# 用OpenRefine清理資料

**TaiBIF 內容經理 劉璟儀**

# 劉璟儀

**TaiBIF 內容經理/ GBIF 亞洲區諮詢顧問**
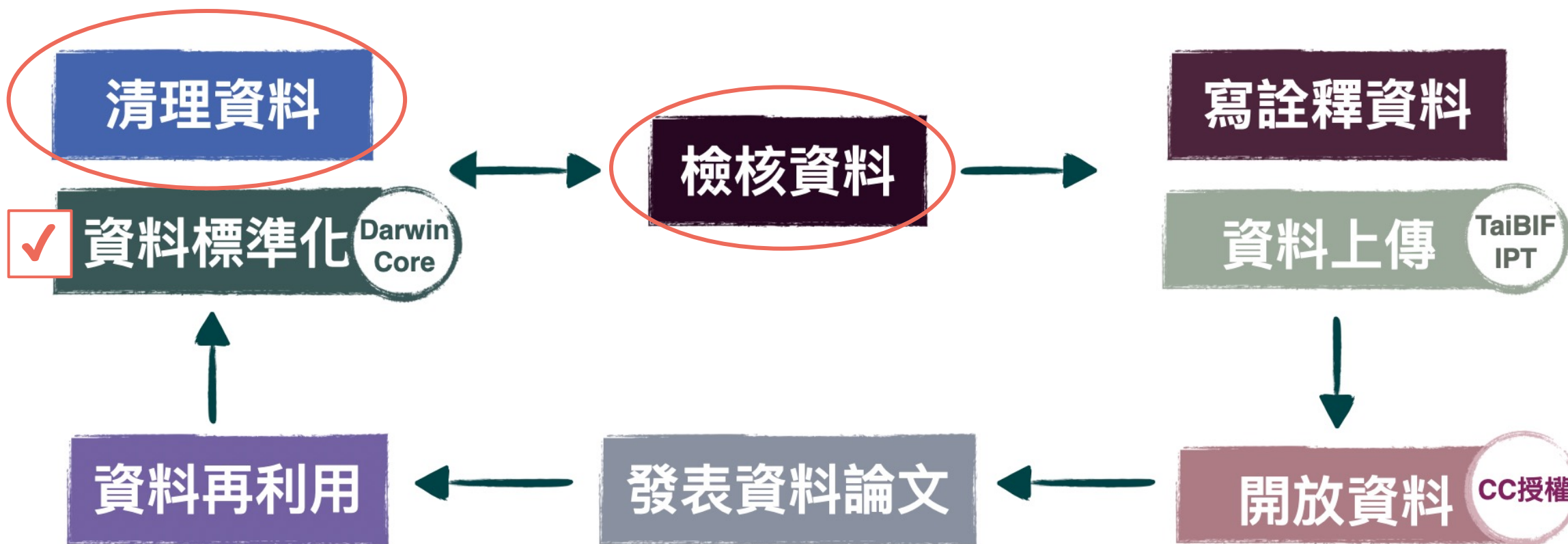
melissaliu0520@gmail.com

大學唸公共衛生，碩士轉領域到海洋科學
主要研究鯨豚的重金屬和碳氮同位素

- 推動及宣傳國內生物多樣性資料標準、資料開放
- 國內生物多樣性資料庫相關單位合作
- TaiBIF 資料庫管理諮詢

- 亞洲計畫團隊開放資料諮詢
- 推動亞洲國家節點合作
- 促進亞洲區域開放資料

# 上傳資料前...
# 你應該準備好這些事



清理資料

資料標準化 Darwin Core ✓

檢核資料

寫詮釋資料

資料上傳 TaiBIF IPT

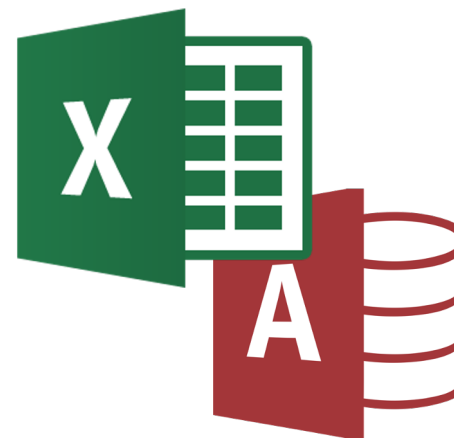開放資料 CC授權

發表資料論文

資料再利用

# 資料清理小工具

## OpenRefine

### 不是資料庫
(無法儲存資料)

### 與 Excel 的使用方式不同
（只能清理資料）

# 用 OpenRefine 清資料

**OpenRefine**
*A power tool for working with messy data.*

Create Project
Open Project
Import Project
Language Settings

**Create a project by importing data. What kinds of data files can I import?**

TSV, CSV, *SV, Excel (.xls and .xlsx), JSON, XML, RDF as XML, and Google Data documents are all supported. Support for other formats can be added with OpenRefine extensions.

Get data from

**This Computer**

Web Addresses (URLs)

Clipboard

Database

Google Data

Locate one or more files on your computer to upload:

選擇檔案　未選擇任何檔案

**Next »**

Version 3.4-beta2 [c67e13b]

使用介面簡單
可一次修正整批資料錯誤/格式
隨時自動暫存且離線操作
可返回任何一步操作
匯入CSV / Excel 不易出現亂碼

# 資料清理小工具
## OpenRefine

**1** Excel
產生資料/ 管理資料

**2** GBIF data validator
**GBIF**
檢核資料

**3** OpenRefine
清理資料

# 清理資料流程

**1** 先產生並
彙整資料

**2** 驗證資料
GBIF Data Validator

**3** 查看資料問題
Validation Issues

**4** 清理資料
OpenRefine

**5** 上傳資料
TaiBIF IPT

**6** 再次確認
資料問題
GBIF dataset 的 Issues & flags

1. 下載 **Data-cleaning-open-refine v20220927**
2. 使用 **GBIF Data Validator** 找出資料錯誤
3. 試著用 **OpenRefine** 找出個別錯誤並修正
4. 利用 **NomenMatch** 比對有效學名/ **Canadensys Coordinate conversion** 轉換座標格式
5. 進階題：使用 **GBIF backbone API** 新增分類階層欄位

