

Optimizing the Frobenius norm under structural constraints.

Chris Salahub

University of Waterloo, csalahub@uwaterloo.ca

Jeffrey Uhlmann

University of Missouri, uhlmannj@missouri.edu

September 28, 2022

Abstract

TODO: Rewrite this to be more general The approximation of a general matrix \mathbf{M} by a circulant matrix \mathbf{C} is explored. Using the discrete Fourier transform matrix \mathbf{F} , the circulant with eigenvalues given by the diagonals of $\mathbf{F}\mathbf{M}\mathbf{F}^*$ is shown to be equivalent to the nearest circulant in the Frobenius norm, \mathbf{C}_M . An intuitive interpretation of this matrix in terms of means and variances of its values is presented.

1 Introduction

TODO: Preamble, justification

Circulant matrices are matrices of the form

$$\mathbf{C} = \begin{bmatrix} c_0 & c_1 & c_2 & \dots & c_{n-2} & c_{n-1} \\ c_{n-1} & c_0 & c_1 & \dots & c_{n-3} & c_{n-2} \\ c_{n-2} & c_{n-1} & c_0 & \dots & c_{n-4} & c_{n-3} \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ c_2 & c_3 & c_4 & \dots & c_0 & c_1 \\ c_1 & c_2 & c_3 & \dots & c_{n-1} & c_0 \end{bmatrix} \quad (1)$$

with $c_0, c_1, \dots, c_{n-1} \in \mathbb{C}$. They see widespread use in signal processing, computation, and physical modelling as matrices with known eigensystems [1, 2]. For \mathbf{C} as in Equation 5, the ordered eigenvalues for $k = 0, 1, \dots, n - 1$ are given by

$$\lambda_k = c_0 + \sum_{l=1}^{n-1} c_l \omega^{lk} \quad (2)$$

and the corresponding k^{th} eigenvector is given by

$$\mathbf{x}_k = \begin{bmatrix} 1 \\ \omega^k \\ \omega^{2k} \\ \vdots \\ \omega^{(M-1)k} \end{bmatrix} \quad (3)$$

where $\omega = \exp(\frac{2\pi i}{n})$ is the complex n^{th} root of unity and $i = \sqrt{-1}$.

Much of the utility of circulant matrices arises from this eigensystem. The $n \times n$ matrix of eigenvectors of \mathbf{C} scaled to be unitary,

$$\mathbf{F} = \frac{1}{\sqrt{n}} [\mathbf{x}_0 | \mathbf{x}_1 | \mathbf{x}_2 | \dots | \mathbf{x}_{n-1}] = \mathbf{F}^T, \quad (4)$$

is simply the discrete Fourier transform (DFT). Circulant matrices therefore have a deep relationship with real and complex analysis [3].

Suppose we would like to take advantage of this depth of theory and practice to approximate the $n \times n$ matrix

$$\mathbf{M} = \begin{bmatrix} m_{11} & m_{12} & \dots & m_{1n} \\ m_{21} & m_{22} & \dots & m_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ m_{n1} & m_{n2} & \dots & m_{nn} \end{bmatrix} \quad (5)$$

by some circulant matrix. This note outlines some results.

2 Approximating \mathbf{M} using \mathbf{F}

Consider the simple approximation algorithm:

1. construct the diagonal matrix \mathbf{D} where $d_{jj} = (\mathbf{F}\mathbf{M}\mathbf{F}^*)_{jj}$ and \mathbf{F}^* is the complex conjugate of \mathbf{F} ,
2. compute $\mathbf{C}_D = \mathbf{F}^*\mathbf{D}\mathbf{F}$.

As \mathbf{D} is diagonal and \mathbf{F} is the matrix of eigenvectors for any circulant matrix, \mathbf{C}_D is a circulant matrix with eigenvalues given by d_{jj} . To determine the elements $(\mathbf{C}_D)_{ij}$ in terms of ω , \mathbf{x} , and \mathbf{M} , first note

$$d_{jj} = \frac{1}{n} \sum_{k=1}^n \sum_{l=1}^n \omega^{(k-l)(j-1)} m_{lk}. \quad (6)$$

Analogously, taking $\mathbf{F}^*\mathbf{D}\mathbf{F}$ gives an i, j element

$$\begin{aligned} (\mathbf{C}_D)_{ij} &= \mathbf{F}^*\mathbf{D}\mathbf{F} \\ &= \frac{1}{n^2} \sum_{l=1}^n \sum_{k=1}^n m_{lk} \mathbf{x}_{k-l}^* \mathbf{x}_{i-j} \\ &= \frac{1}{n} \sum_{l=1}^n \sum_{k=1}^n m_{lk} \delta_{(i-j) \bmod n, (k-l) \bmod n} \end{aligned} \quad (7)$$

where δ_{ij} is the Kronecker delta defined by

$$\delta_{ij} = \begin{cases} 0 & \text{if } i \neq j \\ 1 & \text{if } i = j. \end{cases} \quad (8)$$

