

# Algorithm Recurrence of First Paper

---

Han Liu

July 31, 2019

## 1 Introduction

In the process of training and predicting a gaussian process (GP), we need to get three quantities: the mean of predicted input, the log marginal likelihood function and its derivative. Their expressions are given as following:

$$u_{f|D}(\hat{x}) = u(\hat{x}) + k_{X\hat{x}}^T K^{-1} y \quad (1)$$

$$\log p(y|K) = -\frac{1}{2} \log((2\pi)^k |K|) - \frac{1}{2} y^T K^{-1} y \quad (2)$$

$$\frac{\partial}{\partial \theta_j} \log p(y|X, \theta) = \frac{1}{2} y^T K^{-1} \frac{\partial K}{\partial \theta_j} K^{-1} y - \frac{1}{2} \text{tr}(K^{-1} \frac{\partial K}{\partial \theta_j}) \quad (3)$$

These equation have three operations in common that dominate its time complexity:  $K^{-1}y$ ,  $\log|K|$ ,  $\text{tr}(K^{-1} \frac{\partial K}{\partial \theta_j})$ . Before that, I use Cholesky decomposition of  $K$  to compute all three quantities (see *GP\_code.py*), but it is very computational expensive, thus we need some parallel computational algorithm.

I think the innovative part of this paper is to use matrix-matrix multiply to combine three algorithms (Conjugate Gradient, Lancos Algorithm, Pivoted Cholesky Decomposition) together. The most important formula is as following:

$$[u_0 \ u_1 \ \dots \ u_t] = K^{-1} [y \ z_1 \ \dots \ z_t] \quad (4)$$

The whole algorithm flow chart can be seen in Figure 1. In the following section, I will implement all the three algorithms individually by MATLAB.

## 2 Implementation of mBCG

In this section, I firstly implement standard preconditioned conjugate gradient algorithm, then I try to implement modified conjugate gradient algorithm given in Equation 4.

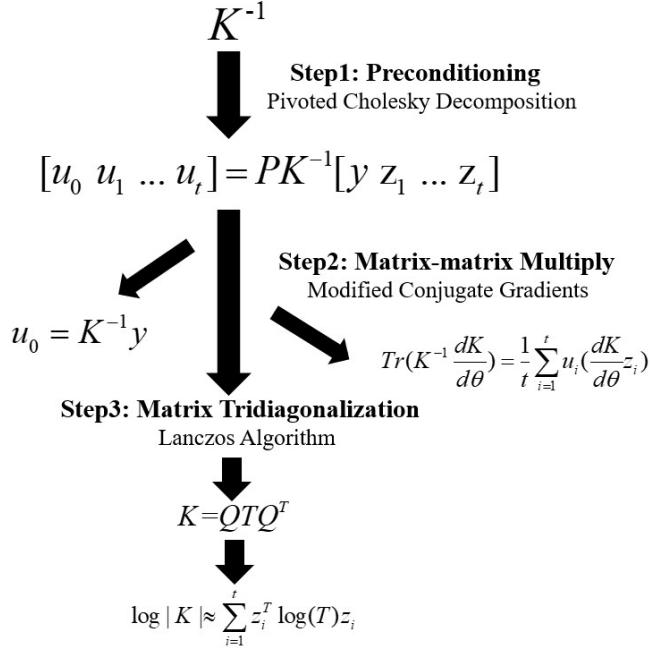


Figure 1: Overall Algorithm Flow Chart

## 2.1 Standard PCG

For the implementation of standard preconditioned conjugate gradient (PCG), I carefully follow the **Algorithm 1** in the paper. But unfortunately, I failed to get the correct results. It bothered me almost entire day, and I refer many literature, and try to implement this algorithm in a different dimension. For example, **Conjugate Gradient Method** [1] from Stanford University specify the process as following:

### CG algorithm

(follows C. T. Kelley)

```

x := 0,   r := b,   ρ₀ := ||r||²
for k = 1, ..., Nmax
  quit if √ρk-1 ≤ ε||b||
  if k = 1 then p := r; else p := r + (ρk-1/ρk-2)p
  w := Ap
  α := ρk-1/pTw
  x := x + αp
  r := r - αw
  ρk := ||r||²
  
```

Figure 2: Conjugate Gradient Method

The MATLAB code is given in *PCG\_stanford.m*. Finally, I found there are two mistakes in the paper, which can be seen in the following figure:

