



计算方法

刘景铖

计算机软件新技术国家重点实验室
南京大学



回顾

- 矩阵多项式的特征空间: $p(A)v = p(\lambda)v$
 - $A \rightarrow p(A)$
 - 特征空间 $\{\lambda_1, \lambda_2, \dots, \lambda_n\} \rightarrow \{p(\lambda_1), p(\lambda_2), \dots, p(\lambda_n)\}$
- **(Cayley-Hamilton)** 存在 n 次多项式 $p(A)$ 使得 $p(A)A = I$
 - 记 $q(x) = 1 - xp(x)$
 - 则 n 次多项式 $q_A(x) = \det(A - xI)$ 满足 $q(A) = 0$
 - 非平凡的: 注意 $\det(A - xI)$ 仅仅为一个数



Cayley-Hamilton

考虑关于 x 的 n 次多项式 $q_A(x) = \det(A - xI)$

定理(Cayley-Hamilton): $q_A(A) = \mathbf{0}$ 为全零矩阵

注意 $q_A(0) = \det(A) \neq \mathbf{0}$

$$A^{-1} = \frac{(-1)^{n-1}}{\det(A)} (A^{n-1} + c_{n-1}A^{n-2} + \cdots + c_1I)$$

$$A^{-1} = p(A)$$

如何计算 $q(x)$? 高斯消元?



Richardson iteration

一个特殊情形: 给定实数 $0 < a < 1$, 如何用 a 表示出 $a^{-1}b$?

$$b + (1 - a)b + (1 - a)^2b + \dots = \frac{b}{1 - (1 - a)} = a^{-1}b$$

等价于 $x_0 = \mathbf{0}$, $x_{k+1} = (1 - a)x_k + b$

Richardson iteration (对于对称正定的矩阵 A):

$$\begin{aligned} x_0 &= \mathbf{0} \\ \overrightarrow{x_{k+1}} &= (I - \alpha A)\overrightarrow{x_k} + \alpha \overrightarrow{b} = \overrightarrow{x_k} - \alpha(A\overrightarrow{x_k} - \overrightarrow{b}) \end{aligned}$$



Richardson iteration

一个特殊情形: 给定实数 $0 < a < 1$, 如何用 a 表示出 $a^{-1}b$?

$$b + (1 - a)b + (1 - a)^2b + \dots = \frac{b}{1 - (1 - a)} = a^{-1}b$$

等价于 $x_0 = \mathbf{0}$, $x_{k+1} = (1 - a)x_k + b$

Richardson iteration:

$$\begin{aligned} x_0 &= \mathbf{0} \\ \vec{x}_{k+1} &= (I - \alpha A)\vec{x}_k + \alpha \vec{b} = \vec{x}_k - \alpha(A\vec{x}_k - \vec{b}) \end{aligned}$$

矩阵多项式的视角: 特征空间: $p(A)v = p(\lambda)v$

- $A \rightarrow p(A)$
- 特征空间 $\{\lambda_1, \lambda_2, \dots, \lambda_n\} \rightarrow \{p(\lambda_1), p(\lambda_2), \dots, p(\lambda_n)\}$

另一个视角: 对于对称的 A , 这正好是目标函数 $\frac{1}{2}x^T Ax - b^T x$ 的梯度下降方法!

$$\nabla \left(\frac{1}{2}x^T Ax - b^T x \right) = \frac{1}{2}(A + A^T)x - b$$



An Optimization Perspective

给定连续可导 $f: \mathbb{R}^n \rightarrow \mathbb{R}$, 最小化 f

例子:

最小二乘法: $f(\mathbf{x}) = \|\mathbf{Ax} - \mathbf{b}\|_2^2$

向量投影: $f(\mathbf{c}) = \|\mathbf{c} \vec{\mathbf{a}} - \vec{\mathbf{b}}\|_2^2$

求伪逆(Pseudoinverse): $f(\mathbf{x}) = \|\mathbf{x}\|_2^2$, subject to $\mathbf{Ax} = \mathbf{b}$

对称矩阵的特征值: $f(\mathbf{x}) = \mathbf{x}^T \mathbf{A} \mathbf{x}$, subject to $\mathbf{x}^T \mathbf{x} = 1$

线性规划: $f(\mathbf{x}) = \mathbf{c}^T \mathbf{x}$, subject to $\mathbf{Ax} \geq \mathbf{b}$

