SN: 2152-7385



# Applied Mathematics

https://www.scirp.org/journal/am

## **Journal Editorial Board**

ISSN Print: 2152-7385 ISSN Online: 2152-7393 https://www.scirp.org/journal/am

**Editorial Board** 

Prof. Tamer Basar University of Illinois at Urbana-Champaign, USA Prof. Leva A. Beklaryan Russian Academy of Sciences, Russia Dr. Aziz Belmiloudi Institut National des Sciences Appliquees de Rennes, France Dr. Anjan Biswas Alabama A&M University, USA Prof. Amares Chattopadhyay Indian School of Mines, India **Prof. Badong Chen** Xi'an Jiaotong University, China Prof. Jose Alberto Cuminato University of Sao Paulo, Spain Prof. Konstantin Dyakonov University of Barcelona, Spain Prof. Rosa Ferrentino University of Salerno, Italy Prof. Elena Guardo University of Catania, Italy Prof. Anwar H. Joarder University of Liberals Arts Bangladesh (ULAB), Bangladesh **Prof. Palle Jorgensen** University of Iowa, USA Dr. Vladimir A. Kuznetsov **Bioinformatics Institute**, Singapore Prof. Kil Hyun Kwon Korea Advanced Institute of Science and Technology, South Korea Prof. Hong-Jian Lai West Virginia University, USA Dr. Goran Lesaja Georgia Southern University, USA Prof. Tao Luo Georgetown University, USA Prof. Hari M. Srivastava University of Victoria, Canada Prof. Addolorata Marasco University of Naples Federico II, Italy Prof. María A. Navascués University of Zaragoza, Spain Prof. Anatolij Prykarpatski AGH University of Science and Technology, Poland Prof. Alexander S. Rabinowitch Moscow State University, Russia Prof. Mohammad Mehdi Rashidi Tongji University, China Prof. Yuriy V. Rogovchenko University of Agder, Norway Prof. Marianna Ruggieri University of Enna "KORE", Italy Prof. Ram Shanmugam Texas State University, USA **Dr. Epaminondas Sidiropoulos** Aristotle University of Thessaloniki, Greece **Prof. Sergei Silvestrov** Mälardalen University, Sweden Prof. Jacob Sturm Rutgers University, USA Prof. Mikhail Sumin Nizhnii Novgorod State University, Russia Dr. Wei Wei Xi'an University of Technology, China Icahn School of Medicine at Mount Sinai, USA Dr. Wen Zhang



## **Table of Contents**

#### 

#### Applied Mathematics (AM) Journal Information

#### **SUBSCRIPTIONS**

The *Applied Mathematics* (Online at Scientific Research Publishing, <u>https://www.scirp.org/</u>) is published monthly by Scientific Research Publishing, Inc., USA.

Subscription rates: Print: \$89 per copy. To subscribe, please contact Journals Subscriptions Department, E-mail: <u>sub@scirp.org</u>

#### SERVICES

Advertisements Advertisement Sales Department, E-mail: service@scirp.org

Reprints (minimum quantity 100 copies) Reprints Co-ordinator, Scientific Research Publishing, Inc., USA. E-mail: <u>sub@scirp.org</u>

#### COPYRIGHT

#### Copyright and reuse rights for the front matter of the journal:

Copyright © 2019 by Scientific Research Publishing Inc. This work is licensed under the Creative Commons Attribution International License (CC BY). http://creativecommons.org/licenses/by/4.0/

#### Copyright for individual papers of the journal:

Copyright © 2019 by author(s) and Scientific Research Publishing Inc.

#### Reuse rights for individual papers:

Note: At SCIRP authors can choose between CC BY and CC BY-NC. Please consult each paper for its reuse rights.

#### Disclaimer of liability

Statements and opinions expressed in the articles and communications are those of the individual contributors and not the statements and opinion of Scientific Research Publishing, Inc. We assume no responsibility or liability for any damage or injury to persons or property arising out of the use of any materials, instructions, methods or ideas contained herein. We expressly disclaim any implied warranties of merchantability or fitness for a particular purpose. If expert assistance is required, the services of a competent professional person should be sought.

#### **PRODUCTION INFORMATION**

For manuscripts that have been accepted for publication, please contact: E-mail: <a href="mailto:am@scirp.org">am@scirp.org</a>



## On the First and Second Locating Zagreb Indices of Graphs

#### Suha A. Wazzan<sup>1\*</sup>, Anwar Saleh<sup>2</sup>

<sup>1</sup>Department of Mathematics, Faculty of Science, KAU King Abdulaziz University, Girls Campus, Jeddah, KSA <sup>2</sup>Department of Mathematics, Faculty of Science, University of Jeddah, Jeddah, KSA Email: \*swazzan@kau.edu.sa, math.msfs@gmail.com

How to cite this paper: Wazzan, S.A. and Saleh, A. (2019) On the First and Second Locating Zagreb Indices of Graphs. *Applied Mathematics*, **10**, 805-816. https://doi.org/10.4236/am.2019.1010057

Received: September 7, 2019 Accepted: September 24, 2019 Published: September 27, 2019

Copyright © 2019 by author(s) and Scientific Research Publishing Inc. This work is licensed under the Creative Commons Attribution International License (CC BY 4.0).

http://creativecommons.org/licenses/by/4.0/

Open Access

#### Abstract

By the distance or degree of vertices of the molecular graph, we can define graph invariant called topological indices. Which are used in chemical graph to describe the structures and predicting some physicochemical properties of chemical compound? In this paper, by introducing two new topological indices under the name first and second Zagreb locating indices of a graph G, we establish the exact values of those indices for some standard families of graphs included the firefly graph.

#### **Keywords**

Zagreb Indices, First and Second Locating Indices, Join Graph, Firefly Graph

#### **1. Introduction**

Topological indices play a significant role mainly in chemistry, pharmacology, etc. (see [1]-[7]). Many of the topological indices of current interest in mathematical chemistry are defined in terms of vertex degrees of the molecular graph. Two of the most famous topological indices of graphs are the first and second Zagreb indices which have been introduced by Gutman and Trinajstic in [8], and defined as  $M_1(G) = \sum_{u \in V(G)} (d(u))^2$  and  $M_2(G) = \sum_{uv \in E(G)} d(u) d(v)$ , respectively. The Zagreb indices have been studied extensively due to their numerous applications in the place of existing chemical methods which need more time and increase the costs. Many new reformulated and extended versions of the Zagreb indices have been introduced for several similar reasons (cf. [9]-[17]).

One of the present authors Saleh [18] has recently introduced a new matrix representation for a graph G by defining the locating matrix Lo(G) over G. We will redefine this representation as in the following.

**Definition 1** ([18]) Let G = (V, E) be a connected graph with vertex set  $V = \{v_1, v_2, \dots, v_n\}$ . A locating function of G denoted by  $\mathbf{L}(G)$  is a function  $\mathbf{L}(G): V(G) \rightarrow (\mathbb{Z}^+ \cup \{0\})^n$  such that

 $\mathbf{L}(v_i) = \mathbf{v}_i = (d(v_1, v_i), d(v_2, v_i), \dots, d(v_n, v_i)), \text{ where } d(v_i, v_j) \text{ is the distance between the vertices } v_i \text{ and } v_j \text{ in } G. \text{ The vector } \mathbf{v}_i \text{ is called the locating vector corresponding to the vertex } v_i, \text{ where } \mathbf{v}_i \cdot \mathbf{v}_j \text{ is actually the dot product of the vectors } \mathbf{v}_i \text{ and } \mathbf{v}_j \text{ in the integers space } (\mathbb{Z}^+ \cup \{0\})^n \text{ such that } v_i \text{ is adjacent to } v_j.$ 

The above locating function and huge applications of Zagreb indices motivated us to introduce two new topological indices, namely *first* and *second locating indices*, based on the locating vectors.

**Definition 2.** Let G = (V, E) be a connected graph with a vertex set  $V = \{v_1, v_2, \dots, v_n\}$  and an edge set E(G). Then we define the first and second locating indices as

$$M_1^{\mathcal{L}}(G) = \sum_{\mathbf{v}_i \in V(G)} (\mathbf{v}_i)^2 \text{ and } M_2^{\mathcal{L}}(G) = \sum_{\mathbf{v}_i \mathbf{v}_j \in E(G)} \mathbf{v}_i \cdot \mathbf{v}_j,$$

respectively.

All graphs in this paper will be assumed simple, undirected and connected unless stated otherwise. For graph theoretical terminologies, we refer [19] to the readers.

#### 2. Some Exact Values in Terms of Locating Indices

In this section, by considering Definition 2, we determine the first and second locating indices for the standard graphs  $K_n$ ,  $C_n$ ,  $K_{n,m}$ ,  $W_n$ ,  $P_n$ , and also for the join graph  $G \cong G_1 + G_2$  such that  $G_1$  and  $G_2$  are both connected graphs with diameter 2 and *G* will be assumed as  $C_3$ ,  $C_5$ -free graphs.

**Theorem 3.** Let  $G \cong K_n$  be the complete graph with a vertex set

$$V(G) = \{v_1, v_2, \dots, v_n\}$$
, where  $n \ge 2$ . Then  $M_1^{\mathcal{L}}(K_n) = n(n-1)$  and  $n(n-1)(n-2)$ 

$$M_2^{\mathcal{L}}(K_n) = \frac{n(n-1)(n-2)}{2}.$$

**Proof.** Let  $\mathbf{v}_i$  be a locating vector corresponding to the vertex  $\mathbf{v}_i \in V(G)$ . Then  $\mathbf{v}_i = (a_1, a_2, \dots, a_n)$  such that  $a_i = 0$  and  $a_{i+1} = 1$ . Thus  $(\mathbf{v}_i)^2 = n-1$ . But we have total *n* vertices in V(G), and so  $M_1^{\mathcal{L}}(K_n) = n(n-1)$ , as required. On the other hand, for any two locating vectors  $\mathbf{v}_i$  and  $\mathbf{v}_j$ , where  $i \neq j$ , we definitely have  $\mathbf{v}_i \cdot \mathbf{v}_j = n-2$ . Hence  $M_2^{\mathcal{L}}(K_n) = \frac{n(n-1)(n-2)}{2}$ .

In the next two Theorems, we investigate the cycle  $C_n$  depends on the status of *n*.

**Theorem 4.** For an even integer  $n \ge 2$ , let  $G \cong C_n$ . Then

$$M_1^{\mathcal{L}}(C_n) = \frac{n(n^2+2)}{12}$$
 and  $M_2^{\mathcal{L}}(C_n) = \frac{n^2(n-2)^2}{12}$ .

**Proof.** By labeling the vertices of the cycle  $C_n$  as  $\{v_1, v_2, \dots, v_n\}$  in the anticlockwise direction, we obtain

$$\begin{split} \mathbf{v}_1 = & \left(0, 1, 2, 3, \cdots, \frac{n}{2}, \frac{n}{2} - 1, \frac{n}{2} - 2, \cdots, 1\right), \\ \mathbf{v}_2 = & \left(1, 0, 1, 2, \cdots, \frac{n}{2} - 1, \frac{n}{2}, \frac{n}{2} - 1, \cdots, 2\right), \\ \mathbf{v}_3 = & \left(2, 1, 0, 1, \cdots, \frac{n}{2} - 2, \frac{n}{2} - 1, \frac{n}{2}, \cdots, 3\right), \\ & \vdots \\ \mathbf{v}_n = & \left(1, 2, 3, \cdots, \frac{n}{2}, \frac{n}{2} - 1, \frac{n}{2} - 2, \frac{n}{2} - 2, \cdots, 0\right), \end{split}$$

and hence  $v_i^2 = 2\sum_{i=1}^{\frac{n}{2}} i^2 - \frac{n^2}{4}$ . It is not difficult to see that each  $v_i$  has the same components within different location, and so each  $v_i^2$  has the same sum as the form of  $v_i^2 = \frac{n(n+1)(n+2)-3n^2}{12}$ . Therefore  $M_1^{\mathcal{L}}(C_n) = \frac{n(n^2+2)}{12}$ . In addition, by the symmetry,

$$\mathbf{v}_{i} \cdot \mathbf{v}_{i+1} = 2\sum_{i=2}^{\frac{n}{2}} i(i-1) = 2\left(\frac{\frac{n}{2}\left(\frac{n}{2}+1\right)(n+1)}{6}-1\right) - 2\left(\frac{\frac{n}{2}\left(\frac{n}{2}+1\right)}{2}-1\right) = \frac{n(n-2)^{2}}{12}$$

which gives  $M_2^{\mathcal{L}}(C_n) = \frac{n^2(n-2)^2}{12}$ .

**Theorem 5.** For an odd integer  $n \ge 3$ , let  $G \cong C_n$ . Then

$$M_1^{\mathcal{L}}(C_n) = \frac{n^2(n^2-1)}{12}$$
 and  $M_2^{\mathcal{L}}(C_n) = \frac{n(n-1)(n-2)(n+3)}{12}$ 

Proof. With a similar procedure as in the proof of Theorem 4, we get

$$\mathbf{v}_{1} = \left(0, 1, 2, 3, \dots, \frac{n-1}{2}, \frac{n-1}{2} - 1, \frac{n-1}{2} - 2, \dots, 1\right),$$

$$\mathbf{v}_{2} = \left(1, 0, 1, 2, \dots, \frac{n-1}{2} - 1, \frac{n-1}{2}, \frac{n-1}{2} - 1, \dots, 2\right),$$

$$\mathbf{v}_{3} = \left(2, 1, 0, 1, \dots, \frac{n-1}{2} - 2, \frac{n-1}{2} - 1, \frac{n-1}{2}, \dots, 3\right),$$

$$\vdots$$

$$\mathbf{v}_{n} = \left(1, 2, 3, \dots, \frac{n-1}{2}, \frac{n-1}{2} - 1, \frac{n-1}{2} - 2, \dots, 0\right)$$

which implies

$$\mathbf{v}_i^2 = 2\sum_{i=1}^{\frac{n-1}{2}} i^2 = \frac{n(n^2-1)}{12},$$

and so 
$$M_1^{\mathcal{L}}(C_n) = \frac{n^2(n^2-1)}{12}$$
. Also, by the symmetry,

DOI: 10.4236/am.2019.1010057

$$\begin{split} \mathbf{v}_{i} \cdot \mathbf{v}_{i+1} &= 2\sum_{i=2}^{\frac{n-1}{2}} i(i-1) + \frac{(n-1)^{2}}{4} \\ &= \left(2\frac{\frac{n-1}{2}\left(\frac{n-1}{2}+1\right)\left(2\frac{n-1}{2}+1\right)}{6}\right) - 1 - \left(2\frac{\frac{n-1}{2}\left(\frac{n-1}{2}+1\right)}{2} - 1\right) + \frac{(n-1)^{2}}{4} \\ &= \frac{(n-1)(n-2)(n+3)}{12} \end{split}$$

which gives the exact value of  $M_2^{\mathcal{L}}(C_n)$  as depicted in the statement of theorem.

Now we will take into account the complete bipartite graphs to determine the locating indices.

**Theorem 6.** Let  $G \cong K_{n,m}$ , where  $1 \le n \le m$ . Then  $M_1^{\mathcal{L}}(K_{n,m}) = 4(n^2 + m^2) - 4(n+m) + 2nm$  and  $M_2^{\mathcal{L}}(K_{n,m}) = 2nm(n+m-2)$ .

**Proof.** For all  $1 \le i \le n$  and  $1 \le j \le m$ , by labeling the adjacent vertices  $v_i$  and  $v_{n+j}$  of  $K_{n,m}$ , the locating vectors  $v_i$  of  $v_i$  are given by:

$$\mathbf{v}_{1} = \left(0, \overline{2, \cdots, 2}, \overline{1, 1, \cdots, 1}\right), \mathbf{v}_{2} = \left(2, 0, \overline{2, \cdots, 2}, \overline{1, 1, \cdots, 1}\right),$$
$$\mathbf{v}_{3} = \left(2, 2, 0, \overline{2, \cdots, 2}, \overline{1, 1, \cdots, 1}\right), \cdots, \mathbf{v}_{n} = \left(\overline{2, \cdots, 2}, 0, \overline{1, 1, \cdots, 1}\right),$$
$$\mathbf{v}_{n+1} = \left(\overline{1, \cdots, 1}, 0, \overline{2, 2, \cdots, 2}\right), \mathbf{v}_{n+2} = \left(\overline{1, \cdots, 1}, 2, 0, \overline{2, \cdots, 2}\right), \cdots,$$
$$\mathbf{v}_{n+m} = \left(\overline{1, \cdots, 1}, \overline{2, \cdots, 2}, 0\right).$$

In here, for any  $i = 1, 2, \dots, n$ , we have  $\mathbf{v}_i^2 = 4(n-1) + m$  and for any  $i = n+1, n+2, \dots, n+m$ , we get  $\mathbf{v}_i^2 = 4(m-1) + n$ . Therefore

$$M_1^{\mathcal{L}}(K_{n,m}) = n(4(n-1)+m) + m(4(m-1)+n)$$
  
= 4(n<sup>2</sup> + m<sup>2</sup>) - 4(n+m) + 2nm.

On the other hand, for any two consecutive locating vertices  $v_i, v_{i+1}$  in  $K_{n,m}$ , since  $v_i \cdot v_{i+1} = 2(n+m-2)$ , we obtain  $M_2^{\mathcal{L}}(K_{n,m}) = 2nm(n+m-2)$ .

Since the following consequences of Theorem 6 are very special cases and clear, we will omit their proofs.

**Corollary 7.** Let  $G \cong K_{n,n}$ , where  $n \ge 1$ . Then  $M_1^{\mathcal{L}}(K_{n,n}) = 2n(5n-4)$  and  $M_2^{\mathcal{L}}(K_{n,n}) = 2n^2(2n-2)$ .

**Corollary 8.** Let  $G \cong K_{1,m}$ . Then  $M_1^{\mathcal{L}}(K_{1,m}) = 2m(2m-1)$  and  $M_2^{\mathcal{L}}(K_{1,m}) = 2m(m-1)$ .

The case for wheel graphs will be investigated in the following result.

**Theorem 9.** Let us consider G as the wheel graph  $W_n$   $(n \ge 4)$  with n+1 vertices. Then we have  $M_1^{\mathcal{L}}(W_n) = 4n(n-2)$  and  $M_2^{\mathcal{L}}(W_n) = n(6n-15)$ .

Proof. With a similar approximation as in the previous results, by labeling the

vertices of V(G) in the anticlockwise direction as  $v_1, v_2, \dots, v_n, v_{n+1}$  such that  $v_{n+1}$  is the center of the wheel, we obtain

$$\mathbf{v}_{1} = \left(0, 1, \overline{2, 2, \cdots, 2}, 1, 1\right), \mathbf{v}_{2} = \left(1, 0, 1, \overline{2, 2, \cdots, 2}, 1\right),$$
$$\mathbf{v}_{3} = \left(2, 1, 0, 1, \overline{2, \cdots, 2}, 1\right), \cdots, \mathbf{v}_{n} = \left(1, \overline{2, 2, \cdots, 2}, 1, 0, 1\right)$$
$$\mathbf{v}_{n+1} = \left(\overline{1, 1, \cdots, 1}, 0\right).$$

Now for any locating vector  $\mathbf{v}_i$  corresponding to a vertex  $v_i$   $(i \in \{1, 2, \dots, n\})$ , we have  $\mathbf{v}_i^2 = 4n + 9$  and  $\mathbf{v}_{n+1}^2 = n$ . Hence  $M_1^{\mathcal{L}}(W_n) = 4n(n-2)$ . For  $M_2^{\mathcal{L}}(W_n)$ , by labeling the vertices as above, we have

$$\mathbf{v}_{1} = \left(\mathbf{0}, \mathbf{1}, \underbrace{2, 2, \cdots, 2}_{n-3}, \mathbf{1}, \mathbf{1}\right), \mathbf{v}_{2} = \left(\mathbf{1}, \mathbf{0}, \mathbf{1}, \underbrace{2, 2, \cdots, 2}_{n-3}, \mathbf{1}\right),$$
$$\mathbf{v}_{3} = \left(2, \mathbf{1}, \mathbf{0}, \mathbf{1}, \underbrace{2, \cdots, 2}_{n-2}, \mathbf{1}\right), \mathbf{v}_{4} = \left(2, 2, \mathbf{1}, \mathbf{0}, \mathbf{1}, \underbrace{2, \cdots, 2}_{n-1}, \mathbf{1}\right),$$
$$\mathbf{v}_{5} = \left(2, 2, 2, \mathbf{1}, \mathbf{0}, \mathbf{1}, \underbrace{2, \cdots, 2}_{n}, \mathbf{1}\right), \cdots, \mathbf{v}_{n} = \left(1, \underbrace{2, 2, \cdots, 2}_{n-3}, \mathbf{1}, \mathbf{0}, \mathbf{1}\right),$$
$$\mathbf{v}_{n+1} = \left(\underbrace{1, 1, \cdots, 1}_{n}, \mathbf{0}\right).$$

Bearing in mind the permutation of components 1, 0, 1 in each vector  $\mathbf{v}_i$ , where  $i = 1, 2, \dots, n$ , it is easy to see that any two adjacent vertices  $v_i$  and  $v_j$   $(i, j \in \{1, 2, \dots, n\})$  satisfy  $\mathbf{v}_i \cdot \mathbf{v}_j = 4n - 11$  and  $\mathbf{v}_i \cdot \mathbf{v}_{i+1} = 2n - 4$  for  $i = 1, 2, \dots, n$ . Hence  $M_2^{\mathcal{L}}(W_n) = n(6n - 15)$ .

The result for determining of locating indices on path graphs can be given as in the following.

