

# 数値計算備忘録

motchy

2022 年 8 月 23 日 ~ 2024 年 9 月 24 日

ver 0.1.0

---

# 目次

第 1 部	固定小数点数	2
第 1.1 章	定義	4
1.1.1	前提	4
1.1.2	構造	4
1.1.3	対応する有理数	4
1.1.4	特性に関する写像	4
1.1.5	積	5
1.1.5.1	定義	5
1.1.5.2	性質	5
第 2 部	線形代数	6
第 2.1 章	Gauss-Seidel 法	7
2.1.1	定義	7
2.1.2	係数行列が狭義優対角ならば厳密解に収束すること	7
第 2.2 章	Cholesky 分解	9
2.2.1	rank-one update	9
第 2.3 章	LDL 分解	11
2.3.1	rank-one update	11

## 第 1 部

# 固定小数点数

## 1. 固定小数点数

---

この部は現時点では充実していないが、将来的に固定小数点数について多く論じる場合に備えて下地を整えておく。

## 第 1.1 章

# 定義

### 1.1.1 前提

1. 「2 の補数表現」など、大学の理工系の教養課程で扱われる内容は前提とし、本書では一々定義しない。
2. 基数は 2 である。

### 1.1.2 構造

固定小数点数とは、1 ビットの記憶領域 bit-cell の有限列と、その各領域に格納される数の組である（乱暴に言えば、容器と中身の両方を合わせたものである）。bit-cell の列は「整数部」と「小数部」という連続した部分配列に二分される。整数部は必ず存在する。小数部は存在しなくてもよい。整数部と小数部の bit-cell の個数（以降は「ビット数」と呼ぶ）をそれぞれ  $N_{\text{int}} (\geq 1)$ ,  $N_{\text{dec}} (\geq 0)$  とする。整数部と小数部の bit-cell の列は添え字が付けられる。整数部の添え字は  $0, 1, \dots, N_{\text{int}}$  である。小数部が存在するとき、その添え字は  $-1, -2, \dots, -N_{\text{dec}}$  である。整数部の最上位ビットを「符号ビット」と呼ぶ。この意味は後述される。整数部のビット数が  $N_{\text{int}}$ 、小数部のビット数が  $N_{\text{dec}}$  である形式を「 $N_{\text{int}}Q_{N_{\text{dec}}}$  形式」と呼ぶ。<sup>\*1</sup>

固定小数点数  $a$  と整数  $k$  について、 $k$  が  $a$  の添え字の範囲に含まれるとき、 $a[k]$  は第  $k$  bit-cell に格納されている数 (0 or 1) を表す。 $k$  が  $a$  の添え字の範囲外であるとき、 $a[k]$  は 0 であると定義する。

### 1.1.3 対応する有理数

$a$  を固定小数点数とする。符号ビットが 0 であるとき、「 $a$  に対応する有理数 (a rational number corresponding to  $a$ )」は次式で定義される。

$$\sum_{k=-N_{\text{dec}}}^{N_{\text{int}}-1} a[k]2^{k-N_{\text{dec}}} = \left( \sum_{k=-N_{\text{dec}}}^{N_{\text{int}}-1} a[k]2^k \right) / 2^{N_{\text{dec}}}$$

符号ビットが 1 であるとき、「 $a$  に対応する有理数」は次式で定義される。

$$-\left( 2^{N_{\text{int}}} - \sum_{k=-N_{\text{dec}}}^{N_{\text{int}}-1} a[k]2^{k-N_{\text{dec}}} \right) = -\left( 2^{N_{\text{int}}+N_{\text{dec}}} - \sum_{k=-N_{\text{dec}}}^{N_{\text{int}}-1} a[k]2^k \right) / 2^{N_{\text{dec}}}$$

以上の変換規則は、 $a$  の整数部と小数部を連結して 2 の補数表現の符号付き整数と見做したのに対応する整数を  $2^{N_{\text{dec}}}$  で割ったものと等しい。

### 1.1.4 特性に関する写像

固定小数点数  $a$  の整数部と小数部のビット数がそれぞれ  $N_{\text{int}}$ ,  $N_{\text{dec}}$  であるとする。次の写像を定義する。

