

共轭梯度下降与图直径估计

本文的主要结论参考自 F. R. K. Chuang, V. Faber, and T. A. Manteuffel. On the diameter of a graph from eigenvalues associated with its Laplacian. *SIAM Journal on Discrete Mathematics*, 7:443-457, 1994, 主要内容的呈现逻辑参考 D. A. Spielman, *Spectral and Algebraic Graph Theory*.

在这里, 我们将呈现优化理论的技巧将如何应用于图论问题当中. 为此, 第一步是将图问题转化成矩阵问题; 第二步则是刻画求解矩阵逆的方法, 并且将其应用到特定的矩阵系统中. 我们假定读者具备基本的线性代数知识.

1 代数图论的一点概念

考虑加权图 $G = (V, E, w)$, 其中 $w : E \rightarrow \mathbb{R}$, 取定线性空间 $\text{span } V$, 并将其基记作 $V = \{v_1, v_2, \dots, v_n\}$. 此时, 定义邻接矩阵为

$$A = (w_{ij}), \quad w_{ij} = w(v_i, v_j) \text{ if } (v_i, v_j) \in E \text{ else } 0$$

这应当是我们在最开始学习图论的时候就已经熟知的. 但是, 这个形式的邻接矩阵显然是最没有意义的——它没有多少性质, 将其视作一个算子或者一个二次型的尝试都是徒劳的. 这意味着, 我们在线性代数所学的技巧对它来说毫无意义, 它无非只是一个数表而已.

那么, 怎么使得这个矩阵变成一个更有意义的矩阵呢? 我们在此展现将其变成一个二次型的方式¹. 考虑一个 \mathbb{R}^n 中的向量 x , 构建二次型:

$$x^T L x = \sum_{i,j \leq n} w_{ij} (x_i - x_j)^2$$

在稍作思考之后, 读者就会发现这是与这个图具备充分的联系的、最自然的二次型: 右边的差纯粹是为了关联两个分量, 如果不做差, 那么顶点之间的连接就无意义. 当然, 稍稍整理之后, 这个二次型的矩阵表示也并不复杂:

$$\begin{aligned} L &= \begin{pmatrix} \sum_{j=1}^n w_{1j} & -w_{12} & \dots & -w_{1n} \\ -w_{21} & \sum_{j=1}^n w_{2j} & \dots & -w_{2n} \\ \vdots & & \ddots & \vdots \\ -w_{n1} & \dots & & \sum_{j=1}^n w_{nj} \end{pmatrix} \\ &= \text{diag} \left(\sum_{j=1}^n w_{1j}, \sum_{j=2}^n w_{2j}, \dots, \sum_{j=1}^n w_{nj} \right) - A =: D - A \end{aligned}$$

容易发现它就是一个实对称矩阵, 而且, 如果这个图具备正边权, 那么它就一定是正定的——这将为我们的讨论奠定基础. 最后是我们讨论的目标:

定义 1 一个图的直径 (diameter) 指的是

$$d(G) := \max_{v_1, v_2 \in V} \max_{\substack{p \text{ a path} \\ \text{from } l_1 \text{ to } l_2}} l(p).$$

¹另一种将其视作算子的方式更加显而易见, 就是扩散矩阵.

2 矩阵（伪）逆的迭代求解法

设 V 是一个 \mathbb{F} 上的线性空间. 对线性代数熟悉的读者应当熟知以下结果:

引理 1 如果 $T \in \text{End } V$ 可逆, 那么存在 $g \in \mathbb{F}[X]$ 使得 $T^{-1} = g(T)$.

这个引理无非是将 $\text{End } V$ 视作线性空间之后的反证, 因此, 它给出的限制是 $\deg g \leq n^2, n = \dim V$. 我们期待, 对于不可逆的矩阵 M , 也能给出类似的多项式用以方便矩阵逆的求解, 或者至少是一种逼近. 下面我们只考虑 Penrose-Moore 逆的情形, 记作 M^\dagger .