**Theorem 10.** Let  $G \cong P_n$   $(n \ge 3)$ . Then

$$M_1^{\mathcal{L}}(P_n) = \sum_{j=1}^{n-1} \frac{(n-j)(n-j+1)(2n-2j+1)}{3},$$

and

$$M_{2}^{\mathcal{L}}(P_{n}) = 2\sum_{j=1}^{n-1} \frac{(n-j)(n-j+1)(n-j-1)}{3}.$$

**Proof.** Assume that *G* is the graph  $P_n$   $(n \ge 3)$ . By labeling the vertices from left to right as  $v_1, v_2, \dots, v_n$  according to the locating function, the corresponding vector for each vertex  $v_i \in V(G)$   $(i = 1, \dots, n)$  will be the form of

$$\mathbf{v}_1 = (0, 1, 2, 3, \dots, n-1), \mathbf{v}_2 = (1, 0, 1, 2, \dots, n-2), \dots,$$
  
$$\mathbf{v}_{n-1} = (n-2, n-1, \dots, 0, 1), \mathbf{v}_n = (n-1, n-2, n-3, \dots, 0)$$

By applying the symmetry on components between the vector pairs  $v_1, v_n$ and  $v_2, v_{n-1}, \cdots$  and so on, we can see that

$$M_1^{\mathcal{L}}(P_n) = 2\sum_{j=1}^{n-1}\sum_{i=1}^{n-j}i^2 = \sum_{j=1}^{n-1}\frac{(n-j)(n-j+1)(2n-2j+1)}{3}.$$

For  $M_2^{\mathcal{L}}(P_n)$ , we see that

$$\mathbf{v}_{1} \cdot \mathbf{v}_{2} = (0 \cdot 1) + (1 \cdot 0) + \dots + (n - 1)(n - 2) = \sum_{i=1}^{n-1} i(i - 1),$$
  
$$\mathbf{v}_{2} \cdot \mathbf{v}_{3} = (1 \cdot 2) + (0 \cdot 1) + \dots + (n - 2)(n - 3) = \sum_{i=1}^{n-1} i(i - 1),$$
  
$$\mathbf{v}_{3} \cdot \mathbf{v}_{4} = \sum_{i=1}^{n-1} i(i - 1),$$
  
$$\vdots$$

However, by the symmetry between the components of the vectors as mentioned above, we get

$$M_{2}^{\mathcal{L}}(P_{n}) = 2\sum_{j=1}^{n-1}\sum_{i=1}^{n-j}i(i-1) = 2\sum_{j=1}^{n-1}\left(\sum_{i=1}^{n-j}i^{2} - \sum_{i=1}^{n-j}i\right)$$

which can be rewritten as in the form

$$M_{2}^{\mathcal{L}}(P_{n}) = 2\sum_{j=1}^{n-1} \left( \frac{(n-j)(n-j+1)(2n-2j+1)}{6} - \frac{(n-j)(n-j+1)}{2} \right)$$
$$= 2\sum_{j=1}^{n-1} \frac{(n-j)(n-j+1)(n-j-1)}{3}.$$

This complete the proof.■

It is known that from the elementary textbooks the join  $G = G_1 + G_2$  of graphs  $G_1$  and  $G_2$  with disjoint vertex sets  $V_1$  and  $V_2$  and edge sets  $E_1$  and  $E_2$  is the graph union  $G_1 \cup G_2$  together with all the edges joining  $V_1$  and  $V_2$ . In the following theorem we find first and second locating indices for the join graph G.

**Theorem 11.** Let  $G \cong G_1 + G_2$  such that  $G_1$  and  $G_2$  are both connected graphs with diameter 2 and G is a  $C_3$  or  $C_5$ -free graph. Assume that  $G_1$  has  $n_1$  vertices and  $m_1$  edges while  $G_2$  has  $n_2$  vertices and  $m_2$  edges. Then

$$M_1^{\mathcal{L}}(G) = 2n_1n_2 + 4(n_1^2 + n_2^2 - n_1 - n_2) - 6(m_1 + m_2),$$

and

$$M_{2}^{\mathcal{L}}(G) = 2(m_{1}n_{1} + m_{2}n_{2}) - 4(m_{1} + m_{2}) + 2n_{1}n_{2}(n_{1} + n_{2} - 2).$$

**Proof.** Assume that G satisfies the conditions in the statement of theorem. Let us label the vertices of the graph G as

$$V_1, V_2, \cdots, V_{n_1}, V_{n_1+1}, V_{n_1+2}, \cdots, V_{n_1+n_2},$$

where  $v_1, v_2, \dots, v_{n_1} \in V(G_1)$  and  $v_{n_1+1}, v_{n_1+2}, \dots, v_{n_1+n_2} \in V(G_2)$ . Also let v be the locating vector corresponding to the vertex v such that  $v \in V(G_1)$ :

$$\boldsymbol{\nu} = \left(0, \underbrace{1, \cdots, 1}_{\deg(\nu)}, \underbrace{2, \cdots, 2}_{n_1 - \deg(\nu) - 1}, \underbrace{1, \cdots, 1}_{n_2}\right).$$

Then  $v^2 = n_2 + 4n_1 - 4 - 3 \deg(v)$ .

Similarly, for any vertex  $w \in V(G_2)$ , the locating vector w corresponding to w:

$$\boldsymbol{w} = \left(\underbrace{1, \dots, 1}_{n_1}, 0, \underbrace{1, \dots, 1}_{\deg(w)}, \underbrace{2, \dots, 2}_{n_2 - \deg(w) - 1}\right)$$

So  $w^2 = n_1 + 4n_2 - 4 - 3 \deg(w)$ . Therefore, by the above equalities on  $v^2$  and  $w^2$ , we obtain

$$M_{1}^{\mathcal{L}}(G) = n_{1}(n_{2} + 4n_{1} - 4) - 6m_{1} + n_{2}(n_{1} + 4n_{2} - 4) - 6m_{2}$$
  
=  $2n_{1}n_{2} + 4(n_{1}^{2} + n_{2}^{2} - n_{1} - n_{2}) - 6(m_{1} + m_{2}).$ 

Now, let us make partition to the set of vertices of G as

$$A = \{u \cdot v : u, v \in V(G_1)\},$$
$$B = \{u \cdot v : u, v \in V(G_2)\},$$
$$C = \{u \cdot v : u \in V(G_1), v \in V(G_2)\}$$

Hence  $M_2^{\mathcal{L}}(G)$  can be written as  $\sum_{u\cdot v\in A} u\cdot v + \sum_{u\cdot v\in B} u\cdot v + \sum_{u\cdot v\in C} u\cdot v$ . To get  $\sum_{u\cdot v\in A} u\cdot v$  for any two adjacent vertices  $u, v \in V(G_1)$ , let us consider

$$\boldsymbol{u} = \left(0, \underbrace{1, \dots, 1}_{\deg(u)}, \underbrace{2, \dots, 2}_{n_1 - \deg(u) - 1}, \underbrace{1, \dots, 1}_{n_2}\right)$$
$$\boldsymbol{v} = \left(1, 0, \underbrace{2, \dots, 2}_{\deg(u) - 1}, \underbrace{1, \dots, 1}_{n_1 - \deg(u) - 1}, \underbrace{1, \dots, 1}_{n_2}\right).$$

We then have

$$u \cdot v = 2(\deg(u)-1) + 2(n_1 - \deg(u)-1) + n_2 = n_2 + 2n_1 - 4$$

which implies  $\sum_{u\cdot v\in A} \boldsymbol{u}\cdot \boldsymbol{v} = m_1(n_2 + 2n_1 - 4)$ . With a similar calculation, we get  $\sum_{u\cdot v\in B} \boldsymbol{u}\cdot \boldsymbol{v} = m_2(n_1 + 2n_2 - 4)$ .

Next, we need to calculate  $\sum_{u \cdot v \in C} u \cdot v$ . To do that let us take  $u \in V(G_1)$  and  $v \in V(G_2)$ , and then labeling as

$$\boldsymbol{u} = \left(0, \underbrace{1, \dots, 1}_{\deg(u)}, \underbrace{2, \dots, 2}_{n_1 - \deg(u) - 1}, \underbrace{1, \dots, 1}_{n_2}\right)$$
$$\boldsymbol{v} = \left(\underbrace{1, \dots, 1}_{\deg(u) + 1}, \underbrace{1, \dots, 1}_{n_1 - \deg(u) - 1}, 0, \underbrace{1, \dots, 1}_{\deg(v)}, \underbrace{2, \dots, 2}_{n_2 - \deg(v) - 1}\right).$$

Hence we get

$$\boldsymbol{u} \cdot \boldsymbol{v} = \deg(u) + 2(n_1 - \deg(u) - 1) + \deg(v) + 2(n_2 - \deg(v) - 1)$$
$$= 2(n_1 + n_2) - 4 - (\deg(u) - \deg(v))$$
and so  $\sum_{u \cdot v \in C} \boldsymbol{u} \cdot \boldsymbol{v} = n_1 n_2 (2(n_1 + n_2) - 4) - 2n_2 m_1 - 2n_1 m_2.$ 

After all above calculations, we finally obtain

$$M_{2}^{\mathcal{L}}(G) = m_{1}(n_{2} + 2n_{1} - 4) + m_{2}(n_{1} + 2n_{2} - 4) + n_{1}n_{2}(2(n_{1} + n_{2}) - 4) - 2n_{2}m_{1} - 2n_{1}m_{2} = 2(m_{1}n_{1} + m_{2}n_{2}) - 4(m_{1} + m_{2}) + 2n_{1}n_{2}(n_{1} + n_{2} - 2)$$

Hence the result.

#### 3. Locating Indices of Firefly Graphs

We recall that a firefly graph  $F_{s,t,n-2s-2t-1}$  ( $s \ge 0, t \ge 0$  and  $n-2s-2t-1\ge 0$ ) is a graph of order *n* that consists of *s* triangles, *t* pendant paths of length 2 and n-2s-2t-1 pendent edges that are sharing a common vertex (cf. [20]). Let  $\mathcal{F}_n$  be the set of all firefly graphs  $F_{s,t,n-2s-2t-1}$ . Note that  $\mathcal{F}_n$  contains the stars  $S_n(\cong F_{0,0,n-1})$ , stretched stars  $(\cong F_{0,t,n-2t-1})$ , friendship graphs  $\left(\cong F_{\frac{n-1}{2},0,0}\right)$  and butterfly graphs  $(\cong F_{s,0,n-2s-1})$ .

In the next theorem we present the first and second locating indices for the firefly graph. In our calculations, for simplicity, we denote n-2s-2t-1 by a single letter *l*.

**Theorem 12.** Let  $G \cong F_{s,t,l}$  ( $s \ge 0, t \ge 0$  and  $l \ge 0$ ) be a firefly graph of order *n*. Then

$$M_{1}^{\mathcal{L}}(G) = 4l^{2} + 16ls + 26lt - 2l + 16s^{2} + 52st - 10s + 38t^{2} - 28t,$$

and

$$M_{2}^{\mathcal{L}}(G) = 2l^{2} + 16ls + 13lt - 2l + 24s^{2} + 52st - 20s + 22t^{2} - 17t.$$

**Proof.** Let  $G \cong F_{s,t,l}$  ( $s \ge 0, t \ge 0$  and  $l \ge 0$ ) is a firefly graph of order *n*. Let us label the vertices with clockwise direction as

$$v_1, v_2, \cdots, v_{2s+1}, v_{2s+2}, v_{2s+3}, \cdots, v_{2s+l+1}, v_{2s+l+2}, v_{2s+l+3}, \cdots, v_{2s+l+t+1}, v_{2s+l+t+2}, v_{2s+l+t+3}, \cdots, v_{2s+2t+l+1},$$

where  $v_1$  is the center of the firefly graph and

$$\underbrace{v_2, v_3, \cdots, v_{2s+1}}_{l}$$
: vertices of triangles,  
$$\underbrace{v_{2s+2}, v_{2s+3}, \cdots, v_{2s+l+1}}_{l}$$
: vertices of pendent edges,  
$$\underbrace{v_{2s+l+2}, v_{2s+l+3}, \cdots, v_{2s+l+t+1}}_{t}$$
: vertices of pendent path of length 1,

 $\underbrace{v_{2s+l+t+2}, v_{2s+l+t+3}, \cdots, v_{2s+2t+l+1}}_{::}$ : vertices of pendent path of length 2.

Now we calculate the corresponding vectors  $v_i$  for each vertex  $v_i \in V(G)$ , where  $i = 1, 2, \dots, 2s + 2t + l + 1$ , as in the following:

$$\boldsymbol{v}_1 = \left(0, \underbrace{1, 1, \cdots, 1}_{2s}, \underbrace{1, 1, \cdots, 1}_{l}, \underbrace{1, 1, \cdots, 1}_{t}, \underbrace{2, 2, \cdots, 2}_{t}\right),$$
$$\boldsymbol{v}_2 = \left(1, 0, 1, \underbrace{2, 2, \cdots, 2}_{2s-2}, \underbrace{2, 2, \cdots, 2}_{l}, \underbrace{2, 2, \cdots, 2}_{t}, \underbrace{3, 3, \cdots, 3}_{t}\right),$$

Suppose that  $A, B, C, D \subset V(G)$  such that

$$A = \{v_2, v_3, \dots, v_{2s+1}\}, C = \{v_{2s+l+2}, v_{2s+l+3}, \dots, v_{2s+l+t+1}\}, B = \{v_{2s+2}, v_{2s+3}, \dots, v_{2s+l+1}\}, D = \{v_{2s+l+t+2}, v_{2s+l+t+3}, \dots, v_{2s+2t+l+1}\}.$$

Therefore we can write

$$M_1^{\mathcal{L}}(G) = \sum_{v \in A} v^2 + \sum_{v \in B} v^2 + \sum_{v \in C} v^2 + \sum_{v \in D} v^2.$$

For the calculation of  $\sum_{v \in A} v^2$ , we have the cases  $v_1^2 = 2s + l + t + 4t = 2s + l + 5t$ and  $v_i^2 = 2 + 4(2s - 2) + 4l + 4t + 9t = 4l + 8s + 13t - 6,$ where  $i = 2, 3, \dots, 2s + 1$ . Hence  $\sum_{v \in A} v^2 = 2s + l + 5t + 2s(4l + 8s + 13t - 6)$  $= 2s + l + 5t + 8sl + 16s^2 + 26st - 12s$  $= l - 10s + 5t + 8ls + 26st + 16s^2.$ On the other hand, for the calculation of  $\sum_{v \in B} v^2$ , we have  $v_i^2 = 1 + 4(2s) + 4(l - 1) + 4t + 9t = 4l + 8s + 13t - 3,$ where for  $i = 2s + 1, 2s + 2, \dots, 2s + l + 1$ . Thus  $\sum_{v \in B} v^2 = l(4l + 8s + 13t - 3) = 4l^2 + 8ls + 13lt - 3l.$ Thirdly to calculate  $\sum_{v \in C} v^2$ , we have

 $\mathbf{v}_{i}^{2} = 2 + 4(2s) + 4l + 4(t-1) + 9(t-1) = 4l + 8s + 13t - 11,$ 

where  $i = 2s + l + 2, 1, 2s + l + 3, \dots, 2s + l + t + 1$ , and so  $\sum_{v \in C} v^2 = t (4l + 8s + 13t - 11) = 4tl + 8ts + 13t^2 - 11t.$ 

Finally, for the case of  $\sum_{\nu \in D} \nu^2$ , we get

$$\mathbf{v}_i^2 = 3 + 9(2s) + 9l + 9(t-1) + 16(t-1) = 9l + 18s + 25t - 22,$$

where  $i = v_{2s+l+t+2}, v_{2s+l+t+3}, \dots, v_{2s+2t+l+1}$ . This gives

$$\sum_{v \in D} v^2 = t \left(9l + 18s + 25t - 22\right) = 9tl + 18ts + 25t^2 - 22t.$$

By collecting all above calculations, we obtain

$$M_{1}^{\mathcal{L}}(G) = \sum_{v \in A} \mathbf{v}^{2} + \sum_{v \in B} \mathbf{v}^{2} + \sum_{v \in C} \mathbf{v}^{2} + \sum_{v \in D} \mathbf{v}^{2}$$
  
=  $l - 10s + 5t + 8ls + 26st + 16s^{2} + 4l^{2} + 8ls + 13lt - 3l$   
+  $4tl + 8ts + 13t^{2} - 11t + 9tl + 18ts + 25t^{2} - 22t$   
=  $4l^{2} + 16ls + 26lt - 2l + 16s^{2} + 52st - 10s + 38t^{2} - 28t$ ,

as required.

Before starting to calculate the index  $M_2^{\mathcal{L}}(G) = \sum_{v_i u_j \in E(G)} v_i \cdot u_j$ , we should remind that for any two adjacent vertices u and v will be denoted by  $u \approx v$ . Now, let us again consider the same subsets A, B, C and D of V(G). Therefore we firstly have

$$\sum_{v_i \in A} \mathbf{v}_i \cdot \mathbf{v}_i = 2s \left( 1 + 2 \left( 2s - 2 + l + t \right) + 3t \right) = 2s \left( 2l + 4s + 5t - 3 \right)$$
  
= 4sl + 8s<sup>2</sup> + 10st - 6s.  
$$\sum_{v_i \in B} \mathbf{v}_1 \cdot \mathbf{v}_i = l \left( 2 \left( 2s + l - 1 + t \right) + 3t \right) = l \left( 2l + 4s + 5t - 2 \right)$$
  
= 2l<sup>2</sup> + 4sl + 5lt - 2l.  
$$\sum_{v_i \in C} \mathbf{v}_1 \cdot \mathbf{v}_i = t \left( 1 + 2 \left( 2s + l + t - 1 \right) + 2t \right) = t \left( 2l + 4s + 4t - 1 \right)$$
  
= 2tl + 4ts + 4t<sup>2</sup> - t.

DOI: 10.4236/am.2019.1010057

Secondly,

$$\sum_{\substack{v_i \in A - \{v_i\}\\v_i \approx v_{i+1}}} \boldsymbol{v}_i \cdot \boldsymbol{v}_{i+1} = 2s(1 + 4(2s - 2 + l + t) + 9t) = 2s(4l + 8s + 13t - 7)$$
  
= 8sl + 16s<sup>2</sup> + 26st - 14s.

Thirdly,

$$\sum_{\substack{v_i \in C, v_{i+t} \in D \\ v_i \approx v_{i+t}}} v_i \cdot v_{i+t} = t \left( 2 + 6 \left( 2s + l + t - 1 \right) + 12 \left( t - 1 \right) \right)$$
$$= t \left( 6l + 12s + 18t - 16 \right) = 6tl + 12ts + 18t^2 - 16t.$$

Again, by collecting all above calculations, we obtain

$$\begin{split} M_{2}^{\mathcal{L}}(G) &= \sum_{v_{i} \in A} \mathbf{v}_{i} \cdot \mathbf{v}_{i} + \sum_{v_{i} \in B} \mathbf{v}_{1} \cdot \mathbf{v}_{i} + \sum_{v_{i} \in C} \mathbf{v}_{1} \cdot \mathbf{v}_{i} + \sum_{v_{i} \in A - \{v_{i}\}} \mathbf{v}_{i} \cdot \mathbf{v}_{i+1} + \sum_{u_{i} \in C, v_{i+t} \in D} \mathbf{v}_{i} \cdot \mathbf{v}_{i+t} \\ &= 4sl + 8s^{2} + 10st - 6s + 2l^{2} + 4sl + 5lt - 2l + 2tl + 4ts + 4t^{2} - t \\ &+ 8sl + 16s^{2} + 26st - 14s + 6tl + 12ts + 18t^{2} - 16t \\ &= 2l^{2} + 16ls + 13lt - 2l + 24s^{2} + 52st - 20s + 22t^{2} - 17t. \end{split}$$

These all above processes complete the proof.■ **Corollary 13.** 1) *For any friendship graph of order n*,

$$M_{1}^{\mathcal{L}}(G) = 4n^{2} - 13n + 9$$
 and  $M_{2}^{\mathcal{L}}(G) = 6n^{2} - 22n + 16$ .

2) For any butterfly graph of order *n*,

 $M_1^{\mathcal{L}}(G) = 4n^2 - 10n - 6s + 6$  and  $M_2^{\mathcal{L}}(G) = 8ns - 24s - 6n + 2n^2 + 4$ .

#### 4. Conclusion

In this paper, two new topological indices based on Zagreb indices are proposed. The exact values of these new topological indices are calculated for some standard graphs and for the firefly graphs. These new indices can be used to investigate the chemical properties for some chemical compound such as drugs, bridge molecular graph etc. For the future work, instead of defining these new topological indices based on the degrees of the vertices, we can redefine them based on the degrees of the edges by defining them on the line graph of any graph. Similar calculations can be computed to indicate different properties of the graph.

#### **Conflicts of Interest**

The authors declare no conflicts of interest regarding the publication of this paper.

#### References

- [1] Diudea, M.V., Florescu, M.S. and Khadikar, P.V. (2006) Molecular Topolgy and Its Applications. Eficon, Bucarest.
- [2] Diudea, M.V. (2010) Nanomolecules and Nanostructures-Poynomial and Indices. Univ. Kragujevac, Kragujevac.
- [3] Gutman, I. and Furula, B. (2012) Distance in Molecular Graphs. Univ. Kragujevac,

Kragujevac.

- [4] Gutman, I. and Furula, B. (2010) Novel Molecular Structure Descriptors-Theory and Applications II. Univ. Kragujevac, Kragujevac.
- [5] Karelson, M. (2000) Molecular Descriptors in QSAR-QSPR. Wiley, New York.
- [6] Todeschini, R. and Consonni, V. (2000) Handbook of Molecular Descriptors. Wiley-VCH, Weinheim. <u>https://doi.org/10.1002/9783527613106</u>
- [7] Todeschini, R. and Consonni, V. (2009) Molecular Descriptors for Chemoinformatics. Wiley-VCH, Weinheim, Vols. I and II. <u>https://doi.org/10.1002/9783527628766</u>
- [8] Gutman, I. and Trinajstic, N. (1972) Graph Theory and Molecular Orbitals, Total π-Electron Energy of Alternant Hydrocarbons. *Chemical Physics Letters*, **17**, 535-538. <u>https://doi.org/10.1016/0009-2614(72)85099-1</u>
- [9] Das, K.C., Yurttas, A., Togan, M., Cevik, A.S. and Cangul, I.N. (2013) The Multiplicative Zagreb Indices of Graph Operation. *Journal of Inequalities and Applications*, 2013, Article No. 90. <u>https://doi.org/10.1186/1029-242X-2013-90</u>
- [10] Alwardi, A., Alqesmah, A., Rangarajan, R. and Cangul, I.N. (2018) Entire Zagreb Indices of Graphs. *Discrete Mathematics, Algorithms and Applications*, **10**, Article ID: 1850037. <u>https://doi.org/10.1142/S1793830918500374</u>
- [11] Braun, J., Kerber, A., Meringer, M. and Rucker, C. (2005) Similarity of Molecular Descriptors: The Equivalence of Zagreb Indices and Walk Counts. *MATCH Communications in Mathematical and in Computer Chemistry*, 54, 163-176.
- [12] Gutman, I. and Das, K.C. (2004) The First Zagreb Index 30 Years after. MATCH Communications in Mathematical and in Computer Chemistry, 50, 83-92.
- Khalifeh, M.H., Youse-Azari, H. and Ashra, A.R. (2009) The First and Second Zagreb Indices of Some Graph Operations. *Discrete Applied Mathematics*, 157, 804-811. <u>https://doi.org/10.1016/j.dam.2008.06.015</u>
- [14] Nikolić, S., Kova č ević, G., Mili č ević, A. and Trinajstić, N. (2003) The Zagreb Indices 30 Years after. *Croatica Chemica Acta*, 76, 113-124.
- [15] Zhou, B. and Gutman, I. (2005) Further Properties of Zagreb Indices. MATCH Communications in Mathematical and in Computer Chemistry, 54, 233-239.
- Zhou, B. and Gutman, I. (2004) Relations between Wiener, Hyper-Wiener and Zagreb Indices. *Chemical Physics Letters*, **394**, 93-95. https://doi.org/10.1016/j.cplett.2004.06.117
- [17] Zhou, B. (2004) Zagreb Indices. MATCH Communications in Mathematical and in Computer Chemistry, 52, 113-118.
- [18] Ramaswamy, H.N., Alwardi, A. and Kumar, N.R. (2017) On the Locating Matrix of a Graph and Its Spectral Analysis. *Computer Science Journal of Moldova*, 25, 260-277.
- [19] Chartand, G. and Lesniak, L. (2005) Graphs and Digraphs. 4th Edition, CRC Press, Boca Raton.
- [20] Li, J.X., Guo, J.M. and Shiu, W.C. (2013) On the Second Largest Laplacian Eignvalues of Graphs. *Linear Algebra and Its Applications*, 438, 2438-2446. https://doi.org/10.1016/j.laa.2012.10.047



## **Approximation of Functions by Quadratic** Mapping in $(\beta, p)$ -Banach Space

#### Xiujiao Chi, Longyin Bao, Liguang Wang\*

School of Mathematical Sciences, Qufu Normal University, Qufu, China Email: chixiujiao5225@163.com, baolongyin1208@163.com, \*wangliguang0510@163.com

How to cite this paper: Chi, X.J., Bao L.Y., and Wang, L.G. (2019) Approximation of Functions by Quadratic Mapping in  $(\beta,$ p)-Banach Space. Applied Mathematics, 10, 817-825.

https://doi.org/10.4236/am.2019.1010058

Received: August 5, 2019 Accepted: September 24, 2019 Published: September 27, 2019

Copyright © 2019 by author(s) and Scientific Research Publishing Inc. This work is licensed under the Creative Commons Attribution International License (CC BY 4.0).

http://creativecommons.org/licenses/by/4.0/ ۲ **Open Access** 

Abstract

In this paper, we study the functions with values in  $(\beta, p)$ -Banach spaces which can be approximated by a quadratic mapping with a given error.