1. 整数部のビット幅： $\bigcup_i a_i := N_{\text{int}}$

<sup>\*1</sup> ARM variant の Q format はこの形式を採用している。

2. 小数部のビット幅： $\underset{\text{d}}{a} := N_{\text{dec}}$
3. 全体のビット幅： $\underset{\text{d}}{a} := N_{\text{int}} + N_{\text{dec}}$
4. 対応する有理数： $\text{Rat}(a) := \text{a rational number corresponding to } a$

次に、有理数から固定小数点数への写像に関わるいくつかの写像を定義する。尚、有理数から固定小数点数への写像の手続きは 1.1.1 で述べたように既知であるとして説明を割愛する。有理数  $p$  を固定小数点数に写像する手続きを実行したものと、整数部のビットを添え字 0 から上に向かって見てゆくとき、ある添え字  $k$  以上では常に 0 となる。このような最小の  $k$  に対して  $k+1$  を以て「整数部の必要ビット数」と定義する。+1 は符号ビットを考慮したものである。また、小数部のビットを添え字 -1 から下に向かって見てゆくとき、次の 2 通りがあり得る：

1. ある添え字  $k$  以下では常にビットが 0 となる。
2. どのような添え字  $k$  をとっても、それより小さい添え字で 1 となるビットが存在する（所謂、無限小数）。

1 の場合、 $-k-1$  を以て「小数部の必要ビット数」と定義する。2 の場合、「小数部の必要ビット数」は無限大であると定義する。

## 1.1.5 積

### 1.1.5.1 定義

固定小数点数  $a, b$  について  $a$  の整数部と小数部のビット数がそれぞれ  $N_{\text{int},a}, N_{\text{dec},a}$  であり、 $b$  の整数部と小数部のビット数がそれぞれ  $N_{\text{int},b}, N_{\text{dec},b}$  であるとする。 $a$  と  $b$  の積  $a \times b$  を「 $\text{Rat}(a) \text{Rat}(b)$  を、整数部が  $N_{\text{int},a} + N_{\text{int},b}$  ビット、小数部が  $N_{\text{dec},a} + N_{\text{dec},b}$  ビットの固定小数点数に写像して得られる数」として定義する。

### 1.1.5.2 性質

前記の定義に従うと、常に  $\text{Rat}(a \times b) = \text{Rat}(a) \text{Rat}(b)$  が成り立つ。

*Proof.*

$a \times b$  の整数部のビット数  $N_{\text{int},a} + N_{\text{int},b}$  が  $\text{Rat}(a) \text{Rat}(b)$  の整数部の必要ビット数以上であり、かつ  $a \times b$  の小数部のビット数  $N_{\text{dec},a} + N_{\text{dec},b}$  が  $\text{Rat}(a) \text{Rat}(b)$  の小数部の必要ビット数以上であることを示せばよい。まず次式が成り立つ。

$$\text{Rat}(a) \text{Rat}(b) = \underbrace{2^{N_{\text{dec},a}} \text{Rat}(a)}_{=:\alpha} \underbrace{2^{N_{\text{dec},b}} \text{Rat}(b)}_{=:\beta} / 2^{N_{\text{dec},a} + N_{\text{dec},b}}$$

$N_a := N_{\text{int},a} + N_{\text{dec},a}$ ,  $N_b := N_{\text{int},b} + N_{\text{dec},b}$  とすると次式が成り立つ。

$$-2^{N_a-1} \leq \alpha \leq 2^{N_a-1} - 1, \quad -2^{N_b-1} \leq \beta \leq 2^{N_b-1} - 1$$

$$\therefore -2^{N_a+N_b-2} + \min\{2^{N_a-1}, 2^{N_b-1}\} \leq \alpha\beta \leq 2^{N_a+N_b-2}$$

$$\therefore -2^{N_{\text{int},a}+N_{\text{int},b}-2} + \min\{2^{N_{\text{int},a}-N_{\text{dec},b}-1}, 2^{N_{\text{int},b}-N_{\text{dec},a}-1}\} \leq \alpha\beta / 2^{N_{\text{dec},a}+N_{\text{dec},b}} \leq \text{Rat}(a) \text{Rat}(b) \leq 2^{N_{\text{int},a}+N_{\text{int},b}-2}$$