# 用**OpenRefine**清理資料

**等等會需要用到的連結**

**練習檔案下載**
**GBIF 資料驗證工具**
**NomenMatch 學名比對**
**Canadensys 座標/日期轉換工具**

# 用 OpenRefine 清資料

下載並安裝在電腦 https://openrefine.org/download.html

# 用 OpenRefine 清資料

下載到電腦打開exe檔　　https://openrefine.org/download.html



在 Windows 開啟openRefine時，會出
現dos視窗，使用時都不要關掉喔！

# 用 **OpenRefine** 清資料

## 要看到這個畫面出現在瀏覽器上才是對的



Windows 下載後打開資料夾直接點這個 ⟶

# 檢核資料—先找出可能的資料錯誤

- **GBIF data validator** https://www.gbif.org/tools/data-validator

# 資料問題

- **找出重複 ID** occurrenceID
- **新增欄位** basisOfRecord
- **內容錯誤或與欄位不符**
  decimalLatitude, decimalLongitude,
  countryCode, country, day, year
- **學名比對&清理** scientificName
- **修正學名格式** ^[A-Z].*\s[A-Z]
- **清除多餘空格** country
- **找出相似文字並合併** County

| | | | |
|---|---|---|---|
| http://rs.tdwg.org/dwc/terms/occurrenceID 🔗 | 100 | ⭕ 100% | 98 |

**Validation Issues**

**GBIF Occurrence Interpretation**

| Basis of record invalid | 98 |
|---|---|
| Continent derived from coordinates | 98 |
| Occurrence status inferred from individual count | 98 |
| Country coordinate mismatch | 13 |
| Presumed negated longitude | 5 |
| Country invalid | 1 |
| Recorded date invalid | 1 |
| Recorded date unlikely | 1 |
| Taxon match fuzzy | 1 |
| Coordinate rounded | 86 |

**Resource Structure**

validation.issueType.OCCURRENCE_NOT_UNIQUELY_IDENTIFIED

# 用 OpenRefine 清資料

**OpenRefine** *A power tool for working with messy data.*

Create Project
Open Project
Import Project
Language Settings

**Create a project by importing data. What kinds of data files can I import?**

TSV, CSV, *SV, Excel (.xls and .xlsx), JSON, XML, RDF as XML, and Google Data documents are all supported. Support for other formats can be added with OpenRefine extensions.

