

データ解析機能

(GakuNin Federated Computing Services)

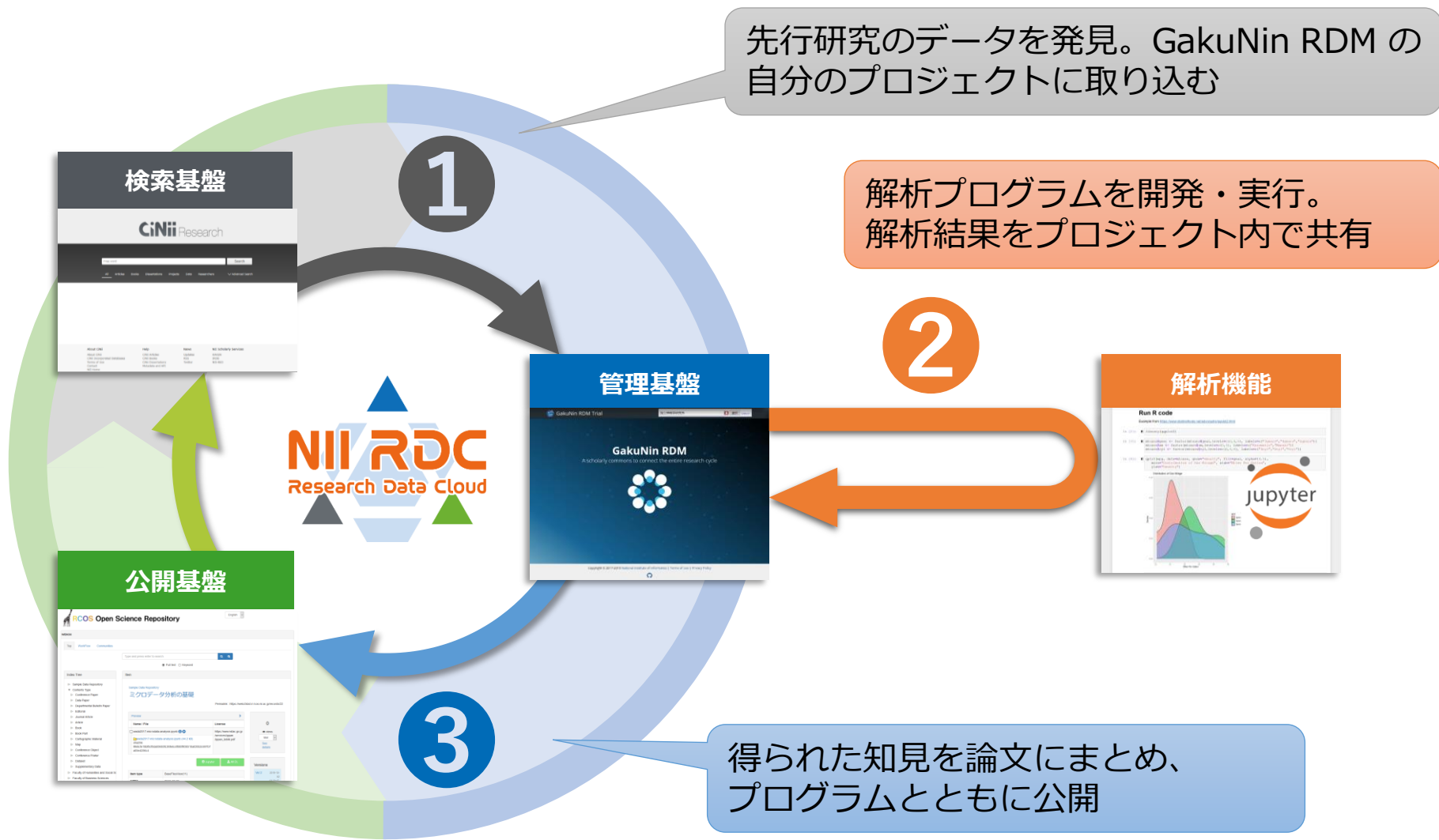
藤原 一毅

国立情報学研究所 オープンサイエンス基盤研究センター

2022-06-01

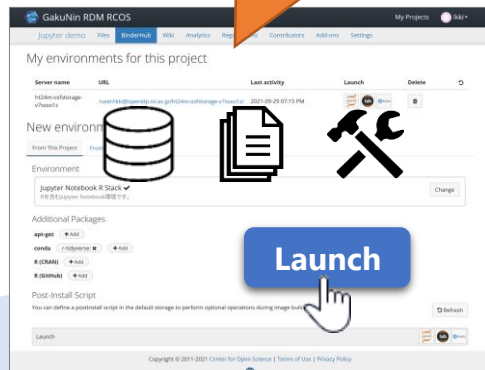
NII学術情報基盤オープンフォーラム

データ（とプログラム）は天下の回りもの



GakuNin RDM データ解析機能 (GakuNin Federated Computing Services)