符号ビットを考慮して  $\text{Rat}(a) \text{Rat}(b)$  の整数部の必要ビット数は  $N_{\text{int},a} + N_{\text{int},b}$  以下である。また、 $\alpha\beta$  は整数であるから  $\text{Rat}(a) \text{Rat}(b)$  の小数部は  $\alpha\beta$  を  $2^{N_{\text{dec},a}+N_{\text{dec},b}}$  で割ったものの小数部と等しい。よって  $\text{Rat}(a) \text{Rat}(b)$  の小数部の必要ビット数は  $N_{\text{dec},a} + N_{\text{dec},b}$  以下である。□

第 2 部

線形代数

## 第 2.1 章

# Gauss-Seidel 法

### 2.1.1 定義

以下に述べる定義は Wikipedia の英語記事 [Gauss Seidel method](#) からの引用である。

$n \in \mathbb{N}$ ,  $A \in \mathbb{C}^{n \times n}$ ,  $\mathbf{b} \in \mathbb{C}^n$  とする。 $A$  は正定値対称、または狭義優対角であるとする。Gauss-Seidel 法とは、線型方程式  $A\mathbf{x} = \mathbf{b}$  の解を求める反復法である。 $\mathbf{x}_1 \in \mathbb{C}^n$  を任意の初期解とし、次の漸化式で解候補を更新してゆく。

$$L_* \mathbf{x}_{k+1} = -U \mathbf{x}_k \quad (k = 1, 2, \dots)$$

ここに  $L_*$  は  $A$  の対角成分およびその下側の要素からなる下三角行列であり、 $U$  は  $A$  の対角成分の上側の要素からなる上三角行列である。

### 2.1.2 係数行列が狭義優対角ならば厳密解に収束すること

*Proof.*

$A$  の次数を  $n$  とする。 $\hat{\mathbf{x}}$  を厳密解とすると  $L_* \hat{\mathbf{x}} = \mathbf{b} - U \hat{\mathbf{x}}$  である。これを解の更新式から減じると次式を得る。

$$L_*(\mathbf{x}_{k+1} - \hat{\mathbf{x}}) = -U(\mathbf{x}_k - \hat{\mathbf{x}}) \quad (1)$$

$\mathbf{v}_k := \mathbf{x}_k - \hat{\mathbf{x}}$  とおくと、式 (1) より次式が成り立つ。

$$L_* \mathbf{v}_{k+1} = -U \mathbf{v}_k \quad (2)$$

$M_k := \max_{i=1, \dots, n} |v_{k,i}|$  ( $v_{k,i}$  は  $\mathbf{v}_k$  の第  $i$  要素) とする。次の 2 つが同時に成り立つことが、 $\mathbf{v}_k$  が  $\mathbf{0}_n$  に収束するための十分条件である。

1. ある  $k \in \mathbb{N}$  に対して  $M_k = 0$  ならば  $M_l = 0$  ( $l = k+1, k+2, \dots$ )
2. 適当な  $0 < \alpha < 1$  が存在して  $M_k > 0 \Rightarrow M_{k+1} < \alpha M_k$

$L_*$  が正則であることと式 (2) より直ちに 1. が成り立つ。次に 2. を数学的帰納法で示す。 $\tilde{\alpha}$  を次式で定義する。

$$\tilde{\alpha} := \min_{i=1,2,\dots,n} \frac{1}{|a_{i,i}|} \sum_{j=1,\dots,n \wedge j \neq i} |a_{i,j}|$$

$A$  は優対角だから  $0 < \tilde{\alpha} < 1$  である。

$$\begin{aligned} a_{1,1} v_{k+1,1} &= - \sum_{j=2}^n a_{1,j} v_{k,j} \\ |a_{1,1}| |v_{k+1,1}| &= \left| \sum_{j=2}^n a_{1,j} v_{k,j} \right| \leq \sum_{j=2}^n |a_{1,j}| |v_{k,j}| \leq M_k \sum_{j=2}^n |a_{1,j}| \\ |v_{k+1,1}| &\leq \frac{M_k}{|a_{1,1}|} \sum_{j=2}^n |a_{1,j}| \leq \tilde{\alpha} M_k \end{aligned}$$