主成分分析 (Principal component analysis): $f(\mathbf{C}) = \|\mathbf{X} - \mathbf{CC}^T \mathbf{X}\|_F$, subject to $\mathbf{C} \mathbf{C}^T = \mathbf{I}_{d \times d}$

离散组合优化: 最小生成树, 最小割、最大流 (网络流), 最短路径, 最大匹配; 最大独立集、团、支配集, 最小覆盖, 背包问题, 最大割, max-SAT, 旅行商问题, 最长路径 (哈密顿路径)



Optimality notion

给定连续可导 $f: \mathbb{R}^n \rightarrow \mathbb{R}$, 最小化 f

x_* 是一个全局最优解 (Global minimum), 如果 $f(x) \geq f(x_*)$, $\forall x$

x_* 是一个局部最优解 (Local minimum), 如果 $\exists \delta > 0$, $\forall x: |x - x_*| < \delta$, 均有 $f(x) \geq f(x_*)$



Optimality condition from calculus

给定连续可导 $f: \mathbb{R}^n \rightarrow \mathbb{R}$, 和一个点 x_0

$$\nabla f(x) := \left(\frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, \dots, \frac{\partial f}{\partial x_n} \right)$$

$$f(x) \approx f(x_0) + \nabla f(x_0) \cdot (x - x_0), \quad \forall x$$

特别地, 可以选取 $x - x_0 = \alpha(\nabla f(x_0))^T$

$$f\left(x_0 + \alpha(\nabla f(x_0))^T\right) - f(x_0) \approx \alpha \|\nabla f(x_0)\|_2^2$$

当 α 足够小的时候, α 的符号决定 f 局部的增减性

当 $\nabla f(x_0) = 0$, 则点 x_0 为驻点(stationary point)

此时需要看二阶导数



Optimality condition from calculus

给定连续可导 $f: \mathbb{R}^n \rightarrow \mathbb{R}$

$$\mathbf{H}_f = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} & \cdots & \frac{\partial^2 f}{\partial x_2 \partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \frac{\partial^2 f}{\partial x_n \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_n^2} \end{bmatrix}.$$

通常也把Hessian矩阵 H_f 记作 $\nabla^2 f$



Optimality condition from calculus

给定连续可导 $f: \mathbb{R}^n \rightarrow \mathbb{R}, \forall x$

$$f(x) \approx f(x_0) + \nabla f(x_0) \cdot (x - x_0) + \frac{1}{2} (x - x_0)^T H_f(x_0) (x - x_0)$$

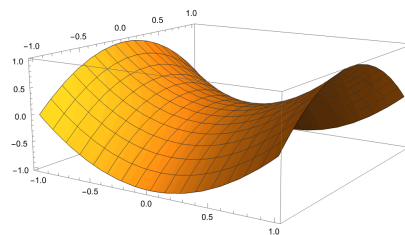
如果 $\nabla f(x_0) = 0$

$H_f(x_0) > 0$: 局部最小

$H_f(x_0) < 0$: 局部最大

$H_f(x_0)$ 同时有正的和负的特征值: saddle point

$H_f(x_0)$ 不可逆: 不一定是局部极值 (Morse theory)