Get data from

**This Computer**

Web Addresses (URLs)

Clipboard

Database

Google Data

Locate one or more files on your computer to upload:

選擇檔案 未選擇任何檔案

Next »

**選擇檔案並按 Next**

Version 3.4-beta2 [c67e13b]

14

# 用 OpenRefine 清資料



選擇檔案後
a. 確認下方文字編碼為 UTF-8
b. 檢視表頭和欄位有沒有抓錯

按下 Create Project 進入使用介面

# 用 **OpenRefine** 清資料



專案列
檔案匯出/ 編輯連結

資料預覽區
資料呈現的地方

資料控制區
顯示選擇的資料
過濾器/查看編輯
歷程

16

# 案例練習 – 進階作業

## 資料問題

- **找出重複 ID** **occurrenceID**
- 新增欄位 basisOfRecord
- 內容錯誤或與欄位不符 decimalLatitude, decimalLongitude, countryCode, country, day, year
- 學名比對&清理 scientificName
- 修正學名格式 ^[A-Z].*\s[A-Z]
- 清除多餘空格 country
- 找出相似文字並合併 County

http://rs.tdwg.org/dwc/terms/occurrenceID 🔗  |  100  |  ○ 100%  |  98

## Validation Issues

### GBIF Occurrence Interpretation

- Basis of record invalid    98
- Continent derived from coordinates    98
- Occurrence status inferred from individual count    98
- Country coordinate mismatch    13
- Presumed negated longitude    5
- Country invalid    1
- Recorded date invalid    1
- Recorded date unlikely    1
- Taxon match fuzzy    1
- Coordinate rounded    86

### Resource Structure

validation.issueType.OCCURRENCE_NOT_UNIQUELY_IDENTIFIED

# 用 OpenRefine 清資料- 找出重複 ID



**Customized facets**
a. 在 **occurrenceID** 那欄
   點選三角形小圖示
b. 選擇 **Facet >>**
   **Customized facet >>**
   **Duplicates facet**

# 案例練習- 進階作業

# 資料問題

- 找出重複 ID **occurrenceID**
- **新增欄位** basisOfRecord
- 內容錯誤或與欄位不符
  decimalLatitude, decimalLongitude,
  countryCode, country, day, year
- 學名比對&清理 scientificName
- 修正學名格式 ^[A-Z].*\s[A-Z]
- 清除多餘空格 country
- 找出相似文字並合併 County

**Validation Issues**

**GBIF Occurrence Interpretation**

Basis of record invalid   98

Continent derived from coordinates   98

Occurrence status inferred from individual count   98

Country coordinate mismatch   13

Presumed negated longitude   5

Country invalid   1

Recorded date invalid   1

Recorded date unlikely   1

Taxon match fuzzy   1

Coordinate rounded   86

**Resource Structure**

validation.issueType.OCCURRENCE_NOT_UNIQUELY_IDENTIFIED

# 用 OpenRefine 清資料- 新增欄位



**Add Column**
a. 在 **occurrenceID** 那欄點選三角形小圖示
b. 選擇 **Edit Column >> Add column based on this column**

# 用 OpenRefine 清資料- 新增欄位



設定內容值
a. 填入新欄位名稱
   **basisOfRecord**
b. 把值都填入
   **"PreservedSpecimen"**

# 資料問題

- 找出重複 ID  **occurrenceID**
- 新增欄位 **basisOfRecord**
- **內容錯誤或與欄位不符**
  **decimalLatitude, decimalLongitude, countryCode, country, day, year**
- 學名比對&清理 **scientificName**
- 修正學名格式 **^[A-Z].*\s[A-Z]**
- 清除多餘空格 **country**
- 找出相似文字並合併 **County**