$|v_{k+1,j}| \leq \tilde{\alpha} M_k$  ( $j = 1, 2, \dots, l$ ) ( $l \in \{1, 2, \dots, n-1\}$ ) ならば  $|v_{k+1,l+1}| \leq \tilde{\alpha} M_k$ であることを示す。式 (2) の  $l+1$  行目を展開すると次式を得る。

$$\begin{aligned}
 \sum_{j=1}^{l+1} a_{l+1,j} v_{k+1,j} &= - \sum_{j=l+2}^n a_{l+1,j} v_{k,j} \\
 a_{l+1,l+1} v_{k+1,l+1} &= - \sum_{j=1}^l a_{l+1,j} v_{k+1,j} - \sum_{j=l+2}^n a_{l+1,j} v_{k,j} \\
 |a_{l+1,l+1}| |v_{k+1,l+1}| &= \left| - \sum_{j=1}^l a_{l+1,j} v_{k+1,j} - \sum_{j=l+2}^n a_{l+1,j} v_{k,j} \right| \leq \sum_{j=1}^l |a_{l+1,j}| |v_{k+1,j}| + \sum_{j=l+2}^n |a_{l+1,j}| |v_{k,j}| \\
 &\leq M_k \sum_{j=1, \dots, n \wedge j \neq l+1} |a_{l+1,j}| \\
 |v_{k+1,l+1}| &\leq \frac{M_k}{|a_{l+1,l+1}|} \sum_{j=1, \dots, n \wedge j \neq l+1} |a_{l+1,j}| \leq \tilde{\alpha} M_k
 \end{aligned}$$

以上より帰納的に  $|v_{k+1,j}| \leq \tilde{\alpha} M_k$  ( $j = 1, 2, \dots, n$ ) が成り立つ。すなわち  $M_{k+1} \leq \tilde{\alpha} M_k$  が成り立つ。 $\tilde{\alpha} < \alpha < 1$  となるように  $\alpha$  を定めることで 2. が示される。  $\square$

## 第 2.2 章

# Cholesky 分解

### 2.2.1 rank-one update

主張

$n \in \mathbb{N}$ ,  $A \in \mathbb{C}^{n \times n}$ ,  $A \succeq O$ ,  $\mathbf{x} \in \mathbb{C}^n$  とし、 $A$  は Hermite 行列であるとする。 $A + \mathbf{x}\mathbf{x}^*$  に対して Cholesky 分解のアルゴリズムを適用すると  $O(n^3)$  の計算量を要する。しかし、 $A$  の Cholesky 分解  $LL^*$  が既に得られているとき、 $A + \mathbf{x}\mathbf{x}^*$  の Cholesky 分解を  $O(n^2)$  で得ることができる。 $\mathbf{x}\mathbf{x}^*$  の階数が 1 以下である (特に 0 となるのは  $\mathbf{x} = \mathbf{0}$  の時かつその時に限る) ことから、この方法は “rank-one update” と呼ばれている。

**導出.** 方針としては、 $n \times n$  行列の rank-one update を  $(n-1) \times (n-1)$  行列の問題に帰着させ、以降同様に逐次的に行列の次数を縮小しながら解を構築する。このアルゴリズムの総計算量が  $O(n^2)$  となるのは明らかであろう。

$A + \mathbf{x}\mathbf{x}^*$  の Cholesky 分解を  $FF^*$  とする。 $L$  の第  $i$  列ベクトルを  $\mathbf{l}_i = [0, \dots, 0, l_{i,i}, \dots, l_{n,i}]^\top \in \mathbb{C}^{n \times n}$  とし、同様に  $F$  の第  $i$  列ベクトルを  $\mathbf{f}_i = [0, \dots, 0, f_{i,i}, \dots, f_{n,i}]^\top \in \mathbb{C}^{n \times n}$  とすると次式が成り立つ。

$$\begin{aligned} \sum_{i=1}^n \mathbf{f}_i \mathbf{f}_i^* &= \mathbf{x}\mathbf{x}^* + \sum_{i=1}^n \mathbf{l}_i \mathbf{l}_i^* \\ \mathbf{f}_1 \mathbf{f}_1^* + \sum_{i=2}^n \mathbf{f}_i \mathbf{f}_i^* &= \mathbf{x}\mathbf{x}^* + \mathbf{l}_1 \mathbf{l}_1^* + \sum_{i=2}^n \mathbf{l}_i \mathbf{l}_i^* \end{aligned} \quad (1)$$