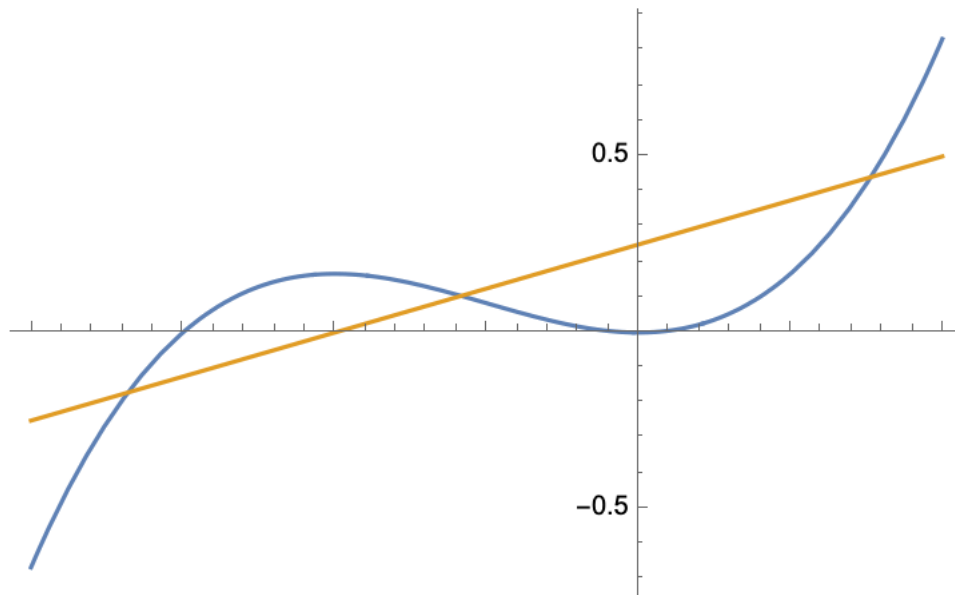
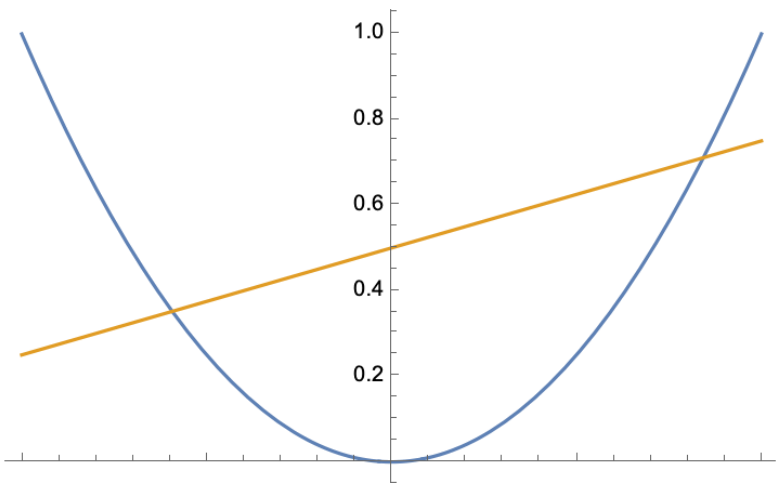




Convexity

如果给定的不仅仅是一个点上的导数，而是一个邻域上导数的信息，那么判断标准可以有所简化。

例子：凸函数(convex function)





Convexity

如果给定的不仅仅是一个点上的导数，而是一个邻域上导数的信息，那么判断标准可以有所简化。

例子：凸函数(convex function)的等价刻画

- Jensen's inequality:

$$\forall t \in [0,1], \forall x, y \in \mathbb{R}^n, f(tx + (1-t)y) \leq tf(x) + (1-t)f(y)$$

- 一阶条件（对可导的 f ）：

$$\forall x, y \in \mathbb{R}^n, f(x) \geq f(y) + \nabla f(y) \cdot (x - y)$$

- 二阶条件（对二阶可导的 f ）：

$$\forall x \in \mathbb{R}^n, H_f(x) \succcurlyeq 0$$



Steepest descent

给定连续可导 $f: \mathbb{R}^n \rightarrow \mathbb{R}$, 最小化 f

使用函数求根的方法, 找极值点

- $\nabla f: \mathbb{R}^n \rightarrow \mathbb{R}^n$, 需要找到 $\nabla f = 0$

梯度下降方法:

- 在当前的点, 找到“最速下降”的方向 (steepest descent)
- $f(x) \approx f(x_0) + \nabla f(x_0) \cdot (x - x_0)$
- 假设 $f(x) = f(x_0) + \nabla f(x_0) \cdot (x - x_0)$
- “最速下降”的方向是什么?
 - 找到单位长度的向量 y , 使得 $\nabla f(x_0) \cdot y$ 长度最大化



Steepest descent

给定向量 \vec{a} ，求单位长度的向量 \vec{b} ，使得 $|\langle \vec{a}, \vec{b} \rangle|$ 最大化。

Cauchy-Schwarz不等式:

$$|\langle \vec{a}, \vec{b} \rangle| \leq \|\vec{a}\|_2 \|\vec{b}\|_2,$$

当且仅当存在乘数(scalar)使得 $\vec{a} = c \vec{b}$ 时，等号可以取到

- 假设 $f(x) = f(x_0) + \nabla f(x_0) \cdot (x - x_0)$
- “最速下降”的方向是什么？
 - 找到单位长度的向量 y ，使得 $\nabla f(x_0) \cdot y$ 长度最大化
 - $y = \alpha \nabla f(x_0)$



回到Richardson iteration

一个特殊情形: 给定实数 $0 < a < 1$, 如何用 a 表示出 $a^{-1}b$?

