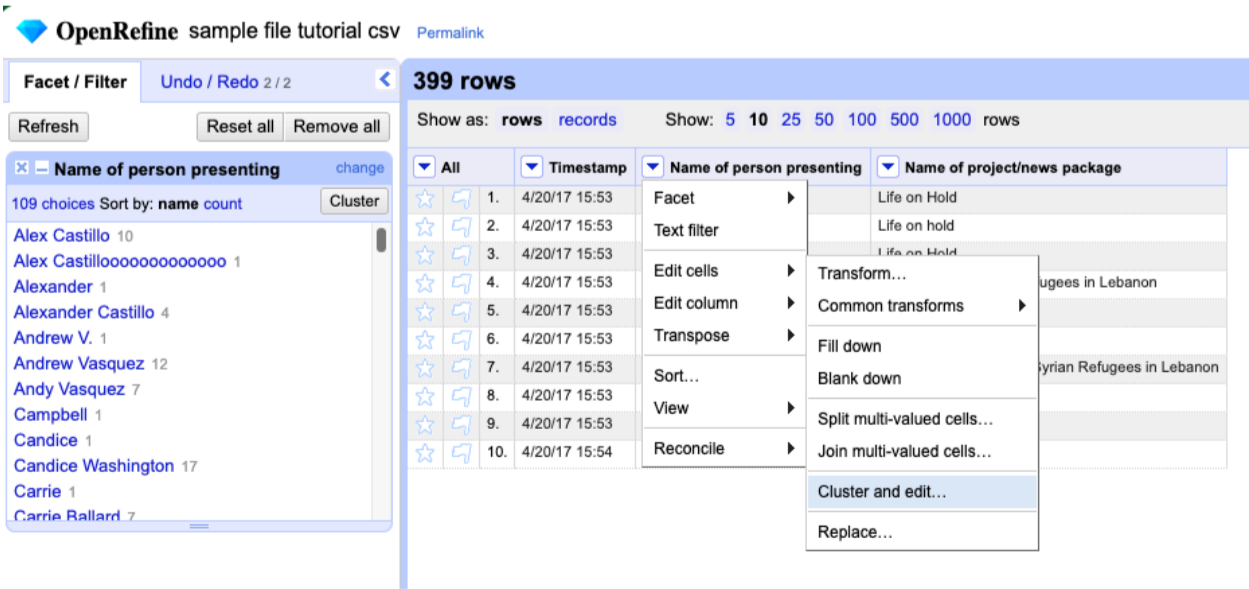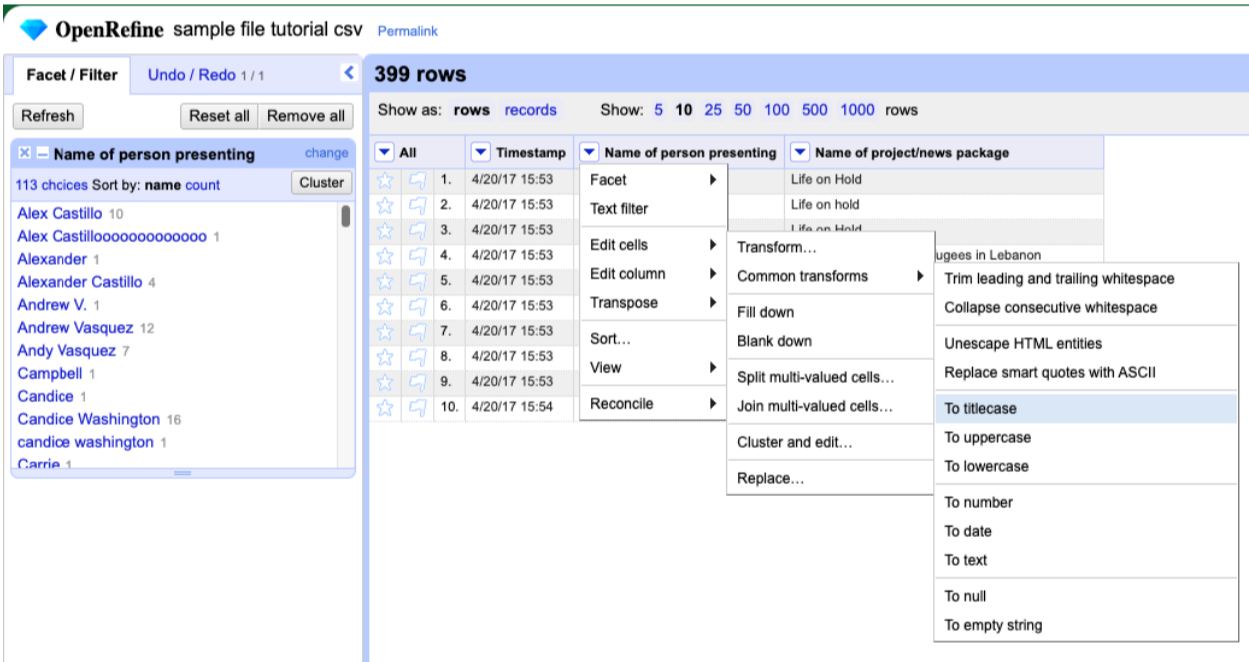