$\mathbf{f}_i \mathbf{f}_i^*$ ,  $\mathbf{l}_i \mathbf{l}_i^*$  ( $i = 2, 3, \dots, n$ ) の第 1 行および第 1 列は 0 であるから、 $\mathbf{f}_1 \mathbf{f}_1^*$  と  $\mathbf{x}\mathbf{x}^* + \mathbf{l}_1 \mathbf{l}_1^*$  の第 1 行および第 1 列が一致する。これより次式が成り立つ。

$$f_{1,1} = \sqrt{l_{1,1}^2 + |x_1|^2} =: r, \quad f_{k,1} = \frac{1}{r} (l_{1,1} l_{k,1} + \overline{x_1} x_k) \quad (k = 2, 3, \dots, n) \quad (2)$$

ただし  $L$  の対角成分が非負の実数であることを前提としている。以上より、 $\tilde{\mathbf{l}}_1 := [0, l_{2,1}, l_{3,1}, \dots, l_{n,1}]^\top$ ,  $\tilde{\mathbf{x}} := [0, x_2, x_3, \dots, x_n]^\top$  とすると次式が成り立つ。

$$\mathbf{f}_1 = r \mathbf{e}_1 + \frac{l_{1,1}}{r} \tilde{\mathbf{l}}_1 + \frac{\overline{x_1}}{r} \tilde{\mathbf{x}}$$

ここに  $\mathbf{e}_1$  は第 1 要素が 1 で他は 0 であるベクトルである。 $\mathbf{f}_1 \mathbf{f}_1^*$  の右下  $(n-1) \times (n-1)$  行列を評価すると次式を得る。

$$\begin{aligned} \frac{1}{r^2} (l_{1,1} \tilde{\mathbf{l}}_1 + \overline{x_1} \tilde{\mathbf{x}}) &= \frac{1}{r^2} (l_{1,1}^2 \tilde{\mathbf{l}}_1 \tilde{\mathbf{l}}_1^* + l_{1,1} x_1 \tilde{\mathbf{l}}_1 \tilde{\mathbf{x}}^* + |x_1|^2 \tilde{\mathbf{x}} \tilde{\mathbf{x}}^* + l_{1,1} \overline{x_1} \tilde{\mathbf{x}} \tilde{\mathbf{l}}_1^*) \\ &= \left(1 - \frac{|x_1|^2}{r^2}\right) \tilde{\mathbf{l}}_1 \tilde{\mathbf{l}}_1^* + \frac{l_{1,1} x_1}{r^2} \tilde{\mathbf{l}}_1 \tilde{\mathbf{x}}^* + \left(1 - \frac{l_{1,1}^2}{r^2}\right) \tilde{\mathbf{x}} \tilde{\mathbf{x}}^* + \frac{\overline{x_1} l_{1,1}}{r^2} \tilde{\mathbf{x}} \tilde{\mathbf{l}}_1^* \\ &= \tilde{\mathbf{l}}_1 \tilde{\mathbf{l}}_1^* + \tilde{\mathbf{x}} \tilde{\mathbf{x}}^* - \frac{1}{r^2} (|x_1|^2 \tilde{\mathbf{l}}_1 \tilde{\mathbf{l}}_1^* + l_{1,1}^2 \tilde{\mathbf{x}} \tilde{\mathbf{x}}^* - x_1 l_{1,1} \tilde{\mathbf{l}}_1 \tilde{\mathbf{x}}^* - \overline{x_1} l_{1,1} \tilde{\mathbf{x}} \tilde{\mathbf{l}}_1^*) \\ &= \tilde{\mathbf{l}}_1 \tilde{\mathbf{l}}_1^* + \tilde{\mathbf{x}} \tilde{\mathbf{x}}^* - \mathbf{y} \mathbf{y}^* \quad \text{where} \quad \mathbf{y} = \frac{1}{r} (l_{1,1} \tilde{\mathbf{x}} - x_1 \tilde{\mathbf{l}}_1) \end{aligned}$$

上式の  $\tilde{l}_1 \tilde{l}_1^* + \tilde{x} \tilde{x}^*$  は  $\mathbf{x} \mathbf{x}^* + \mathbf{l}_1 \mathbf{l}_1^*$  の右下  $(n-1) \times (n-1)$  行列である。以上より次式が成り立つ。

$$\mathbf{f}_1 \mathbf{f}_1^* = \mathbf{x} \mathbf{x}^* + \mathbf{l}_1 \mathbf{l}_1^* - \mathbf{y} \mathbf{y}^*$$