#### **Keywords**

Hyers-Ulam-Rassias Stability, Quadratic Mapping,  $(\beta, p)$ -Banach Space

#### **1. Introduction**

The stability problem of functional equations originated from a question of Ulam [1] in 1940 concerning the stability of group homomorphisms.

Give a group  $(G_1,*)$  and a metric group  $(G_2,\cdot,d)$  with the metric  $d(\cdot,\cdot)$ . Given  $\varepsilon > 0$ , does there exist a  $\delta > 0$  such that if  $f: G_1 \to G_2$  satisfies  $d(f(x*y), f(x) \cdot f(y)) < \delta$  for all  $x, y \in G_1$ , then there is a homomorphism  $g: G_1 \to G_2$  with  $d(f(x), g(x)) < \varepsilon$  for all  $x \in G_1$ ?

Hyers [2] gave the first affirmative partial answer to the question of Ulam for Banach spaces. Hyers's Theorem was generalized by Aoki [3] for additive mappings and by Rassias [4] for linear mappings by considering an unbounded Cauchy difference. The paper of Th. M. Rassias has provided a lot of influence in the development of what we call generalized Hyers-Ulam-Rassias stability of functional equations. Beginning around 1980, the stability problems of several functional equations and approximate homomorphisms have been extensively investigated by a number of authors and there are many interesting results concerning this problem (see [5]-[18]).

The functional equation

$$f(x+y)+f(x-y)=2f(x)+2f(y)$$

is called the quadratic functional equation. Every solution of the quadratic func-

tional equation is said to be a quadratic mapping. The Hyers-Ulam stability for quadratic functional equation was first proved by Skof [5] for mappings acting between a normed space and a Banach space. P. W. Cholewa [6] showed that Skof's Theorem is also valid if the normed space is replaced with an abelian group.

Now we recall some basic facts concerning  $(\beta, p)$ -Banach spaces. We fixed real numbers  $\beta$  with  $0 < \beta \le 1$  and p with  $0 . Let <math>\mathbb{K} = \mathbb{R}$  or  $\mathbb{C}$ . Let X be linear space over  $\mathbb{K}$ . A quasi- $\beta$ -norm  $\|\cdot\|$  is a real-valued function on Xsatisfying the following conditions:

- (i)  $||x|| \ge 0, \forall x \in X$ ; ||x|| = 0 if and only if x = 0;
- (ii)  $\|\lambda x\| = |\lambda|^{\beta} \|x\|, \forall x \in X, \beta \in K;$
- (iii) There is a constant  $K \ge 1$  such that  $||x + y|| \le K(||x|| + ||y||), \forall x, y \in X$ .

The pair  $(X, \|\cdot\|)$  is called a quasi- $\beta$ -normed space if  $\|\cdot\|$  is a quasi- $\beta$ -norm on X. The smallest possible K is called the module of concavity of  $\|\cdot\|$ . A quasi- $\beta$ -Banach space is a complete quasi- $\beta$ -normed space.

A quasi- $\beta$ -norm  $\|\cdot\|$  is called a  $(\beta, p)$ -norm if  $\|x+y\|^p \le \|x\|^p + \|y\|^p$  for all  $x \in X$ . In this case, a quasi- $(\beta, p)$ -Banach space is called a  $(\beta, p)$ -Banach space. For more details and related stability results on  $(\beta, p)$ -Banach spaces, we refer to [19] [20]. Recently, L. Găvruta and P. Găvruta [21] studied the approximation of functions in Banach space. In this paper, we will consider this problem in  $(\beta, p)$ -Banach spaces and extend previous result for quadratic functional equations.

#### 2. Main Results

Given  $0 < \beta \le 1$  and  $0 . Throughout this paper we always assume that X is a linear space, Y is a <math>(\beta, p)$ -Banach space and  $f: X \to Y$  is a mapping.

**Definition 2.1.** Let  $f: X \to Y$  be a mapping. We say f is  $\Phi$ -approximable by a quadratic map if there exists a quadratic mapping  $Q: X \to Y$  such that

$$\left\|f\left(x\right) - Q\left(x\right)\right\| \le \Phi\left(x\right) \tag{1}$$

for all  $x \in X$ . In this case, we say that Q is the quadratic  $\Phi$ -approximation of f.

The following result is our main result in this paper.

**Theorem 2.2.** Let 
$$V_1 = \left\{ \Phi : X \to \mathbb{R}_+ : \lim_{n \to \infty} 4^{n\beta p} \Phi^p \left( \frac{1}{2^n} x \right) = 0, \forall x \in X \right\}$$
 and

suppose  $\Phi \in V_1$ . Then *f* is  $\Phi$ -approximable by a quadratic map if and only if the following two condition hold:

(i) 
$$\lim_{n \to \infty} 4^{n\beta p} \left\| f\left(\frac{1}{2^n}x + \frac{1}{2^n}y\right) + f\left(\frac{1}{2^n}x - \frac{1}{2^n}y\right) - 2f\left(\frac{1}{2^n}x\right) - 2f\left(\frac{1}{2^n}y\right) \right\|^p = 0$$
,  
 $x, y \in X$ ;

(ii) There exists  $\Psi \in V_1$  such that

$$\left\|f\left(\frac{1}{2^n}x\right) - \frac{1}{4^n}f\left(x\right)\right\|^p \le \Psi^p\left(\frac{1}{2^n}x\right) + \frac{1}{4^{n\beta p}}\Phi^p\left(x\right), x \in X.$$

In this case, the quadratic  $\Phi$ -approximation of *f* is unique and is given by

$$Q(x) = \lim_{n \to \infty} 4^n f\left(\frac{1}{2^n}x\right)$$

for all  $x \in X$ .

**Proof.** We first assume that *f* is  $\Phi$ -approximable by a quadratic map. Then for  $x, y \in X$ , we have

$$\left|f(x+y) - Q(x+y)\right| \le \Phi(x+y)$$

and

$$\left\|f(x-y)-Q(x-y)\right\| \leq \Phi(x-y).$$

It follows that

$$\begin{aligned} & \left\| f(x+y) + f(x-y) - 2f(x) - 2f(y) \right\|^{p} \\ & \leq \left\| f(x+y) - Q(x+y) \right\|^{p} + \left\| f(x-y) - Q(x-y) \right\|^{p} \\ & + \left\| 2f(x) - 2Q(x) \right\|^{p} + \left\| 2f(y) - 2Q(y) \right\|^{p} \\ & \leq \Phi^{p}(x+y) + \Phi^{p}(x-y) + 2^{\beta p} \Phi^{p}(x) + 2^{\beta p} \Phi^{p}(y) \end{aligned}$$

for all  $x, y \in X$ . Hence

$$4^{n\beta p} \left\| f\left(\frac{1}{2^{n}}x + \frac{1}{2^{n}}y\right) + f\left(\frac{1}{2^{n}}x - \frac{1}{2^{n}}y\right) - 2f\left(\frac{1}{2^{n}}x\right) - 2f\left(\frac{1}{2^{n}}y\right) \right\|^{p} \\ \leq 4^{n\beta p} \Phi^{p} \left(\frac{1}{2^{n}}x + \frac{1}{2^{n}}y\right) + 4^{n\beta p} \Phi^{p} \left(\frac{1}{2^{n}}x - \frac{1}{2^{n}}y\right) \\ + 4^{n\beta p} \cdot 2^{\beta p} \Phi^{p} \left(\frac{1}{2^{n}}x\right) + 4^{n\beta p} \cdot 2^{\beta p} \Phi^{p} \left(\frac{1}{2^{n}}y\right)$$

for all  $x, y \in X$ . By letting  $n \to \infty$ , we obtain condition (i) since  $\Phi \in V_1$ . Since Q is quadratic, we have

$$\begin{split} \left\| f\left(\frac{1}{2^{n}}x\right) - \frac{1}{4^{n}}f\left(x\right) \right\|^{p} &\leq \left\| f\left(\frac{1}{2^{n}}x\right) - Q\left(\frac{1}{2^{n}}x\right) \right\|^{p} + \left\| \frac{1}{4^{n}}Q\left(x\right) - \frac{1}{4^{n}}f\left(x\right) \right\|^{p} \\ &\leq \Phi^{p}\left(\frac{1}{2^{n}}x\right) + \frac{1}{4^{n\beta p}}\Phi^{p}\left(x\right) \end{split}$$

for all  $x \in X$  . We take  $\Phi = \Psi \in V_1$  in the first position, then for all  $x \in X$  , we have

$$\left\|f\left(\frac{1}{2^{n}}x\right)-\frac{1}{4^{n}}f\left(x\right)\right\|^{p} \leq \Psi^{p}\left(\frac{1}{2^{n}}x\right)+\frac{1}{4^{n\beta p}}\Phi^{p}\left(x\right)$$

and the condition (ii) holds.

Conversely we suppose that (i) and (ii) hold. It follows from condition (ii) that for all  $x \in X$ , we have

$$\left\|4^{n}f\left(\frac{1}{2^{n}}x\right)-f\left(x\right)\right\|^{p} \leq 4^{n\beta p}\Psi^{p}\left(\frac{1}{2^{n}}x\right)+\Phi^{p}\left(x\right).$$
(2)

Then  $\left\{4^n f\left(\frac{1}{2^n}x\right)\right\}$  is a Cauchy sequence. Indeed, by using  $\frac{1}{2^m}x$  replace x,

we get

$$\left\| 4^n f\left(\frac{1}{2^{n+m}}x\right) - f\left(\frac{1}{2^m}x\right) \right\|^p \le 4^{n\beta p} \Psi^p\left(\frac{1}{2^{n+m}}x\right) + \Phi^p\left(\frac{1}{2^m}x\right),$$

and by multipling  $4^{m\beta p}$ , for all  $x \in X$ , we have

$$\left\| 4^{n+m} f\left(\frac{1}{2^{n+m}} x\right) - 4^m f\left(\frac{1}{2^m} x\right) \right\|^p \le 4^{(n+m)\beta p} \Psi^p\left(\frac{1}{2^{n+m}} x\right) + 4^m \Phi^p\left(\frac{1}{2^m} x\right).$$

Hence, for all  $x \in X$ ,

$$\left\|4^{n+m}f\left(\frac{1}{2^{n+m}}x\right)-4^mf\left(\frac{1}{2^m}x\right)\right\|^p\to 0$$

as  $m, n \to \infty$ . Since *Y* is a  $(\beta, p)$ -Banach space, the limit

$$Q(x) := \lim_{n \to \infty} 4^n f\left(\frac{1}{2^n}x\right)$$
 exists. Let  $n \to \infty$  in relation (2), we get condition (1).

Now we show that Q satisfies the required conditions. From the hypothesis, for all  $x,y\in X$  ,

$$\lim_{n \to \infty} 4^{n\beta p} \left\| f\left(\frac{1}{2^n} x + \frac{1}{2^n} y\right) + f\left(\frac{1}{2^n} x - \frac{1}{2^n} y\right) - 2f\left(\frac{1}{2^n} x\right) - 2f\left(\frac{1}{2^n} y\right) \right\|_p^p = 0.$$

Hence for all  $x, y \in X$ ,

$$||Q(x+y)+Q(x-y)-2Q(x)-2Q(y)||=0.$$

Therefore

$$Q(x+y)+Q(x-y)=2Q(x)+2Q(y)$$

and Q is a quadratic map. Now we show the uniqueness of Q. We suppose that Q satisfies

$$\left\|f(x) - Q(x)\right\| \le \Phi(x)$$

for all  $x \in X$  and there exists a Q' satisfying

$$\left\|f(x)-Q'(x)\right\| \leq \Phi(x).$$

Since Q and Q' are quadratic mappings, we have

$$\left\| f\left(\frac{1}{2^n}x\right) - Q\left(\frac{1}{2^n}x\right) \right\| = \left\| f\left(\frac{1}{2^n}x\right) - \frac{1}{4^n}Q\left(x\right) \right\| \le \Phi\left(\frac{1}{2^n}x\right)$$

for all  $x \in X$ . Hence for all  $x, y \in X$ ,

$$\begin{split} \left\| \mathcal{Q}(x) - \mathcal{Q}'(x) \right\|^p &\leq \left\| \mathcal{Q}(x) - 4^n f\left(\frac{1}{2^n} x\right) \right\|^p + \left\| 4^n f\left(\frac{1}{2^n} x\right) - \mathcal{Q}'(x) \right\|^p \\ &\leq 2 \cdot 4^{n\beta p} \Phi^p \left(\frac{1}{2^n} x\right). \end{split}$$

Since  $\Phi \in V_1$ , for all  $x \in X$ , we have

$$\left\|Q(x)-Q'(x)\right\|^{p} \leq 2\lim_{n\to\infty}4^{n\beta p}\Phi^{p}\left(\frac{1}{2^{n}}x\right)=0.$$

Hence for all  $x \in X$ , Q(x) = Q'(x). This completes the proof. **Corollary 2.3.** Let  $\varphi: X \times X \to [0,\infty)$  be a mapping satisfying

$$\Phi_1^p(x, y) = \sum_{n=0}^{\infty} 4^{n\beta p} \varphi^p\left(\frac{1}{2^{n+1}}x, \frac{1}{2^{n+1}}y\right) < \infty$$

and

$$\lim_{n \to \infty} 4^{n\beta p} \Phi^p \left( \frac{1}{2^n} x \right) = 0$$

for all  $x, y \in X$  where  $\Phi(x) = \Phi_1(x, x)$ . Suppose  $f: X \to Y$  a function with f(0) = 0 and satisfying

$$\left\|f\left(x+y\right)+f\left(x-y\right)-2f\left(x\right)-2f\left(y\right)\right\|^{p} \le \varphi^{p}\left(x,y\right)$$
(3)

for all  $x, y \in X$ . Then there exists a unique quadratic function  $Q: X \to Y$  such that

$$\left\|f\left(x\right) - Q\left(x\right)\right\| \le \Phi\left(x\right), x \in X$$

which is defined

$$Q(x) = \lim_{n \to \infty} 4^n f\left(\frac{1}{2^n}x\right)$$

for all  $x \in X$ .

**Proof.** Replace x and y by  $\frac{1}{2}x$  in (3), we have

$$f(x)-4f\left(\frac{x}{2}\right)^{p} \leq \varphi^{p}\left(\frac{x}{2},\frac{x}{2}\right).$$

Dividing by  $4^{\beta p}$ , we have

$$\left\|\frac{1}{4}f\left(x\right) - f\left(\frac{x}{2}\right)\right\|^{p} \le \frac{1}{4^{\beta p}}\varphi^{p}\left(\frac{x}{2}, \frac{x}{2}\right).$$
(4)

Replacing *x* by  $\frac{1}{2}x$  in (4), we get

$$\left\|\frac{1}{4}f\left(\frac{x}{2}\right) - f\left(\frac{x}{4}\right)\right\|^p \le \frac{1}{4^{\beta p}}\varphi^p\left(\frac{x}{4}, \frac{x}{4}\right).$$
(5)

Then we have

$$\begin{split} \left\| \frac{1}{4^2} f(x) - f\left(\frac{1}{2^2} x\right) \right\|^p &= \left\| \frac{1}{4^2} f(x) - \frac{1}{4} f\left(\frac{x}{2}\right) \right\|^p + \left\| \frac{1}{4} f\left(\frac{x}{2}\right) - f\left(\frac{1}{2^2} x\right) \right\|^p \\ &\leq \frac{1}{4^{2\beta_p}} \varphi^p \left(\frac{x}{2}, \frac{x}{2}\right) + \frac{1}{4^{\beta_p}} \varphi^p \left(\frac{x}{4}, \frac{x}{4}\right) \\ &= \frac{1}{4^{2\beta_p}} \left[ \varphi^p \left(\frac{x}{2}, \frac{x}{2}\right) + 4^{\beta_p} \varphi^p \left(\frac{x}{4}, \frac{x}{4}\right) \right] \\ &\leq \frac{1}{4^{2\beta_p}} \Phi^p \left(x\right) \end{split}$$

for all  $x \in X$ . We claim that

$$\left\|\frac{1}{4^m}f(x) - f\left(\frac{1}{2^m}x\right)\right\|^p \le \frac{1}{4^{m\beta p}}\Phi^p(x).$$
(6)

holds for all  $m \ge 1$  and  $x \in X$ . When m = 1, this is obviously by (4). Suppose (6) holds when m = k, *i.e.* for all  $x \in X$ ,

$$\left\|\frac{1}{4^{k}}f(x)-f\left(\frac{1}{2^{k}}x\right)\right\|^{p} \leq \frac{1}{4^{k\beta_{p}}}\Phi^{p}(x).$$

Then for m = k + 1, we have

$$\begin{split} & \left\| \frac{1}{4^{k+1}} f(x) - f\left(\frac{1}{2^{k+1}} x\right) \right\|^{p} \\ & \leq \left\| \frac{1}{4^{k+1}} f(x) - \frac{1}{4^{k}} f\left(\frac{x}{2}\right) \right\|^{p} + \left\| \frac{1}{4^{k}} f\left(\frac{x}{2}\right) - f\left(\frac{1}{2^{k+1}} x\right) \right\|^{p} \\ & \leq \frac{1}{4^{(k+1)\beta p}} \left[ \varphi^{p}\left(\frac{x}{2}, \frac{x}{2}\right) + 4^{\beta p} \Phi^{p}\left(\frac{x}{2}\right) \right] \\ & \leq \frac{1}{4^{(k+1)\beta p}} \Phi^{p}(x) \end{split}$$

for all  $x \in X$ . By induction, (6) is true for all  $m \ge 1$  and  $x \in X$ . Replacing (x, y) by  $\left(\frac{1}{2^n}x, \frac{1}{2^n}y\right)$  in (3) and multiplying both side by  $4^{n\beta p}$ , we have

$$4^{n\beta p} \left\| f\left(\frac{1}{2^{n}}x + \frac{1}{2^{n}}y\right) + f\left(\frac{1}{2^{n}}x - \frac{1}{2^{n}}y\right) - 2f\left(\frac{1}{2^{n}}x\right) - 2f\left(\frac{1}{2^{n}}y\right) \right\|^{p} \\ \leq 4^{n\beta p} \varphi^{p} \left(\frac{1}{2^{n}}x, \frac{1}{2^{n}}y\right).$$

Since

$$\Phi_1^p(x, y) = \sum_{n=0}^{\infty} 4^{n\beta p} \varphi^p\left(\frac{1}{2^{n+1}}x, \frac{1}{2^{n+1}}y\right) < \infty,$$

we have

$$\lim_{n \to \infty} 4^{n\beta p} \varphi^p \left( \frac{1}{2^{n+1}} x, \frac{1}{2^{n+1}} y \right) = 0$$

for all  $x, y \in X$ . Hence for all  $x, y \in X$ ,

$$\lim_{n \to \infty} 4^{n\beta p} \left\| f\left(\frac{1}{2^n} x + \frac{1}{2^n} y\right) + f\left(\frac{1}{2^n} x - \frac{1}{2^n} y\right) - 2f\left(\frac{1}{2^n} x\right) - 2f\left(\frac{1}{2^n} y\right) \right\|^p = 0.$$

It follows from Theorem 2.2 (with  $\Psi = 0$  there) that there exists a unique quadratic function Q such that

$$\left\|f(x) - Q(x)\right\| \le \Phi(x)$$

for all  $x \in X$ .

**Theorem 2.4.** Let  $V_2 = \left\{ \Phi : X \to \mathbb{R}_+ : \lim_{n \to \infty} \frac{1}{4^{n\beta p}} \Phi^p \left( 2^n x \right) = 0, \forall x \in X \right\}$ . Sup-

pose  $\Phi \in V_2$ . Then f is  $\Phi$ -approximable by a quadratic map if and only if the following two condition

(i) 
$$\lim_{n \to \infty} \frac{1}{4^{n\beta p}} \left\| f\left(2^n x + 2^n y\right) + f\left(2^n x - 2^n y\right) - 2f\left(2^n x\right) - 2f\left(2^n y\right) \right\|^p = 0;$$

#### (ii) There exists a $\Psi \in V_2$ such that

$$\left\|f\left(2^{n}x\right)-4^{n}f\left(x\right)\right\|^{p} \leq \Psi^{p}\left(2^{n}x\right)+4^{n\beta p}\Phi^{p}\left(x\right)$$

hold for all  $x, y \in X$ . In this case, the quadratic  $\Phi$ -approximation of f is unique and is given by

$$Q(x) = \lim_{n \to \infty} \frac{1}{4^n} f(2^n x), x \in X.$$

**Proof.** The proof is similar to that of Theorem 2.2 and we omit it. **Corollary 2.5.** Let  $\varphi: X \times X \rightarrow [0, \infty)$  be a mapping such that

$$\Phi_{1}^{p}(x, y) = \sum_{n=0}^{\infty} 4^{-(n+1)\beta p} \varphi^{p}(2^{n} x, 2^{n} y) < \infty$$

for all  $x, y \in X$ . Let  $\Phi(x) = \Phi_1(x, x)$ . Suppose  $\lim_{n \to \infty} \frac{1}{4^{n\beta p}} \Phi^p(2^n x) = 0$  all

 $x \in X$ . Let  $f: X \to Y$  a function with f(0) = 0 and satisfying

$$\left\|f\left(x+y\right)+f\left(x-y\right)-2f\left(x\right)-2f\left(y\right)\right\|^{p} \leq \varphi^{p}\left(x,y\right)$$

for all  $x, y \in X$ . Then there exists a unique quadratic function  $Q: X \to Y$  such that

$$\left\|f(x) - Q(x)\right\| \le \Phi(x)$$

for all  $x \in X$ .