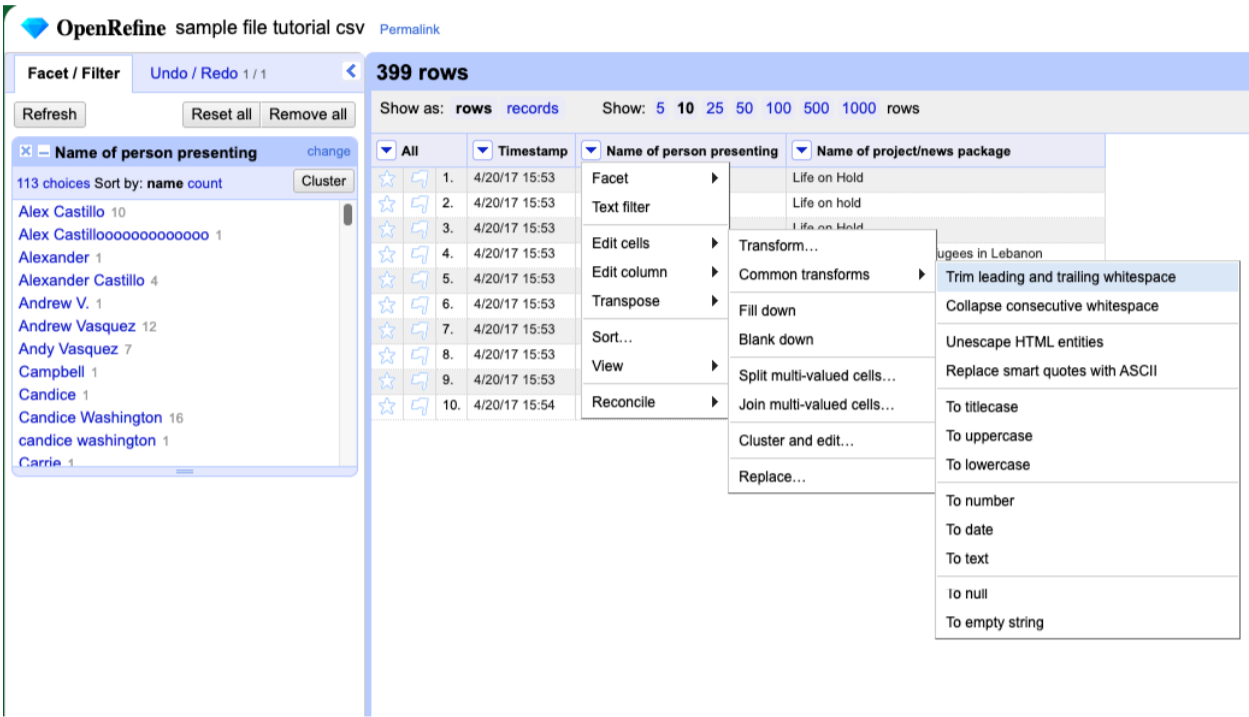# Business Intelligence - E07 - ETL (Open Refine)