$$b + (1-a)b + (1-a)^2b + \dots = \frac{b}{1-(1-a)} = a^{-1}b$$

等价于 $x_0 = \mathbf{0}$, $x_{k+1} = (1-a)x_k + b$

Richardson iteration:

$$\begin{aligned} x_0 &= \mathbf{0} \\ \overrightarrow{x_{k+1}} &= (I - \alpha A)\overrightarrow{x_k} + \alpha \overrightarrow{b} = \overrightarrow{x_k} - \alpha(A\overrightarrow{x_k} - \overrightarrow{b}) \end{aligned}$$

矩阵多项式的视角: 特征空间: $p(A)v = p(\lambda)v$

- $A \rightarrow p(A)$
- 特征空间 $\{\lambda_1, \lambda_2, \dots, \lambda_n\} \rightarrow \{p(\lambda_1), p(\lambda_2), \dots, p(\lambda_n)\}$

另一个视角: 对于对称的 A , 这正好是目标函数 $\frac{1}{2}x^T Ax - b^T x$ 的梯度下降方法!

$$\nabla \left(\frac{1}{2}x^T Ax - b^T x \right) = \frac{1}{2}(A + A^T)x - b = A\overrightarrow{x_k} - \overrightarrow{b}$$



Richardson iteration

$$A^{-1} = p(A)$$

可否近似 $p(A)$?

等价地, 记 $q(x) = 1 - xp(x)$; 需要寻找 $q(0) = 1, q(x) \approx 0, \forall x > 0$

要解 $Ax = b$, 即找出 $x_* = p(A)b \in \text{span}\{b, Ab, A^2b, \dots\}$

Richardson iteration: 设 $x_0 = 0$,

$$x_{k+1} = (I - \alpha A)x_k + \alpha b$$

误差 $e_k = x_k - x_*$ 满足 $e_k = (I - \alpha A)e_{k-1} = (I - \alpha A)^k e_0$

收敛性: 当且仅当 $\rho(I - \alpha A) = \max\{|1 - \alpha\lambda_1|, |1 - \alpha\lambda_n|\} < 1$

可以令 $1 - \alpha\lambda_1 = -(1 - \alpha\lambda_n)$ 即 $\alpha = \frac{2}{\lambda_1 + \lambda_n}$

此时 $\rho(I - \alpha A) = \frac{\lambda_n - \lambda_1}{\lambda_1 + \lambda_n}$, 收敛需要 $O\left(\left(1 + \frac{\lambda_n}{\lambda_1}\right) \cdot \log \frac{1}{\epsilon}\right)$ 步

回顾: $\frac{\lambda_n}{\lambda_1}$ 正是条件数!



Richardson iteration

$$A^{-1} = p(A)$$

可否近似 $p(A)$?

等价地, 记 $q(x) = 1 - xp(x)$; 需要寻找 $q(0) = 1, q(x) \approx 0, \forall x > 0$

要解 $Ax = b$, 即找出 $x_* = p(A)b \in \text{span}\{b, Ab, A^2b, \dots\}$

Richardson iteration: 设 $x_0 = 0$,

$$x_{k+1} = (I - \alpha A)x_k + \alpha b$$

注意, 这里的 $x_k = p_k(A)b \in \text{span}\{b, Ab, A^2b, \dots, A^{k-1}b\}$; Krylov子空间

可以证明: 相当于考虑 $q(x) = 1 - xp(x) = (1 - \alpha x)^k$

$$-Ae_k = b - Ax_k = (I - Ap_k(A))b = q_k(A)b$$

$$-e_k = x_* - x_k = (I - p_k(A)A)x_* = q_k(A)x_*$$



Chebyshev iteration (选讲)

$$A^{-1} = p(A)$$

目标: 近似 $p(A)$

等价地, 记 $q(x) = 1 - xp(x)$; 需要寻找 $q(0) = 1, q(x) \approx 0, \forall x > 0$

要解 $Ax = b$, 即找出 $x_* = p(A)b \in \text{span}\{b, Ab, A^2b, \dots\}$

Chebyshev iteration:

$$x_{k+1} = (I - \alpha_k A)x_k + \alpha_k b$$

误差 $e_k = x_k - x_*$ 满足 $e_k = \prod_i (I - \alpha_i A) e_0$

$$\|e_k\| \leq \left\| \prod_i (I - \alpha_i A) \right\| \|e_0\|$$

要最小化 $\|\prod_i (I - \alpha_i A)\|$: Chebyshev 多项式!