これを式 (1) に適用して次式を得る。

$$\sum_{i=2}^n \mathbf{f}_i \mathbf{f}_i^* = \mathbf{y} \mathbf{y}^* + \sum_{i=2}^n \mathbf{l}_i \mathbf{l}_i^*$$

これは  $(n-1) \times (n-1)$  行列の rank-one update である。このようにして行列の次数を逐次的に縮小し、最後はスカラーの計算に帰着する。次数  $k$  の問題に対し式 (2) の計算量は  $O(k)$  であるから、このアルゴリズムの総計算量は  $n(n+1)/2$  に比例する。  $\square$

このアルゴリズムの Julia 1.8.0 による実装例を `Cholesky-decomposition_rank-one_update.ipynb` に記した。本文書の Git リポジトリ内でファイル検索すれば見つかる。

## 第 2.3 章

# LDL 分解

### 2.3.1 rank-one update

主張

$n \in \mathbb{N}$ ,  $A \in \mathbb{C}^{n \times n}$ ,  $A \succeq O$ ,  $\mathbf{x} \in \mathbb{C}^n$  とし、 $A$  は Hermite 行列であるとする。 $A + \mathbf{x}\mathbf{x}^*$  に対して LDL 分解のアルゴリズムを適用すると  $O(n^3)$  の計算量を要する。しかし、 $A$  の LDL 分解  $LDL^*$  が既に得られているとき、 $A + \mathbf{x}\mathbf{x}^*$  の LDL 分解を  $O(n^2)$  で得ることができる。 $\mathbf{x}\mathbf{x}^*$  の階数が 1 以下である (特に 0 となるのは  $\mathbf{x} = \mathbf{0}$  の時かつその時に限る) ことから、この方法は “rank-one update” と呼ばれている。

**導出.** 方針は Cholesky 分解の rank-one update と同様である。 $A + \mathbf{x}\mathbf{x}^*$  の LDL 分解を  $FGF^*$  とする。 $D, G$  の第  $i$  対角成分をそれぞれ  $d_i, g_i$  とする。但し  $d_i \geq 0$  を前提とする。 $L$  の第  $i$  列ベクトルを  $\mathbf{l}_i = [0, \dots, 0, 1, l_{i+1,i}, \dots, l_{n,i}]^\top \in \mathbb{C}^{n \times n}$  とし、同様に  $F$  の第  $i$  列ベクトルを  $\mathbf{f}_i = [0, \dots, 0, 1, f_{i+1,i}, \dots, f_{n,i}]^\top \in \mathbb{C}^{n \times n}$  とすると次式が成り立つ。

$$\begin{aligned} \sum_{i=1}^n \mathbf{f}_i g_i \mathbf{f}_i^* &= \mathbf{x}\mathbf{x}^* + \sum_{i=1}^n \mathbf{l}_i d_i \mathbf{l}_i^* \\ \mathbf{f}_1 g_1 \mathbf{f}_1^* + \sum_{i=2}^n \mathbf{f}_i g_i \mathbf{f}_i^* &= \mathbf{x}\mathbf{x}^* + \mathbf{l}_1 d_1 \mathbf{l}_1^* + \sum_{i=2}^n \mathbf{l}_i d_i \mathbf{l}_i^* \end{aligned} \quad (1)$$

$\mathbf{f}_i g_i \mathbf{f}_i^*$ ,  $\mathbf{l}_i d_i \mathbf{l}_i^*$  ( $i = 2, 3, \dots, n$ ) の第 1 行および第 1 列は 0 であるから、 $\mathbf{f}_1 g_1 \mathbf{f}_1^*$  と  $\mathbf{x}\mathbf{x}^* + \mathbf{l}_1 d_1 \mathbf{l}_1^*$  の第 1 行および第 1 列が一致する。これより次式が成り立つ。

$$g_1 = d_1 + |x_1|^2 =: g, \quad f_{k,1} = \frac{1}{g} (d_1 l_{k,1} + \overline{x_1} x_k) \quad (k = 2, 3, \dots, n) \quad (2)$$

以上より、 $\tilde{\mathbf{l}}_1 := [0, l_{2,1}, l_{3,1}, \dots, l_{n,1}]^\top$ ,  $\tilde{\mathbf{x}} := [0, x_2, x_3, \dots, x_n]^\top$  とすると次式が成り立つ。

$$\mathbf{f}_1 = \mathbf{e}_1 + \frac{d_1}{g} \tilde{\mathbf{l}}_1 + \frac{\overline{x_1}}{g} \tilde{\mathbf{x}}$$