**Validation Issues**

**GBIF Occurrence Interpretation**

Basis of record invalid    98

Continent derived from coordinates    98

Occurrence status inferred from individual count    98

Country coordinate mismatch    13

Presumed negated longitude    5

Country invalid    1

Recorded date invalid    1

Recorded date unlikely    1

Taxon match fuzzy    1

Coordinate rounded    86

**Resource Structure**

validation.issueType.OCCURRENCE_NOT_UNIQUELY_IDENTIFIED

# 案例練習– 進階作業

# 資料問題

座標和國家不符

Country coordinate mismatch  13  ×

| recordId | dwc:decimalLatitude | dwc:decimalLongitude | dwc:geodeticDatum | dwc:country | dwc:cour |
|---|---|---|---|---|---|
| c4da5630-3da3-11ed-b878-0242ac120002 | 17.563668 | 0.10294211 | WGS84 | Guatemala | GT |
| c4da05a4-3da3-11ed-b878-0242ac120002 | 5° 35' 12" N | 75° 46' 18" W | WGS84 | Guatemala | GT |
| c4da20d4-3da3-11ed-b878-0242ac120002 | 7° 18' 10.12" N | 75° 04' 25.03" W | WGS84 | Guatemala | GT |

# 用 OpenRefine 清資料- 內容錯誤(座標)



Facet / Filter　Undo / Redo 1 / 1

Refresh　Reset All　Remove All

**2**　✕ ─ **decimalLongitude**　invert reset

^[0-9]

☐ case sensitive　☑ regular expression

此部分無法批次複製修改，僅能個別修正

20 matching rows (100 total)

Show as: **rows** records　Show: 5 10 25 50 **100** 500 1000 rows　Sort ▾　« first ‹ previo

| ovince | ▼ decimalLatitude | ▼ decimalLongitude | ▼ county | ▼ locality | ▼ verbatimElevation | ▼ geodeticDatum |
|---|---|---|---|---|---|---|
| | 17.56 | Facet　▶ | San Andres | Río Inírida. caño Nabuquen | 220 | WGS84 |
| | 17.7783778 | **1** Text filter | San | Hacienda | 292 | WGS84 |
| | 17.7783778 | Edit cells　▶ | | | | |
| | | Edit column | | | | |
| | 17.7783778 | Transpose | | | | |
| | | Sort... | | | | |
| 5° 35' 12" N | | View | | | | |
| 5° 35' 12" N | 75° 46' 18" W | Reconcile | | | | |
| 5° 35' 12" N | 75° 46' 18" W | | Flores | Los Hornitos | 150 | WGS84 |
| 7° 18' 10.12" N | 75° 04' 25.03" W | | Flores | Lote en las afueras del pueblo | 150 | WGS84 |
| 6° 4' 20.210" N | 75° 38' 20.440" W | | La Libertad | Barrio Oneti | 514 | WGS84 |

## Text Filter
a. 利用正規表示式 **^[0-9]** 篩選出第一個字是數字的資料
b. 找出非十進位座標並修正成十進位

24

# 用 OpenRefine 清資料- 內容錯誤(座標)



**Canadensys Coordinate conversion**

利用座標轉換工具，將度分秒的座標格式換成十進位

**1** 貼上座標並按Submit

# 案例練習– 進階作業

## 資料問題

推定經度應為負值

Presumed negated longitude | 5 | ✕

?

| recordId | dwc:decimalLatitude | dwc:decimalLongitude |
|---|---|---|
| c4da1594-3da3-11ed-b878-0242ac120002 | 17.7783778 | 90.84424953 |
| c4da4f50-3da3-11ed-b878-0242ac120002 | 17.7783778 | 90.84424953 |
| c4da21a6-3da3-11ed-b878-0242ac120002 | 17.7783778 | 90.84424953 |
| c4da5b26-3da3-11ed-b878-0242ac120002 | 17.2160555 | 89.50767314 |
| c4da499c-3da3-11ed-b878-0242ac120002 | 17.4114231 | 90.18308898 |