**Proof.** The proof is similar to that of Corollary 2.3 and we omit it.

#### Funding

This article is partially supported by NSFC (11871303 and 11671133) and NSF of Shandong Province (ZR2019MA039).

#### **Conflicts of Interest**

The authors declare no conflicts of interest regarding the publication of this paper.

#### References

- Ulam, S.M. (1960) A Collection of Mathematical Problems. Interscience Publishers, New York.
- [2] Hyers, D.H. (1941) On the Stability of the Linear Functional Equation. Proceedings of the National Academy of Sciences of the United States of America, 27, 222-224. https://doi.org/10.1073/pnas.27.4.222
- [3] Aoki, T. (1950) On the Stability of the Linear Transformation in Banach Spaces. Journal of the Mathematical Society of Japan, 2, 64-66. <u>https://doi.org/10.2969/jmsj/00210064</u>
- [4] Rassias, Th.M. (1978) On the Stability of the Linear Mapping in Banach Space. Proceedings of the American Mathematical Society, 72, 297-300.

https://doi.org/10.2307/2042795

- [5] Skof, F. (1983) Proprieta locali e approssimazione di operatori. *Rendiconti del Seminario Matematico e Fisico di Milano*, 53, 113-129. https://doi.org/10.1007/BF02924890
- [6] Cholewa, P.W. (1984) Remark on the Stability of Functional Equations. Aequationes Mathematicae, 27, 67-86. <u>https://doi.org/10.1007/BF02192660</u>
- [7] Wang, L.G., Liu, B. and Bai, R. (2010) Stability of a Mixed Type Functional Equation on Multi-Banach Spaces: A Fixed Point Approach. *Fixed Point Theory and Applications*, 2010, Article ID: 283827. <u>https://doi.org/10.1155/2010/283827</u>
- [8] Wang, L.G. (2010) The Fixed Point Method for Intuitionistic Fuzzy Stability of a Quadratic Functional Equation. *Fixed Point Theory and Applications*, 2010, Article ID: 107182. https://doi.org/10.1155/2010/107182
- [9] Wang, L.G. and Li, J. (2012) On the Stability of a Functional Equation Deriving from Additive and Quadratic Functions. *Advances in Difference Equations*, 2012, Article No. 98. <u>https://doi.org/10.1186/1687-1847-2012-98</u>
- [10] Wang, L.G., Xu, K.P. and Liu, Q.W. (2014) On the Stability a Mixed Functional Equation Deriving from Additive, Quadratic and Cubic Mappings. *Acta Mathematica Sinica*, **30**, 1033-1049. https://doi.org/10.1007/s10114-014-3335-9
- [11] Badea, C. (1994) On the Hyers-Ulam Stability of Mappings: The Direct Method. In: Rassias, Th.M. and Tabor, J., Eds., *Stability of Mapping of Hyers-Ulam Type*, Hadronic Press, Palm Harbor, 7-13.
- [12] Czerwik, S. (1992) On the Stability of Quadratic Mapping in Normed Spaces. Abhandlungen aus dem Mathematischen Seminar der Universität Hamburg, 62, 59-64. https://doi.org/10.1007/BF02941618
- [13] Lee, S.H., Im, S.M. and Hawng, I.S. (2005) Quadratic Functional Equation. *Journal of Mathematical Analysis and Applications*, **307**, 387-394. https://doi.org/10.1016/j.jmaa.2004.12.062
- [14] Czerwik, S. (2003) Stability of Functional Equations of Ulam-Hyers-Rassias Type. Hadronic Press, Palm Harbor.
- [15] Jung, S.-M., Popa, D. and Rassias, M.Th. (2014) On the Stability of the Linear Functional Equation in a Single Variable on Complete Metric Groups. *Journal of Global Optimization*, **59**, 165-171. <u>https://doi.org/10.1007/s10898-013-0083-9</u>
- [16] Jung, S.-M. (1996) On the Hyers-Ulam-Rassias Stability of Approximately Additive Mappings. *Journal of Mathematical Analysis and Applications*, **204**, 221-226. <u>https://doi.org/10.1006/jmaa.1996.0433</u>
- [17] Lee, Y.-H., Jung, S.-M. and Rassias, M.Th. (2014) On an n-Dimensional Mixed Type Additive and Quadratic Functional Equation. *Applied Mathematics and Computation*, 228, 13-16. <u>https://doi.org/10.1016/j.amc.2013.11.091</u>
- [18] Mortici, C., Rassias, M.Th. and Jung, S.-M. (2014) On the Stability of a Functional Equation Associated with the Fibonacci Numbers. *Abstract and Applied Analysis*, 2014, Article ID: 546046. <u>https://doi.org/10.1155/2014/546046</u>
- [19] Rassias, J.M. and Kim, H.M. (2009) Generalized Hyers-Ulam Stability for General Additive Functional Equation in Quasi-β-Normed Space. *Journal of Mathematical Analysis and Applications*, 356, 302-309. <u>https://doi.org/10.1016/j.jmaa.2009.03.005</u>
- [20] Wang, L.G. and Liu, B. (2010) The Hyers-Ulam Stability of a Functional Equation Deriving from Quadratic and Cubic Functions in Quasi-β-Normed Spaces. Acta Mathematica Sinica, English Series, 26, 2335-2348. https://doi.org/10.1007/s10114-010-9330-x

[21] Găvruta, L. and Găvruta, P. (2016) Approximation of Functions by Additive and by Quadratic Mappings. In: Rassias, T.M. and Gupta, V., Eds., *Mathematical Analysis, Approximation Theorem and Their Applications*, Springer Optimization and Its Applications, Berlin, 281-292. <u>https://doi.org/10.1007/978-3-319-31281-1\_12</u>



## Optimization in Transition between Two Dynamic Systems Governed by a Class of Weakly Singular Integro-Differential Equations

#### **Shihchung Chiang**

Department of Finance, Chung Hua University, Taiwan Email: Chiang@chu.edu.tw

How to cite this paper: Chiang, S. (2019) Optimization in Transition between Two Dynamic Systems Governed by a Class of Weakly Singular Integro-Differential Equations. *Applied Mathematics*, **10**, 826-835. https://doi.org/10.4236/am.2019.1010059

Received: September 2, 2019 Accepted: October 7, 2019 Published: October 10, 2019

Copyright © 2019 by author(s) and Scientific Research Publishing Inc. This work is licensed under the Creative Commons Attribution International License (CC BY 4.0).

http://creativecommons.org/licenses/by/4.0/

CC O Open Access

#### Abstract

This study presents numerical methods for solving the minimum energies that satisfy typical optimal requirements in the transition between two dynamic systems where each system is governed by a different kind of weakly singular integro-differential equation. The class of weakly singular integro-differential equations originates from mathematical models in aeroelasticity. The proposed numerical methods are based on earlier reported approximation schemes for the equations of the first kind and the second kind. The main result of this study is the development of numerical techniques for determining the stability between two dynamic systems in the minimum energy sense.

#### **Keywords**

Optimal Requirement, Transition, Weakly Singular Integro-Differential Equations, Stability

#### **1. Introduction**

The minimum energy problem and the associated optimal control problem have been investigated for more than half a century. The system constraints can be ordinary differential equations, partial differential equations, or functional differential equations. This study introduces a numerical method for finding the minimum energy to satisfy the general criterion that can be adjusted to minimize various requirements through the selection of appropriate parameters. One system constraint is the class of equations of the first kind, which originates from an aeroelasticity problem where the mathematical model consists of eight integro-differential equations [1]. In the model, the most determinate equation is a scalar weakly singular integro-differential equation of the first kind [2] [3]. Furthermore, because of the natural facts of transition between liquid water and solid ice [4] or the aviation transition between vertical take-off and horizontal flight of an unmanned aerial vehicle [5], we were interested in the energy issue in the transition between two basically different (but related) dynamic systems. For the setting, the second dynamic system was constructed from the first system using finite derivative delay terms that included the boundary points of the considered interval. This study followed the structure of other relevant studies [6] in assuming that the forcing terms of the system are the control forces. This study is organized as follows: Section 2 presents the criteria for the optimal issues. Section 3 presents the approach for determining the minimum energy for the transition procedure. Section 4 presents the numerical results attained by choosing different parameters for various cost requirements. Section 5 presents the summary of this study.

#### 2. The Model

Consider the class of weakly singular integro-differential equations of the first kind

$$\frac{\mathrm{d}}{\mathrm{d}t}Dx_t = u\left(t\right) \tag{1}$$

with initial data

$$x(s) = \phi(s), \quad -b \le s \le 0.$$
<sup>(2)</sup>

The difference operator D is defined as

$$Dx_{t} = \int_{-b}^{0} g\left(s\right) x_{t}\left(s\right) \mathrm{d}s, \qquad (3)$$

where

$$x_t(s) = x(t+s). \tag{4}$$

The weighting kernel g is integrable, positive, nondecreasing, and weakly singular at s = 0. The control force u(t) is assumed to be locally integrable for t > 0. Although a more general kernel g also works, this study focused on the Abel-type kernel (*i.e.*,  $g(s) = |s|^{-p}$ , where  $s \in [-b, 0]$  and p = 0.5 from the original aeroelastic model).

The initial condition  $\phi(s), -b \le s \le 0$  is in  $L_{1,g}$ , which is a weighted  $L_1$  space with weight  $g(\cdot)$ . Note that the initial value problem in Equations (1)-(2) can be written as

$$Dx_t = Dx_0 + \int_0^t u(\tau) \mathrm{d}\tau, \qquad (5)$$

provided that the function

$$Dx_{t} = \int_{-b}^{0} g\left(s\right) x\left(t+s\right) \mathrm{d}s \tag{6}$$

is absolutely continuous for t > 0 and the function  $g(\cdot)\phi(\cdot)$  belongs to  $L_1[-b,0]$ . Without a loss of generality, we assume that b = 1.

The second system is a class of weakly singular integro-differential equations of the second kind

$$\sum_{i=1}^{l} \frac{\mathrm{d}}{\mathrm{d}t} x \left( t - \sigma_i \right) + \frac{\mathrm{d}}{\mathrm{d}t} D x_i = u \left( t \right), \tag{7}$$

where *l* is a positive integer and  $0 \le \sigma_i \le 1, i = 1, \dots, l$ . The initial condition is

$$x(s) = \phi(s), -1 \le s \le 0.$$
 (8)

For the partition between systems (2) and (3), a parameter  $\lambda \in [0,1]$  is assumed. Therefore, the combined system can be written as

$$\sum_{i=1}^{l} \frac{\mathrm{d}}{\mathrm{d}t} \lambda x (t - \sigma_i) + \frac{\mathrm{d}}{\mathrm{d}t} D \lambda x_i = u(t)$$

$$\left| \frac{\mathrm{d}}{\mathrm{d}t} D (1 - \lambda) x_i = v(t) \right|$$
(9)

with initial data

$$x(s) = \phi(s), -1 \le s \le 0.$$
 (8)

Although the proposed methods can be applied to more general cost functions, this study primarily considered the typical cost function for comparison:

$$\Phi(\lambda) = \Phi_1(\lambda) + \Phi_2(\lambda), \tag{10}$$

and

$$\Phi_{1}(\lambda) = \alpha_{1}(\lambda x(1) - h)^{2} + \alpha_{2} \int_{0}^{1} (\lambda x(t) - \eta(t))^{2} dt + \alpha_{3} \int_{0}^{1} u(t)^{2} dt, \qquad (11)$$

$$\Phi_{2}(\lambda) = \alpha_{1}((1-\lambda)x(1)-h)^{2} + \alpha_{2}\int_{0}^{1}((1-\lambda)x(t)-\eta(t))^{2} dt + \alpha_{3}\int_{0}^{1}\nu(t)^{2} dt, \quad (12)$$

where *h* is a constant of final target state,  $\eta(t)$  is a target function, and parameters  $\alpha_1, \alpha_2$  and  $\alpha_3$  are nonnegative constants with a total sum of 1.

#### 3. The Numerical Method

This procedure is proposed to discretize system (9) and the cost function (10) simultaneously to construct two corresponding linear systems with unknowns as states and controls. The space mesh points (corresponding to the *s* variable) are discretized as  $-1 = \tau_n < \tau_{n-1} < \cdots < \tau_1 < \tau_0 = 0$ , and a new variable  $\xi$  is defined as

$$\xi(t,s) = x(t+s), \ -1 \le s \le 0, \ t > 0.$$
(13)

System (9) can then be reformulated as a first-order hyperbolic equation

$$\frac{\partial}{\partial t}\xi(t,s) = \frac{\partial}{\partial s}\xi(t,s), \quad -1 \le s \le 0, \tag{14}$$

with the condition

$$\begin{cases} \lambda \sum_{i=1}^{l} \frac{\mathrm{d}}{\mathrm{d}t} \xi(t, -\sigma_{i}) + \lambda \int_{-1}^{0} |s|^{-p} \frac{\partial}{\partial s} \xi(t, s) \mathrm{d}s = u(t), \\ (1-\lambda) \int_{-1}^{0} |s|^{-p} \frac{\partial}{\partial s} \xi(t, s) \mathrm{d}s = v(t). \end{cases}$$
(15)

Next, assume that the solution to Equation (8) has the form

$$\xi(t,s) = \sum_{i=0}^{n} \kappa_i(t) B_i(s), \qquad (16)$$

where the basis,  $B_i(s), i = 0, \dots, n$  is given by

$$B_{i}(s) = \begin{cases} \frac{1}{(\tau_{i} - \tau_{i+1})} (s - \tau_{i+1}) & s \in [\tau_{i+1}, \tau_{i}], \\ \frac{1}{(\tau_{i-1} - \tau_{i})} (\tau_{i-1} - s) & s \in [\tau_{i}, \tau_{i-1}], \\ 0 & \text{otherwise.} \end{cases}$$
(17)

Namely,  $B_i(s), i = 0, \dots, n$  are piecewise linear functions. After substituting the special form of  $\xi$  in Equation (16) into Equations (14)-(15), the governing equations for  $\kappa_i(t), i = 0, \dots, n$  become the following:

$$\frac{\mathrm{d}}{\mathrm{d}t}\kappa_{i}\left(t\right) = \frac{1}{\delta_{i}}\left(\kappa_{i-1}\left(t\right) - \kappa_{i}\left(t\right)\right), i = 1, \cdots, n,$$
(18)

$$\begin{cases} \lambda \sum_{i=1}^{l} \frac{\mathrm{d}}{\mathrm{d}t} \kappa_{\bar{\sigma}_{i}}\left(t\right) + \lambda \int_{-1}^{0} \left|s\right|^{-p} \sum_{i=0}^{n} \kappa_{i}\left(t\right) \frac{\mathrm{d}}{\mathrm{d}s} B_{i}\left(s\right) \mathrm{d}s = u\left(t\right), \\ \left(1 - \lambda\right) \int_{-1}^{0} \left|s\right|^{-p} \sum_{i=0}^{n} \kappa_{i}\left(t\right) \frac{\mathrm{d}}{\mathrm{d}s} B_{i}\left(s\right) \mathrm{d}s = v\left(t\right), \end{cases}$$

$$\tag{19}$$

where  $\delta_i = \tau_{i-1} - \tau_i > 0$ , for  $i = 1, \dots, n$ . For time *t*, discretization contains  $T^0, T^1, \dots, T^m$ , for  $0 = T^0 < T^1 < \dots < T^m = 1$ . Define  $\Delta^k = T^{k+1} - T^k$ , for  $k = 0, \dots, m-1$ . By assuming  $\alpha_i^k = \kappa_i (T^k)$ , for  $i = 0, 1, \dots, n$ , and  $k = 0, \dots, m$ , and without losing generality, we assume l = 2,  $\overline{\sigma_1} = 0$ ,  $\overline{\sigma_2} = n$ , and Equations (18)-(19) can now be written as

$$\frac{1}{\Delta^k} \left( \alpha_i^{k+1} - \alpha_i^k \right) = \frac{1}{\delta_i} \left( \alpha_{i-1}^k - \alpha_i^k \right), \tag{20}$$

$$\begin{cases} \frac{\lambda}{\delta_{1}}\alpha_{0}^{k+1} - \frac{\lambda}{\delta_{1}}\alpha_{1}^{k+1} + \frac{\lambda}{\delta_{n-1}}\alpha_{n-1}^{k+1} - \frac{\lambda}{\delta_{n-1}}\alpha_{n}^{k+1} + \lambda\sum_{i=1}^{n}\frac{g_{i}}{\delta_{i}}\left(\alpha_{i-1}^{k+1} - \alpha_{i}^{k+1}\right) = u\left(T^{k+1}\right),\\ \left(1 - \lambda\right)\sum_{i=1}^{n}\frac{g_{i}}{\delta_{i}}\left(\alpha_{i-1}^{k+1} - \alpha_{i}^{k+1}\right) = v\left(T^{k+1}\right),\end{cases}$$
(21)

for  $i = 1, \dots, n$ ,  $k = 0, \dots, m-1$ , and  $g_i = \int_{\tau_i}^{\tau_{i-1}} |s|^{-p} ds$ .

Furthermore, we assume a uniform mesh for both space and time, and the mesh points are  $\tau_i$ ,  $i = 0, \dots, n$  and  $T^k$ ,  $k = 0, \dots, m$ . Specifically, we have  $\tau_i = -\frac{i}{n}$ ,  $T^k = \frac{k}{m}$ , for some positive integers n and m. The associated differences are defined as  $\Delta^k = T^{k+1} - T^k$ ,  $k = 0, \dots, m-1$ , for the time variable and  $\delta_i = \tau_{i-1} - \tau_i$ ,  $i = 1, \dots, n$ , for the space variable. Thus, we obtain  $\Delta^k = 1/m$  and  $\delta_i = 1/n$ , for  $k = 0, \dots, m-1$ , and  $i = 1, \dots, n$ . Setting m = n produces the relation

 $\Delta^k = \delta_i = 1/n$  for  $k = 0, \dots, n-1$ , and  $i = 1, \dots, n$ , and deriving Equations (20)-(21) lead to the following system:

$$\alpha_i^{k+1} = \alpha_{i-1}^k, \tag{22}$$

and

$$\begin{bmatrix}
\frac{\lambda}{\delta_{1}}\alpha_{0}^{k+1} - \frac{\lambda}{\delta_{1}}\alpha_{1}^{k+1} + \frac{\lambda}{\delta_{n-1}}\alpha_{n-1}^{k+1} - \frac{\lambda}{\delta_{n-1}}\alpha_{n}^{k+1} + \lambda\sum_{i=1}^{n}\frac{1}{\delta_{i}}\left(\alpha_{i-1}^{k+1} - \alpha_{i}^{k+1}\right) \cdot \frac{1}{1-p}\left[-\left(-\tau_{i-1}\right)^{1-p} + \left(-\tau_{i}\right)^{1-p}\right] = u\left(T^{k+1}\right) = u_{k+1} \\
\left(1-\lambda\right)\sum_{i=1}^{n}\frac{1}{\delta_{i}}\left(\alpha_{i-1}^{k+1} - \alpha_{i}^{k+1}\right) \cdot \frac{1}{1-p}\left[-\left(-\tau_{i-1}\right)^{1-p} + \left(-\tau_{i}\right)^{1-p}\right] = v\left(T^{k+1}\right) = v_{k+1}
\end{cases}$$
(23)

for  $i = 1, \dots, n$ , and  $k = 0, \dots, n-1$ .