Ömer Furkan Çoban

## Tasks - Part I

1.In order to get familiar with the tool follow this tutorial on OpenRefine"Using OpenRefine to Clean Your Data"–https://multimedia.report/tutorials/openrefine/

## Cluster and edit column "Name of person presenting"

Find groups of different cell values that might be other representations of the same thing. For example, "New York" and "new york" likely refer to the same concept and just differ by capitalization, and "Gödel" and "Godel" probably refer to the same person. Find out more…

Method [Key collision ▾]   Keying function [Fingerprint ▾]   [Manage clustering fuctions]
☐ Auto-update   **3** clusters found

| Merge? | Values in cluster | New cell value | Cluster size | Row count |
|---|---|---|---|---|
| ☑ | ☑ Sheila Rhodes & Jake Wheeler (2 rows)<br>☑ Sheila Rhodes, Jake Wheeler | Sheila Rhodes & Jake Wheeler | 2 | 3 |
| ☑ | ☑ Candice Washington (17 rows)<br>☑ Washington Candice | Candice Washington | 2 | 18 |
| ☑ | ☑ Jay And Sheila<br>☑ Sheila And Jay | Jay And Sheila | 2 | 2 |

**# Rows in cluster**
2 — 18

**Average length of choices**
14 — 28

**Length variance of choices**
0 — 0.5

[Select all] [Deselect all]          [Export clusters] [**Merge selected & re-cluster**] [Merge selected & Close] [Close]

---

## Cluster and edit column "Name of person presenting"

Find groups of different cell values that might be other representations of the same thing. For example, "New York" and "new york" likely refer to the same concept and just differ by capitalization, and "Gödel" and "Godel" probably refer to the same person. Find out more…

Method [Key collision ▾]   Keying function [n-Gram fingerprint ▾]   [Manage clustering fuctions]   n-Gram size [2]
☐ Auto-update   **3** clusters found

| Merge? | Values in cluster | New cell value | Cluster size | Row count |
|---|---|---|---|---|
| ☑ | ☑ Hattie Garett (2 rows)<br>☑ Hatie Garett | Hattie Garett | 2 | 3 |
| ☑ | ☑ Sheila Rhodes, Jacob Wheeler (3 rows)<br>☑ Sheila Rhodes/jacob Wheeler | Sheila Rhodes, Jacob Wheeler | 2 | 4 |
| ☑ | ☑ Hattie Garrett (10 rows)<br>☑ Hattie Garret (3 rows) | Hattie Garrett | 2 | 13 |

**# Rows in cluster**
3 — 13

**Average length of choices**
12 — 28

[Select all] [Deselect all]          [Export clusters] [**Merge selected & re-cluster**] [Merge selected & Close] [Close]

---

## Cluster and edit column "Name of person presenting"

Find groups of different cell values that might be other representations of the same thing. For example, "New York" and "new york" likely refer to the same concept and just differ by capitalization, and "Gödel" and "Godel" probably refer to the same person. Find out more…