# 用 OpenRefine 清資料- 內容錯誤(座標)



**Text Filter**
a. 利用正規表示式 **^[0-9]** 篩選出第一個字是數字的資料
b. 再從此篩選結果點選 Text Facet，找出那幾筆錯誤的十進位座標並修正成負值

# 案例練習– 進階作業

# 資料問題

國家代碼無效

Country invalid  1  ✕

| recordId | dwc:country | dwc:countryCode |
| --- | --- | --- |
| c4da28ea-3da3-11ed-b878-0242ac120002 | Guatemala | 17.3857972 |

# 用 OpenRefine 清資料- 內容錯誤(countryCode)



**Text Facet**
將錯誤的值修改成GT

# 案例練習– 進階作業

## 資料問題

- 找出重複 ID  occurrenceID
- 新增欄位 basisOfRecord
- 內容錯誤或與欄位不符
  decimalLatitude, decimalLongitude,
  countryCode, country, day, year
- **學名比對&清理** scientificName
- 修正學名格式 ^[A-Z].*\s[A-Z]
- 清除多餘空格 country
- 找出相似文字並合併 County

**Validation Issues**

**GBIF Occurrence Interpretation**

- Basis of record invalid    98
- Continent derived from coordinates    98
- Occurrence status inferred from individual count    98
- Country coordinate mismatch    13
- Presumed negated longitude    5
- Country invalid    1
- Recorded date invalid    1
- Recorded date unlikely    1
- Taxon match fuzzy    1
- Coordinate rounded    86

**Resource Structure**

- validation.issueType.OCCURRENCE_NOT_UNIQUELY_IDENTIFIED

# 用 OpenRefine 清資料- 學名比對

Taxon match fuzzy　1　✕

分類未對應 GBIF backbone

| recordId | dwc:genus | dwc:class | dwc:phylum | dwc:scientificNameAuthorship | |
|---|---|---|---|---|---|
| c4da38bc-3da3-11ed-b878-0242ac120002 | Paepalanthus | Equisetopsida | Magnoliophyta | (Körn.) Tissot-Squalli | |

# 用 OpenRefine 清資料- 學名比對



**NomenMatch**
將有問題的學名貼上按 Check names

結果會顯示與有效學名差異之處，以及比對吻合度的分數

# 用 OpenRefine 清資料- 學名比對



**Global Names Resolver**
如果NomenMatch找不到，也可以用這個比對看看

# 用 OpenRefine 清資料- 學名清理



清除多餘空格
將連續空格清除成一個

# 案例練習- 進階作業

# 資料問題

- 找出重複 ID **occurrenceID**
- 新增欄位 **basisOfRecord**
- 內容錯誤或與欄位不符
  **decimalLatitude, decimalLongitude, countryCode, country, day, year**
- 學名比對**&**清理 **scientificName**
- **修正學名格式 ^[A-Z].*\s[A-Z]**
- 清除多餘空格 **country**
- 找出相似文字並合併 **County**

**Validation Issues**

**GBIF Occurrence Interpretation**

| Basis of record invalid | 98 |
| Continent derived from coordinates | 98 |
| Occurrence status inferred from individual count | 98 |
| Country coordinate mismatch | 13 |
| Presumed negated longitude | 5 |
| Country invalid | 1 |
| Recorded date invalid | 1 |
| Recorded date unlikely | 1 |
| Taxon match fuzzy | 1 |
| Coordinate rounded | 86 |

