# Introducing the Official Thailand R Local Community: R x TH

Nathakhun Wiroonsri

Organizer

# 🔴 **Materials**

Slides

# Outline

- Our mission and aims

- Our team

- Past contributions and awards

- How to join us

- Q&A

- Alumni Talk

- Workshop

# 🔴 Our Mission

R x TH มีพันธกิจหลักในการรวบรวมผู้ใช้ และ ขยายฐานของผู้ใช้งาน R ในประเทศไทย และสร้างความรู้ความเข้าใจถึงประสิทธิภาพและประโยชน์ของภาษา R ทั้งในภาคการศึกษา และภาคธุรกิจ
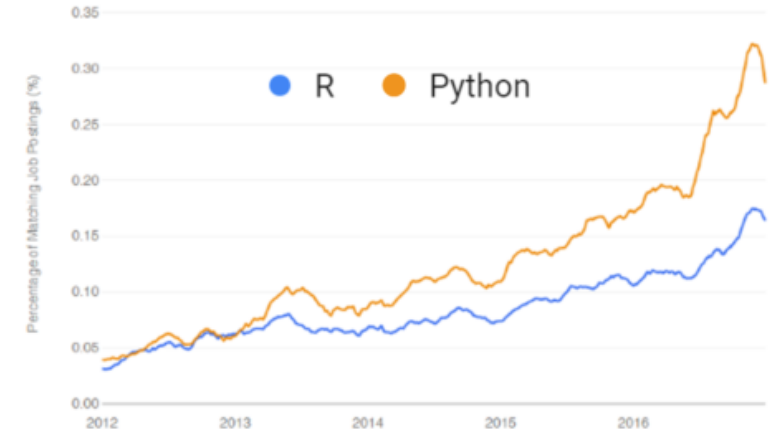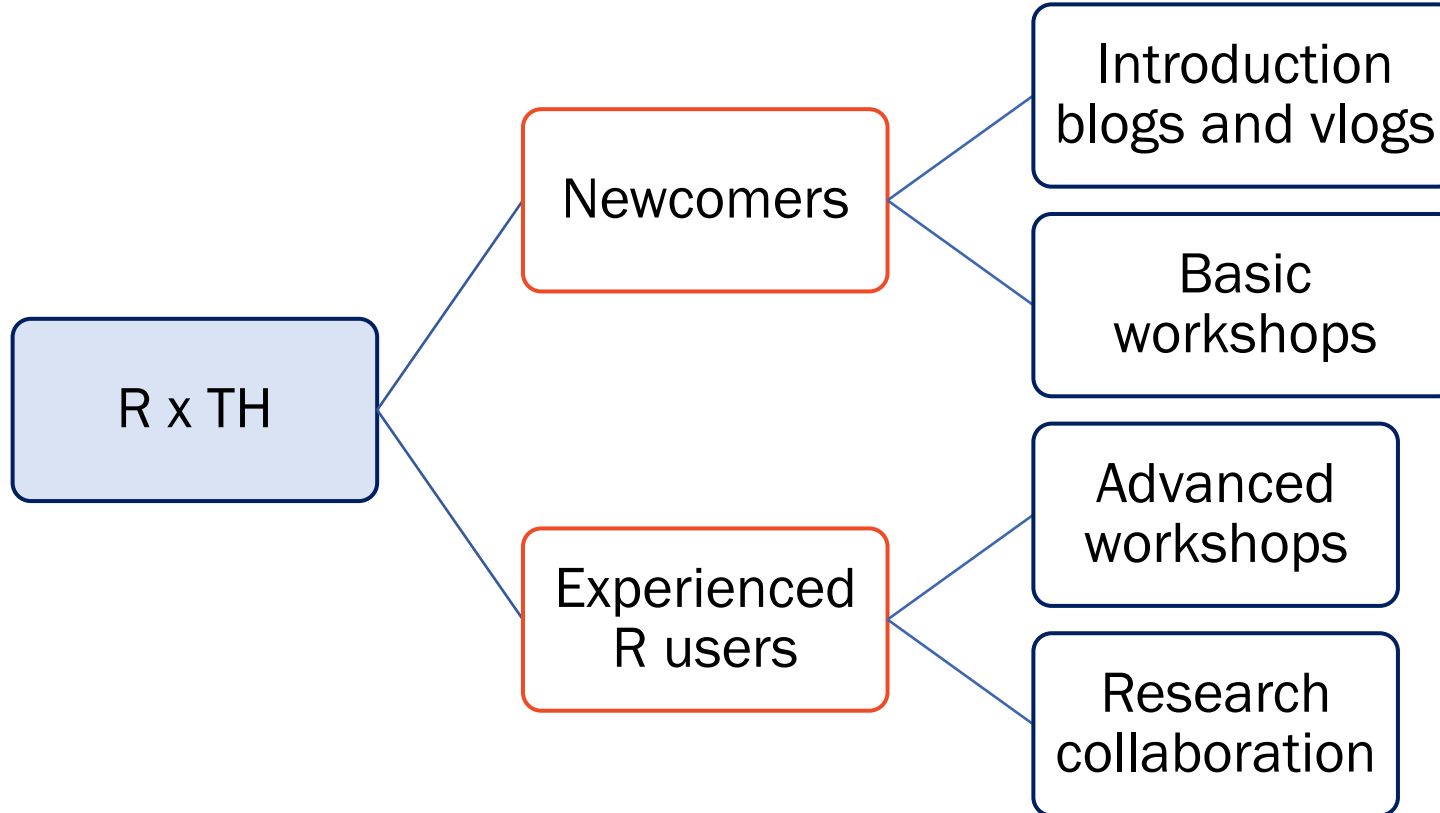