After defining corresponding constants  $c_0, c_1, \dots, c_n$ , and  $d_0, d_1, \dots, d_n$ , Equation (23) can be written in the following simplified form:

$$\begin{cases} \lambda \left( \alpha_0^{k+1} c_0 + \alpha_0^k c_1 + \dots + \alpha_0^0 c_{k+1} + \dots + \alpha_{n-k-1}^0 c_n \right) = u_{k+1} \\ \left( 1 - \lambda \right) \left( \alpha_0^{k+1} d_0 + \alpha_0^k d_1 + \dots + \alpha_0^0 d_{k+1} + \dots + \alpha_{n-k-1}^0 d_n \right) = v_{k+1} \end{cases}, \quad k = 0, \dots, n-1 \quad (24)$$

The connection between the solution x(t) and a's is as follows: Because  $\xi(t,s) = x(t+s)$ , for  $-1 \le s \le 0$ , t > 0, and  $\xi(t,s) = \sum_{i=0}^{n} \kappa_i(t) B_i(s)$ , it follows that  $x(\overline{t})$ , for  $\overline{t} > 0$  can be obtained in the following case:

$$x(T^{j}) = \sum_{l=0}^{n} \kappa_{l}(T^{j}) B_{l}(0) = \kappa_{0}(T^{j}) = \alpha_{0}^{j}, \text{ for } j = 1, \cdots, n.$$

$$(25)$$

For the cost function

$$\Phi_1(\lambda) = \alpha_1 \left(\lambda x(1) - h\right)^2 + \alpha_2 \int_0^1 \left(\lambda x(t) - \eta(t)\right)^2 dt + \alpha_3 \int_0^1 u(t)^2 dt,$$
  
ratigad form is:

the discretized form is:

$$\Phi_{1}(\lambda) = \alpha_{1}(\lambda\alpha_{0}^{n} - h)^{2} + \alpha_{2}\frac{1}{n}\sum_{i=1}^{n}(\lambda\alpha_{0}^{i} - \eta(T^{i}))^{2} + \alpha_{3}\frac{1}{n}\sum_{i=1}^{n}u_{i}^{2}.$$
 (26)

Taking the first derivatives of  $\Phi_1(\lambda)$  with respect to  $u_i, i = 1, \dots, n$ , and setting them to zero yields the following equations:

$$n\lambda^{2}\alpha_{1} \cdot \alpha_{0}^{n} \cdot aa(n) + \lambda^{2}\alpha_{2} \Big[ aa(1) \cdot \alpha_{0}^{1} + aa(2) \cdot \alpha_{0}^{2} + \dots + aa(n) \cdot \alpha_{0}^{n} \Big] + \alpha_{3} \cdot u_{1}$$

$$= n\lambda\alpha_{1}h \cdot aa(n) + \lambda\alpha_{2} \Big[ \eta(t_{1}) \cdot aa(1) + \dots + \eta(t_{n}) \cdot aa(n) \Big],$$

$$\vdots$$

$$n\lambda^{2}\alpha_{1} \cdot \alpha_{0}^{n} \cdot aa(n-j+1) + \lambda^{2}\alpha_{2} \Big[ aa(1) \cdot \alpha_{0}^{j} + aa(2) \cdot \alpha_{0}^{j+1} + \dots + aa(n-j+1) \cdot \alpha_{0}^{n} \Big] + \alpha_{3} \cdot u_{j}$$

$$= n\lambda\alpha_{1}h \cdot aa(n-j+1) + \lambda\alpha_{2} \Big[ \eta(t_{j}) \cdot aa(1) + \dots + \eta(t_{n}) \cdot aa(n-j+1) \Big],$$

$$\vdots$$

$$n\lambda^{2}\alpha_{1} \cdot \alpha_{0}^{n} \cdot aa(1) + \lambda^{2}\alpha_{2} \cdot \alpha_{0}^{n} \cdot aa(1) + \alpha_{3} \cdot u_{n}$$

$$= n\lambda\alpha_{1}h \cdot aa(1) + \lambda\alpha_{2} \cdot \eta(t_{n}) \cdot aa(1),$$
(27)

where

$$aa(1) = \frac{1}{\lambda c_0},$$
  

$$aa(2) = -\frac{c_1}{c_0} \cdot aa(1),$$
  

$$\vdots$$
  

$$aa(j) = -\frac{c_1}{c_0} aa(j-1) - \frac{c_2}{c_0} \cdot aa(j-2) - \dots - \frac{c_{j-1}}{c_0} \cdot aa(1),$$

DOI: 10.4236/am.2019.1010059

$$aa(n) = -\frac{c_1}{c_0}aa(n-1) - \frac{c_2}{c_0} \cdot aa(n-2) - \dots - \frac{c_{n-1}}{c_0} \cdot aa(1).$$

:

Systems (24) with  $\lambda$  and (27) can be set up as [A][x] = [b], where the vector [x] consists of the unknowns  $\alpha_0^j$ ,  $j = 1, \dots, n$ , and  $u_k, k = 1, \dots, n$ . The structure of matrix [A] is

$$\begin{bmatrix} \lambda c_{0} & 0 & \cdots & 0 & -1 & 0 & \cdots & 0 \\ \lambda c_{1} & \lambda c_{0} & \cdots & 0 & 0 & -1 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \ddots & \vdots \\ \lambda c_{n-1} & \lambda c_{n-2} & \cdots & c_{0} & 0 & 0 & \cdots & -1 \\ \lambda^{2} \alpha_{2} aa(1) & \lambda^{2} \alpha_{2} aa(2) & \cdots & \lambda^{2} (\alpha_{2} + n\alpha_{1}) aa(n) & \alpha_{3} & 0 & \cdots & 0 \\ 0 & \lambda^{2} \alpha_{2} aa(1) & \cdots & \vdots & 0 & \alpha_{3} & \cdots & 0 \\ \vdots & \vdots & \ddots & \lambda^{2} (\alpha_{2} + n\alpha_{1}) aa(3) & \vdots & \vdots & \ddots & \vdots \\ \vdots & \vdots & \lambda^{2} \alpha_{2} aa(1) & \lambda^{2} (\alpha_{2} + n\alpha_{1}) aa(2) & & \alpha_{3} & 0 \\ 0 & 0 & \cdots & \lambda^{2} (\alpha_{2} + n\alpha_{1}) aa(1) & 0 & \cdots & 0 & \alpha_{3} \end{bmatrix}_{2n \times 2n}$$

and vector [b] is given by

$$\begin{split} & \lambda \begin{bmatrix} -\alpha_{0}^{0}c_{1} - \alpha_{1}^{0}c_{2} - \dots - \alpha_{n-1}^{0}c_{n} - b(t_{1}) \\ -\alpha_{0}^{0}c_{2} - \alpha_{1}^{0}c_{3} - \dots - \alpha_{n-2}^{0}c_{n} - b(t_{2}) \\ -\alpha_{0}^{0}c_{3} - \alpha_{1}^{0}c_{4} - \dots - \alpha_{n-3}^{0}c_{n} - b(t_{3}) \\ \vdots \\ & \ddots \\ & -\alpha_{0}^{0}c_{n} - b(t_{n}) \\ \alpha_{2} \Big[ \eta(t_{1})aa(1) + \dots + \eta(t_{n-1})aa(n-1) \Big] + \Big[ \alpha_{2}\eta(t_{n}) + n\alpha_{1}h \Big] aa(n) \\ \alpha_{2} \Big[ \eta(t_{2})aa(1) + \dots + \eta(t_{n-1})aa(n-2) \Big] + \Big[ \alpha_{2}\eta(t_{n}) + n\alpha_{1}h \Big] aa(n-1) \\ \vdots \\ & \left[ \alpha_{2}\eta(t_{n}) + n\alpha_{1}h \Big] aa(1) \end{bmatrix}_{2n \times 1} \end{split}$$

For the cost function

 $\Phi_2(\lambda) = \alpha_1 \left( (1-\lambda) x(1) - h \right)^2 + \alpha_2 \int_0^1 \left( (1-\lambda) x(t) - \eta(t) \right)^2 dt + \alpha_3 \int_0^1 v(t)^2 dt,$ the discretized form is:

$$\Phi_{2}(\lambda) = \alpha_{1}((1-\lambda)\alpha_{0}^{n}-h)^{2} + \alpha_{2}\frac{1}{n}\sum_{i=1}^{n}((1-\lambda)\alpha_{0}^{i}-\eta(T^{i}))^{2} + \alpha_{3}\frac{1}{n}\sum_{i=1}^{n}v_{i}^{2}.$$
 (28)

Taking first derivatives of  $\Phi_2(\lambda)$  with respect to  $v_i, i = 1, \dots, n$ , and setting them to zero produces the following equations:

$$n(1-\lambda)^{2} \alpha_{1} \cdot \alpha_{0}^{n} \cdot aa(n) + (1-\lambda)^{2} \alpha_{2} \Big[ aa(1) \cdot \alpha_{0}^{1} + aa(2) \cdot \alpha_{0}^{2} + \cdots \\ + aa(n) \cdot \alpha_{0}^{n} \Big] + \alpha_{3} \cdot v_{1}$$

$$= n(1-\lambda)\alpha_{1}h \cdot aa(n) + (1-\lambda)\alpha_{2} \Big[ \eta(t_{1}) \cdot aa(1) + \cdots + \eta(t_{n}) \cdot aa(n) \Big],$$

$$\vdots$$

$$n(1-\lambda)^{2} \alpha_{1} \cdot \alpha_{0}^{n} \cdot aa(n-j+1) + (1-\lambda)^{2} \alpha_{2} \Big[ aa(1) \cdot \alpha_{0}^{j} + aa(2) \cdot \alpha_{0}^{j+1} + \cdots \\ + aa(n-j+1) \cdot \alpha_{0}^{n} \Big] + \alpha_{3} \cdot v_{j}$$

$$= n(1-\lambda)\alpha_{1}h \cdot aa(n-j+1) + \lambda\alpha_{2} \Big[ \eta(t_{j}) \cdot aa(1) + \cdots + \eta(t_{n}) \cdot aa(n-j+1) \Big],$$

$$(1-\lambda)^{2} \alpha_{0}^{n} (n \cdot \alpha_{1} \cdot aa(1)) + \alpha_{2} \cdot aa(1) + \alpha_{3} \cdot \nu_{n}$$
  
=  $(1-\lambda) aa(1) [n\alpha_{1}h + \alpha_{2} \cdot \eta(t_{n})].$  (29)

where

$$aa(1) = \frac{1}{(1-\lambda)c_0},$$

$$aa(2) = -\frac{c_1}{c_0} \cdot aa(1),$$

$$\vdots$$

$$aa(j) = -\frac{c_1}{c_0}aa(j-1) - \frac{c_2}{c_0} \cdot aa(j-2) - \dots - \frac{c_{j-1}}{c_0} \cdot aa(1),$$

$$\vdots$$

$$aa(n) = -\frac{c_1}{c_0}aa(n-1) - \frac{c_2}{c_0} \cdot aa(n-2) - \dots - \frac{c_{n-1}}{c_0} \cdot aa(1).$$

Systems (24) with  $1-\lambda$  and (29) can be set up as [A][x] = [b], where the vector [x] consists of the unknowns  $\alpha_0^j$ ,  $j = 1, \dots, n$ , and  $\nu_k, k = 1, \dots, n$ . The structure of matrix [A] is

$$\begin{bmatrix} (1-\lambda)d_0 & 0 & \cdots & 0 & -1 & 0 & \cdots & 0\\ (1-\lambda)d_1 & (1-\lambda)d_0 & \cdots & 0 & 0 & -1 & \cdots & 0\\ \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \ddots & \vdots\\ (1-\lambda)d_{n-1} & (1-\lambda)d_{n-2} & \cdots & d_0 & 0 & 0 & \cdots & -1\\ (1-\lambda)^2\alpha_2aa(1) & (1-\lambda)^2\alpha_2aa(2) & \cdots & (1-\lambda)^2(\alpha_2+n\alpha_1)aa(n) & \alpha_3 & 0 & \cdots & 0\\ 0 & (1-\lambda)^2\alpha_2aa(1) & \cdots & \vdots & 0 & \alpha_3 & \cdots & 0\\ \vdots & \vdots & \ddots & (1-\lambda)^2(\alpha_2+n\alpha_1)aa(3) & \vdots & \vdots & \ddots & \vdots\\ \vdots & \vdots & (1-\lambda)^2\alpha_2aa(1) & (1-\lambda)^2(\alpha_2+n\alpha_1)aa(2) & \alpha_3 & 0\\ 0 & 0 & \cdots & (1-\lambda)^2(\alpha_2+n\alpha_1)aa(1) & 0 & \cdots & 0 & \alpha_3 \end{bmatrix}_{2n\times 2n}$$

and vector [b] is given by

$$(1-\lambda) \begin{bmatrix} -\alpha_{0}^{0}d_{1} - \alpha_{1}^{0}d_{2} - \dots - \alpha_{n-1}^{0}d_{n} - b(t_{1}) \\ -\alpha_{0}^{0}d_{2} - \alpha_{1}^{0}d_{3} - \dots - \alpha_{n-2}^{0}d_{n} - b(t_{2}) \\ -\alpha_{0}^{0}d_{3} - \alpha_{1}^{0}d_{4} - \dots - \alpha_{n-3}^{0}d_{n} - b(t_{3}) \\ \vdots \\ -\alpha_{0}^{0}d_{n} - b(t_{n}) \\ \alpha_{2} [\eta(t_{1})aa(1) + \dots + \eta(t_{n-1})aa(n-1)] + [\alpha_{2}\eta(t_{n}) + n\alpha_{1}h]aa(n) \\ \alpha_{2} [\eta(t_{2})aa(1) + \dots + \eta(t_{n-1})aa(n-2)] + [\alpha_{2}\eta(t_{n}) + n\alpha_{1}h]aa(n-1) \\ \vdots \\ [\alpha_{2}\eta(t_{n}) + n\alpha_{1}h]aa(1) \end{bmatrix}_{2n \times 1}$$

DOI: 10.4236/am.2019.1010059

#### **4. Numerical Examples**

Consider examples involving p = 0.5,  $\lambda \in [0,1]$ , initial conditions  $\phi(s) = 0, -1 \le s \le 0$ , different target final state *h*, and different target functions  $\eta(t), 0 \le t \le 1$ . For different criteria, the combinations of constants *a*'s in the cost functions are changed accordingly.

For the case  $(\alpha_1, \alpha_2, \alpha_3) = (0, 1, 0)$ , the problem is the "tracking problem".

Typical cost distribution is as the following two graphs (Figure 1 and Figure 2).

Example 1: n = 100,  $(\alpha_1, \alpha_2, \alpha_3) = (0.3, 0.5, 0.2)$ 

h = 1	$\eta(t) = 1$	mincost $\Phi = 0.5951$	when $\lambda = 0$
h = 1	$\eta(t) = t$	mincost $\Phi = 0.3313$	when $\lambda = 0$
h = 0	$\eta(t) = 1 - t$	mincost $\Phi = 0.1517$	when $\lambda = 0$

Example 2: n = 100,  $(\alpha_1, \alpha_2, \alpha_3) = (0, 0, 1)$ 

h = 1	$\eta(t) = 1$	mincost $\Phi = 0$	when $\lambda = 0$
h = 1	$\eta(t) = t$	mincost $\Phi = 0$	when $\lambda = 0$
h = 0	$\eta(t) = 1 - t$	mincost $\Phi = 0$	when $\lambda = 0$

Example 3: n = 100,  $(\alpha_1, \alpha_2, \alpha_3) = (1, 0, 0)$ 

h = 1	$\eta(t) = 1$	mincost $\Phi = 2.2132e - 28$	when $\lambda = 0$
h = 1	$\eta(t) = t$	mincost $\Phi = 2.2132e - 28$	when $\lambda = 0$
h = 0	$\eta(t) = 1 - t$	mincost $\Phi = 0$	when $\lambda = 0$

Example 4: n = 100,  $(\alpha_1, \alpha_2, \alpha_3) = (0, 1, 0)$ 

h = 1	$\eta(t) = 1$	mincost $\Phi = 5.8587e - 29$	when $\lambda = 0$
h = 1	$\eta(t) = t$	mincost $\Phi = 2.5009e - 29$	when $\lambda = 0$
h = 0	$\eta(t) = 1 - t$	mincost $\Phi = 2.9966e - 29$	when $\lambda = 0$

#### Example 5: n = 100, $(\alpha_1, \alpha_2, \alpha_3) = (0.9, 0, 0.1)$

h = 1	$\eta(t) = 1$	mincost $\Phi = 0.3424$	when $\lambda = 0$
h = 1	$\eta(t) = t$	mincost $\Phi = 0.3424$	when $\lambda = 0$
h = 0	$\eta(t) = 1 - t$	mincost $\Phi = 0$	when $\lambda = 0$

Example 6: n = 100,  $(\alpha_1, \alpha_2, \alpha_3) = (0, 0.9, 0.1)$ 

h = 1	$\eta(t) = 1$	mincost $\Phi = 0.5712$	when $\lambda = 0$
h = 1	$\eta(t) = t$	mincost $\Phi = 0.1785$	when $\lambda = 0.5$
h = 0	$\eta(t) = 1 - t$	mincost $\Phi = 0.2321$	when $\lambda = 0$



**Figure 1.** Total cost for  $\lambda$  from 0 to 1.



**Figure 2.** Total cost for  $\lambda$  from 0 to 1.

#### **5.** Conclusion

This study presented a numerical method for finding the minimum of the total cost when it contains two partial costs from two dynamic systems, and each cost contains three weights to adjust for different considerations of energy and different combinations of the measurable parameter  $\lambda$  between two systems. The effectiveness of the proposed method was tested by examples. The numerical re-

sults indicated that the most stable situations are  $\lambda = 0$ . In other words, dynamic system with the first kind integro-differential equation is the most stable system in the minimum cost sense.

#### Acknowledgements

Author would like to thank MOST (Ministry of Science and Technology) under grant No.108-2914-I-216-002-A1 to partially support this project.

#### **Conflicts of Interest**

The author declares no conflicts of interest regarding the publication of this paper.

#### References

- Burns, J.A., Cliff, E.M. and Herdman, T.L. (1983) A State-Space Model for an Aeroelastic System. *Proceedings of 22nd IEEE Conference on Decision and Control*, San Antonio, December 1983, 1074-1077. <u>https://doi.org/10.1109/CDC.1983.269685</u>
- [2] Burns, J.A., Herdman, T.L. and Stech, H.W. (1983) Linear Functional Differential Equations as Semigroups on Product Spaces. *SIAM Journal on Mathematical Analysis*, 14, 98-116. <u>https://doi.org/10.1137/0514007</u>
- [3] Kappel, F. and Zhang, K.P. (1986) Equivalence of Functional Equations of Neutral Type and Abstract Cauchy Problems. *Monatshefte für Mathematik*, 101, 115-133. <u>https://doi.org/10.1007/BF01298925</u>
- [4] Vrbka, L. and Jungwirth, P. (2006) Homogeneous Freezing of Water Starts in the Subsurface. *The Journal of Physical Chemistry B*, **110**, 18126-18129. <u>https://doi.org/10.1021/jp064021c</u>
- [5] Vorsin, D. and Arogeti, S. (2017) Flight Transition Control of a Multipurpose UAV.
   2017 13 th IEEE International Conference on Control & Automation (ICCA), Ohrid,
   3-6 July 2017, 507-512. https://doi.org/10.1109/ICCA.2017.8003112
- [6] Chiang, S.C. and Herdman, T.L. (2013) Revised Numerical Methods for Optimal Control of a Class of Singular Integro-Differential Equations. *Mathematics in Engineering, Science and Aerospace*, 4, 171-187.



## Inequality of Realization of a Stochastic Dynamics Based on the Erdös Discrepancy Problem

#### Hiroyuki Kato

Faculty of Management and Economics, Kaetsu University, Tokyo, Japan Email: hiroyuki-kat0@kaetsu.ac.jp

How to cite this paper: Kato, H. (2019) Inequality of Realization of a Stochastic Dynamics Based on the Erdös Discrepancy Problem. *Applied Mathematics*, **10**, 836-847. https://doi.org/10.4236/am.2019.1010060

Received: September 8, 2019 Accepted: October 9, 2019 Published: October 12, 2019

Copyright © 2019 by author(s) and Scientific Research Publishing Inc. This work is licensed under the Creative Commons Attribution International License (CC BY 4.0). http://creativecommons.org/licenses/by/4.0/

co 🛈 Open Access

#### Abstract

This paper proposes a stochastic dynamics model in which people who are endowed with different discount factors chose to buy the capital stock periodically with different periodicities and are exposed to randomness at arithmetic progression times. We prove that the realization of a stochastic equilibrium may render to the people quite unequal benefits. Its proof is based on Erdös Discrepancy Problem that an arithmetic progression sum of any sign sequence goes to infinity, which is recently solved by Terence Tao [1]. The result in this paper implies that in some cases, the sources of inequality come from pure luck.

#### **Keywords**

Erdös Discrepancy Problem, Arithmetic Progression, Inequality, Economic Dynamics

#### **1. Introduction**

The existence of inequality of wages, assets, and other incomes in a society has been gaining wide attention recently especially since Pikkety [2] (Atkinson *et al.* [3], Gabaix *et al.* [4], Grossman and Helpman [5], Jones [6], Jones and Kim [7], Kasa and Lei [8], Mankiw [9] to name only a few). Many researchers tackled this problem by providing models that explain the empirical data, say, the large gap between capital income and labor one, or inequality among labor incomes, and its extent of that inequality. They employ growth models that endogenously induce the inequality underlining the market mechanism. However, whether the inequality is the problem that needs some remedy or should be taken as mere phenomena depends on the sources of inequality. If the inequality arises from the pure market forces, some people think that interference must be as little as possible and the inequality is not a serious problem. If the inequality is born beyond the individual capacity (e.g. inheritance or pure luck), governmental or nongovernmental policies are considered to be required in many respects (tax, wage control, nationalization of institutions and so on) and the inequality is an important problem we must grapple with.

What this paper concerns is the sources of inequality and especially we focus on the possibility that the inequality arises from pure luck. We provide a simple stochastic model in which the ex-post realization of the equilibrium stochastic process is quite biased among people.

To complete this purpose, we have to investigate the existence of some regularity within randomness. Intuitively, the realization of randomness from uniform distribution offers quite equal benefit among people in the long run, for example, in throwing dices or flipping coins, the same numbers realize in almost the same times as experiments continue infinitely. However, from a different mathematical viewpoint, it is possibly said that the same number arises in a regular manner so that the same numbers fall upon almost the same people. To support this aspect, we employ a monumental mathematical theorem which is recently solved. That theorem is the so called Erdös Discrepancy Problem, long time being conjecture from around 1932, which is proved by Terence Tao in [1]. This theorem roughly states that for any random sequence, the realization of which contains almost the same number periodically.

In this paper, we construct a stochastic equilibrium model in which consumers who have different discount factors buy periodically the capital stock so that they are exposed to randomness at arithmetic progression times. Therefore according to the Erdös Discrepancy Problem, there are some people who obtain high wages arbitrary larger times than low wages or who get low wages arbitrary larger times than high wages corresponding to their discount factors.<sup>1</sup>

The main feature in this paper is its approach to elucidating the inequality. The existent models (such as [1]-[9]) basically attribute the inequality to intrinsic character such as productivity, ability and income resource. Since we aim to investigate the other resources that give rise to inequality, the model developed in this paper is in a class of its own though based on standard economics notions such as utility, production and equilibrium, and we draw the distinctive conclusion that the pure randomness possibly causes inequality. The underlying mathematics is the Erdös Discrepancy Problem which is deep and new theorem in the number theory. After Tao's proof [1], some papers clarify the substance of this problem (such as Soundararajan [10]).

The next section describes the stochastic model in which people who have different discount factors select capital stock with different periodicity. The third

<sup>&</sup>lt;sup>1</sup>The claim that the possession of capital becomes biased among people according to heterogeneous discount factors is apparently related to the Ramsey's conjecture, which says that the people who have the lowest discount factors own all the capital and is solved by many authors in various settings (e.g. Becker [11], Mitra and Sorger [12]). However, in our paper, the discount factor endowed by people who have much capital depends on the realization of stochastic processes and it is not necessarily the lowest discount factor's people who have the large capital.

section explains the Erdös Discrepancy Problem and applies it to prove the realization of stochastic equilibrium. The last section offers concluding remark.

#### 2. The Model

Let  $(\Omega, \mathcal{F}, P)$  be a probability space and define a two point valued stochastic process  $a_t \in \{\overline{a}, \underline{a}\}$  for  $t \in \mathbb{N} := \{1, 2, \cdots\}$  with  $\overline{a} > \underline{a} > 0$ . The producers' behavior is described as the following maximization problem.

$$\max_{t} a_t L_t - w_t L_t$$

where  $L_t$  means the aggregate labor and  $w_t$  is the wage rate.

We normalize  $E_P[a_t] = 1$ .

Consumers buy the capital stock and directly obtain the utility from it and supply labors that yield disutility. Let  $x_t$  be the quantity of capital and denote the labor supply by  $l_t \in [0,1]$  at *t*. The quantity of initial capital  $x_1 > 0$  decays at the depreciation rate of  $0 < \delta < 1$ . So the stock remains like  $x_1, \delta x_1, \delta^2 x_1, \cdots$  as the time passes until the period written by  $t = t_1$ . Consumers buy the new capital and replace the old one at  $t = t_1 + 1$ . We assume that in the period  $t = t_1 + 1$  no capital is available because buying and replacement are assumed to take a time. Next, the new capital is installed after one period at  $t_1 + 2$ . Then by the same manner  $x_{t_n+2}$ ,  $n \in \mathbb{N}$ , decays as  $x_{t_n+2}, \delta^2 x_{t_n+2}, \cdots$  until  $t = t_{n+1}$ . Consumers buy the new capital and replace the old one at  $t_{n+1} + 2$ . So we need  $t_{n+1} \ge t_n + 2$  and the period  $[t_n + 2, t_{n+1}]$  represents the length of time during which the capital is available. Define for  $n \in \mathbb{N}$ ,

$$\hat{x}_t := \begin{cases} x_{t_{n-1}+2} & t_{n-1}+2 \le t \le t_n \\ 0 & t = t_n+1, \end{cases} \quad g(t) := \begin{cases} t - (t_{n-1}+2) & t_{n-1}+2 \le t \le t_n \\ 0 & t = t_n+1 \end{cases}$$

where  $t_0 + 2 := 1$ . Consumers' objective function can be described by

$$E_{P}\left[\sum_{t=1}^{\infty}\rho^{t-1}\left\{u\left(\delta^{g(t)}\hat{x}_{t}\right)-v\left(l_{t}\right)\right\}\right]$$

with  $t_0 + 2 = 1$  where *u* and *v* stand for the utility function and disutility one respectively. In what follows, we assume that the utility and disutility functions are linear.