①環境定義・共有



GakuNin RDM

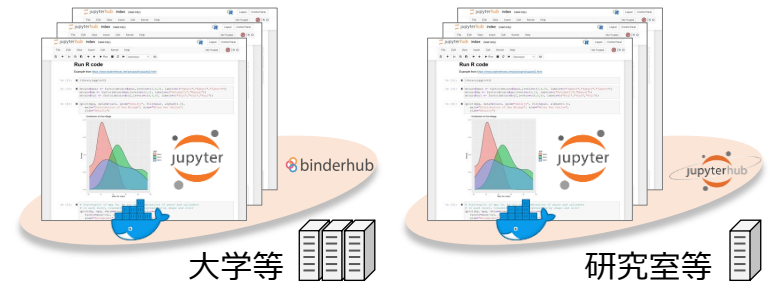
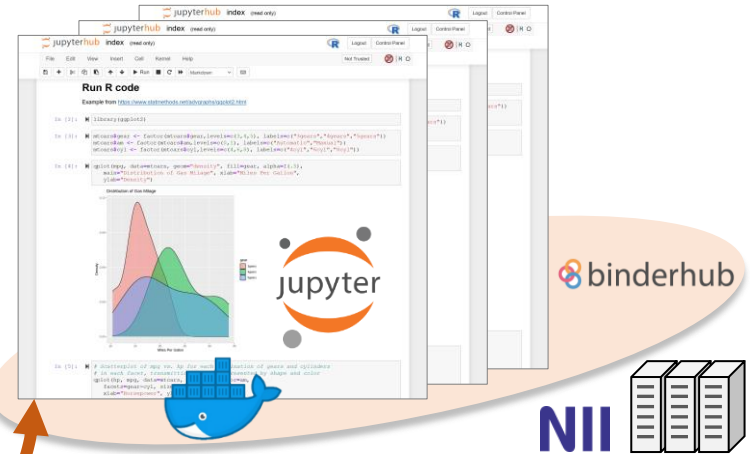
標準ストレージ

機関ストレージ

②取り込み

③書き戻し

④読み書き



デモ (1/4)

ストレージプロバイダーをクリックするか、ドラッグ&ドロップしてファイルをアップロードします

名前	サイズ	バージョン	ダウンロ...	最終更新日時
Reanalysis Example				
- NII Storage				
analyses.ipynb	3.2 kB	1	0	2022-05-27 10:34 AM
supplementary materials.xlsx	150.7 kB	1	0	2022-05-27 10:34 AM

解析プログラム (Jupyter Notebook)

データファイル

デモ (2/4)



新しい解析環境

① 解析環境を構成

基本イメージ

Python 3.9 + R 4.1.3 ✓
Jupyter Notebook, JupyterLab, RStudio, Shinyが使えます。

変更

追加パッケージ

apt-get fonts-noto-cjk: ✕ + 追加

conda seaborn: ✕ openpyxl: ✕ + 追加

pip + 追加

R (MRAN) + 追加

自動実行スクリプト

```
#!/bin/bash
set -x
...
```

保存

② 計算機を選択して起動!

環境作成

このプロジェクトのデフォルトストレージの内容がコピーされます。

新しい解析環境を作成: <https://binder.cs.rcos.nii.ac.jp>

デモ (3/4)

③ ファイルがGakuNin RDMからコピーされている

② 書き戻しボタン

④ ファイルを読み込んで解析

⑤ 解析結果を ~/result/ に保存

```
[1]: df = pd.read_excel("supplementary_materials.xlsx", index_col=0)
df = df.rename(columns={'day(1=3/31, 2=4/30, 3=5/31, 4=6/10)': 'day'})
df['day'] = df['day'].map({'1': '3/31', 2: '4/30', 3: '5/31', 4: '6/10'})
df

[3]: top10 = pd.pivot_table(df, index='country').nlargest(10, 'Infections')
df1 = pd.pivot_table(df[df['country'].isin(top10)], index='country', columns='day')
ax = df1['Infections'].plot(xlabel='調査日', ylabel='百万人あたり感染者数')
# ax.legend(loc='center left', bbox_to_anchor=(1, 0, 0.5))
plt.savefig("result/graph1.png")
```

調査日	Belgium	Chile	Ireland	Kuwait	Luxembourg	Peru	Qatar	Singapore	Spain	United States of America
3/31	~1000	~500	~500	~500	~500	~500	~500	~500	~500	~500
4/30	~1500	~1000	~1000	~1000	~1000	~1000	~1000	~1000	~1000	~1000
5/31	~2000	~1500	~1500	~1500	~1500	~1500	~1500	~1500	~1500	~1500
6/10	~2500	~2000	~2000	~2000	~2000	~2000	~2000	~2000	~2000	~2000

```
[4]: df2 = pd.read_excel("supplementary_materials.xlsx", sheet_name=4,
index_col=0, header=5, skipfooter=3,
usecols=[1,2,3,5,6,8,9,11,12], skiprows=[6])
df2 = df2.dropna()
df2 = df2.set_axis(['CF1', 'SE1', 'CF2', 'SE2', 'CF3', 'SE3', 'CF4', 'SE4'], axis=1)
df2['3/31'] = df2['CF1'] / df2['SE1']
df2['4/30'] = df2['CF2'] / df2['SE2']
df2['5/31'] = df2['CF3'] / df2['SE3']
df2['6/10'] = df2['CF4'] / df2['SE4']
```

デモ (4/4)

The screenshot displays the GakuNin RDM RCOS web interface. The browser address bar shows the URL <https://rcos.rdm.nii.ac.jp/yfbz/>. The page title is "Reanalysis of COVID-19 Infect...". The navigation menu includes "ファイル", "Wiki", "解析", "メンバー", "アドオン", "設定", and "証跡管理". The user profile "Ikki" is visible in the top right.

