# 開放資料的前置準備

**TaiBIF 內容經理 劉璟儀**
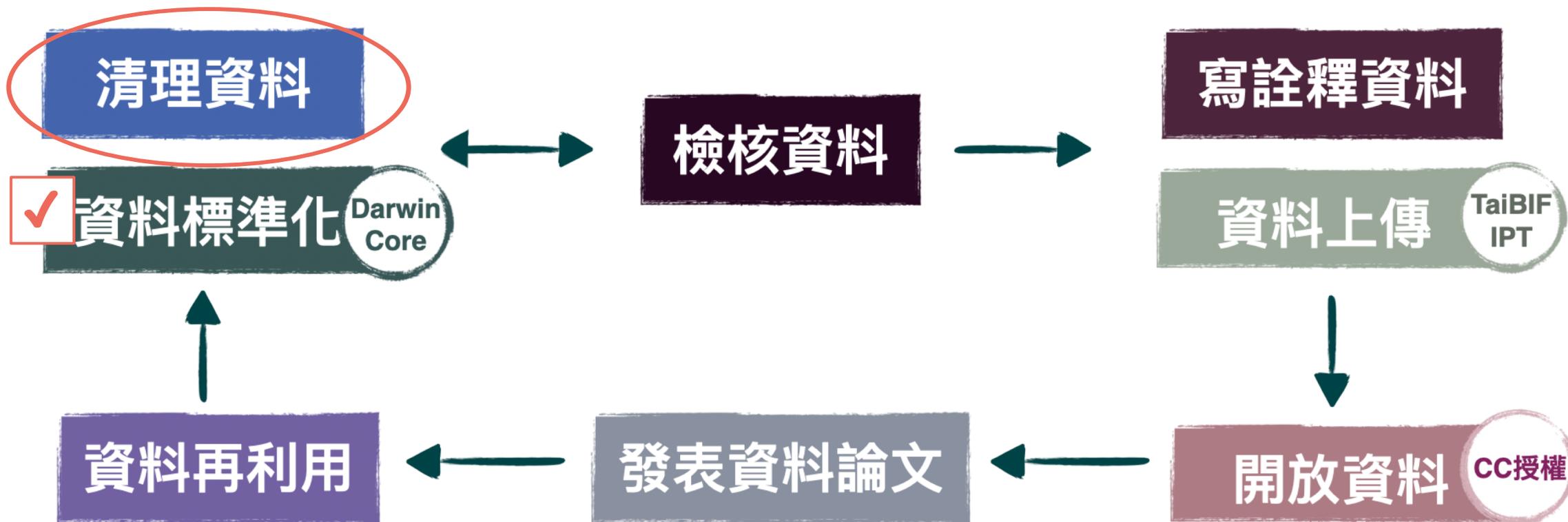
# 上傳資料前...
# 你應該準備好這些事

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

# 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

資料再利用

發表資料論文

開放資料 CC授權

# 什麼是資料清理？

## 你可能以為的資料清理…

把別人的資料弄成自己看得懂的樣子

把資料表 A 貼到資料表 B

把不顯著的資料刪掉 刪除明顯錯誤的資料

Ctrl + C

Ctrl + V

# 什麼是資料清理？

透過找出資料的錯誤、缺漏並修正，進而提升資料品質的過程



你敢保證你的資料品質完美無缺？

對啊，怎樣？

你這裡多了一個空格

# 什麼是資料清理？

**透過找出資料的錯誤、缺漏並修正，進而提升資料品質的過程**

**讓資料適合被利用**

# 常見的資料錯誤

- **格式不一致/錯誤**
- **拼字錯誤**
- **資料缺漏**
- **範圍錯誤**
- **ID/編碼重複**

什麼！那樣也算？！

# 資料清理小工具

- **地理分布線上座標系統轉換** https://portal.taibif.tw/coordinateConverter.php



**地理分布線上座標轉換** 批次轉換

**原始坐標**

投影座標系統　TWD97/ 臺灣地區

座標資料

**TWD 97**
**(生態調查常用)**

**轉換坐標**

投影座標系統　WGS 84 經緯度

轉換結果

↓

**WGS 84**
**(國際通用/ GBIF 預設)**

送出

**批次轉換說明**

將經度與緯度的坐標（經度在前，緯度在後）利用逗號或空格隔開，貼於原始坐標的輸入方塊中，並選擇原始輸入的坐標系統與欲輸出的坐標系統。再按下送出的按鈕，即可完成多筆坐標的轉換。

**網路服務使用說明**

參數說明：

- **source:** 來源坐標系統
  - 1: WGS 84 經緯度
  - 2: TWD 67 經緯度

# 資料清理小工具

- **座標轉換（度分秒- 十進位）** https://data.canadensys.net/tools/coordinates

# 資料清理小工具

- ## 有效學名比對 NomenMatch　　http://match.taibif.tw/

NomenMatch (code name: MyMatch): a scientific-name checking tool

### Query settings

| Result format | Sources | Version | Best results only? | Solr endpoints |
|---|---|---|---|---|
| table | ALL | #N/A | Yes (fast and simple) | DEFAULT (http://solr:8983/solr/taxa) |

### Scientific names

You can input one scientific name per line without or with authors, such as *Taiwania cryptomerioides* or *Taiwania cryptomerioides* Hayata