---

**Algorithm 1:** Standard preconditioned conjugate gradients (PCG).

---

**Input :** `mvm_A()` – function for matrix-vector multiplication (MVM) with matrix  $A$   
 $b$  – vector to solve against  
 $P^{-1}()$  – function for preconditioner

**Output :**  $A^{-1}b$ .

```

u0 ← 0 // Current solution
r0 ← mvm_A(u0) - b // Current error → should be b-mvm_A(u0)
z0 ← P⁻¹(r0) // Preconditioned error
d0 ← z0 // "Search" direction for next solution
for j ← 0 to T do
    vj ← mvm_A(dj⁻¹)
    αj ← (rj⁻¹ᵀ zj⁻¹) / (dj⁻¹ᵀ vj)
    uj ← uj⁻¹ + αj dj⁻¹
    rj ← rj⁻¹ - αj vj
    if ||rj||₂ < tolerance then return uj ;
    zj ← P⁻¹(rj)
    βj ← (zjᵀ zj) / (zj⁻¹ᵀ zj⁻¹)
    dj ← zj - βj dj⁻¹ → should be zj + beta_j d_j-1
end
return uj+1

```

---

Figure 3: Correction of PCG Algorithm

A simple test for the following input is given in Figure 4:

$$A = \begin{bmatrix} 6 & 3 & 0 \\ 3 & 6 & -6 \\ 0 & -6 & 11 \end{bmatrix}, b = \begin{bmatrix} 3 \\ 1 \\ 2 \end{bmatrix} \quad (5)$$

```

ans =
-3.3704
8.0741
5.2222

>> inv(mvm_A)*b

ans =
-3.3704
8.0741
5.2222

```

Figure 4: Matlab Test Results

## 2.2 Modified PCG

To implement a modified (batched) PCG, we need to convert MVM to MMM. There are several modifications of original algorithm, and some proofs are given in appendix. The MATLAB code is as following:

```

1 %Modified preconditioned conjugate gradients
2
3 mmm_A= [-3.3704    0.2593   -3.6296;
4 8.0741     0.4815    7.9259;
5 5.2222     0.4444   4.7778];
6 B=[4 3 2;7 1 8;9 2 5];
7 col=size(B,2);
8 U=zeros(3,3);
9 R=B-mmm_A*U;
10 Z=R;
11 D=Z;
12 alpha=zeros(3,1);
13 beta=zeros(3,1);
14 t=3;          %number of iteration
15 e=0.0001;
16
17 for j=1:3
18     V=mmm_A*D;
19     alpha=(dot(R,Z)./dot(D,V))';
20     U= (U'+diag(alpha)*D)';
21     R= (R'-diag(alpha)*V)';
22     norm_matrix=zeros(1,col);
23     for i=1:col
24         norm_matrix(i)=norm(R(:,i),2);
25     end
26     if any(any(norm_matrix<e))
27         break
28     end
29     Z_old=dot(Z,Z);
30     Z=R;
31     beta=(dot(Z,Z)./Z_old)';
32     D=(Z'+diag(beta)*D)';
33 end
34 display(U)

```

## 3 Implementation of Lanczos Algorithm

To implement Lanczos Algorithm for tridiagonalizing a symmetric matrix A, I refer the alogorithm given in **Algorithm 2**. Specifically, we have:

$$\begin{aligned} \forall i \quad [T_i]_{j,j} &= 1/[\alpha_j]_i + [\beta_{j-1}]_i/[\alpha_{j-1}]_i \\ [T_i]_{j-1,j}, [T_i]_{j,j-1} &= \sqrt{[\beta_{j-1}]/[\alpha_j]}_i \end{aligned} \tag{6}$$

I have implemented this algorithm in MATLAB as following:

```

1 function [ T ] = Lanczos_Algorithm( alpha,beta )
2     col=size(alpha,2); %iteration
3     row=size(alpha,1); %third-dimension of T
4     T=zeros(col,col,row);
5     for i=1:row
6         for j=1:col

```

```

7      if j==1
8          T(1,1,i)=1/alpha(i,1);
9      else
10         T(j,j,i)=1/alpha(i,j)+beta(i,j-1)/alpha(i,j-1);
11         T(j-1,j,i)=sqrt(beta(i,j-1))/alpha(i,j-1);
12         T(j,j-1,i)= T(j-1,j,i);
13     end
14   end
15 end
16 end

```

To test the correctness of this algorithm, I have choose two matrixes:

$$\alpha = \begin{bmatrix} 6 & 3 & 1 \\ 3 & 2 & -6 \\ 1 & -6 & 11 \end{bmatrix}, \beta = \begin{bmatrix} 1 & 2 & 3 \\ 3 & 4 & 5 \\ 6 & 7 & 8 \end{bmatrix} \quad (7)$$