The main content area shows a file upload interface for "graph1.png (バージョン: 1)". Action buttons include "チェックアウト", "タイムスタンプを打つ", "削除", "ダウンロード", "プレビュー", and "バージョン管理".

On the left, a file explorer sidebar shows the file structure. The file "graph1.png" is highlighted in blue and enclosed in a red box. An orange callout bubble points to this file with the text: "⑦解析結果がGakuNin RDMに書き戻される".

The main content area features a line graph titled "country" showing the number of cases per 100,000 people (百万人あたり感染数) over time (調査日). The x-axis ranges from 3/31 to 6/10, and the y-axis ranges from 0 to 30,000. The legend includes: Belgium, Chile, Ireland, Kuwait, Luxembourg, Peru, Qatar, Singapore, Spain, and United States of America. The Qatar line shows a sharp increase starting around 5/31, reaching approximately 30,000 cases per 100,000 people by 6/10.

調査日	Belgium	Chile	Ireland	Kuwait	Luxembourg	Peru	Qatar	Singapore	Spain	United States of America
3/31	~1000	~1000	~1000	~1000	~1000	~1000	~1000	~1000	~1000	~1000
4/30	~2000	~2000	~2000	~2000	~2000	~2000	~2000	~2000	~2000	~2000
5/31	~3000	~3000	~3000	~3000	~3000	~3000	~10000	~3000	~3000	~3000
6/10	~4000	~4000	~4000	~4000	~4000	~4000	~30000	~4000	~4000	~4000

こんな用途に使えます

研究

- ご自身の研究のためのデータ分析



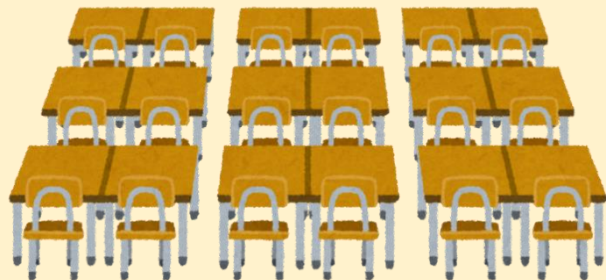
公開・共有

- 他の研究者の二次分析に資するデータとプログラムの公開



教育・学習

- 学生たちにデータ分析をさせるゼミ・講義・演習など



引き継ぎ

- 先輩の研究環境を後輩が再現し、研究を継続する



ご期待ください

データ解析機能

(開発完了)

- GakuNin RDM ユーザーに機関ごとに提供されます。

外部計算機連携機能

(テスト中)

- 機関の計算機上に GakuNin RDM と連携する解析環境を構築できます。

計算再現パッケージ

機能 (設計中)

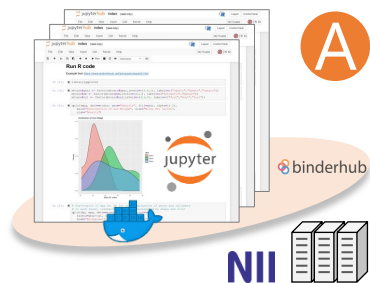
- 論文・データ・プログラムを含む一連の研究成果を公開・発見・再利用しやすくなります。

秘密計算機能

(構想中)

- 暗号化と秘密分散により、機密性の高いデータ分析にも利用できるようになります。

外部計算機連携のバリエーション



	システム	運用主体	認証方法	同時起動数	ドメイン名+サーバ証明書	バックエンド
A	BinderHub	NII	学認	10個/ユーザ	○	Kubernetes
B	BinderHub	機関	学認, OAuth, LDAP, etc.	任意設定	必要	Kubernetes
C	JupyterHub	研究室等	OAuth, LDAP, ローカル	1個/ユーザ	不要	Linux VM

- バッチスケジューラやワークフローエンジンを介して利用するスパコン等との連携も模索中。

詳しい情報

導入手続き

- データ解析機能は、GakuNin RDM のオプション機能として、機関単位で提供されます。
- 利用機関の情報基盤センター等で初期設定を行うと、その機関に所属するユーザーが利用可能となります。
- 現在 GakuNin RDM を正式利用されている機関の担当者様に、解析機能の追加についてご案内する予定です。

マニュアル

- <https://support.rdm.nii.ac.jp/>

お問い合わせ

- 開発・活用に関すること → cs-support@nii.ac.jp
- 導入手続きに関すること → rdm_support@nii.ac.jp