Equation 7 indicates that \mathbf{C}_D is generated by replacing the values of M along each circulant diagonal by the corresponding diagonal mean. Expressed as in Equation 1, \mathbf{C}_D is the circulant matrix with

$$c_k = \frac{1}{n} \sum_{\{i,j|(i-j) \bmod n=k\}} m_{ij} := \overline{m}_k \quad (9)$$

Remarkably, though this was not the original motivation, \mathbf{C}_D is also optimal in the sense of [1]: it minimizes the Frobenius norm.

Theorem 1 (\mathbf{C}_D is Frobenius optimal). \mathbf{C}_D minimizes $\|\mathbf{C} - \mathbf{M}\|_F$ for circulant \mathbf{C} , where $\|\mathbf{A}\|_F$ is the Frobenius norm of \mathbf{A} .

Proof. We can write $\|\mathbf{C} - \mathbf{M}\|_F$ as

$$\sqrt{\text{trace}((\mathbf{C} - \mathbf{M})^*(\mathbf{C} - \mathbf{M}))}. \quad (10)$$

Any \mathbf{C} which minimizes Equation 10 will also minimize $\|\mathbf{C} - \mathbf{M}\|_F^2$. Therefore we seek to minimize

$$\text{trace}((\mathbf{C} - \mathbf{M})^*(\mathbf{C} - \mathbf{M})) = \text{trace} \mathbf{M}^* \mathbf{M} - \text{trace} \mathbf{M}^* \mathbf{C} - \text{trace} \mathbf{C}^* \mathbf{M} + \text{trace} \mathbf{C}^* \mathbf{C}. \quad (11)$$

$\mathbf{M}^* \mathbf{M}$ is constant in \mathbf{C} , so this term can be ignored in the optimization. The latter three terms can be considered individually to express them as terms of the c_i . $\text{trace} \mathbf{C}^* \mathbf{C}$ is the simplest, as

$$\text{trace} \mathbf{C}^* \mathbf{C} = n \sum_{i=0}^{n-1} c_i^* c_i. \quad (12)$$

The negative terms can be expressed

$$\begin{aligned} \text{trace} \mathbf{M}^* \mathbf{C} &= \sum_{i=1}^n \sum_{j=1}^n m_{ji}^* c_{(i-j) \bmod n} \\ &= \sum_{i=0}^{n-1} c_i \left(\sum_{j=1}^{n-i} m_{j,j+i}^* + \sum_{j=1}^i m_{n-i+j,j}^* \right) \\ &= n \sum_{i=0}^{n-1} c_i \bar{m}_i^* \end{aligned} \quad (13)$$

and

$$\text{trace} \mathbf{C}^* \mathbf{M} = n \sum_{i=0}^{n-1} c_i^* \bar{m}_i. \quad (14)$$

So we seek to minimize

$$\begin{aligned} F(\mathbf{c}) &= n \sum_{i=0}^{n-1} c_i^* c_i - n \sum_{i=0}^{n-1} c_i^* \bar{m}_i - n \sum_{i=0}^{n-1} c_i \bar{m}_i^* \\ &= n (\langle \mathbf{c}, \mathbf{c} \rangle - \langle \bar{\mathbf{m}}, \mathbf{c} \rangle - \langle \mathbf{c}, \bar{\mathbf{m}} \rangle) \end{aligned} \quad (15)$$

where $\langle \mathbf{x}, \mathbf{y} \rangle \geq 0$ is the Hermitian inner product of $\mathbf{x}, \mathbf{y} \in \mathbb{C}$, $\mathbf{c} = (c_0, c_1, \dots, c_{n-1})^\top$ is the first row of \mathbf{C} , and $\bar{\mathbf{m}} = (\bar{m}_0, \bar{m}_1, \dots, \bar{m}_{n-1})^\top$ is the vector of diagonal means of \mathbf{M} .