**R User Groups**

| Members | Groups | Countries |
|---------|--------|-----------|
| 75,570  | 91     | 39        |

# 🔴 **Our Aims**

- เพื่อรวบรวมผู้ใช้งาน R ให้รับรู้ถึงการใช้งานที่กว้างขวางในประเทศไทย
- เพื่อพัฒนาความรู้ให้กับผู้ใช้งาน R
- เพื่อแนะนำ R และขยายฐานผู้งานใหม่
- เพื่อร่วมกันพัฒนางานวิจัย และ ซอฟต์แวร์แพ็คเกจ

# 🔴 Our Plan

```
                    ┌──────────────────┐
              ┌─────│   Introduction   │
              │     │ blogs and vlogs  │
   ┌──────────┐     └──────────────────┘
   │ Newcomers│
   └──────────┘     ┌──────────────────┐
              └─────│      Basic       │
┌────────┐          │    workshops     │
│ R x TH │          └──────────────────┘
└────────┘
              ┌──────────────────┐
        ┌─────│     Advanced     │
┌──────────────┐   │    workshops     │
│ Experienced  │   └──────────────────┘
│   R users    │
└──────────────┘   ┌──────────────────┐
        └─────│     Research     │
              │  collaboration   │
              └──────────────────┘
```

# 🔴 Our Team

**Nathakhun Wiroonsri (Nat)**
Assistant Professor
Organizer

**Wasamon Jantai (Mam)**
Lecturer
Chula Co-organizer

**Onthada Preedasawakul (WW)**
Master student
Co-organizer

**Noppanon Teangthae (Pure)**
Master student
Co-organizer

🔴 **Our Team**

**Raywat Tanadkithirun**
Member, Assistant Professor

**Porntip Dechpichai**
Member, Lecturer

**Thanet Chitsuphaphan**
Member, Lecturer

**Chutiphan Charoensuk**
Member, Master student

**Hafsah Tabassum**
Member, Post-doc

**Natapon Aeimsri**
Member, Master student

# 🔴 **Sponsors**

# 🔴 Past contributions: Packages



**UniversalCVI: Hard and Soft Cluster Validity Indices**

Algorithms for checking the accuracy of a clustering result with known classes, computing cluster validity indices, and generating plots for comparing them. The package is compatible with K-means, fuzzy C means, EM clustering, and hierarchical clustering (single, average, and complete linkage). The details of the indices in this package can be found in: C. Alok. (2010) <https://hdl.handle.net/10603/93443>, J. C. Bezdek, M. Moshtaghi, T. Runkler, C. Leckie (2016) <doi:10.1109/TFUZZ.2016.2540063>, T. Calinski, J. Harabasz (1974) <doi:10.1080/03610927408827101>, C. H. Chou, M. C. Su, E. Lai (2004) <doi:10.1007/s10044-004-0218-1>, D. L. Davies, D. W. Bouldin (1979) <doi:10.1109/TPAMI.1979.4766909>, J. C. Dunn (1973) <doi:10.1080/01969727308546046>, F. Haouas, Z. Ben Dhiaf, A. Hammouda, B. Solaiman (2017) <doi:10.1109/FUZZ-IEEE.2017.8015651>, M. Kim, R. S. Ramakrishna (2005) <doi:10.1016/j.patrec.2005.04.007>, S. H. Kwon (1998) <doi:10.1049/EL:19981523>, S. H. Kwon, J. Kim, S. H. Son (2021) <doi:10.1049/ell2.12249>, G. W. Miligan (1980) <doi:10.1007/BF02293907>, M. K. Pakhira, S. Bandyopadhyay, U. Maulik (2004) <doi:10.1016/j.patcog.2003.06.005>, M. Popescu, J. C. Bezdek, T. C. Havens, J. M. Keller (2013) <doi:10.1109/TSMCB.2012.2205679>, S. Saitta, B. Raphael, I. Smith (2007) <doi:10.1007/978-3-540-73499-4_14>, A. Starczewski (2017) <doi:10.1007/s10044-015-0525-8>, Y. Tang, F. Sun, Z. Sun (2005) <doi:10.1109/ACC.2005.1470111>, N. Wiroonsri (2024) <doi:10.1016/j.patcog.2023.109910>, N. Wiroonsri, O. Preedasawakul (2023) <doi:10.48550/arXiv.2308.14785>, C. H. Wu, C. S. Ouyang, L. W. Chen, L. W. Lu (2015) <doi:10.1109/TFUZZ.2014.2322495> and X. Xie, G. Beni (1991) <doi:10.1109/34.85677>.