注意, 这里的 $x_k \in \text{span}\{b, Ab, A^2b, \dots\}$



回到Richardson iteration: Steepest descent

Richardson iteration:

$$\begin{aligned} \mathbf{x}_0 &= \mathbf{0} \\ \overrightarrow{\mathbf{x}}_{k+1} &= (I - \alpha A)\overrightarrow{\mathbf{x}}_k + \alpha \overrightarrow{\mathbf{b}} = \overrightarrow{\mathbf{x}}_k - \alpha(A\overrightarrow{\mathbf{x}}_k - \overrightarrow{\mathbf{b}}) \end{aligned}$$

对于对称的 A , 这正好是目标函数 $\frac{1}{2}\mathbf{x}^T A\mathbf{x} - \mathbf{b}^T \mathbf{x}$ 的梯度下降方法!

$$\nabla \left(\frac{1}{2}\mathbf{x}^T A\mathbf{x} - \mathbf{b}^T \mathbf{x} \right) = \frac{1}{2}(A + A^T)\mathbf{x} - \mathbf{b} = A\overrightarrow{\mathbf{x}}_k - \overrightarrow{\mathbf{b}}$$

目标函数 $\frac{1}{2}\mathbf{x}^T A\mathbf{x} - \mathbf{b}^T \mathbf{x}$ 是convex的: A 是正定的

利用strong convexity, 也可以证明梯度下降的收敛时间正比于条件数 $\frac{\lambda_n}{\lambda_1}$

参考: Convex Optimization by Stephen Boyd and Lieven Vandenberghe



Conjugate gradients

给定正定矩阵 A ，可定义内积：

$$\langle x, y \rangle_A := x^T A y$$

- 线性
- 对称性
- 正定性 $\langle x, x \rangle_A > 0 \quad \forall x \neq 0$

如果 $\langle x, y \rangle_A = 0$ 则称 x, y 关于 A 共轭

内积自带的范数： $\|x\|_A^2 := x^T A x$

令 x_* 满足 $Ax_* = b$ ，考虑关于 x 的函数：

$$\frac{1}{2} x^T A x - b^T x = \frac{1}{2} \|x - x_*\|_A^2 + \text{constant}$$



Conjugate gradients

给定正定矩阵 A ，可定义内积：

$$\langle x, y \rangle_A := x^T A y$$

性质：共轭向量一定是线性无关的。因此可以组成一组基。

假设 p_1, p_2, \dots, p_k 互相共轭，并且是线性相关的：存在 $\alpha_1, \dots, \alpha_k$ 使得：

$$\alpha_1 p_1 + \dots + \alpha_k p_k = 0$$

另一方面

$$\begin{aligned} & (\alpha_1 p_1 + \dots + \alpha_k p_k)^T A (\alpha_1 p_1 + \dots + \alpha_k p_k) \\ &= \alpha_1^2 \|p_1\|_A^2 + \alpha_2^2 \|p_2\|_A^2 + \dots + \alpha_k^2 \|p_k\|_A^2 > 0 \end{aligned}$$

与上式矛盾



Conjugate gradients

$$A^{-1} = p(A)$$

目标：近似 $p(A)$

等价地，记 $q(x) = 1 - xp(x)$ ；需要寻找 $q(0) = 1$ ， $q(x) \approx 0, \forall x > 0$

要解 $Ax = b$ ，即找出 $x_* = p(A)b \in \text{span}\{b, Ab, A^2b, \dots\}$

考虑所有这样的迭代算法：

记 Krylov 子空间 $K_0 = \{0\}$, $K_i = \text{span}\{b, Ab, \dots, A^{i-1}b\}$

设 $x_i = \arg \min_{x \in K_i} \|x - x_*\|_A^2$ ，其中 x_* 满足 $Ax_* = b$

引理：记 $v_i := x_i - x_{i-1}$ ，则 $\{v_i\}$ 两两共轭： $v_i^T Av_j = 0, \forall i \neq j$



Conjugate gradients

记 Krylov 子空间 $K_0 = \{0\}$, $K_i = \text{span}\{b, Ab, \dots, A^{i-1}b\}$

设 $x_i = \arg \min_{x \in K_i} \|x - x_*\|_A^2$, 其中 x_* 满足 $Ax_* = b$

引理: 记 $v_i := x_i - x_{i-1}$, 则 $\{v_i\}$ 两两共轭: $v_i^T Av_j = 0, \forall i \neq j$