Taiwania cryptomerioides

Check names

Data Citations
- GBIF Secretariat: GBIF Backbone Taxonomy. doi:10.15468/39omei Accessed via http://www.gbif.org/dataset/d7dddbf4-2cf0-4f39-9b2a-bb099caae36c on 2016-09-06
- Roskov Y., Abucay L., Orrell T., Nicolson D., Kunze T., Culham A., Bailly N., Kirk P., Bourgoin T., DeWalt R.E., Decock W., De Wever A., eds. (2015). Species 2000 & ITIS Catalogue

14

# 資料清理小工具

- **有效學名比對 NomenMatch**

# 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）

# 資料清理小工具
## OpenRefine

| 1 | 2 | 3 |
|---|---|---|
| Excel | GBIF data validator | OpenRefine |

產生資料/ 管理資料　　　　　　檢核資料　　　　　　清理資料

# 資料清理小工具
## 用 OpenRefine 清資料 https://openrefine.org/download.html

# 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 »

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

Version 3.4-beta2 [c67e13b]

19

# 用 OpenRefine 清資料



**1** 選擇檔案後
**a. 確認下方文字編碼為 UTF-8**

**2** 按下 Create Project 進入使用介面

# 用 OpenRefine 清資料 範例1



善用每個欄位中的**Facet 功能**
可作基本的文字內容歸類
找出重複或文字錯誤

# 用 OpenRefine 清資料 範例2



使用正規表示式來搜尋格式不符的學名
正規表示式 Regular expression 說明

# 上傳資料前...
# 你應該準備好這些事

# 檢核小工具

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

# 檢核小工具

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

## Halieutichthys-上傳資料_melissa.xlsx

**Row type:** Darwin Core Occurrence

**Number of lines:** 823

**Number of interpreted dates:** 54

**Number of rows with interpreted taxon:** 823

### Term Frequency

| Term | Count | Percentage | Interpreted |
|---|---|---|---|
| dcterms:identifier 🔗 | 591 | 72% | |
| dwc:occurrenceID 🔗 | 591 | 72% | |
| dcterms:rightsHolder 🔗 | 591 | 72% | |
| dwc:institutionCode 🔗 | 591 | 72% | |
| dwc:basisOfRecord 🔗 | 591 | 72% | |
| dwc:catalogNumber 🔗 | 591 | 72% | |
| dwc:identifiedBy 🔗 | 591 | 72% | |
| dwc:scientificName 🔗 | 591 | 72% | |
| dwc:individualCount 🔗 | 589 | 72% | |
| dwc:lifeStage 🔗 | 591 | 72% | |
| dwc:establishmentMeans 🔗 | 591 | 72% | |
| dwc:preparations 🔗 | 591 | 72% | |
| dwc:disposition 🔗 | 591 | 72% | |
| dwc:associatedReferences 🔗 | 591 | 72% | |
| dwc:recordedBy 🔗 | 364 | 44% | |
| dwc:eventDate 🔗 | 575 | 70% | |
| dwc:occurrenceRemarks 🔗 | 30 | 4% | |

### Validation Issues

**Resource Structure**

Record not uniquely identified    231

**GBIF Occurrence Interpretation**

Recorded date invalid    521
Taxon match none    232
Basis of record invalid    232
Taxon match higherrank    46

26

# 用**OpenRefine**清理資料

**TaiBIF 內容經理 劉璟儀**

# 用 **OpenRefine** 清資料



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
檢核資料

**3**  OpenRefine
清理資料

4

# 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

# 用OpenRefine清理資料
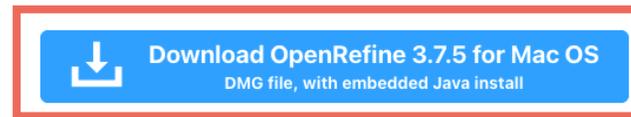
## 會需要用到的連結

練習檔案下載
GBIF Data Validator
NomenMatch 學名比對
Global Names Resolver

# 用 OpenRefine 清資料

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

# 用 OpenRefine 清資料

**下載並安裝在電腦** 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！

# 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

12

# 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]

# 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 Creat 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

15

# 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 ⇕

# 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
"PresevedSpecimen"

# 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

# 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)



**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

此部分無法批次複製修改，僅能個別修正

# 用 **OpenRefine** 清資料- 內容錯誤(座標)



## Coordinate conversion

Use this tool to convert geographic coordinates from DDMMSS to decimal degrees. Type coordinate pairs on separate lines or paste la
longitude columns from a spreadsheet. Each row may be optionally preceded by an identifier followed by a pipe or tab.

5° 35' 12" N, 75° 46' 18" W

**1**

貼上座標並按Submit

**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

## Coordinate conversion results

| original | decimalLatitude | decimalLongitude |
|---|---|---|
| 5° 35' 12" N, 75° 46' 18" W | 5.5866667 | −75.7716667 |

24

# 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)

# 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

# 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



**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

記得下面兩個選項要打勾

# 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? |
| 2. | | |
| 3. | | |
| 4. | | |

OK

**語法在下一頁，請整串複製貼上**

## 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

43

# 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 colimnus
   >> Re-order/ remove columns
b. 拖曳左邊不想要的欄位到右邊區
   域並按 OK

# Thank you!