**Resource Structure**

validation.issueType.OCCURRENCE_NOT_UNIQUELY_IDENTIFIED

# 用 OpenRefine 清資料- 修正學名格式



記得下面兩個選項要打勾

## Text Filter

a. 利用正規表示式 **^[A-Z].*\s[A-Z]** 篩選出第一個字開頭是大寫字母，同時第二個字開頭也是大寫字母的資料

# 用 OpenRefine 清資料- 修正學名格式



**Text Facet**
修正學名格式，第二個字開頭應為小寫字母

可以批次修改

# 用 **OpenRefine** 清資料**-** 修正學名格式



記得下面兩個選項要打勾

**Text Filter**
**1.** 利用正規表示式 **^[a-z].*\s[a-z]** 篩選出第一個字開頭是小寫字母，同時第二個字開頭也是小寫字母的資料
**2.** 將第一個字開頭修正為大寫

# 案例練習- 進階作業

# 資料問題

- 找出重複 ID  occurrenceID
- 新增欄位 basisOfRecord
- 內容錯誤或與欄位不符
  decimalLatitude, decimalLongitude,
  countryCode, country, day, year
- 學名比對&清理 scientificName
- 修正學名格式 ^[A-Z].*\s[A-Z]
- **清除多餘空格** country
- 找出相似文字並合併 County

**Validation Issues**

**GBIF Occurrence Interpretation**

| | |
|---|---|
| Basis of record invalid | 98 |
| Continent derived from coordinates | 98 |
| Occurrence status inferred from individual count | 98 |
| Country coordinate mismatch | 13 |
| Presumed negated longitude | 5 |
| Country invalid | 1 |
| Recorded date invalid | 1 |
| Recorded date unlikely | 1 |
| Taxon match fuzzy | 1 |
| Coordinate rounded | 86 |

**Resource Structure**

validation.issueType.OCCURRENCE_NOT_UNIQUELY_IDENTIFIED

# 用 OpenRefine 清資料- 清除多餘空格2



**清除多餘空格**

a. 選擇 Country 那欄
b. 點選 Edit cells >>
Common transforms >>
Trim leading and trailing
whitespace
c. 將文字前後的多餘空格去除

# 案例練習 – 進階作業

## 資料問題

- 找出重複 ID **occurrenceID**
- 新增欄位 **basisOfRecord**
- 內容錯誤或與欄位不符
  **decimalLatitude, decimalLongitude, countryCode, country, day, year**
- 學名比對**&**清理 **scientificName**
- 修正學名格式 **^[A-Z].*\s[A-Z]**
- 清除多餘空格 **country**
- **找出相似文字並合併** county

**Validation Issues**

**GBIF Occurrence Interpretation**

| Basis of record invalid | 98 |
| Continent derived from coordinates | 98 |
| Occurrence status inferred from individual count | 98 |
| Country coordinate mismatch | 13 |
| Presumed negated longitude | 5 |
| Country invalid | 1 |
| Recorded date invalid | 1 |
| Recorded date unlikely | 1 |
| Taxon match fuzzy | 1 |
| Coordinate rounded | 86 |

**Resource Structure**

validation.issueType.OCCURRENCE_NOT_UNIQUELY_IDENTIFIED

# 用 OpenRefine 清資料- 統一資料格式



## Cluster 比對相似資料及合併

a. 選擇 Text Facet
b. 點選 Cluster
c. 結果找出可能是一樣但格式不一致的值
d. 勾選要合併的值，按 Merge Selected & Re-cluster

# 進階題-自動匯入高階層分類欄位



**連接 GBIF backbone API**
a. 選擇 **scientificName**
b. 點選 **Edit column >> Add column by fetching URLs**

# 進階題-自動匯入高階層分類欄位

**Add column by fetching URLs based on column scientificName**



**2** New column name `Api_name`    Throttle delay `250` millise **3**

On error    ● set to blank  ○ store error        ☑ Cache responses

HTTP headers to be used when fetching URLs: Show

**Formulate the URLs to fetch:**

Expression            Language [ General Refine Expression Language (GREL) ∨ ]

**4**
```
"http://api.gbif.org/v1/species/match?
verbose=true&name="+escape(value,'url')
```
        No syntax error.

| Preview | History | Starred | Help |

| row | value | "http://api.gbif.org/v1/specie ... |
|-----|-------|-------------------------------------|
| 1. | Vriesea drewii | http://api.gbif.org/v1/species/match? |