ここに  $\mathbf{e}_1$  は第 1 要素が 1 で他は 0 であるベクトルである。 $\mathbf{f}_1 g_1 \mathbf{f}_1^*$  の右下  $(n-1) \times (n-1)$  行列を評価すると次式を得る。

$$\begin{aligned} & \frac{1}{g} \left( d_1 \tilde{\mathbf{l}}_1 + \tilde{\mathbf{x}} \tilde{\mathbf{x}}^* \right) \left( d_1 \tilde{\mathbf{l}}_1 + \tilde{\mathbf{x}} \tilde{\mathbf{x}}^* \right)^* = \frac{d_1}{g} \tilde{\mathbf{l}}_1 d_1 \tilde{\mathbf{l}}_1^* + \frac{|x_1|^2}{g} \tilde{\mathbf{x}} \tilde{\mathbf{x}}^* + \frac{d_1}{g} \left( x_1 \tilde{\mathbf{l}}_1 \tilde{\mathbf{x}}^* + \overline{x_1} \tilde{\mathbf{x}} \tilde{\mathbf{l}}_1^* \right) \\ &= \frac{g - |x_1|^2}{g} \tilde{\mathbf{l}}_1 d_1 \tilde{\mathbf{l}}_1^* + \frac{g - d_1}{g} \tilde{\mathbf{x}} \tilde{\mathbf{x}}^* + \frac{d_1}{g} \left( x_1 \tilde{\mathbf{l}}_1 \tilde{\mathbf{x}}^* + \overline{x_1} \tilde{\mathbf{x}} \tilde{\mathbf{l}}_1^* \right) \\ &= \tilde{\mathbf{l}}_1 d_1 \tilde{\mathbf{l}}_1^* + \tilde{\mathbf{x}} \tilde{\mathbf{x}}^* - \frac{d_1}{g} \left[ |x_1|^2 \tilde{\mathbf{l}}_1 \tilde{\mathbf{l}}_1^* + \tilde{\mathbf{x}} \tilde{\mathbf{x}}^* - x_1 \tilde{\mathbf{l}}_1 \tilde{\mathbf{x}}^* - \overline{x_1} \tilde{\mathbf{x}} \tilde{\mathbf{l}}_1^* \right] \\ &= \tilde{\mathbf{l}}_1 d_1 \tilde{\mathbf{l}}_1^* + \tilde{\mathbf{x}} \tilde{\mathbf{x}}^* - \mathbf{y} \frac{d_1}{g} \mathbf{y}^* \quad \text{where} \quad \mathbf{y} = x_1 \tilde{\mathbf{l}}_1 - \tilde{\mathbf{x}} \end{aligned}$$

上式の  $\tilde{\mathbf{l}}_1 d_1 \tilde{\mathbf{l}}_1^* + \tilde{\mathbf{x}} \tilde{\mathbf{x}}^*$  は  $\mathbf{x}\mathbf{x}^* + \mathbf{l}_1 d_1 \mathbf{l}_1^*$  の右下  $(n-1) \times (n-1)$  行列である。以上より次式が成り立つ。

$$\mathbf{f}_1 g_1 \mathbf{f}_1^* = \mathbf{x}\mathbf{x}^* + \mathbf{l}_1 d_1 \mathbf{l}_1^* - \mathbf{y} \frac{d_1}{g} \mathbf{y}^*$$

これを式 (1) に適用して次式を得る。

$$\sum_{i=2}^n f_i g_i f_i^* = \mathbf{y} \frac{d_1}{g} \mathbf{y}^* + \sum_{i=2}^n l_i d_i l_i^*$$

これは  $(n-1) \times (n-1)$  行列の rank-one update である。このようにして行列の次数を逐次的に縮小し、最後はスカラーの計算に帰着する。次数  $k$  の問題に対し式 (2) の計算量は  $O(k)$  であるから、このアルゴリズムの総計算量は  $n(n+1)/2$  に比例する。  $\square$

このアルゴリズムの Julia 1.8.0 による実装例を `LDL-decomposition_rank-one_update.ipynb` に記した。本文書の Git リポジトリ内でファイル検索すれば見つかる。