Method [Key collision ▾]   Keying function [Metaphone3 ▾]   [Manage clustering fuctions]
☐ Auto-update   **19** clusters found

| Merge? | Values in cluster | New cell value | Cluster size | Row count |
|---|---|---|---|---|
| ☑ | ☑ Deanna Zimmerman (6 rows)<br>☑ Diana Zimmerman (4 rows)<br>☑ Deanna Zimerman (2 rows)<br>☑ Deana Zimerman<br>☑ Deana Zimmerman<br>☑ Diana Zimerman | Deanna Zimmerman | 6 | 15 |
| ☑ | ☑ Gordon Barry (6 rows)<br>☑ Gordon Berry (6 rows)<br>☑ Gordon Bary (4 rows)<br>☑ Gordon Bery | Gordon Barry | 4 | 17 |
| ☑ | ☑ Shana Lindsey (10 rows)<br>☑ Shana Lindsay (4 rows)<br>☑ Shawna Lindsay<br>☑ Shawna Lindsey | Shana Lindsey | 4 | 16 |
| ☑ | ☑ Sonya Mann (6 rows)<br>☑ Sonia Mann (5 rows) | Sonya Mann | 3 | 14 |

**# Choices in cluster**
2 — 6

**# Rows in cluster**
2 — 17

**Average length of choices**
4 — 27

**Length variance of choices**
0 — 6

[Select all] [Deselect all]          [Export clusters] [**Merge selected & re-cluster**] [Merge selected & Close] [Close]

**Cluster and edit column "Name of person presenting"**

Find groups of different cell values that might be other representations of the same thing. For example, "New York" and "new york" likely refer to the same concept and just differ by capitalization, and "Gödel" and "Godel" probably refer to the same person. Find out more…

Method | Key collision ▾    Keying function | Cologne phonetic ▾    [ Manage clustering fuctions ]
☐ Auto-update    **4 clusters found**

| Merge? | Values in cluster | New cell value | Cluster size | Row count |
|---|---|---|---|---|
| ☑ | ☑ Sonya Mann (14 rows)<br>☑ Sonja Mann | Sonya Mann | 2 | 15 |
| ☑ | ☑ Jacqueline (4 rows)<br>☑ Jacqueline J | Jacqueline | 2 | 5 |
| ☑ | ☑ Deanna (6 rows)<br>☑ Tim (5 rows) | Deanna | 2 | 11 |
| ☑ | ☑ Natasha Roberson (10 rows)<br>☑ Natasha Robertson | Natasha Roberson | 2 | 11 |

**# Rows in cluster**
5 — 15

**Average length of choices**
4 — 17

**Length variance of choices**
0 — 1.5

[ Select all ] [ Deselect all ]    [ Export clusters ] [ **Merge selected & re-cluster** ] [ Merge selected & Close ] [ Close ]

---

**Cluster and edit column "Name of person presenting"**

Find groups of different cell values that might be other representations of the same thing. For example, "New York" and "new york" likely refer to the same concept and just differ by capitalization, and "Gödel" and "Godel" probably refer to the same person. Find out more…

Method | Key collision ▾    Keying function | Daitch-Mokotoff ▾    [ Manage clustering fuctions ]
☐ Auto-update    **3 clusters found**

| Merge? | Values in cluster | New cell value | Cluster size | Row count |
|---|---|---|---|---|
| ☑ | ☑ Sheila Rhodes And Jacob Wheeler (6 rows)<br>☑ Sheila Rhodes, Jacob Wheeler (4 rows)<br>☑ Sheila Rhodes & Jake Wheeler (3 rows)<br>☑ Shelia Rhodes (2 rows) | Sheila Rhodes And Jacob Whe | 4 | 15 |
| ☑ | ☑ Alexander Castillo (4 rows)<br>☑ Alexander | Alexander Castillo | 2 | 5 |
| ☑ | ☑ Jacqueline Estrada (4 rows)<br>☑ Jacquiline Es | Jacqueline Estrada | 2 | 5 |

**# Choices in cluster**
2 — 4

**# Rows in cluster**
5 — 15

**Average length of choices**
13 — 25

**Length variance of choices**
2.5 — 7.04

[ Select all ] [ Deselect all ]    [ Export clusters ] [ **Merge selected & re-cluster** ] [ Merge selected & Close ] [ Close ]

---

2.In addition to the tutorial perform following manipulations on the dataset provided to you:–Create a new column that contains length of the "Name of project/news package"

**Add column based on column Name of person presenting**

New column name    project_name_length

On error    ● set to blank  ○ store error  ○ copy value from original column

Expression    Language  General Refine Expression Language (GREL) ▾

value.trim().length()                                    No syntax error.