#### Assumption 1.

$$u(x) = x, \quad v(l) = \eta l, \quad \eta > 0.$$

For convenience, we write down the consumers' maximization problem by setting the length of the remaining period of stock,  $t_n - (t_{n-1} + 2) =: k_n$ , namely the period between the beginning of newly installed capital,  $t_{n-1} + 2$  and the end of it,  $t_n$ . Note that the period at which no capital is not yet available is written by  $t_n + 1 = \sum_{i=1}^{n} (k_i + 2)$ . Denote the set of time at which the capital exists by

$$\mathbb{T} := \mathbb{N} \setminus \left\{ \sum_{i=1}^{n} \left( k_i + 2 \right) \mid n = 1, 2, \cdots \right\}.$$

Denote the set of nonnegative integers by  $\mathbb{Z}_+$  (namely  $\mathbb{N} \bigcup \{0\}$ ). Then the consumers' maximization problem is written as follows.

$$\begin{aligned} \max_{\{k_i\}_{i=1,2,\cdots} \subset \mathbb{Z}_+, \{(x_i, l_i)\}_i} \\ &E_P \Big[ x_1 \Big( 1 + \rho \delta + \rho^2 \delta^2 + \dots + \rho^{k_1} \delta^{k_1} \Big) - \eta \Big( l_1 + \rho l_2 + \dots + \rho^{k_1} l_{k_1+1} \Big) \\ &- \rho^{k_1+1} \eta l_{k_1+2} + \rho^{k_1+2} x_{k_1+3} \Big( 1 + \rho \delta + \rho^2 \delta^2 + \dots + \rho^{k_2} \delta^{k_2} \Big) \\ &- \rho^{k_1+2} \eta \Big( l_{k_1+3} + \rho l_{k_1+4} + \dots + \rho^{k_2} l_{k_1+3+k_2} \Big) - \rho^{k_1+2+k_2+1} \eta l_{k_1+2+k_2+2} \\ &+ \rho^{k_1+2+k_2+2} x_{k_1+2+k_2+3} \Big( 1 + \rho \delta + \rho^2 \delta^2 + \dots + \rho^{k_3} \delta^{k_3} \Big) \\ &- \rho^{k_1+2+k_2+2} \eta \Big( l_{k_1+2+k_2+3} + \dots + \rho^{k_3} l_{k_1+2+k_2+3+k_3} \Big) \\ &- \rho^{k_1+2+k_2+2+k_3+1} \eta l_{k_1+2+k_2+2+k_3+2} \\ &+ \dots + \rho^{\sum_{i=1}^{n} (k_i+2)} x_{\sum_{i=1}^{n} (k_i+2)+1} \Big( 1 + \rho \delta + \rho^2 \delta^2 + \dots + \rho^{k_{n+1}} \delta^{k_{n+1}} \Big) \\ &- \rho^{\sum_{i=1}^{n} (k_i+2)} \eta \Big( l_{\sum_{i=1}^{n} (k_i+2)+1} + \dots + \rho^{k_{n+1}} l_{\sum_{i=1}^{n} (k_i+2)+1+k_{n+1}} \Big) \\ &- \rho^{\sum_{i=1}^{n} (k_i+2)+k_{n+1}+1} \eta l_{\sum_{i=1}^{n+1} (k_i+2)} + \dots \Big] \end{aligned}$$

subject to

$$S_t \Delta \theta_t = w_t l_t, \ t \in \mathbb{T}$$
$$x_{t+1} + S_t \Delta \theta_t = w_t l_t, \ t \notin \mathbb{T}$$

where  $S_t$  is the price of stock, which is used for financing the capital or saving, and  $\Delta \theta_t$  means the increment of quantity of the stock at *t*. Notice that  $x_{t+1}$  is bought at *t*.

We assume that the price of stock has no trend.

Assumption 2.

$$E_{P}[S_{t}] = S_{1} > 0$$

for some  $S_1 > 0$ .

Thus consumers prefer buying at most capital to saving something at the periods other than  $\mathbb{T}$  due to the linearity of utility, presence of discounting  $\rho$  and no trend of stock prices. They save only when being in  $\mathbb{T}$  and buy the capital using all the savings and current wages while being in other than  $\mathbb{T}$ . Hence we can express as

$$x_{k_{1}+3} = S_{k_{1}+2} \sum_{t=1}^{k_{1}+1} \Delta \theta_{t} + w_{k_{1}+2} l_{k_{1}+2} = S_{k_{1}+2} \sum_{t=1}^{k_{1}+1} \left\{ \frac{w_{t} l_{t}}{S_{t}} \right\} + w_{k_{1}+2} l_{k_{1}+2},$$
(2)

and

$$x_{\sum_{i=1}^{n+1}(k_i+2)+1} = S_{\sum_{i=1}^{n+1}(k_i+2)} \sum_{t=\sum_{i=1}^{n}(k_i+2)+1}^{\sum_{i=1}^{n}(k_i+2)+1} \left\{\frac{w_t l_t}{S_t}\right\} + w_{\sum_{i=1}^{n+1}(k_i+2)} l_{\sum_{i=1}^{n+1}(k_i+2)}$$
(3)

for  $n = 1, 2, \dots$ . Set the price process  $S_t$  and  $w_t$  by for some  $\varepsilon > 0$ ,

$$eS_t = w_t$$
 and  $E_P[w_t] = 1.$  (4)

Since it needs to hold  $a_t = w_t$  in equilibrium due to the linearity of produc-

tion function, putting  $E_p[w_t] = 1$  is required for equilibrium. Next we impose parametric assumptions, which lead to the situation where consumers postpone working as late as possible but cannot help but work when buying capital in the time of the form  $\sum_{i=1}^{n} (k_i + 2)$ .

Let  $\overline{\rho}$  be the positive solution to  $\overline{\rho}^2 + (\eta \delta) \overline{\rho} - \eta = 0$ , namely,

$$\overline{\rho} = \left(-\eta \delta + \sqrt{\left(\eta \delta\right)^2 + 4\eta}\right) / 2$$

Assumption 3.

$$\eta \leq \frac{1}{2(2-\delta)}$$
 and  $\overline{\rho} > \rho > \frac{\eta}{1+\eta\delta}$ .

Since  $\delta < 1$ , we have  $\eta \le 1/(4-2\delta) < 1/(1+\delta)$ . Due to  $\eta^2 + (\eta\delta)\eta - \eta < 0$ , we have  $\overline{\rho} > \eta > \eta/(1+\eta\delta)$ . Thus the assumptions are consistent. From the latter part of Assumption 3 and due to  $0 > \rho^2 + (\eta\delta)\rho - \eta$ , we obtain

$$\frac{\rho}{1-\rho\delta} > \eta > \frac{\rho^2}{1-\rho\delta}.$$
(5)

Note from (1) and (2) that for  $1 \le t \le k_1 + 1$ , the marginal utility of  $l_t$  that contributes to  $x_{k_1+3}$  is  $\rho^{k_1+2}S_{k_1+2}w_t \left(1 + \rho\delta + \rho^2\delta^2 + \dots + \rho^{k_2}\delta^{k_2}\right)/S_t$ , and marginal disutility is  $\rho^{t-1}\eta$ .

We calculate as follows; for  $1 \le t \le k_1 + 1$ ,

$$E_{P}\left[\left(\rho^{k_{1}+2}S_{k_{1}+2}\frac{w_{t}}{S_{t}}\left(1+\rho\delta+\rho^{2}\delta^{2}+\cdots\rho^{k_{2}}\delta^{k_{2}}\right)-\eta\rho^{t-1}\right)l_{t}\right]$$

$$\leq E_{P}\left[\left(\rho^{k_{1}+2}S_{k_{1}+2}\frac{w_{t}}{S_{t}}\left(1+\rho\delta+\rho^{2}\delta^{2}+\cdots\rho^{k_{2}}\delta^{k_{2}}\right)-\eta\rho^{k_{1}}\right)l_{t}\right]$$

$$= E_{P}\left[\rho^{k_{1}}\left(\rho^{2}S_{k_{1}+2}\varepsilon\left(1+\rho\delta+\rho^{2}\delta^{2}+\cdots\rho^{k_{2}}\delta^{k_{2}}\right)-\eta\right)l_{t}\right]$$

$$< E_{P}\left[\rho^{k_{1}}\left(\rho^{2}S_{k_{1}+2}\varepsilon\frac{1}{1-\rho\delta}-\eta\right)l_{t}\right] = \rho^{k_{1}}\left(\rho^{2}\frac{1}{1-\rho\delta}-\eta\right)l_{t} < 0$$

The second and fourth equalities come from (4) and the last inequality is obtained by (5). Hence consumers select  $l_i = 0$  for  $1 \le t \le k_1 + 1$ . For  $\sum_{i=1}^{n} (k_i + 2) + 1 \le t \le \sum_{i=1}^{n} (k_i + 2) + k_{n+1} + 1$ ,  $n = 1, 2, 3, \cdots$ , the same arguments apply. Hence we conclude that

$$l_t = 0 \quad \text{for } t \in \mathbb{T}. \tag{6}$$

Consider  $t = k_1 + 2$ . We see from (1) and (2) that the marginal utility of  $l_{k_1+2}$  that contributes to  $x_{k_1+3}$  is  $\rho^{k_1+2} w_{k_1+2} (1 + \rho \delta + \rho^2 \delta^2 + \dots + \rho^{k_2} \delta^{k_2})$ , and marginal disutility is  $\rho^{k_1+1} \eta$ . We calculate as

$$E_{P}\left[\left(\rho^{k_{1}+2}w_{k_{1}+2}\left(1+\rho\delta+\rho^{2}\delta^{2}+\dots+\rho^{k_{2}}\delta^{k_{2}}\right)-\eta\rho^{k_{1}+1}\right)l_{k_{1}+2}\right]$$
  
=  $E_{P}\left[\rho^{k_{1}+1}\left(\rho w_{k_{1}+2}\left(1+\rho\delta+\rho^{2}\delta^{2}+\dots+\rho^{k_{2}}\delta^{k_{2}}\right)-\eta\right)l_{k_{1}+2}\right]$   
=  $\rho^{k_{1}+1}\left(\rho\left(1+\rho\delta+\rho^{2}\delta^{2}+\dots+\rho^{k_{2}}\delta^{k_{2}}\right)-\eta\right)l_{k_{1}+2}.$ 

Therefore if  $\rho(1 + \rho\delta + \dots + (\rho\delta)^{k_2}) > \eta$ , it follows that  $l_{k_1+2} = 1$ , if otherwise,

 $l_{k_1+2} = 0$  holds. But if  $l_{k_1+2} = 0$ , consumers cannot replace the capital so  $k_1 = \infty$ , whose case can be neglected because the period  $k_2$  arising from new capital is not selected at the outset. For  $k_n + 2$ ,  $n = 2, 3, \cdots$ , the same arguments apply. Thus we see that

$$l_t = 1 \text{ for } t \notin \mathbb{T}. \tag{7}$$

Hence from (2) and (3) it holds that for  $n = 1, 2, \dots, n$ 

$$x_{\sum_{i=1}^{n}(k_i+2)+1} = w_{\sum_{i=1}^{n}(k_i+2)}.$$
(8)

Thus we have  $E_p\left[x_{\sum_{i=1}^n(k_i+2)+1}\right] = 1$  for all *n*. So we can rewrite the objective function in (1) as follows.

Objective function

$$= \left(1 + \rho\delta + \rho^{2}\delta^{2} + \dots + \rho^{k_{1}}\delta^{k_{1}} - \rho^{k_{1}+1}\eta\right) \\ + \rho^{k_{1}+2}\left(1 + \rho\delta + \rho^{2}\delta^{2} + \dots + \rho^{k_{2}}\delta^{k_{2}} - \rho^{k_{2}+1}\eta\right) \\ + \rho^{k_{1}+2+k_{2}+2}\left(1 + \rho\delta + \rho^{2}\delta^{2} + \dots + \rho^{k_{3}}\delta^{k_{3}} - \rho^{k_{3}+1}\eta\right) + \dots \\ + \rho^{\sum_{l=1}^{n}(k_{l}+2)}\left(1 + \rho\delta + \rho^{2}\delta^{2} + \dots + \rho^{k_{n+1}}\delta^{k_{n+1}} - \rho^{k_{n+1}+1}\eta\right) + \dots \\ = \left(1 + \rho\delta + \rho^{2}\delta^{2} + \dots + \rho^{k_{1}}\delta^{k_{1}} - \rho^{k_{1}+1}\eta\right) \\ + \rho^{k_{1}+2}\left\{\left(1 + \rho\delta + \rho^{2}\delta^{2} + \dots + \rho^{k_{2}}\delta^{k_{2}} - \rho^{k_{2}+1}\eta\right) \\ + \rho^{k_{2}+2}\left(1 + \rho\delta + \rho^{2}\delta^{2} + \dots + \rho^{k_{3}}\delta^{k_{3}} - \rho^{k_{3}+1}\eta\right) + \dots \\ + \rho^{\sum_{l=2}^{n}(k_{l}+2)}\left(1 + \rho\delta + \rho^{2}\delta^{2} + \dots + \rho^{k_{n+1}}\delta^{k_{n+1}} - \rho^{k_{n+1}+1}\eta\right) + \dots\right\},$$

Define the value function by

$$V := \sup_{\{k_i\}_{i=1,2,\dots} \subset \mathbb{Z}_+} [\text{objective function}].$$

We can write

$$V = \max_{k_1} \left[ \left( 1 + \rho \delta + \rho^2 \delta^2 + \dots \rho^{k_1} \delta^{k_1} - \rho^{k_1 + 1} \eta \right) + \rho^{k_1 + 2} V \right]$$

From the recursive character seen above, we see that at the optimal, all  $k_i$ ,  $i = 1, 2, \cdots$  are the same. Write the optimal  $k_i$  as  $k^*$ . Then the maximal value takes the form

$$V = \left(1 + \rho\delta + (\rho\delta)^{2} + \dots + (\rho\delta)^{k^{*}} - \rho^{k^{*}+1}\eta\right) \left[1 + \rho^{k^{*}+2} + \rho^{2(k^{*}+2)} + \rho^{3(k^{*}+2)} + \dots\right]$$
$$= \frac{1 + \rho\delta + (\rho\delta)^{2} + \dots + (\rho\delta)^{k^{*}} - \rho^{k^{*}+1}\eta}{1 - \rho^{k^{*}+2}}.$$

In what follows, we aim to determine the concrete number of  $k^*$ . Let us put

$$v(k) \coloneqq \frac{1 + \rho \delta + (\rho \delta)^2 + \dots + (\rho \delta)^k - \rho^{k+1} \eta}{1 - \rho^{k+2}}$$

Then the optimization problem boils down to

 $V = \max_{k} v(k).$ 

We investigate the difference of the above v(k) with respect to k. One has v(k+1)-v(k)

$$=\frac{(1-\rho^{k+3})(\rho\delta)^{k+1}+\rho^{k+1}\eta(1-\rho)-\rho^{k+2}(1-\rho)(1+\rho\delta+(\rho\delta)^{2}+\dots+(\rho\delta)^{k+1})}{(1-\rho^{k+3})(1-\rho^{k+2})}.$$
(9)

It suffices to know the sign of the numerator in (9) to determine the sign of the fraction (9). Note that

$$sign(v(k+1)-v(k)) = sign\{(1-\rho^{k+3})(\rho\delta)^{k+1}+\rho^{k+1}\eta(1-\rho) - \rho^{k+2}(1-\rho)(1+\rho\delta+(\rho\delta)^{2}+\dots+(\rho\delta)^{k+1})\}$$
$$= sign\{(1-\rho^{k+3})\delta^{k+1}+\eta(1-\rho)-\rho(1-\rho)(1+\rho\delta+(\rho\delta)^{2}+\dots+(\rho\delta)^{k+1})\}.$$

Put

$$\Delta(k, k+1) := (1-\rho^{k+3})\delta^{k+1} + \eta(1-\rho) - \rho(1-\rho)(1+\rho\delta + (\rho\delta)^2 + \dots + (\rho\delta)^{k+1}).$$

Then we can write as

$$\operatorname{sign}(v(k+1)-v(k)) = \operatorname{sign}\Delta(k,k+1).$$

Now we further put the following assumption on the parameters. **Assumption 4**.

$$\delta < \frac{1-\rho^3}{1-\rho^4}$$
 and  $\overline{\rho} > \frac{\eta+\delta}{1-\delta}$ .

Both inequalities in Assumption 4 are satisfied when  $\delta$  is sufficiently small because if  $\delta = 0$ , all inequalities hold consistently with Assumption 3.

It follows from former part of Assumption 4 that  $\delta - 1 + \rho^3 (1 - \rho \delta) < 0$ , which leads to  $\delta - 1 + \rho^{k+3} (1 - \rho \delta) < 0$  for  $k = 0, 1, 2, \cdots$ , and further we see that

$$\begin{split} \delta^{k+1} \left\{ \delta - 1 + \rho^{k+3} \left( 1 - \rho \delta \right) \right\} \\ &= \delta^{k+1} \left\{ \left( 1 - \rho^{k+4} \right) \delta - \left( 1 - \rho^{k+3} \right) \right\} = \left( 1 - \rho^{k+4} \right) \delta^{k+2} - \left( 1 - \rho^{k+3} \right) \delta^{k+1} < 0 \end{split}$$

for  $k = 0, 1, 2, \cdots$ . Thus we find that  $(1 - \rho^{k+3})\delta^{k+1}$  is strictly decreasing function in k. So equivalently  $(1 - \rho^{k+3})\delta^{k+1} + \eta(1 - \rho)$  is strictly decreasing function in k. Together with the fact that  $\rho(1 - \rho)(1 + \rho\delta + (\rho\delta)^2 + \cdots + (\rho\delta)^{k+1})$  is strictly increasing in k, we see that

$$\Delta(k, k+1)$$
 is strictly decreasing in k. (10)

Thus we process the following arguments.

If  $v(1)-v(0) \le 0$ , it holds that v(k+1)-v(k) < 0 for  $k = 1, 2, \cdots$ . In this case, we have  $k^* = 0$  since v is decreasing entirely. It requires  $\rho > \eta$ , which induces  $l_{i(k^*+2)} = 1$ ,  $i = 1, 2, \cdots$ .

If v(1)-v(0) > 0 and  $v(2)-v(1) \le 0$ , it holds that v(k+1)-v(k) < 0 for  $k = 2, 3, \cdots$ . In this case, we have  $k^* = 1$  if  $\rho(1+\rho\delta) > \eta$  which leads to  $l_{i(k^*+2)} = 1$ ,  $i = 1, 2, \cdots$ .

If v(2)-v(1) > 0 and  $v(3)-v(2) \le 0$ , it holds that v(1)-v(0) > 0 and v(k+1)-v(k) < 0 for  $k = 3, 4, \cdots$ . In this case, we have  $k^* = 2$  if  $\rho(1+\rho\delta+(\rho\delta)^2) > \eta$  which implies  $l_{i(k^*+2)} = 1$ ,  $i = 1, 2, \cdots$ . And so forth  $\cdots$ . Thus we see that for  $k = 1, 2, \cdots$ ,

if v(k)-v(k-1) > 0 and  $v(k+1)-v(k) \le 0$ , then we have  $k = k^*$ . if  $\rho(1+\rho\delta+(\rho\delta)^2+\dots+(\rho\delta)^k) > \eta$  which leads to  $l_{k^*+2} = 1$ ,  $i = 1, 2, \dots$ .

Note that  $\Delta(k, k+1) < 0$  for large k because  $(1-\rho^{k+3})\delta^{k+1} + \eta(1-\rho)$  converges to  $\eta(1-\rho)$  and  $\rho(1-\rho)(1+\rho\delta+(\rho\delta)^2+\dots+(\rho\delta)^{k+1})$  converges to  $\rho(1-\rho)/(1-\rho\delta)$  as  $k \to \infty$ , and because  $\eta(1-\rho) < \rho(1-\rho)/(1-\rho\delta)$  by (5).

Since  $\eta \leq 1/(4-2\delta)$  in Assumption 3, we have  $(\eta\delta)^2 + 4\eta \leq 1 + 2\eta\delta + (\eta\delta)^2$ , which leads to  $-\eta\delta + \sqrt{(\eta\delta)^2 + 4\eta} \leq 1$ . Thus it holds that  $\overline{\rho} \leq \frac{1}{2}$  (recall the definition of  $\overline{\rho}$  before Assumption 3). We can confirm for  $\rho \in (0, \overline{\rho}]$  that

 $\Delta(k, k+1)$  is strictly increasing as  $\rho \to 0$ , (11)

and that for any *k*,

$$\Delta(k, k+1) > 0 \text{ holds for small } \rho.$$
(12)

Consider the case of  $\overline{\rho} > \rho \ge (\eta + \delta)/(1 - \delta)$ . Then

$$\Delta(0,1) = (1-\rho^3)\delta + \eta(1-\rho) - \rho(1-\rho)(1+\rho\delta)$$
  
=  $(1-\rho)(1+\rho+\rho^2)\delta + \eta(1-\rho) - \rho(1-\rho)(1+\rho\delta)$   
=  $(1-\rho)\{\eta+\delta-(1-\delta)\rho\} \le 0.$ 

Since  $\Delta(k, k+1) < 0$  for  $k = 1, 2, \cdots$  and  $\rho > \eta$ , we obtain  $k^* = 0$ .

Although  $\Delta(1,2) < 0$  at  $\rho = (\eta + \delta)/(1-\delta)$  (since k rises and (10)), we see from (11) and (12) that  $\Delta(1,2) = 0$  for some  $\rho_1 < (\eta + \delta)/(1-\delta)$ . Since

$$\Delta(1,2) = (1-\rho_1^4)\delta^2 + \eta(1-\rho_1) - \rho_1(1-\rho_1)(1+\rho_1\delta + (\rho_1\delta)^2)$$
$$= (1-\rho_1^2)\delta^2 + \eta(1-\rho_1) - \rho_1(1-\rho_1)(1+\rho_1\delta) = 0,$$

we obtain  $\eta - \rho_1(1+\rho_1\delta) < 0$ . Therefore for  $(\eta + \delta)/(1-\delta) > \rho \ge \rho_1$ , we see  $\Delta(0,1) > 0$  and  $\Delta(1,2) \le 0$  with  $\rho(1+\rho\delta) > \eta$ . So it follows that  $k^* = 1$  for  $(\eta + \delta)/(1-\delta) > \rho \ge \rho_1$ .