The calculation results from calculator and program operation results are as following:

$T(:,:,1) =$

$$\begin{bmatrix} 0.1667 & 0.1667 & 0 \\ 0.1667 & 0.5000 & 0.4714 \\ 0 & 0.4714 & 1.6667 \end{bmatrix}$$

$T(:,:,2) =$

$$\begin{bmatrix} 0.3333 & 0.5774 & 0 \\ 0.5774 & 1.5000 & 1.0000 \\ 0 & 1.0000 & 1.8333 \end{bmatrix}$$

$T(:,:,3) =$

$$\begin{bmatrix} 1.0000 & 2.4495 & 0 \\ 2.4495 & 5.8333 & -0.4410 \\ 0 & -0.4410 & -1.0758 \end{bmatrix}$$

simple test:  
 $\alpha = \begin{bmatrix} 6 & 3 & 1 \\ 3 & 2 & -6 \\ 1 & -6 & 11 \end{bmatrix}, \beta = \begin{bmatrix} 1 & 2 & 3 \\ 3 & 4 & 5 \\ 6 & 7 & 8 \end{bmatrix}$

$$T_{(1)} = \begin{bmatrix} \frac{1}{\alpha_1} & \frac{\sqrt{\beta_1}}{\alpha_1} & \frac{\beta_1}{\alpha_1} \\ \frac{\sqrt{\beta_1}}{\alpha_1} & \frac{1}{\alpha_2 + \alpha_1} & \frac{\sqrt{\beta_2}}{\alpha_2 + \alpha_1} \\ \frac{\beta_1}{\alpha_1} & \frac{\sqrt{\beta_2}}{\alpha_2 + \alpha_1} & \frac{1}{\alpha_3 + \alpha_2 + \alpha_1} \end{bmatrix}$$

for  $\alpha_{(1)} = [6, 3, 1]$ ,  $\beta_{(1)} = [1, 2, 3]$

$$T_{(1)} = \begin{bmatrix} 0.1667 & 0.1667 & 0 \\ 0.1667 & 0.5 & 0.4714 \\ 0 & 0.4714 & 1.6667 \end{bmatrix}$$

for  $\alpha_{(2)} = [3, 2, -6]$ ,  $\beta_{(2)} = [3, 4, 5]$

$$T_{(2)} = \begin{bmatrix} 0.3333 & 0.5774 & 0 \\ 0.5774 & 1.5 & 1 \\ 0 & 1 & 1.8333 \end{bmatrix}$$

for  $\alpha_{(3)} = [1, -6, 11]$ ,  $\beta_{(3)} = [6, 7, 8]$

$$T_{(3)} = \begin{bmatrix} 1 & 2.4495 & 0 \\ 2.4495 & 5.8333 & -0.441 \\ 0 & -0.441 & -1.0758 \end{bmatrix}$$

(a) MATLAB Results

(b) Calculator Results

Figure 5: Comparison of Two Implementations

## 4 Implementation of Pivoted Cholesky Decomposition

The preconditioning is an operation of utilizing hardware more efficiently. In this paper, the author use **Pivoted Cholesky Decomposition (PCG)**. Specifically, PCG algorithm will produce a low-

rank approximation of a positive definite matrix, which is:

$$K_{XX} = L_k L_k^T \quad (8)$$

The pseudo-code of PCG algorithm is as following [2]:

### **Algorithm: cost $O(nm^2)$**

---

**Data:** matrix  $\mathbf{A} = [a_{i,j}] \in \mathbb{R}^{n \times n}$  and error tolerance  $\varepsilon > 0$   
**Result:** low-rank approximation  $\mathbf{A}_m = \sum_{i=1}^m \ell_i \ell_i^T$  such that  
 $\text{trace}(\mathbf{A} - \mathbf{A}_m) \leq \varepsilon$

```

begin
    set  $m := 1$ ;
    set  $\mathbf{d} := \text{diag}(\mathbf{A})$  and  $\text{error} := \|\mathbf{d}\|_1$ ;
    initialize  $\boldsymbol{\pi} := (1, 2, \dots, n)$ ;
    while  $\text{error} > \varepsilon$  do
        set  $i := \arg \max\{d_{\pi_j} : j = m, m+1, \dots, n\}$ ;
        swap  $\pi_m$  and  $\pi_i$ ;
        set  $\ell_{m,\pi_m} := \sqrt{d_{\pi_m}}$ ;
        for  $m+1 \leq i \leq n$  do
            compute  $\ell_{m,\pi_i} := \left( a_{\pi_m, \pi_i} - \sum_{j=1}^{m-1} \ell_{j,\pi_m} \ell_{j,\pi_i} \right) / \ell_{m,\pi_m}$ ;
            update  $d_{\pi_i} := d_{\pi_i} - \ell_{m,\pi_m} \ell_{m,\pi_i}$ ;
            compute  $\text{error} := \sum_{i=m+1}^n d_{\pi_i}$ ;
        increase  $m := m + 1$ ;
end