**Preview**    History    Starred    Help

| row | value | value.trim().length() |
| --- | --- | --- |
| 1. | Tim Campbell | 12 |
| 2. | Tim Campbell | 12 |
| 3. | Timothy Campbell | 16 |
| 4. | Timothy Campbell | 16 |
| 5. | Deanna | 6 |
| 6. | Campbell | 8 |

OK    Cancel

–Create a new column that contains length of the "given name"

| given_name |
| --- |
| 12 |
| 12 |
| 16 |
| 16 |
| 6 |
| 8 |
| 12 |
| 6 |
| 12 |
| 16 |

–Create a new column that contains length of the "surname"

**Add column based on column Name of person presenting**

New column name    surname_length

On error    ● set to blank  ○ store error  ○ copy value from original column

Expression    Language  General Refine Expression Language (GREL) ▾

value.split(" ").reverse()[0].length()                   No syntax error.

**Preview**    History    Starred    Help

| row | value | value.split(" ").reverse()[0]. ... |
| --- | --- | --- |
| 1. | Tim Campbell | 8 |
| 2. | Tim Campbell | 8 |
| 3. | Timothy Campbell | 8 |
| 4. | Timothy Campbell | 8 |
| 5. | Deanna | 6 |
| 6. | Campbell | 8 |

OK    Cancel

3.Answer following questions

–What is the longest given name in this dataset?



–What is the longest "Name of project/news package"



–What is the average length of "Name of project/news package"



# Tasks –Part II

1.Import data from the "dates_synthetic.csv" file into your OpenRefineproject



2.Create a new column that contains the first 2 octets of the provided IP address

## Add column based on column internal_ip

New column name: IP_First_Two_Octets

On error: ◉ set to blank ○ store error ○ copy value from original column

Expression: Language: General Refine Expression Language (GREL) ▾

```
value.split(".")[0,2].join(".")
```

No syntax error.

**Preview** | History | Starred | Help

| row | value | value.split(".")[0,2].join("." ... |
|-----|-------|-----------------------------------|
| 1. | 10.234.114.114 | 10.234 |
| 2. | 10.121.236.103 | 10.121 |
| 3. | 192.168.241.110 | 192.168 |
| 4. | 10.42.131.205 | 10.42 |
| 5. | 172.22.243.247 | 172.22 |
| 6. | 10.120.188.48 | 10.120 |

OK | Cancel

3.Create multiple new columns with the following content for each date in the column "datetime_eu"

–day of the week

## Add column based on column datetime_eu

New column name: day_of_week

On error: ◉ set to blank ○ store error ○ copy value from original column

Expression: Language: General Refine Expression Language (GREL) ▾

```
value.toString("EEEE")
```

No syntax error.

**Preview** | History | Starred | Help

| row | value | value.toString("EEEE") |
|-----|-------|------------------------|
| 1. | 2019-10-13T12:25:04Z | Sunday |
| 2. | 2020-04-25T18:21:30Z | Saturday |
| 3. | 2019-07-13T08:17:39Z | Saturday |
| 4. | 2019-08-14T14:05:43Z | Wednesday |
| 5. | 2020-06-28T03:06:00Z | Sunday |
| 6. | 2019-12-29T21:55:58Z | Sunday |

OK | Cancel

–week of the year

## Add column based on column datetime_eu

New column name: week_of_the_year

On error: ● set to blank  ○ store error  ○ copy value from original column

Expression    Language: General Refine Expression Language (GREL) ▾

```
value.toString("w")
```
No syntax error.