In the same way, we have  $\Delta(2,3) = 0$  for some  $\rho_2 < \rho_1$ . Then it holds that  $k^* = 2$  for  $\rho_1 > \rho \ge \rho_2$ . We have  $\Delta(3,4) = 0$  for some  $\rho_3 < \rho_2$ . Then it holds that  $k^* = 3$  for  $\rho_2 > \rho \ge \rho_3$ , and so on.

If  $\rho \leq \eta/(1+\eta\delta)$ , which is out of concern due to Assumption 3, it holds  $\rho(1+\rho\delta+\dots+(\rho\delta)^k) < \eta$  for all *k*, which means  $k^* = \infty$ , in other words, this consumer wants to hold the initial stock forever.

To summarize we conclude that  $k^* = 0$  for  $\rho \in [(\eta + \delta)/(1 - \delta), \overline{\rho}), k^* = 1$ 

for 
$$\rho \in [\rho_1, (\eta + \delta)/(1 - \delta))$$
,  $k^* = 2$  for  $\rho \in [\rho_2, \rho_1)$ , ...,  $k^* = i$  for  
 $\rho \in [\rho_i, \rho_{i-1})$ , .... Since  

$$\lim_{k \to \infty} \Delta(k, k + 1)$$

$$= \lim_{k \to \infty} [(1 - \rho^{k+3})\delta^{k+1} + \eta(1 - \rho) - \rho(1 - \rho)(1 + \rho\delta + (\rho\delta)^2 + \dots + (\rho\delta)^{k+1})]$$

$$= \eta(1 - \rho) - \rho(1 - \rho)/(1 - \rho\delta) < 0 \text{ for } \rho > \eta/(1 + \eta\delta),$$

it follows for some  $\underline{\rho}$  that  $\rho \downarrow \underline{\rho} \ge \eta/(1+\eta\delta)$  in order for  $\Delta(k,k+1) = 0$  to hold as  $k \to \infty$  (note  $\rho = \eta/(1+\eta\delta)$  is equivalent to  $\rho/(1-\rho\delta) = \eta$ ). Hence we have

$$(\underline{\rho}, \overline{\rho}) = \bigcup_{i=0}^{\infty} [\rho_i, \rho_{i-1})$$

where  $\rho_0 = (\eta + \delta)/(1 - \delta)$  and  $\rho_{-1} = \overline{\rho}$ .

Let  $\varphi_{i+2} := [\rho_i, \rho_{i-1}]$ . A consumer who has a discount factor in  $\varphi_{i+2}$  selects  $k^* = i$ ,  $i = 0, 1, 2, \cdots$ . For people who belong to  $\varphi_{i+2}$ , the supply of labor is one when  $t = (i+2)n, n = 1, 2, \cdots$ , which is the unique opportunity of receiving wages and being exposed by uncertainty, for example,  $\varphi_2$  -people who select  $k^* = 0$  supply one labor at  $t = 2, 4, \cdots, 2n, \cdots$ ,  $\varphi_3$  -people who select  $k^* = 1$  provide one labor at  $t = 3, 6, \cdots, 3n, \cdots$ ,  $\varphi_4$  -people who select  $k^* = 2$  supply one labor at  $t = 4, 8, \cdots, 4n, \cdots$ , and so on. For example, in the case of t = 12, the prime factorization is  $t = 2 \times 2 \times 3$  and people who supply one labor are represented by  $\{2, 3, 4, 6, 12\}$ , namely,  $6^{\text{th}}$  time of  $\varphi_2$  -people,  $4^{\text{th}}$  time of  $\varphi_3$ -people,  $3^{\text{rd}}$  time of  $\varphi_4$  -people,  $2^{\text{nd}}$  time of  $\varphi_6$  -people and  $1^{\text{st}}$  time of t by

$$t = \prod_{j=1}^{\omega(t)} p_j^{\alpha}$$

where  $p_j$  is a prime number and  $\alpha_j$  means its multiplicity. Notation  $\omega(t)$  obeys the convention in the number theory, which means the number of distinct primes and approximately follows normal distribution (Erdös and Kac Theorem). We write the following expansion as

$$(1 + p_1 + p_1^2 + \dots + p_1^{\alpha_1}) (1 + p_2 + p_2^2 + \dots + p_2^{\alpha_2}) \cdots (1 + p_{\omega(t)} + p_{\omega(t)}^2 + \dots + p_{\omega(t)}^{\alpha_{\omega(t)}})$$
  
=: 1 + x<sub>1</sub> + x<sub>2</sub> + \dots + x<sub>M(t)</sub>

where  $M(t) := (1 + \alpha_1)(1 + \alpha_2) \cdots (1 + \omega(t)) - 1$ . Denote the set of label of people who supply one labor at t by

$$J(t) \coloneqq \left\{ x_1, x_2, \cdots, x_{M(t)} \right\}$$

In the case of a forementioned example t = 12, one see that  $(1+2+2^2)(1+3) = 1+2+3+2^2+6+2^2 \times 3$ . So we get  $J(t) = \{2,3,4,6,12\}$  that stands for the set who supply one labor as before. Therefore we have

$$L_t = \#J(t) = M(t).$$

Since labor demand is arbitrary from linearity, the supply of labor #J(t) = M(t) is always in equilibrium of the labor market. The goods equilibrium condition is

denoted by

supply 
$$a_{t}M(t) = \text{demand } x_{t+1}M(t)$$
.

Because  $a_t = w_t = x_{t+1}$  (note  $x_{t+1}$  is available from t+1 but bought at t), its condition holds. The stock market equilibrium condition is written as

supply 
$$\frac{w_t - x_{t+1}}{S_t} M(t)$$
  
= demand  $\frac{w_t}{S_t} \cdot [$  the total number of labor other than  $J(t)].$ 

Because  $w_t = x_{t+1}$  and the labor other than J(t) equals 0 (for each consumer, labor supply is zero in  $\mathbb{T}$ ), the above condition follows.

#### 3. Realization of Stochastic Capital Process

This section concerns the realization of stochastic capital process. The aggregate capital process in equilibrium is described by

$$a_t M(t)$$

as in the previous section where  $a_t (= w_t)$  is the exogenous productivity stochastic process taking value  $\overline{a}$  or  $\underline{a}$ , and M(t) stands for a deterministic one endogenously determined in equilibrium. However, each individual consumer potentially face and really encounter at arithmetic progression times the exogenous stochastic productivity process (equivalently wages)  $\{a_1, a_2, a_3, \cdots\}$ , which is realized as  $(\overline{a}, \underline{a}, \underline{a}, \overline{a}, \underline{a}, \cdots)$ , or  $(\overline{a}, \overline{a}, \underline{a}, \overline{a}, \underline{a}, \cdots)$ , or  $(\underline{a}, \underline{a}, \overline{a}, \underline{a}, \overline{a}, \cdots)$ , and so on.

At this point, we introduce the monumental mathematical theorem, known as Erdös Discrepancy Problem, long time being conjecture from around 1932, proved by Terence Tao in (2016). It states that for any sign sequence  $f : \mathbb{N} \to \{-1, 1\}$ ,

$$\sup_{n,d\in\mathbb{N}}\left|\sum_{j=1}^{n}f\left(jd\right)\right|$$

is infinite. Formally, for any C > 0 and f, there exist n and d such that

$$\left|\sum_{j=1}^{n} f\left(jd\right)\right| \geq C.$$

Roughly speaking, given infinite sign sequence, say,  $\{-1, -1, +1, -1, -1, +1, \cdots\}$ , pick up each number skipping d-1 times (e.g. pick up  $+1, +1, \cdots$  skipping 2 times (avoiding -1, -1)), which adds up to sufficiently large for sufficiently large length of numbers. This topic generally concerns the problem as to whether there exists some regularity within random sequences. Van der Waerden's theorem (1927) asserts that for any f and  $k \in \mathbb{N}$ , there exist a and  $r \in \mathbb{N}$  such that

$$f(a) = f(a+r) = f(a+2r) = \dots = f(a+(k-1)r)$$

namely, for any sign sequence there exists any long arithmetic progression with the same number. Erdös Discrepancy Problem says the similar statements that taking a *homogeneous* arithmetic progression, the either sign outnumbers the other one by arbitrary large extent. We apply this Erdös Discrepancy Problem to the exogenous  $\{a_i\}$ , which is equal to  $\{w_i\}$  in equilibrium. In the model in the previous section, consumers are exposed to randomness periodically, namely,  $\varphi_{i+2}$ -people encounter the randomness at (i+2)n,  $n \in \mathbb{N}$  periods for  $i \in \mathbb{Z}_+$ . By redefining  $f: \mathbb{N} \to \{\overline{a}, \underline{a}\}$  and reinterpret (i+2) = d, we can say the following theorem.

**Theorem**. For arbitrary large number, C > 0, and any realization of  $\{a_i\}$ , there exists a long period of time,  $N \in \mathbb{N}$ , and  $\varphi_{i+2}$ -people who face the high wages or low wages for periods that outnumber the other ones by difference C > 0, namely,

$$\varphi_{i+2} \text{-people encounter} \\ \left\{ \# \overline{a} - \# \underline{a} \ge C ( \text{lucky people case} ) \\ \text{or} \\ \# \underline{a} - \# \overline{a} \ge C ( \text{unlucky people case} ) \end{cases}$$

for  $t = 1, 2, \dots, N$ .

Roughly speaking, even under random environment, there may be a fixed member in a society who is almost always lucky or unlucky for large period of time. Note that in the case of d = 0 that attains the given *C*, we take periods, say,  $t = 2, 4, 6, \dots, 2n, \dots$  that  $\varphi_2$ -people encounter and reinterpret it as original sequence, then we can take subsequence that attains the given *C*.

#### 4. Conclusions

This paper proposes a stochastic dynamics in which people who are endowed with different discount factors buy the capital stock periodically and are exposed to randomness at arithmetic progression times. We prove that the realization of the stochastic equilibrium may render to the people quite unequal benefits. Its proof is based on Erdös Discrepancy Problem that an arithmetic progression sum of any sign sequence goes to infinity, which is recently solved by Terence Tao (2016). There are some people who obtain high wages arbitrary larger times than low wages or who get low wages arbitrary larger times than high wages corresponding to their distinct discount factors. The result in this paper implies that in a certain society, the sources of inequality come from pure luck.

Finally we note the topics that remain in future research. Inequality arising from realization of stochastic processes only identifies the most lucky or the least one and does not explain the distribution of various income realization. In addition, whether people face the fortunate case or not reflects observation of the finite time and we cannot say anything about what occurs beyond the periods. The type of phenomena that is in this paper out of scope may be explained by other approach or more generalized mathematical theorem on the number theory or stochastic analysis.

#### **Conflicts of Interest**

The author declares no conflicts of interest regarding the publication of this paper.

#### References

- [1] Piketty, T. (2014) Capital in The Twenty-First Century. Translated by Arthur Goldhammer, The Belknap Press of Harvard University Press, Cambridge, MA.
- Tao, T. (2016) The Erdös Discrepancy Problem. *Discrete Analysis*, 1-27. https://doi.org/10.19086/da.609
- [3] Atkinson, A.B., Piketty, T. and Saez, E. (2011) Top Incomes in the Long Run of History. *Journal of Economic Literature*, 49, 3-71. <u>https://doi.org/10.1257/jel.49.1.3</u>
- [4] Gabaix, X., Lasry, J.M., Lions, P.L., Moll, B. and Qu, Z. (2016) The Dynamics of Inequality. *Econometrica*, 84, 2071-2111. <u>https://doi.org/10.3982/ECTA13569</u>
- [5] Grosman, G.M. and Helpman, E. (2018) Growth, Trade, and Inequality. *Econometrica*, 86, 37-83. <u>https://doi.org/10.3982/ECTA14518</u>
- [6] Jones, C.I. (2015) Pareto and Piketty: The Macroeconomics of Top Income and Wealth Inequality. *Journal of Economic Perspectives*, 29, 29-46. <u>https://doi.org/10.1257/jep.29.1.29</u>
- Jones, C.I. and Kim, J. (2018) A Schumpeterian Model of Top Income Inequality. Journal of Political Economy, 126, 1785-1826. <u>https://doi.org/10.1086/699190</u>
- [8] Kasa, K. and Lei, X. (2018) Risk, Uncertainty and the Dynamics of Inequality. *Journal of Monetary Economics*, 94, 60-78. https://doi.org/10.1016/j.jmoneco.2017.11.008
- [9] Mankiw, N.G. (2015) Yes, *r>g.* So What? *American Economic Review*, 105, 43-47. https://doi.org/10.1257/aer.p20151059
- Soundararajan, K. (2018) Tao's Resolution of the Erdös Discrepancy Problem. Bulletin of the American Mathematical Society, 55, 81-92. https://doi.org/10.1090/bull/1598
- Becker, R.A. (1980) On the Long-Run Steady State in a Simple Dynamic Model of Equilibrium with Heterogeneous Households. *The Quarterly Journal of Economics*, 95, 375-382. <u>https://doi.org/10.2307/1885506</u>
- [12] Mitra, T. and Sorger, G. (2013) On Ramsey's Conjecture. *Journal of Economic Theory*, **148**, 1953-1976. <u>https://doi.org/10.1016/j.jet.2013.05.003</u>



## **An Improved Randomized Circle Detection Algorithm Using in Printed Circuit Board Locating Mark**

#### Jingkun Liu<sup>1</sup>, Qi Fan<sup>2\*</sup>

<sup>1</sup>Chengyi University College, Jimei University, Xiamen, China <sup>2</sup>Xiamen Sinic-Tek Intelligent Technology Co., Ltd., Xiamen, China Email: liujingkun@126.com, \*qi.fan@sinictek.com

How to cite this paper: Liu, J.K. and Fan, Q. (2019) An Improved Randomized Circle Detection Algorithm Using in Printed Circuit Board Locating Mark. Applied Mathematics, 10, 848-861.

https://doi.org/10.4236/am.2019.1010061

Received: September 14, 2019 Accepted: October 20, 2019 Published: October 23, 2019

Copyright © 2019 by author(s) and Scientific Research Publishing Inc. This work is licensed under the Creative Commons Attribution International License (CC BY 4.0).

http://creativecommons.org/licenses/by/4.0/  $\odot$ 

#### Abstract

This paper presents an improved Randomized Circle Detection (RCD) algorithm with the characteristic of circularity to detect randomized circle in images with complex background, which is not based on the Hough Transform. The experimental results denote that this algorithm can locate the circular mark of Printed Circuit Board (PCB).

#### **Keywords**

Circle Detection, Randomized Algorithm, Characteristic of Circularity, Printed Circuit Board

#### **1. Introduction Open Access**

Detecting circles from a digital image is very important in image processing [1] and computer vision [2]. This problem has extensive application value in engineering such as product inspection and assembly, traffic monitoring, robot vision, face recognition, vectorization of hand-sketched drawing. For example, it plays an important role in the production of Printed Circuit Board (PCB) to achieve automatic positioning of PCB and to locate the reference point named as fiducial mark. Fiducial marks are used to locate the position of all features on PCB. In recent years, two main research directions are how to improve the accuracy and reduce the computation performance.

As the most well-known approach for circle detection, Hough Transform (TH) [3] has gained the widespread interest from researchers. Standard Circle Hough Transform (CHT) [4] has been shown with possessing robustness in dealing with noisy images. Set (x, y) be an edge pixel on a circle with center coordinates (a,b) and radius r, the circle can be represented by

$$(x-a)^{2} + (y-b)^{2} = r^{2}.$$
 (1)

Each edge pixel (x, y) in the image can be mapped into a conic surface in the three-dimensional (a,b,r)-parameter space. Although the technique can regularly achieve a relatively high degree of accuracy, there are some influencing factors such as huge demand of memory space for the accumulator and high computational complexity due to the time-consuming voting procedure. Several HT-based methods have been developed to overcome these problems. Randomized Hough Transform (RHT) [5] relative to CHT has faster processing speed and smaller storage requirement. RHT uses three noncollinear edges to deal with the three parameters (a,b,r) of Equation (1). This method randomly selects three edge pixels in the image with equiprobability every time, and the corresponding mapped pixel in the three-dimensional parameter space is collected. RHT performs well in high quality image, but low performance in noisy image.

In place of creating an accumulator array for mapping the extracted edge pixels in images to the circle parameters in HT-based method, Randomized Circle Detection (RCD) [6] does not use an accumulator for saving the information of related parameters in Randomized Sample Consensus (RANSAC) [7] based method. The main concept is that the algorithm randomly chooses four edge pixels from the image first, and then uses a distance criterion to determine whether they belong to a possible circle in the image. After finding a possible circle, RCD uses an evidence-collecting step to further determine whether the candidate circle is a real-circle. Since RCD does not need extra accumulator storage, the memory requirements needed in RCD are only a few variables, and the method has some other advantages such as real-time speed and more robust to noise. However, sampling for RCD randomly happens on all edge pixels of the whole image and verification of the hypothetical circles also use all the edge pixels, which both occupy a mass of time and obtain uncertainty of results [8]. To solve these problems, an improved randomized circle detection algorithm in the complex background image is proposed in the paper, which uses improved RCD algorithm and the characteristic of circularity. Firstly, the improved RCD based on connected contours is applied to detect possible circles. Then, the characteristic of circularity is used to discard some inaccurate possible circles. The algorithm is faster than RCD for it samples only on the connected curve [8]. The experimental results show that the refined algorithm has good detection performance.

In this paper, we search circles with an improved RCD algorithm in the complex back ground image, and use the characteristic of circularity to eliminate incorrect circles. Simulation results show that the proposed algorithm can locate circle mark effectively in PCB. The rest of this paper is organized as follows: Section 2 discusses the theory about normal Randomized Circle Detection (RCD). Section 3 introduces the improved Randomized Circle Detection in detail. Section 4 states the experimental results. Finally, some conclusions are given in Section 5.

#### 2. Randomized Circle Detection (RCD)

This section describes the algorithm of standard RCD. Store all edge pixels  $v_i = (x_i, y_i)$  to the set *V*. Any three noncollinear pixels  $v_i = (x_i, y_i), i = 1, 2, 3$  can exactly determine one circle  $C_{123}$ . By Equation (1), we have

$$2x_i a_{123} + 2y_i b_{123} + r_{123}^2 - a_{123}^2 - b_{123}^2 = x_i^2 + y_i^2, \ i = 1, 2, 3.$$

A representation of Equation (2) in terms of matrix form

L.

$$\begin{pmatrix} 2x_1 & 2y_1 & 1 \\ 2x_2 & 2y_2 & 1 \\ 2x_3 & 2y_3 & 1 \end{pmatrix} \begin{pmatrix} a_{123} \\ b_{123} \\ r_{123}^2 - a_{123}^2 - b_{123}^2 \end{pmatrix} = \begin{pmatrix} x_1^2 + y_1^2 \\ x_2^2 + y_2^2 \\ x_3^2 + y_3^2 \end{pmatrix}$$
(3)

Applying Gaussian elimination and Cramer's rule, the center  $(a_{123}, b_{123})$  can be obtained by

$$a_{123} = \frac{\begin{vmatrix} x_2^2 + y_2^2 - x_1^2 - y_1^2 & y_2 - y_1 \\ x_3^2 + y_3^2 - x_1^2 - y_1^2 & y_3 - y_1 \end{vmatrix}}{2(x_2 - x_1)(y_3 - y_1) - 2(x_3 - x_1)(y_2 - y_1)},$$
(4)

$$b_{123} = \frac{\begin{vmatrix} x_2 - x_1 & x_2^2 + y_2^2 - x_1^2 - y_1^2 \\ x_3 - x_1 & x_3^2 + y_3^2 - x_1^2 - y_1^2 \end{vmatrix}}{2(x_2 - x_1)(y_3 - y_1) - 2(x_3 - x_1)(y_2 - y_1)}.$$
(5)

After obtaining the center  $(a_{123}, b_{123})$ , the radius can be calculated by

$$r_{123} = \sqrt{\left(x_i - a_{123}\right)^2 + \left(y_i - b_{123}\right)^2}, \ i = 1, 2, 3.$$
(6)

Let  $v_4 = (x_4, y_4)$  be the fourth edge pixel, then the distance between  $v_4$  and the boundary of the circle  $C_{123}$  denoted by

$$d_{4\to 123} := \left| \sqrt{\left( x_4 - a_{123} \right)^2 + \left( y_4 - b_{123} \right)^2} - r_{123} \right|, \tag{7}$$

where |z| denotes the absolute value of *z*. The four edge pixels

 $v_i = (x_i, y_i), i = 1, 2, 3, 4$  can obtain four circles  $C_{123}$ ,  $C_{124}$ ,  $C_{134}$ ,  $C_{234}$  with respect to four distances  $d_{4\to 123}$ ,  $d_{3\to 124}$ ,  $d_{2\to 134}$ ,  $d_{1\to 234}$ . If  $d_{l\to ijk}$  is smaller than the given threshold  $T_d$ , these four edge pixels are co-circular and the circle  $C_{ijk}$  is a possible circle. If the four distances are all larger than  $T_d$ , RCD algorithm chooses the other four edge pixels.

Set a counter C = 0 for this possible circle in order to count how many edge pixels lie on the possible circle. For each edge pixel  $v_n$  in V, if  $d_{n \rightarrow ijk} \leq T_d$ , the counter C will be incremented by one and  $v_n$  should be taken out of V; otherwise the next edge pixel will be proceed to. Continue above process until all the edge pixels in V have been examined. In the evidence-collecting process, the final value of C is equal to  $n_p$  which is the number of edge pixels on the possible circle. If  $n_p$  is larger than the given global threshold  $T_g$ , the possible circle is claimed to be true. Otherwise, we discard the false circle and return those  $n_p$  edge pixels back into V.

In practical application, the following thresholds are used to search circle in

image:

$T_{f}$	The number of failures that can be tolerated. The running time must be finite, regardless of whether the correct results are obtained.
$T_r$	The ratio threshold of co-circle pixels. It is used to replace the global threshold
$T_{g}$	The number of pixels at the edge of circle, $T_s = 2\pi r T_r$ , <i>r</i> is the radius of possible circle.
$T_a$	The threshold of distance. The distance between any two agent pixels of the possible circle should be larger than $~T_{_a}$ .
$T_d$	The distance threshold for co-circle. It means that the fourth selected pixel is closed to the possible circle, which is determined by the other three pixels.
$T_{c}$	The number of circles that we want to detect in image.
$T_{rL}$	The minimal radius that the possible circle should fit.
$T_{rU}$	The maximal radius that the possible circle should fit.

The standard RCD algorithm has some advantages such as fewer memory requirements, faster speed and simpler algorithm than CHT or RHT.

#### 3. The Improved RCD Algorithm

This section presents the improved algorithm consisting of the following. The general idea is to reduce the computational complexity and enhance the accuracy by extracting connected contours. It can discard some inaccurate possible circles by using characteristic of circularity.