2.1 一阶 Richardson 迭代和条件数

让我们依旧从 $Ax = b$ 的求解开始. 注意到, 此时我们有:

$$\alpha Ax = \alpha b \Rightarrow x + (\alpha A - I)x = \alpha b \Rightarrow x = (I - \alpha A)x + \alpha b, \forall \alpha \in \mathbb{F}$$

出于简化问题的考虑, 下面直接设 A 是对称正定的, 那么求解 $Ax = b$ 的问题无非就是寻找二次函数

$$f(x) = x^T Ax - bx$$

的极小值点的问题, 而上面的式子给出的无非就是最速下降法给出的速度. 因此, 分析一阶 Richardson 迭代的敛性其实等价于分析这种最速下降法的敛性. 我们同样熟知的结果是, 这种最速下降法的敛性与条件数有关:

定理 2 设 A 的特征值为 $0 < \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$, 则对于任意 ε 和

$$t > \frac{1}{2} \left(\frac{\lambda_n}{\lambda_1} + 1 \right) \ln \left(\frac{1}{\varepsilon} \right)$$

都有 $\|x^* - x_t\| \leq \varepsilon \|x^* - x_0\|$, 其中 x^* 为真实解, x_t 为第 t 次迭代之后的解, x_0 为预设的迭代起始值. $\kappa(A) := \lambda_n/\lambda_1$ 被称作矩阵 A 的条件数 (condition number)².

证明 直接计算

$$\begin{aligned} x^* - x_t &= ((I - \alpha A)x^* + \alpha b) - ((I - \alpha A)x_{t-1} + \alpha b) \\ &= (I - \alpha A)(x^* - x_{t-1}) = (I - \alpha A)^t(x^* - x_0). \end{aligned}$$

因此, 我们需要研究的是 $I - \alpha A$ 的谱. 不难发现 $I - \alpha A$ 的特征值就是 $\{1 - \alpha\lambda_i\}, i = 1, 2, \dots, n$, 因此, 它的最大特征值就是

$$|\max(1 - \alpha\lambda_1, 1 - \alpha\lambda_n)| \geq \left| \frac{\lambda_n - \lambda_1}{\lambda_n + \lambda_1} \right|$$

极值当 $\alpha = 2/(\lambda_n + \lambda_1)$ 时取得. 因此

$$\|x^* - x_t\| \leq \left(\frac{\lambda_n - \lambda_1}{\lambda_n + \lambda_1} \right)^t \|x^* - x_0\| \leq e^{-2t \frac{\lambda_1}{\lambda_n + \lambda_1}} \|x^* - x_0\|$$

因此

²对实对称矩阵来说如此, 否则它定义为最大奇异值和最小奇异值的商.

$$t > \frac{1}{2} \left(\frac{\lambda_n}{\lambda_1} + 1 \right) \ln \left(\frac{1}{\varepsilon} \right)$$

时原不等式一定成立。

■

这个算法称不上是什么相当好的方法，但是它给了我们一个相当有趣的多项式近似手段。稍微多写几项，我们发现：

$$x_t = \sum_{i=0}^t (I - \alpha A)^i \alpha b$$

而且，下面的式子在收敛的情况下也是显然成立的，这为我们所熟知：

$$\alpha \sum_{i=0}^{\infty} (I - \alpha A)^i = \alpha (I - (I - \alpha A))^{-1} = \alpha (\alpha A)^{-1} = A^{-1}$$

