T \quad (Q \text{ orthogonal}, \quad \Lambda = \text{diag}(\lambda_1, \dots, \lambda_n))$$

$$\text{define } y = Q^T x \quad b = Q^T b \quad y^* = Q^T x^*$$

$$\begin{aligned} f(x) &= \frac{1}{2} x^T A x + b^T x \\ &= \frac{1}{2} x^T Q \Lambda Q^T x + b^T Q \Lambda Q^T x \\ &= \frac{1}{2} y^T \Lambda y - \bar{b}^T y = \bar{f}(y) \\ &= \frac{1}{2} \sum_i \lambda_i y_i^2 - \bar{b}_i y_i \end{aligned}$$

$$y_i^* = \frac{\bar{b}_i}{\lambda_i} \quad f^* = -\frac{1}{2} \sum_i \left(\frac{\bar{b}_i^2}{\lambda_i} \right)$$

Krylov Sequence in terms of y

$$y^{(k)} = \arg \min_{y \in E_k} \bar{f}(y) = \arg \min_{y \in E_k} \frac{1}{2} y^T \Lambda y - \bar{b}^T y \quad E_k = \{ \bar{b}, \Lambda \bar{b} - \lambda^{k-1} \bar{b} \}$$

$$y^{(k)} = P_k(\lambda_i) \bar{b}; \quad \deg(P_k) \leq k$$

$$\begin{aligned} P_k &= \arg \min_{\substack{\deg(P) \leq k}} \sum_i^n \frac{1}{2} \bar{b}_i^2 - \lambda_i \cdot P(\lambda_i)^T - \bar{b}_i^T P(\lambda_i) \\ &= \arg \min_{\deg(P) \leq k} \sum_i^n \bar{b}_i^2 \left(\frac{1}{2} \lambda_i P(\lambda_i)^T - P(\lambda_i) \right) \end{aligned}$$

$$\begin{aligned} f(x^{(k)}) - f^* &= \bar{f}(y^{(k)}) - f^* \\ &= \min_{\deg(P) \leq k} \sum_i^n \bar{b}_i^2 \left(\frac{1}{2} \lambda_i P(\lambda_i)^T - P(\lambda_i) \right) - -\frac{1}{2} \sum_i^n \left(\bar{b}_i^2 / \lambda_i \right) \\ &= \min_{\deg(P) \leq k} \frac{1}{2} \sum_i^n \bar{b}_i^2 \left(\frac{(\lambda_i P(\lambda_i) - 1)^2}{\lambda_i} \right) \\ &= \min_{\deg(P) \leq k} \frac{1}{2} \sum_i^n y_i^{*2} \lambda_i \underbrace{\left(\frac{(\lambda_i P(\lambda_i) - 1)^2}{\lambda_i} \right)}_{\text{new polynomial } q} \\ &= \min_{\deg(q) \leq k} q(\lambda_i) = \frac{1}{2} \sum_i^n y_i^{*2} \lambda_i q(\lambda_i) \\ &= \min_{\deg(q) \leq k} q(\lambda_i) = \frac{1}{2} \sum_i^n \bar{b}_i^2 \frac{q(\lambda_i)^2}{\lambda_i} \end{aligned}$$

$$\gamma^{(k)} = \frac{f(x^{(k)}) - f^*}{\|x^{(k)} - x^*\|_2} = \frac{f(x^{(k)}) - f^*}{\|y^*\|_2} = \frac{f(x^{(k)}) - f^*}{\|y^*\|_2}$$

$$= \underbrace{\min_{\deg(g) \leq k, g(0)=1} \sum_i^n y_i^* \lambda_i g(\lambda_i)^2}_{\sum_i^n y_i^* \lambda_i} \quad \text{a weighted sum of } g(\lambda_i)$$

$$\leq \min_{\deg(g) \leq k, g(0)=1} \max_{i=1 \dots n} |g(\lambda_i)|$$

if $\exists g$ of degree k , $g(0)=1$, that is small on spectrum (eigenvalue) of A , then $f(x^{(k)}) - f^*$ is small

if eigenvalue are clustered in k groups, then $y^{(k)}$ is a good approximate solution