```

---

Figure 6: Pivoted Cholesky Decomposition Algorithm

I implement this algorithm by MATLAB:

```

1 %Implementation of Pivited Cholesky Composition
2 A=[ 6 3 0;
3 3 6 -6;
4 0 -6 11];      %semi-positive matrix
5
6 e=0.1;           %Error
7 n=size(A,1);
8
9 m=1;
10 d=diag(A);
11 error=norm(d,1);
12 v=randperm(n);
13 Pi = sort(v);
14 l=zeros(n,n);   %triangular matrix we want
15
16 while error>e
17     [argvalue, i] = max(d(Pi(m:end)));
18     a=Pi(m);
19     Pi(m)=Pi(i);
20     Pi(i)=a;

```

```

21     l(m,Pi(m))=sqrt(d(Pi(m)));
22     for i = (m+1:n)
23         S=0;
24         for j =(1:m-1)
25             l(j,Pi(m))
26             l(j,Pi(i))
27             S=S+l(j,Pi(m)).*l(j,Pi(i));
28         end
29         l(m,Pi(i))=(A(Pi(m),Pi(i))-S)./l(m,Pi(m));
30         d(Pi(i))=d(Pi(i))-l(m,Pi(m)).*l(m,Pi(i));
31     end
32 %compute error
33 error_s=0;
34 for i = (m+1:n)
35     error_s=error_s+d(Pi(i));
36 end
37 error=error_s;
38 m=m+1;
39 end
40
41 S_Am=zeros(3,3);
42 for i=(1:m-1)
43     S_Am=S_Am+l(:,i)*l(:,i)';
44 end
45 Am=S_Am;

```

## References

- [1] [https://stanford.edu/class/ee364b/lectures/conj\\_grad\\_slides.pdf](https://stanford.edu/class/ee364b/lectures/conj_grad_slides.pdf)
- [2] Helmut Harbrecht, Michael Peters, and Reinhold Schneider, "The pivoted Cholesky decomposition and its application to stochastic PDEs", p7. <http://www.mathe.tu-freiberg.de/naspde2010/sites/default/files/harbrecht.pdf>

## A Derivation of mBCG

$$\begin{bmatrix} u_0 \\ u_1 \\ u_2 \\ \vdots \\ u_k \end{bmatrix} = \begin{bmatrix} -y \\ z_1 \\ z_2 \\ \vdots \\ z_k \end{bmatrix} \underbrace{k_{xx}^{-1}}_{\text{对角}} \Rightarrow u_i = z_i^T k_{xx}^{-1}$$

Lanczos tridiagonal matrix:

purpose: To get a tridiagonalization of a symmetric matrix  $A \rightarrow A = Q T Q^T \rightarrow$  tridiagonal

$A: n \times n$       orthonormal:  $Q^T = Q^{-1}$

① 分解有无数种，每一种唯一被 probe vector 确定 (Q的第1列)

② 次迭代:  $A = Q T Q^T$  (exact)  
P次迭代: (低秩近似):

$$A \approx \tilde{Q} T \tilde{Q}^T, Q \in \mathbb{R}^{n \times p}, T \in \mathbb{R}^{p \times p}$$

$\left[ [U_j]_1, \dots, [U_j]_t \right] = U_j = U_{j-1} + \text{diag}(c_j) D_{j-1}$

$\text{diag}(c_j) D_{j-1} \Rightarrow \begin{bmatrix} c_1 & & & \\ & c_2 & & \\ & & c_3 & \\ & & & \ddots \end{bmatrix} \begin{bmatrix} [D_{j-1}]_1 \\ [D_{j-1}]_2 \\ [D_{j-1}]_3 \\ \vdots \end{bmatrix}$

$= \begin{bmatrix} c_1 [D_{j-1}]_1 \\ c_2 [D_{j-1}]_2 \\ c_3 [D_{j-1}]_3 \\ \vdots \end{bmatrix}$

Figure 7: Derivation of mBCG