也就是说，本质上我们只是找到了一个 Taylor 展开，然后利用这个 Taylor 级数作截断，最后实现了一个多项式水平的近似。

2.2 初见端倪：Chebyshev 下降法

接下来的问题是，怎么找到更好的近似结果。为了获得更好的近似结果，就要找到更好的多项式。注意到，我们希望获得：

$$\|p(A)b - x\| \leq \varepsilon \Rightarrow \|p(A)A - I\| \leq \varepsilon$$

由于 A 是对称正定的，所以 $p(A)A - I$ 也是对称的，因此它的范数无非就是它的最大特征值的绝对值。于是，我们希望找到一个多项式使得

$$|\lambda_i p(\lambda_i) - 1| \leq \varepsilon, \forall i = 1, 2, \dots, n$$

进一步地，定义 $q(x) = 1 - xp(x)$ ，则我们可以彻底将这个问题转变成为一个多项式零点的分布问题：

$$q(0) = 1, \quad |q(\lambda_i)| \leq \varepsilon, \quad \forall i = 1, 2, \dots, n$$

更粗糙的假设是，干脆让它在 $[\lambda_1, \lambda_n]$ 上整体的绝对值受限。这时，我们在小学二年级就学过的 Chebyshev 多项式就给了我们一个很好的直观，定义

$$T_i(x) = \frac{1}{2} \left[\left(x + \sqrt{x^2 - 1} \right)^i + \left(x - \sqrt{x^2 - 1} \right)^i \right]$$

它在 $[-1, 1]$ 上的行为是高度受控的，且在这个范围之外的增长飞快。因此，我们用它来构造多项式：

$$q_i(x) = T_i \left(\frac{\lambda_n + \lambda_1 - 2x}{\lambda_n - \lambda_1} \right) / T_i \left(\frac{\lambda_n + \lambda_1}{\lambda_n - \lambda_1} \right)$$

于是，此时的 ε 无非就是讨论分母项有多大。一个简单的估计告诉我们：

$$T_i \left(\frac{\lambda_n + \lambda_1}{\lambda_n - \lambda_1} \right) = T_i \left(\frac{\kappa + 1}{\kappa - 1} \right) = T_i(\cosh \theta)$$

这里反双曲函数的技巧在处理 Chebyshev 多项式相关的式子时是常规的. 我们求解

$$\frac{\kappa + 1}{\kappa - 1} = \frac{e^\theta + e^{-\theta}}{2} \Rightarrow e^\theta = \frac{\sqrt{\kappa} \pm 1}{\sqrt{\kappa} \mp 1}$$

于是

$$T_i(\cosh \theta) = \cosh i\theta = \frac{1}{2} \left[\left(\frac{\sqrt{\kappa} + 1}{\sqrt{\kappa} - 1} \right)^i + \left(\frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^i \right] \geq \frac{1}{2} \left(\frac{\sqrt{\kappa} + 1}{\sqrt{\kappa} - 1} \right)^i.$$

因此, 我们指出, 利用 Chebyshev 多项式 T_i 做逼近时, 有

$$\varepsilon \leq 2 \left(\frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^i.$$

因为 Chebyshev 多项式可以迭代计算, 所以实际上, 它就给了我们一个与 $\sqrt{\kappa}$ 相关的收敛速率. 对比一下一阶 Richardson 方法, 它的收敛速率无疑与之相比获得了二次的提升.

2.3 共轭梯度下降

当然, 对于一阶 Richardson 方法的讨论当中, 我们实际上构建了下面的空间, 这将带给我们一些额外的直观:

定义 2 取 $x \in \mathbb{R}^n, A \in \text{Mat}_n \mathbb{R}$, 定义

$$K_m(x, A) = \text{span}\{x, Ax, \dots, A^{m-1}x\}$$

为由 A 和 x 张成的 m 阶 Krylov 子空间.

在我们的证明当中, 不难发现, 所有的 x_t 都落在 $x_0 + K_t(x^* - x_0, I - \alpha A) = x_0 + K_t(x^* - x_0, A)$ 当中. 因此, 接下来需要探讨的就是怎么选取这个仿射子空间中的元素, 使得我们在其中获得最小的残差 $r_t = b - Ax_t$.