**Preview**  History  Starred  Help

| row | value | value.toString("w") |
|-----|-------|---------------------|
| 1. | 2019-10-13T12:25:04Z | 42 |
| 2. | 2020-04-25T18:21:30Z | 17 |
| 3. | 2019-07-13T08:17:39Z | 28 |
| 4. | 2019-08-14T14:05:43Z | 33 |
| 5. | 2020-06-28T03:06:00Z | 27 |
| 6. | 2019-12-29T21:55:58Z | 1 |

OK  Cancel

–month of the year (as number)

## Add column based on column datetime_eu

New column name: month_of_year

On error: ● set to blank  ○ store error  ○ copy value from original column

Expression    Language: General Refine Expression Language (GREL) ▾

```
value.datePart("month") + 1
```
No syntax error.

**Preview**  History  Starred  Help

| row | value | value.datePart("month") + 1 |
|-----|-------|------------------------------|
| 1. | 2019-10-13T12:25:04Z | 11 |
| 2. | 2020-04-25T18:21:30Z | 5 |
| 3. | 2019-07-13T08:17:39Z | 8 |
| 4. | 2019-08-14T14:05:43Z | 9 |
| 5. | 2020-06-28T03:06:00Z | 7 |
| 6. | 2019-12-29T21:55:58Z | 13 |

OK  Cancel

–month of the year ( as string -'January')

**Add column based on column datetime_eu**

New column name  month_name

On error  ◉ set to blank  ○ store error  ○ copy value from original column

Expression  Language  General Refine Expression Language (GREL) ⌄

```
toString(value, "MMMM")
```
No syntax error.

**Preview**  History  Starred  Help

| row | value | toString(value, "MMMM") |
|-----|-------|-------------------------|
| 1. | 2019-10-13T12:25:04Z | October |
| 2. | 2020-04-25T18:21:30Z | April |
| 3. | 2019-07-13T08:17:39Z | July |
| 4. | 2019-08-14T14:05:43Z | August |
| 5. | 2020-06-28T03:06:00Z | June |
| 6. | 2019-12-29T21:55:58Z | December |

OK  Cancel

4.Homogenize given date and time for all columns with datetime value, use the following format "%Y-%m-%d %H:%M"

**Custom text transform on column datetime_eu**

Expression  Language  General Refine Expression Language (GREL) ⌄

```
toString(value, "yyyy-MM-dd HH:mm")
```
No syntax error.

**Preview**  History  Starred  Help

| row | value | toString(value, "yyyy-MM-dd HH ... |
|-----|-------|------------------------------------|
| 1. | 2019-10-13T12:25:04Z | 2019-10-13 12:25 |
| 2. | 2020-04-25T18:21:30Z | 2020-04-25 18:21 |
| 3. | 2019-07-13T08:17:39Z | 2019-07-13 08:17 |
| 4. | 2019-08-14T14:05:43Z | 2019-08-14 14:05 |
| 5. | 2020-06-28T03:06:00Z | 2020-06-28 03:06 |
| 6. | 2019-12-29T21:55:58Z | 2019-12-29 21:55 |

On error  ◉ keep original  ☐ Re-transform up to [10] times until no change
          ○ set to blank
          ○ store error

OK  Cancel

## Custom text transform on column datetime_usa1

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

```
toString(value, "yyyy-MM-dd HH:mm")
```
No syntax error.

**Preview**    History    Starred    Help

| row | value | toString(value, "yyyy-MM-dd HH ... |
|-----|-------|-----------------------------------|
| 1. | 2019-12-21T03:11:35Z | 2019-12-21 03:11 |
| 2. | 2019-05-02T19:26:00Z | 2019-05-02 19:26 |
| 3. | 2019-09-21T02:50:03Z | 2019-09-21 02:50 |
| 4. | 2019-01-03T07:34:47Z | 2019-01-03 07:34 |
| 5. | 2019-12-12T21:16:44Z | 2019-12-12 21:16 |
| 6. | 2019-09-08T12:39:07Z | 2019-09-08 12:39 |

On error    ◉ keep original    ☐ Re-transform up to [ 10 ] times until no change
            ○ set to blank
            ○ store error

**OK**   Cancel

---

## Custom text transform on column datetime_usa2

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

```
toString(value, "yyyy-MM-dd HH:mm")
```
No syntax error.