#### 3.1. Noise Suppression

In order to obtain higher recognition quality, an image preprocessing is used to suppress and eliminate noise from image. A Gaussian smoothing is a widely used to reduce image noise, which is represented by Gaussian filter [9]. Gaussian filter is a windowed filter of linear class with weighted nature, which is calculated according to Gaussian distribution. Gaussian filter in the most general function form of

$$G_{\sigma}\left(x\right) = \frac{1}{\sqrt{2\pi\sigma}} e^{-\frac{x^{2}}{2\sigma^{2}}},$$
(8)

in one-dimensional is able to fulfill this criterion optimally, where  $\sigma$  is the standard deviation of the Gaussian curve. In two-dimensions the Gaussian filter is defined by

$$G_{\sigma}(x, y) = \frac{1}{2\sigma^{2}\pi} e^{-\frac{x^{2}+y^{2}}{2\sigma^{2}}}$$
(9)

This filter can be separable by

$$G_{\sigma}(x, y) = G_{\sigma}(x)G_{\sigma}(y).$$
<sup>(10)</sup>

In order to remove small details from the image and fill small gaps in lines or

curves, the Gaussian filter should be used before other operations. The above mentioned Gaussian filter is used in continuous domain. In one discrete digital image, a discrete Gaussian low-pass filter should be used. The value of every pixel in the image will be replaced by the average of the intensity levels in the neighborhood defined by the filter mask. Because the smoothing process leads to "un-sharp" transition in intensities, smoothing filters have the undesirable side effect, which blur edges of object. In order to keep the balance between eliminating random noise and preserving sharpness of edges, the Gaussian filter with mask size  $3 \times 3$ 

$$\frac{1}{9} \times \begin{pmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{pmatrix}$$
(11)

should be used in our experiment. It can be best seen by substituting the coefficient of the mask into the characteristic response R [10] with the sum of products as

$$R = \sum_{k=1}^{9} w_k z_k = W^{\mathrm{T}} Z, \qquad (12)$$

where W and Z are nine-dimensional vectors formed from the coefficients of a  $3\times 3$  filter mark and the image intensities are encompassed by the mask respectively. Then we obtain

$$R = \frac{1}{9} \sum_{i=1}^{9} z_i.$$
 (13)

**Figure 1** shows the results of applying a Gaussian smoothing in digital image. As expected from the Gaussian filter with mask size  $3 \times 3$ , the detail of noise is blurred.

#### **3.2. Edge Detection**

The purpose of edge detection is to identify points in digital image at which the image brightness changes sharply. In the past 30 years, there are many methods for edge detection, but most of them can be classified into two categories, the Template Matching (TM) and the Differential Gradient (DG) [11]. In either



**Figure 1.** Gaussian smoothing. (a) Original image of size  $307 \times 307$  pixels; (b) Result of smoothing using Gaussian filter.

case, the aim is to find where the intensity gradient magnitude g is large enough to be taken as a reliable indicator of the edge. In the DG approach, the gradient information is estimated vectorially. There are some well-known edge operators, such as Sobel edge operator, Canny edge operator, Robbers edge operator, Prewitt edge operator and so on. In practice, Sobel gradient operator is most frequently used according to its simplicity and effectiveness. The local edge magnitude may be calculated vectorially using the nonlinear transformation [11]

$$g = \sqrt{g_x^2 + g_y^2},$$
 (14)

In order to save computational effort, the approximate formula

$$g = \left| g_x \right| + \left| g_y \right|, \tag{15}$$

could be used in practice. Figure 2 illustrates the edge segmentation, which are obtained from input original image Figure 1(a) and smoothed image Figure 1(b). It is clear that Figure 2(b) shows fewer noisy in the segmentation.

#### **3.3. Extracting Contours**

In order to reduce the computation running time and enhance the accuracy, the contours are obtained by using software OpenCV. RCD will be executed in the connected contour to find possible circles. The computational complexity is reduced from all pixels to contours.

As illustrated in Figure 3(a), an image with all contours from Figure 2(b) is generated. Figure 3(b) shows another more complex image with a different fiducial mark. All contours extracted from Figure 3(b) are shown in Figure 3(c). The circle is divided into some broken edges. Every connected contour is drawn with its own particular color in images.

#### 3.4. The Improved RCD Algorithm

As above mentioned, we describe how to determine possible circles according to



Figure 2. Edge segmentation. (a) Segmentation of Figure 1(a); (b) Segmentation of Figure 1(b).



Figure 3. Connected contours. (a) Contours from Figure 2(b); (b) Original image with size  $195 \times 195$  pixels; (c) Contours obtained from (b).

the extracted contours. Let *V* denote the set of all edge contours in the image. The procedure randomly picks four pixels from each contour to make a decision whether the contour belongs to a possible circle based on the RCD algorithm. If the four pixels meet Equation (6) and Equation (7), all pixels from all contours will be collected. We continue the above process until the contour is determined as a possible circle or the global threshold  $T_f$  has been reached. The procedure proceeds to next contour and the global Threshold  $T_f$  is reset at the beginning.

Because the proposed algorithm samples only on the connected contours, it can overcome the interference between different objects and improve the accuracy of detection. The computational complexities are reduced and the execution-time is saved.

#### 3.5. The Characteristic of Circularity

The above technique can detect all possible circles around the true fiducial mark. There exists a bias problem, because of noise and nonstandard circular circle. The circularity will be calculated to determine the best matching circle.

Usually the used form features can be grouped into two categories, the contour characteristic and the regional characteristic. The contour characteristic of image is mainly aimed at the boundary of image, and the regional characteristic of image is related to the whole shape region. The circular degree could be used to describe the circle. There are four circular degree measures [12], such as circularity, boundary energy, density and ratio of area to average distance quadratic sum. The characteristic of circularity is easy to be realized and more suitable for engineering application. It is needed to explain that the extraction of characteristic of circularity is based on image pre-processing and edge detection.

The characteristic of circularity is defined by all boundary pixels of the region D. An image may be defined as one two-dimensional function f(x, y), and the amplitude of f at any pair of coordinates (x, y) is called the gray level of the image at the pixel. Suppose that the image is of size  $N \times M$  elements. Set a pair of coordinates  $(x_k, y_k)$  of any pixel in the image. The barycentric coordinates

are defined as

$$\overline{x} := \frac{\sum_{i=1}^{N} \sum_{j=1}^{M} i \cdot f(i, j)}{\sum_{i=1}^{N} \sum_{j=1}^{M} f(i, j)}, \quad \overline{y} := \frac{\sum_{i=1}^{N} \sum_{j=1}^{M} j \cdot f(i, j)}{\sum_{i=1}^{N} \sum_{j=1}^{M} f(i, j)}.$$
(16)

The average distance from the regional barycentric to the boundary pixel is in the form of

$$\mu_D \coloneqq \frac{1}{K} \sum_{k=0}^{K-1} \left\| \left( x_k, y_k \right) - \left( \overline{x}, \overline{y} \right) \right\|, \tag{17}$$

where  $\|(x_1, y_1) - (x_2, y_2)\| \coloneqq \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$ , and the mean square deviation of distance from the regional barycentric to the boundary pixel is in the form of

$$\delta_D \coloneqq \frac{1}{K} \sum_{k=0}^{K-1} \left[ \left\| \left( x_k, y_k \right) - \left( \overline{x}, \overline{y} \right) \right\| - \mu_D \right]^2$$
(18)

can be simplified to

$$\delta_D \coloneqq \frac{1}{K} \sum_{k=0}^{K-1} \left[ \left( x_k - \overline{x} \right)^2 + \left( y_k - \overline{y} \right)^2 \right]^2 - \mu_D^2.$$
(19)

The calculation formula for the characteristic of circularity is defined by

$$C \coloneqq \frac{\delta_D}{\mu_D}.$$
 (20)

When the regional D tends to be a circle, the characteristic of circularity C is monotonically decreasing and tends to be infinitesimal. It's not affected by regional translation, rotation and variation of scales. It means that we can determine the position of the circle in the image by the minimum point of characteristics of circularity in local region.

#### 3.6. Stability, Accuracy and Time Speed

There are still some shortcomings in the standard RCD. Because the standard RCD randomly picks four edge pixels each time. The search process ended when it found the true circle or reached the number of failures  $T_f$ . It means that the searching result may be quite different according to the different four starting pixels. If the staring pixels are close to the target circle, the circle may be detected rapidly. Sometimes, the start four pixels are too far from the target circle and the input image contains too many useful and useless pixels. The RCD algorithm may end without finding the true circle. In order to find possible circle rapidly, these thresholds  $T_d$ ,  $T_r$  and  $T_a$  should not be too small. Then many incorrect circles would be found.

By suppressing noise and finding of the connected contours, which contains at least a certain number of points, it will reduces the calculation time and find fewer possible circles. At last, after verifying the circularity, the real circle will be selected. So the improved RCD provides the better stability and accuracy than standard RCD. As the image grows larger and more complex, the time will become more longer. But the improved RCD will significantly reduce the time. The following experiment will be used to verify these conclusions.

#### 4. Experimental Results

In this section, some experimental results are demonstrated to show the execution-time and accuracy advantages of our improved RCD algorithm, compared with standard RCD algorithm. Four original images are used to evaluate the performance of the proposed algorithm. The first three images with the same mark are shown in **Figures 4(a)-(c)**, which have different size  $131 \times 124$ ,  $262 \times 153$  and  $307 \times 307$  pixels. The last image **Figure 4(d)** has another different mark. The radius of the first mark *R* is 26.5 pixels and the radius of the second mark *R* is 36.5 pixels. **Table 1** describes all parameters used in our experiments. Some





**Figure 4.** Four original images with different size and mark. (a)-(c) Original images with the same mark; (d) Original image with another different mark.

Table 1. All thresholds used in experiments.

	$T_{_f}$	$T_r$	$T_a$	$T_{d}$	$T_{c}$	$T_{rL}$	$T_{_{rU}}$	$T_{cc}$
Pure RCD	100,000	0.6	0.25R	1.5	20	<i>R</i> – 5	R + 5	1
Improved RCD	Square of the length of contour, which is different in searching circle in each chain	0.6	0.25R	1.5	20	<i>R</i> – 5	<i>R</i> + 5	1

additional variables are defined as following: P denotes the center pixel (x, y) of detected circle. R is the pre-input radius and r is the radius of detected circle. PN is the number of pixels on the contour of detected circle. C is the circularity of detected circle. Avg T (ms) represents the average running time in millise-cond.  $T_{cc}$  is the threshold of circularity.

#### 4.1. Locating Mark with Standard RCD Algorithm

In this part the marks will be detected by using standard RCD algorithm. As shown in **Figure 5** there are many correct and incorrect circles in images. In **Figure 5(c)** the true mark circle is not included in the detected possible circles. There is no way to discard the incorrect circles and find the true mark. According to the threshold  $T_c$ , there are twenty possible detected circles at most, which are drawn in **Figure 5**. In **Tables 2-5** only ten possible circles are described.



(c)

(d)

**Figure 5.** Locating results using standard RCD algorithm. (a)-(c) The detected results, obtained from **Figures 4(a)-(c)**; (d) The result obtained from **Figure 4(d)**.

Table 2.	Ten results shown	in	Figure	5(	(a)	).
----------	-------------------	----	--------	----	-----	----

Р	r	PN	Р	r	PN
(71.039, 81.479)	25.911	404	(70.902, 80.576)	26.225	369
(90.787, 30.734)	28.790	133	(75.500, 83.700)	28.452	138
(70.930, 79.679)	26.190	302	(45.153, 11.418)	27.853	58
(72.391, 81.334)	25.750	383	(71.071, 80.071)	25.014	320.000
(71.982, 81.718)	25.382	379	(62.561, 16.912)	23.763	66

Р	r	PN	Р	r	PN
(189.352, 66.432)	26.473	175	(157.550, 47.532)	26.822	170
(194.103, 72.685)	24.486	116	(195.472, 52.751)	24.591	166
(157.036, 36.878)	24.956	151	(223.636, 50.818)	26.623	106
(201.467, 86.570)	23.407	89	(185.626, 41.030)	23.896	143
(169.820, 65.472)	27.184	158	(206.284, 80.524)	27.742	67
(71.982, 81.718)	25.382	379	(62.561, 16.912)	23.763	66

Table 3. Ten results shown in Figure 5(b).

#### Table 4. Ten results shown in Figure 5(c).

Р	r	PN	Р	r	PN
(34.786, 163.67)	27.257	150	(252.645, 48.271)	28.241	122
(177.100, 43.500)	28.500	148	(103.150, 122.58)	23.413	74
(69.098, 179.443)	25.349	103	(205.563, 55.500)	26.300	137
(145.883, 32.260)	27.740	151	(167.003, 42.453)	26.043	204
(216.647, 44.956)	25.610	97	(148.500, 13.571)	27.472	119
(71.982, 81.718)	25.382	379	(62.561, 16.912)	23.763	66

#### Table 5. Ten results shown in Figure 5(d).

Р	r	PN	Р	r	PN
(116.345, 107.928)	33.594	286	(115.135, 107.68)	34.296	318
(116.459, 106.745)	34.481	348	(137.832, 112.68)	33.834	102
(116.271, 107.794)	34.310	337	(116.256, 106.40)	33.951	309
(132.030, 109.909)	35.091	83	(107.256, 129.41)	34.519	101
(116.398, 101.856)	35.538	108	(116.355, 106.31)	34.913	325

#### 4.2. Locating Mark with Improved RCD Algorithm

In the following experiments, the fiducial marks will be detected by our improved RCD algorithm and corresponding circularity. As illustrated in Figure 6, the left four images (a)-(d) denote the results, which are obtained without calculating minimal circularity. Because of unsharpness of edge, the correct and bias circles are detected, which have similar radius and position. It is difficult to determine the best true circle without additional decision condition. After calculating the minimal circularity, the true marks are located in the right four images (e)-(h). The detected possible circles are shown in Tables 6-9. Compared to the former standard RCD algorithm, the number of possible circles is fewer and the results are close to the true mark. By calculating minimal circularity, in Table 6 the true circle with center pixel (71.500, 80.500) and radius 26.163 pixels is detected. In Table 7 the true circle with center pixel (71.369, 81.622) and radius 27.132 pixels is detected. In Table 8 the true circle with center pixel (75.601, 81.013) and



**Figure 6.** Results obtained by using our improved RCD algorithm and circularity. (a)-(d) Results obtained without calculating minimal circularity; (e)-(h) Results obtained by calculating the minimal circularity.

#### Table 6. Results shown in Figure 6(a).

Р	r	РВ	С	
(72.313, 81.737)	26.795	309	0.8077	
(71.500, 80.500)	26.163	392	0.1288	Selected As true

#### Table 7. Results shown in Figure 6(b).

	С	PN	r	Р
	0.8077	356	26.291	(71.404, 82.005)
Selected As true	0.1288	253	27.132	(71.369, 81.622)
	0.2732	272	27.016	(71.194, 81.403)
	0.1548	290	26.903	(71.570, 80.779)

#### Table 8. Results shown in Figure 6(c).

Р	r	PN	С	
(74.985, 81.435)	27.104	257	0.6047	
(75.601, 81.013)	26.951	253	0.0339	Selected As true

_					
	Р	r	PN	С	
	(114.992, 108.471)	34.629	356	0.0595	Selected As true
	(14.192, 107.192)	34.529	376	0.0985	
	(116.826, 106.174)	34.699	461	0.2000	
	(116.500, 108.500)	34.821	383	0.1547	
-					

Table 9. Results shown in Figure 6(d).

**Table 10.** Time performance comparison between standard RCD and the improved RCD in terms of milliseconds.

Original Images	Avg T (RCD)	Avg T (improved RCD)
Figure 4(a)	77.8 ms	121 ms
Figure 4(b)	466.8 ms	478.2 ms
Figure 4(c)	1172.8 ms	710 ms
Figure 4(d)	148.2 ms	145 ms

radius 26.951 pixels is detected. In **Table 9** the true circle with center pixel (114.992, 108.471) and radius 34.629 pixels is detected.

#### 4.3. Comparing the Execution-Time

In order to compare the execution-time, all concerned experiments are performed on the Intel i5-5200U CPU with 2.20 GHz and 8 GB RAM. The adopted operating system is MS-Windows 7 and the programming environment is VS2010. In order to accurately evaluate the execution-time, we run each image 100 times and calculate the average time in **Table 10**. For locating mark in small images, the execution-time is almost the same for both algorithms. The time will be significantly reduced, when the mark is detected in a larger and complex image.

#### **5.** Conclusions

This paper has presented the proposed improved RCD strategy to improve the performance of execution-time and the accuracy of detection. The refined method can suppress the interference between different objects significantly. After calculating minimal circularity, the true mark will be located efficiently and accurately. The experimental results demonstrate that the proposed improved RCD improves significantly the accuracy of detection. Experimental results also demonstrate that the proposed algorithm provides a considerable execution-time improvement. The algorithm has significant execution-time superiority in large and complex image.

At the current stage, the quality of the image, such as the stability of the light source, the degree of image damage, clarity, etc., has a great influence on the calculation results. In the future, some normalization methods will be studied to reduce these interferences, such as combining with pattern recognition.

#### Founding

This article is supported by Science and Technology Project of Fujian Provincial Department of Education under contract JAT170917 and Youth Science and Research Foundation of Chengyi College Jimei University under contract C16005.

#### **Conflicts of Interest**

The authors declare no conflicts of interest regarding the publication of this paper.

#### **References**

- [1] Gonzalez, R.C. and Woods, R.E. (1992) Digital Image Processing. Addison Wesley, New York.
- [2] Forsyth, D.A. (2002) Computer Vision: A Modern Approach. Prentice-Hall, New Jersey.
- [3] Leavers, V.F. (1993) Survey: Which Hough Transform. CVGIP: Image Understanding, 58, 250-264. <u>https://doi.org/10.1006/ciun.1993.1041</u>
- [4] Duda, R.O. and Hart, P.E. (1972) Use of the Hough Transformation to Detect Lines and Curves in Pictures. *Communications of the ACM*, 15, 11-15. https://doi.org/10.1145/361237.361242
- [5] Xu, L. and Oja, E. (1993) Randomized Hough Transform (RHT): Basic Mechanisms, Algorithms, and Computational Complexities. *CVGIP*. *Image Understanding*, 57, 131-154. <u>https://doi.org/10.1006/ciun.1993.1009</u>
- [6] Chen, T.C. and Chung, K.L. (2001) An Efficient Randomized Algorithm for Detecting Circles. *Computer Vision and Image Understanding*, 83, 172-191. <u>https://doi.org/10.1006/cviu.2001.0923</u>
- [7] Fischer, M.A. and Bolles, R.C. (1981) Random Sample Consensus: A Paradigm to Model Fitting with Applications to Image Analysis and Automated Cartography. *Communications of the ACM*, 24, 381-395. <u>https://doi.org/10.1145/358669.358692</u>
- [8] Jia, L.Q., Peng, C.Z., Liu, H.M. and Wang, Z.H. (2011) A Fast Randomized Circle Detection Algorithm. *International Congress on Image and Signal Processing*, 4, 835-838. <u>https://doi.org/10.1109/CISP.2011.6100372</u>
- [9] Gomes, J. and Velho, L. (1997) Image Processing for Computer Graphics. Springer-Verlag, New York. <u>https://doi.org/10.1007/978-1-4757-2745-6</u>
- [10] Gonzalez, R.C. and Woods, R.E. (2008) Digital Image Processing. Pearson Prentice Hall, New Jersey.
- [11] Davies, E.R. (2005) Machine Vision: Theory Algorithms Practicalities. Morgan Kaufmann, San Francisco.
- [12] Fan, Y. (2007) Digital Image Processing and Analysis. Beijing University of Aeronautics and Astronautics, Beijing.

## Call for Papers

![](_page_61_Picture_1.jpeg)

## **Applied Mathematics (AM)**

### ISSN Print: 2152-7385 ISSN Online: 2152-7393 https://www.scirp.org/journal/am

**Applied Mathematics (AM)** is an international journal dedicated to the latest advancement of applied mathematics. The goal of this journal is to provide a platform for scientists and academicians all over the world to promote, share, and discuss various new issues and developments in different areas of applied mathematics.

#### Subject Coverage

All manuscripts must be prepared in English, and are subject to a rigorous and fair peer-review process. Accepted papers will immediately appear online followed by printed hard copy. The journal publishes original papers including but not limited to the following fields:

- Applied Probability
- Applied Statistics
- Approximation Theory
- Chaos Theory
- Combinatorics
- Complexity Theory
- Computability Theory
- Computational Methods in Mechanics and Physics
- Continuum Mechanics
- Control Theory
- Cryptography
- Discrete Geometry
- Dynamical Systems
- Elastodynamics

- Evolutionary Computation
- Financial Mathematics
- Fuzzy Logic
- Game Theory
- Graph Theory
- Information Theory
- Inverse Problems
- Linear Programming
- Mathematical Biology
- Mathematical Dielogy
   Mathematical Chemistry
- Mathematical Economics
- Mathematical Physics
- Mathematical Daysh
- Mathematical PsychologyMathematical Sociology

- Matrix Computations
- Neural Networks
- Nonlinear Processes in Physics
- Numerical Analysis
- Operations Research
- Optimal Control
- Optimization
- Ordinary Differential Equations
- Partial Differential Equations
- Probability Theory
- Statistical Finance
- Stochastic Processes
- Theoretical Statistics

We are also interested in: 1) Short Reports—2-5 page papers where an author can either present an idea with theoretical background but has not yet completed the research needed for a complete paper or preliminary data; 2) Book Reviews—Comments and critiques.

#### Notes for Intending Authors

Submitted papers should not have been previously published nor be currently under consideration for publication elsewhere. Paper submission will be handled electronically through the website. All papers are refereed through a peer review process. For more details about the submissions, please access the website.

#### Website and E-mail

https://www.scirp.org/journal/am

E-mail: am@scirp.org

#### What is SCIRP?

Scientific Research Publishing (SCIRP) is one of the largest Open Access journal publishers. It is currently publishing more than 200 open access, online, peer-reviewed journals covering a wide range of academic disciplines. SCIRP serves the worldwide academic communities and contributes to the progress and application of science with its publication.

#### What is Open Access?

Art and Design Review

Advances in

idvances in Biological

Entomolog

Applied Mathematics

Engineering

entill a

All original research papers published by SCIRP are made freely and permanently accessible online immediately upon publication. To be able to provide open access journals, SCIRP defrays operation costs from authors and subscription charges only for its printed version. Open access publishing allows an immediate, worldwide, barrier-free, open access to the full text of research papers, which is in the best interests of the scientific community.

• High visibility for maximum global exposure with open access publishing model

Soft

- Rigorous peer review of research papers
- Prompt faster publication with less cost
- Guaranteed targeted, multidisciplinary audience

![](_page_62_Picture_8.jpeg)

Website: https://www.scirp.org Subscription: sub@scirp.org Advertisement: service@scirp.org