## 貼上語法串接API
**a.** 將新欄位名稱設定為 **Api_name**
**b. Throttle delay** 設定為 **250**
**c.** 在 **Expression** 貼上語法

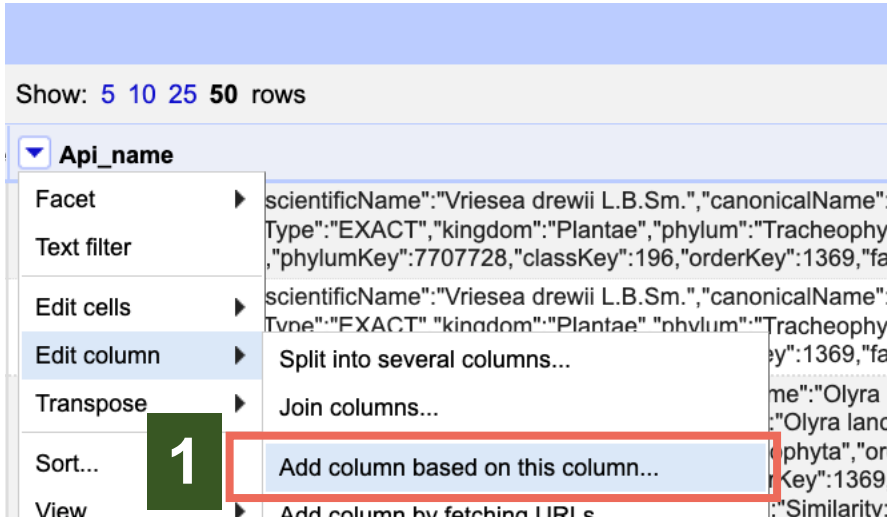**語法在下一頁，請整串複製貼上**

44

# 進階題-自動匯入高階層分類欄位

## 語法在此，請整串複製貼上

⬇

"http://api.gbif.org/v1/species/match?verbose=true&name="+escape(value,'url')

# 進階題-自動匯入高階層分類欄位



呼叫各分類階層的值

a. 到 Api_name 欄位並選擇 Edit column >> Add column based on this column
b. 將新欄位名稱寫為 higherClassification
c. 貼上語法按 OK

語法在下一頁，請整串複製貼上

# 進階題-自動匯入高階層分類欄位

## 語法在此，請整串複製貼上

⬇

value.parseJson().get("kingdom")+", "+value.parseJson().get("phylum")+", "+value.parseJson().get("class")+", "+value.parseJson().get("order")+", "+value.parseJson().get("family")

**複製貼上請注意語法是否有空格和空行，請刪除**

# 進階題-自動匯入高階層分類欄位

將一個欄位中的值分成不同欄位

a. 到 **higherClassification** 欄位並選擇 **Edit column >> Split into several columns**

b. 確認該欄位的分隔符號是逗號並按 **OK**

c. 一一將欄位名稱改為界、門、綱…

# 進階題-自動匯入高階層分類欄位



**將不要的欄位刪除**
a. 到 All欄位並選擇Edit columns
   >> Re-order/ remove columns
b. 拖曳左邊不想要的欄位到右邊區
   域並按 OK

# TaiBIF Open Data Toolkit 開放資料整合工具

## TaiBIF 內容經理 劉璟儀

# 開發緣由與背景

## 大家開放資料挫敗感太重

- DwC 那麼多要怎麼選
- 資料太多錯誤很難一次找到修正

## 很多工具可以用但尚未整合

- 如果有一個可以從資料欄位產生到編輯到清理都可以在同一個地方做完的工具就好了…

# **TaiBIF Open Data Toolkit**
# 開放資料整合工具出爐！

位產生資料 ➜ 編輯資料 ➜ 驗證資料 ➜ 清理資料 ➜ 打包成資料模

直接可以上傳IPT～

# 第一步：建立資料模板

# 第二步：編輯資料

# 第三步：驗證資料

# 第四步：清理資料與打包

# TaiBIF Open Data Toolkit
# 開放資料整合工具出爐！

位產生資料 → 編輯資料 → 驗證資料 → 清理資料 → 打包成資料檔

預計今年優化完成後，年底上線

# 下個月記得交作業喔～



**Open Refine / DwC 作業繳交 - https://forms.gle/zVCr1Ujeq3971TSs6**

**需 要 認 證 證 書 的 人 一 定 要 交 喔 ！**