| | |
|---|---|
| Version: | 1.1.2 |
| Depends: | R (≥ 2.10) |
| Imports: | e1071, mclust |
| Published: | 2024-03-31 |
| DOI: | 10.32614/CRAN.package.UniversalCVI |
| Author: | Nathakhun Wiroonsri [cre, aut], Onthada Preedasawakul [aut] |
| Maintainer: | Nathakhun Wiroonsri <nathakhun.wir at kmutt.ac.th> |
| License: | GPL (> 3) |

downloads 3134

**BayesCVI: Bayesian Cluster Validity Index**

Algorithms for computing and generating plots with and without error bars for Bayesian cluster validity index (BCVI) (O. Preedasawakul, and N. Wiroonsri, A Bayesian Cluster Validity Index, Computational Statistics & Data Analysis, 202, 108053, 2025. <doi:10.1016/j.csda.2024.108053>) based on several underlying cluster validity indexes (CVIs) including Calinski-Harabasz, Chou-Su-Lai, Davies-Bouldin, Dunn, Pakhira-Bandyopadhyay-Maulik, Point biserial correlation, the score function, Starczewski, and Wiroonsri indices for hard clustering, and Correlation Cluster Validity, the generalized C, HF, KWON, KWON2, Modified Pakhira-Bandyopadhyay-Maulik, Pakhira-Bandyopadhyay-Maulik, Tang, Wiroonsri-Preedasawakul, Wu-Li, and Xie-Beni indices for soft clustering. The package is compatible with K-means, fuzzy C means, EM clustering, and hierarchical clustering (single, average, and complete linkage). Though BCVI is compatible with any underlying existing CVIs, we recommend users to use either WI or WP as the underlying CVI.

| | |
|---|---|
| Version: | 1.0.1 |
| Depends: | R (≥ 2.10) |
| Imports: | e1071, mclust, ggplot2, UniversalCVI |
| Published: | 2024-09-04 |
| DOI: | 10.32614/CRAN.package.BayesCVI |
| Author: | Nathakhun Wiroonsri [aut], Onthada Preedasawakul [cre, aut] |
| Maintainer: | Onthada Preedasawakul <o.preedasawakul at gmail.com> |
| License: | GPL (> 3) |

downloads 1841

# Past contributions: Talks

# 🔴 Past contributions: Conferences



- ICSDS 2023 by IMS, Lisbon, Portugal
  - useR! 2024, Salzburg, Austria
  - ICSDS 2024 by IMS, Nice, France

# Past contributions: research papers

# 🔴 **Past contributions: research papers**

# Past contributions: Awards

# 🔴 Please join us

https://www.meetup.com/r-x-th/



ส่วนหนึ่งของ **R User Groups - 91 กลุ่ม** ⓘ

## R x TH

📍 Bangkok, ประเทศไทย

👥 สมาชิก 69 คน · กลุ่มสาธารณะ ⓘ

👤 จัดโดย **Nathakhun Wiroonsri** and **2 others**

แชร์: 🅽 📘 🐦 💼 ✉️

เกี่ยวกับ   กิจกรรม   สมาชิก   รูปภาพ   การสนทนา

**เข้าร่วมกลุ่มนี้**  **...**

### สิ่งที่เรากำลังจะ

Welcome to the R x TH (R local user group in Thailand).

Our aim is to gather and expand R users in both academic and industry around Thailand. We also intend to organize meetings, workshops, and talks to encourage local R users to collaborate and develop some special projects and research together.

อ่านเพิ่มเติม

### Organizers

Nathakhun Wiroonsri and **2 others**
✉️ ข้อความ

### Members (69)

ดูทั้งหมด

# Evaluation

# THANK YOU.

# ANY QUESTIONS?