$\langle \bar{\mathbf{m}}, \mathbf{c} \rangle$ and $\langle \mathbf{c}, \bar{\mathbf{m}} \rangle$ are maximized when $\mathbf{c} = t\bar{\mathbf{m}}$ for $t \in \mathbb{R}$ and $\langle \mathbf{c}, \mathbf{c} \rangle$ simply gives the squared magnitude of \mathbf{c} . Therefore, the minimizer of Equation 15 must be $\mathbf{c} = t\bar{\mathbf{m}}$ for some $t \in \mathbb{R}$. Substituting this into Equation 15:

$$\begin{aligned} F(t\bar{\mathbf{m}}) &= n(t^2\langle \bar{\mathbf{m}}, \bar{\mathbf{m}} \rangle - t\langle \bar{\mathbf{m}}, \bar{\mathbf{m}} \rangle - t\langle \bar{\mathbf{m}}, \bar{\mathbf{m}} \rangle) \\ &= n\|\bar{\mathbf{m}}\|^2(t^2 - 2t), \end{aligned} \tag{16}$$

which has a minimum of $-n\|\bar{\mathbf{m}}\|^2$ when $t = 1$. Therefore, $\mathbf{c} = \bar{\mathbf{m}}$ minimizes $F(\mathbf{c})$ and so the optimal circulant matrix \mathbf{C} to approximate \mathbf{M} in the Frobenius norm satisfies Equation 9. \square

Recognizing that the value of the squared Frobenius norm is

$$\|\mathbf{C} - \mathbf{M}\|_F^2 = \text{trace } \mathbf{M}^* \mathbf{M} + F(\mathbf{c}) = \sum_{i=1}^n \sum_{j=1}^n \|m_{ij}\|^2 + F(\mathbf{c}),$$

substituting $\mathbf{c} = \bar{\mathbf{m}}$ gives a minimum

$$\begin{aligned} \|\mathbf{C}_D - \mathbf{M}\|_F^2 &= \sum_{i=1}^n \sum_{j=1}^n \|m_{ij}\|^2 - n\|\bar{\mathbf{m}}\|^2 \\ &= n \left(\sum_{k=0}^{n-1} \sum_{\{1 \leq i, j \leq n \mid (i-j) \bmod n = k\}} \frac{\|m_{ij}\|^2}{n} - \sum_{k=0}^{n-1} \|\bar{m}_k\|^2 \right) \\ &= n \sum_{k=0}^{n-1} \sigma_k^2 \end{aligned} \tag{17}$$

where σ_k^2 is the variance of values along the k^{th} diagonal of \mathbf{M} . Therefore the value of the Frobenius norm $\|\mathbf{C}_D - \mathbf{M}\|_F$ is given by the total standard deviation of the values of \mathbf{M} from their respective circulant diagonals!

3 General structured matrices

The proof provided for Theorem 1 does not apply to circulant matrices alone, as the grouping of terms in the sums is arbitrary. First, we provide a general definition of a structured matrix.

Definition 1 (Structured matrix). Suppose we have a matrix \mathbf{A} with entries a_{ij} following a regular pattern in i and j , that is

$$a_{ij} = a_{f(i,j)} \quad (18)$$

where $f : \{0, 1, \dots, n-1\}^2 \mapsto \{0, 1, 2, \dots, K\}$ is the index function defining the membership of the index pair i, j to a constant index set indexed by k . Then we say that \mathbf{A} is structured. Additionally define the k^{th} index set

$$\mathcal{A}_k = \{(i, j) | f(i, j) = k\}$$

with cardinality $|\mathcal{A}_k| = n_k > 0$.

With this definition, we can prove that means following the structured pattern are Frobenius-optimal.

Theorem 2 (Means optimize the structured approximation of \mathbf{M} in the Frobenius norm.). The optimal structured matrix \mathbf{A} for index function $f(i, j)$ and index sets $\mathcal{A}_0, \mathcal{A}_1, \dots, \mathcal{A}_K$ to approximate a given matrix \mathbf{M} is given by the matrix \mathbf{A}_M with

$$a_{ij} = a_{f(i,j)} = \bar{m}_k \quad (19)$$

where

$$\bar{m}_k := \frac{1}{n_k} \sum_{\mathcal{A}_k} m_{ij}. \quad (20)$$

is the mean of entries of \mathbf{M} over the corresponding index set. Furthermore, $\|\mathbf{M} - \mathbf{A}_M\|_F$ is given by the total standard deviation of the values of \mathbf{M} within each index set from the corresponding mean.