证明: 假设 $i < j$. 注意到 x_j 的定义是 K_j 中最小化 $\|x - x_*\|_A^2$ 的, 所以 $\nabla \frac{1}{2} \|x_j - x_*\|_A^2 = Ax_j - b$ 必须与 K_j 正交, 因此与 K_{j-1} 正交。类似地 $Ax_{j-1} - b$ 必须与 K_{j-1} 正交。因此,

$$Av_j = Ax_j - Ax_{j-1} \in K_{j-1}^\perp$$

而 $v_i \in K_i \subset K_{j-1}$, 因此 $v_i^T Av_j = 0$

一个推论: $K_i = \text{span}\{v_1, v_2, \dots, v_i\}$



Conjugate gradients

记 Krylov 子空间 $K_0 = \{0\}$, $K_i = \text{span}\{b, Ab, \dots, A^{i-1}b\}$

设 $x_i = \arg \min_{x \in K_i} \|x - x_*\|_A^2$, 其中 x_* 满足 $Ax_* = b$

引理: 记 $v_i := x_i - x_{i-1}$, $r_i := b - Ax_i$, 则

$$v_i = \frac{v_i^T r_{i-1}}{\|r_{i-1}\|^2} \left(r_{i-1} - \frac{r_{i-1}^T A v_{i-1}}{v_{i-1}^T A v_{i-1}} v_{i-1} \right)$$

证明(sketch): 首先注意到 $Ax_{i-1} - b \in K_{i-1}^\perp$, 但是 $Ax_{i-1} - b \in K_i$, 因此 $K_i = \text{span}\{v_1, v_2, \dots, v_{i-1}, r_{i-1}\}$

换言之, 可以写出 $v_i = c_0 r_{i-1} + \sum_{j=1}^{i-1} c_j v_j$

要确定 c_0 , 考虑 $v_i^T r_{i-1} = c_0 \|r_{i-1}\|^2$

要确定 c_{i-1} , 注意到 $v_j^T A v_i = 0, \forall j < i - 1$, 令 $v_{i-1}^T A v_i = 0$ 即可

类似地, 令 $v_j^T A v_i = 0$ 可得到其它的 $c_j = 0$



Conjugate gradients

记 Krylov 子空间 $K_0 = \{0\}, K_i = \text{span}\{b, Ab, \dots, A^{i-1}b\}$

设 $x_i = \arg \min_{x \in K_i} \|x - x_*\|_A^2$, 其中 x_* 满足 $Ax_* = b$

引理: 记 $v_i := x_i - x_{i-1}, r_i := b - Ax_i$, 则

$$v_i = \frac{v_i^T r_{i-1}}{\|r_{i-1}\|^2} (r_{i-1} - \frac{r_{i-1}^T A v_{i-1}}{v_{i-1}^T A v_{i-1}} v_{i-1})$$

简化: 令 $d_i = \frac{\|r_{i-1}\|^2}{v_i^T r_{i-1}} v_i$,

即有 $x_i = x_{i-1} + \frac{\|r_{i-1}\|^2}{d_i^T A d_i} d_i, d_i = r_{i-1} + \frac{\|r_{i-1}\|^2}{\|r_{i-2}\|^2} d_{i-1}$

可以证明: conjugate gradients 在运行 k 步之后, 误差最多是正比于从所有 $\deg(q) < k$ 且 $q(0) = 1$ 的多项式中, 最小的 $\|q\|_\infty$

事实上:

$$\|x_k - x_*\|_A^2 \leq \inf_{q(0)=1, \deg(q) \leq k} \max_i q(\lambda_i)^2 \cdot \|b\|_{A^{-1}}^2$$

与高斯消元相比: 对于稀疏矩阵, 恰当运用迭代法可以收敛更快
最多迭代 $n+1$ 次

Conjugate Gradient Method

$x_0 =$ initial guess

$d_0 = r_0 = b - Ax_0$

for $k = 0, 1, 2, \dots, n - 1$

if $r_k = 0$, stop, end

$$\alpha_k = \frac{r_k^T r_k}{d_k^T A d_k}$$

$$x_{k+1} = x_k + \alpha_k d_k$$

$$r_{k+1} = r_k - \alpha_k A d_k$$

$$\beta_k = \frac{r_{k+1}^T r_{k+1}}{r_k^T r_k}$$

$$d_{k+1} = r_{k+1} + \beta_k d_k$$

end