注记 1 观察上面的 Richardson 方法, 不难发现其中我们每次下降给出的残差是在直线

$$l: x_{t-1} + \alpha(b - Ax_{t-1})$$

上的残差之 2-范数最小者, 但是, 它不一定在 Krylov 子空间当中全局最小.

现在, 鉴于 A 是一个对称正定的矩阵, 我们有一个非常轻松的度量残差大小的方法, 就是 A -范数 $\|x\|_A := \sqrt{x^T A x}$. 对称性和正定性告诉我们它确实是一个合理的范数, 而且它自然地让我们能够给出内积结构 $\langle x, y \rangle_A = x^T A y$ ——众所周知, 正交性在任何 Hilbert 空间的优化问题中都是一个相当好的性质, 它几乎能够保证最佳的收敛性.

基于上面的观察, 我们尝试在 Krylov 子空间中寻找全局最小的、以 A -范数度量的残差. 取 $t+1$ 阶 Krylov 子空间的一组 A -正交基 $\{v_0, v_1, \dots, v_t\}$, 令 $x_t = \sum_{i=0}^t c_i v_i$, 则优化问题就是:

$$\frac{1}{2} x_t^T A x_t - b^T x_t = \frac{1}{2} \left(\sum_{i=0}^t c_i v_i \right)^T A \left(\sum_{i=0}^t c_i v_i \right) - b^T \left(\sum_{i=0}^t c_i v_i \right)$$

$$= \frac{1}{2} \sum_{i=0}^t c_i^2 v_i^T A v_i - \sum_{i=0}^t c_i b^T v_i$$

注意到，取 A -正交基在此的益处就在于把交叉项消除了，从而我们可以只考察如何极小化每一个 i 单独对应的项，只需将它关于 c_i 的导数置为 0，从而就有

$$c_i = \frac{b^T v_i}{v_i^T A v_i}$$

这就意味着，只要取出这样一组正交基，梯度下降的方向就是显然的。而取得这样的正交基的过程就是我们熟知的 Gram-Schmidt 正交化：

$$v_{t+1} = A v_t - \sum_{i=0}^t v_i \frac{(A v_t)^T A v_i}{v_i^T A v_i} = A v_t - v_t \frac{(A v_t)^T (A v_t)}{v_t^T A v_t} - v_{t-1} \frac{(A v_t)^T (A v_{t-1})}{v_{t-1}^T A v_{t-1}}$$

其中其余项都因为正交性显然地消失了。于是我们的迭代项就已经被完美解决了。

2.4 收敛性分析

注意到，我们是在 Krylov 子空间中寻找到的 A -范数度量的最小残差，而 Krylov 子空间的起始项正好是 $x^* - x_0$ ，于是，根据 Krylov 子空间的定义，应当有

$$x^* - x_t = \left(I + \sum_{i=1}^t c_i A^i \right) (x^* - x_0) = p_t(A) (x^* - x_0)$$

其中 p_j 是一个 j 次多项式。接下来，我们考虑以 A 的特征向量作为一组标准正交基 $\{\eta_1, \eta_2, \dots, \eta_n\}$ ，对应的特征值为 $0 < \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$ ，这时我们有 $p_t(A)\eta_j = p_t(\lambda_j)\eta_j$ ，

$$\|x^* - x_t\|_A^2 = \sum_{i=1}^n \xi_j^2 p_t(\lambda_i)^2 \lambda_i \leq \min_{p_t} \max_{\lambda} p_t(\lambda)^2 \|x^* - x_0\|_A^2$$

如此，这个优化问题化规到了 Chebyshev 下降的情形，上面我们给出的 q_i 恰好就构成了这样的优化问题的解。实际上，读者不难发现，有这么一个多项式：

$$q(x) = \frac{\prod_{i=1}^n (\lambda_i - x)}{\prod_{i=1}^n \lambda_i}$$

正好满足了我们在 Chebyshev 下降当中定义的 $q(x)$ 的要求，且 $\varepsilon = 0$ ，也就是说，利用这个多项式获得的解应当是完美的解——这无非就是 Cayley-Hamilton 定理的一种重述。也因为这个多项式是 n 次的，所以我们搜索到第 n 次时， x^* 一定已经处于 n 阶 Krylov 子空间中，这时，我们已经给出了一个完美的解。