Proof. Suppose we have a matrix \mathbf{M} and would like to approximate it with a matrix \mathbf{A} structured as in Definition 1 to minimize the squared Frobenius norm of the difference:

$$\|\mathbf{M} - \mathbf{A}\|_F^2.$$

Taking \bar{m}_k to be the mean of entries in \mathbf{M} for the k^{th} index set as in Equation 20, define the vector of all such means

$$\bar{\mathbf{m}} = (\bar{m}_0, \bar{m}_1, \dots, \bar{m}_K)^\top.$$

Further, denote the vector of unique a_k as

$$\mathbf{a} = (a_0, a_1, \dots, a_K)^\top$$

and the diagonal matrix of n_k as

$$\mathbf{N} = \text{diag}(n_0, n_1, \dots, n_K).$$

Then we can generalize Equation 15 as

$$\begin{aligned} F(\mathbf{a}) &= \mathbf{a}^* \mathbf{N} \mathbf{a} - \overline{\mathbf{m}}^* \mathbf{N} \mathbf{a} - \mathbf{a}^* \mathbf{N} \overline{\mathbf{m}} \\ &= (\mathbf{a} - \overline{\mathbf{m}})^* \mathbf{N} (\mathbf{a} - \overline{\mathbf{m}}) - \overline{\mathbf{m}}^* \mathbf{N} \overline{\mathbf{m}}. \end{aligned} \quad (21)$$

As $n_k > 0$ for all $k = 0, 1, \dots, K$, \mathbf{N} is positive definite, and so the quadratic form $\mathbf{x}^* \mathbf{N} \mathbf{x}$ has a minimum of zero when $\mathbf{x} = \mathbf{0}$. Therefore $F(\mathbf{a})$ is minimized for $\mathbf{a} = \overline{\mathbf{m}}$ and has a minimum of

$$F(\overline{\mathbf{m}}) = -\overline{\mathbf{m}}^* \mathbf{N} \overline{\mathbf{m}} = -\sum_{k=0}^K n_k \|\overline{m}_k\|^2. \quad (22)$$

So \mathbf{A}_M is the Frobenius-optimal structured matrix \mathbf{A} approximating \mathbf{M} . It has a squared Frobenius norm of

$$\begin{aligned} \|\mathbf{A}_M - \mathbf{M}\|_F^2 &= \sum_{k=0}^K \sum_{\mathcal{A}_k} \|m_{ij}\|^2 - \sum_{k=0}^K n_k \|\overline{m}_k\|^2 \\ &= \sum_{k=0}^K n_k \left(\sum_{\mathcal{A}_k} \frac{\|m_{ij}\|^2}{n_k} - \|\overline{m}_k\|^2 \right) \\ &= \sum_{k=0}^K n_k \sigma_k^2 \end{aligned} \quad (23)$$

where

$$\sigma_k^2 = \frac{1}{n_k} \sum_{\mathcal{A}_k} (m_{ij} - \overline{m}_k)^2$$

is the variance of the m_{ij} for the index set \mathcal{A}_k . □

3.1 Examples

Two well-known examples of structural matrices are Toeplitz matrices and Hankel matrices. Toeplitz matrices have an index function

$$f(i, j) = j - i$$

Hankel matrices take

$$f(i, j) = i + j$$

Another approach to Toeplitz approximation is to embed a Toeplitz matrix in a larger circulant one.

Toeplitz matrices are matrices of the form

$$\mathbf{T} = \begin{bmatrix} t_0 & t_1 & t_2 & \dots & t_{n-2} & t_{n-1} \\ t_{-1} & t_0 & t_1 & \dots & t_{n-3} & t_{n-2} \\ t_{-2} & t_{-1} & t_0 & \dots & t_{n-4} & t_{n-3} \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ t_{2-n} & t_{3-n} & t_{4-n} & \dots & t_0 & t_1 \\ t_{1-n} & t_{2-n} & t_{3-n} & \dots & t_{-1} & t_0 \end{bmatrix} \quad (24)$$

Hankel matrices are matrices of the form

$$\mathbf{H} = \begin{bmatrix} h_0 & h_1 & h_2 & \dots & h_{n-2} & h_{n-1} \\ h_1 & h_2 & h_3 & \dots & h_{n-1} & h_n \\ h_2 & h_3 & h_4 & \dots & h_n & h_{n+1} \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ h_{n-2} & h_{n-1} & h_n & \dots & h_{2n-4} & h_{2n-3} \\ h_{n-1} & h_n & h_{n+1} & \dots & h_{2n-3} & h_{2n-2} \end{bmatrix} \quad (25)$$

References

- [1] Tony F Chan. An optimal circulant preconditioner for toeplitz systems. *SIAM Journal on Scientific and Statistical Computing*, 9(4):766–771, 1988.
- [2] Robert M Gray. *Toeplitz and circulant matrices: A review*. now Publishers Inc., 2006.
- [3] Ulf Grenander and Gabor Szegő. *Toeplitz forms and their applications*. University of California Press, 1958.