注记 2 读者同样不难注意到，如果特征值的分布是均匀的，那么 Chebyshev 多项式实际上恰好就给出了这样一个完美解，因此，特征值的均匀性对优化算法实际上也起到了重要的作用，我们在后面讨论预处理时还会再次碰到这一点。

3 图半径估计和多项式

下面我们将这种视角应用于一个图的 Laplacian 矩阵 L . 当然, 这个矩阵是奇异的, 所以它的特征值中有 0. 但是, 这时我们只要考虑它的 Penrose-Moore 伪逆即可, 0 特征值对应的特征向量是 $(1, 1, \dots, 1)^T$, 所以我们考虑:

$$\|LL^\dagger - P\| \leq \varepsilon, \quad P \text{ is the projection to the complement of the span of } (1, 1, \dots, 1)^T$$

此时当然 $\lambda_1 = 0$ 无意义, 所以只要考虑 λ_2 来代替原来的 λ_1 即可, 我们注意到, 根据共轭梯度下降的构造 (或者说, 根据 Cayley-Hamilton 定理), 存在一个多项式 p , 使得 $\deg p = k - 1$, k 为 L 中互不相同的特征值的个数, 满足 $p(L) = L^\dagger$.

直观上讲, 直径越小的图当然是越简单的, 而它所对应的 Laplacian 构成的线性系统也当然应当更容易求解, 实际上:

引理 3 设 $G = (V, E)$ 是一个具备正边权的连通无向图, 其对应的 Laplacian 具备 k 个不同的特征值, 则这个图的直径不大于 k .

证明 首先我们刻画直径. 注意到, 直径最简单的刻画方式是用邻接矩阵:

$$d(G) = \min\{n : \forall i \neq j, \exists A^m \text{ with } 0 \leq m \leq n, \text{ s.t. } A_{ij}^m \neq 0\}$$

其中, 存在一个 A^m 使得 $A_{ij}^m \neq 0$ 无非就是说 i, j 之间存在一个长为 m 的通路. 不难发现, 这个条件等价于对于任意 $i \neq j$, 存在一个 m 阶多项式 p 满足 $0 \leq m \leq n$, 使得 $p(A)_{ij} \neq 0$. 通过对所有这样的多项式求和, 我们表明:

$$d(G) = \min\{n : \forall i \neq j, \exists p(x) \text{ with degree } n, \text{ s.t. } p(A)_{ij} \neq 0\}.$$

另外, 因为对角阵的幂次仍然是对角阵, 所以上面的刻画当中, 我们可以把 A 换成 L , 使得这个结论成立. 而注意到, 如果 s, t 是距离 d 的两个顶点, 那么 $P_{st} = -\frac{1}{n} \neq 0$, 而我们有 $LL^\dagger = P = Lp(L)$, 如果 p 的次数小于 d 的话, $(Lp(L))_{st}$ 就只能为 0, 因为二者之间没有长度小于 d 的通路. 这与前面的假设矛盾, 所以 $p \geq d$ 对于任意的 s, t 都成立, 即 $p \geq d(G)$. ■

基于类似的想法, 因为我们有 ε 估计, 所以我们有以下定理:

定理 4 (Chung-Faber-Manteuffel) 设 $G = (V, E)$ 是一个具备正边权的连通无向图, 其对应的 Laplacian 特征值为 $0 \leq \lambda_2 \leq \dots \leq \lambda_n$, 则

$$d(G) \leq \left(\frac{1}{2} \sqrt{\frac{\lambda_n}{\lambda_2}} + 1 \right) \ln 2n$$

证明 在 Chebyshev 下降当中, 我们要求多项式幂次满足

$$\frac{1}{n} \leq 2 \left(\frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^k$$

在这里，取 $\kappa = \lambda_n/\lambda_2$ ，然后直接利用估计

$$k \leq \frac{\ln(2n)}{\ln\left(\frac{\sqrt{\kappa}-1}{\sqrt{\kappa}+1}\right)} \leq \left(\frac{1}{2}\sqrt{\frac{\lambda_n}{\lambda_2}} + 1\right) \ln 2n.$$

其中右侧的估计应当是我们在高中就熟知的。 ■

注记 3 这个结论实际上可以推广到可以一笔画的 Euler 图，因为它的每个顶点入度和出度相等，所以处理手法基本上是类似的。

4 餐后甜点：预处理

最后，作为餐后甜点，我们讨论优化的预处理技巧 (preconditioning)，以及预处理技巧从图论中可以学到什么。首先，我们定义：

定义 3 记矩阵 A 和矩阵 B 的相对条件数

$$\kappa(A, B) = \frac{\beta}{\alpha}$$

其中 β 是使得 $\beta B - A$ 半正定的最小正数， α 是使得 $A - \alpha B$ 半正定的最大正数。

为了简明起见，记 $A \preceq B$ ，如果 $B - A$ 半正定。注意到上面的定义实际上就等价于

引理 5 $\kappa(A, B) = \kappa(B^{-1}A)$ 。

证明 ... ■

下面我们考虑一个比较合适的预处理器，能够使得

$$(1 - \varepsilon)B \preceq A \preceq (1 + \varepsilon)B$$

此时我们注意到在 A -范数的意义下，有：

$$\begin{aligned} \|B^{-1}b - x\|_A &= \|A^{\frac{1}{2}}B^{-1}b - A^{\frac{1}{2}}x\| \\ &= \|A^{\frac{1}{2}}B^{-1}Ax - A^{\frac{1}{2}}x\| \\ &= \|A^{\frac{1}{2}}B^{-1}A^{\frac{1}{2}}(A^{\frac{1}{2}}x) - A^{\frac{1}{2}}x\| \\ &\leq \|A^{\frac{1}{2}}B^{-1}A^{\frac{1}{2}} - I\| \|A^{\frac{1}{2}}x\| \\ &\leq \varepsilon \|A^{\frac{1}{2}}x\| = \varepsilon \|x\|_A \end{aligned}$$

因此，只要 B 的选取足够好，我们就可以借助对 B 对应的线性方程的求解来解出方程的解。同时，如果我们想进一步优化这个结果，只要注意到

$$x_1 = B^{-1}b, \quad r_1 = A(x - x_1) = b - Ax_1$$

再次迭代求解，就能利用 B 解出 $x_2 = B^{-1}r_1$ ，这时：

$$\|x - x_1 - x_2\|_A \leq \varepsilon \|x - x_1\|_A \leq \varepsilon^2 \|x\|_A$$

也就是说，每次迭代都实现了指数级别的优化。那么，怎么给出这样一个比较好的预处理器就成了我们需要讨论的问题。对于 Laplacian 矩阵，Vaidya 给出了一个洞见：对于任意的子图 $H \leq G$ ，都有 $L_H \preceq L_G$ 。于是，我们讨论的问题就变成了：

如何找到一个子图 H 使得 L_H 很容易取逆，且 $L_H^{-1}L(G)$ 不会太大？

一种回答是树，自然地我们也就会想到这个图的最小生成树。这会让我们讨论图的有效阻抗 (effective resistance)，并且估计树对应的展宽 (stretch)。基于此，我们能够诞生非常多有趣的子图预处理器，见 Michael B. Cohen et al., *Solving sdd linear systems in nearly $m \log 1/2n$ time*, STOC'14.

另一种回答是稀疏化 (sparsification)，...