**Preview**    History    Starred    Help

| row | value | toString(value, "yyyy-MM-dd HH ... |
|-----|-------|-----------------------------------|
| 1. | 2019-12-21T03:11:00Z | 2019-12-21 03:11 |
| 2. | 2019-05-02T19:26:00Z | 2019-05-02 19:26 |
| 3. | 2019-09-21T02:50:00Z | 2019-09-21 02:50 |
| 4. | 2019-01-03T07:34:00Z | 2019-01-03 07:34 |
| 5. | 2019-12-12T21:16:00Z | 2019-12-12 21:16 |
| 6. | 2019-09-08T12:39:00Z | 2019-09-08 12:39 |

On error    ◉ keep original    ☐ Re-transform up to [ 10 ] times until no change
            ○ set to blank
            ○ store error

**OK**   Cancel

5.Calculate "time delta" with today's date for the column with the date "datetime_eu"

## Add column based on column datetime_eu

New column name: `Time_Delta_Days`

On error: ● set to blank ○ store error ○ copy value from original column

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

```
diff(now(), value, "days")
```
No syntax error.

**Preview**    History    Starred    Help

| row | value | diff(now(), value, "days") |
|-----|-------|----------------------------|
| 1. | 2019-10-13T12:25:00Z | 2274 |
| 2. | 2020-04-25T18:21:00Z | 2079 |
| 3. | 2019-07-13T08:17:00Z | 2366 |
| 4. | 2019-08-14T14:05:00Z | 2334 |
| 5. | 2020-06-28T03:06:00Z | 2015 |
| 6. | 2019-12-29T21:55:00Z | 2197 |

[ OK ] [ Cancel ]

# Tasks –Part III

1.Import data from the "UOL News.csv" file into your OpenRefineproject

2.Check, if provided within the dataset URLs are still available

## Add column based on column URL_UOL_News

New column name: `URL_Status`

On error: ● set to blank ○ store error ○ copy value from original column

Expression    Language [Python / Jython ▾]

```
import urllib2
try:
    response = urllib2.urlopen(value)
    return response.getcode()
except Exception as e:
```
No syntax error.

**Preview**    History    Starred    Help

| row | value | import urllib2 try: respon ... |
|-----|-------|-------------------------------|
| 1. | https://uol.de/aktuelles/artikel/digitale-frachtpapiere-und-quantenbiologie-4637 | 200 |
| 2. | https://uol.de/aktuelles/artikel/demokratisches-defizit-das-parlament-in-der-pandemie-4605 | 200 |
| 3. | https://uol.de/aktuelles/artikel/kurze-clips-ueber-neue-buecher-4597 | 200 |
| 4. | https://uol.de/aktuelles/artikel/hoerwahrnehmung-besser-verstehen-durch-virtuelle-realitaet-4587 | 200 |

[ OK ] [ Cancel ]

3.Fetch data from a URL using "Add Column By Fetching URL" feature

**10 rows**

Show as: **rows** records    Show: 5 **10** 25 50 100 500 1000 rows

| ▼ All | ▼ Column | ▼ URL_UOL_News | ▼ Page_content | ▼ URL_Statu |
|---|---|---|---|---|
| ☆ ⟲ 1. | 1 | https://uol.de/aktuelles/artikel/digitale-frachtpapiere-und-quantenbiologie-4637 | <!DOCTYPE html><br><html class="uol" lang="de"><br><head><br><br><meta charset="utf-8"><br><!--<br>    This website is powered by TYPO3 - inspiring people to share!<br>    TYPO3 is a free open source Content Management Framework initially created by Kasper Skaarhoj and licensed under GNU/GPL.<br>    TYPO3 is copyright 1998-2026 of Kasper Skaarhoj.<br>    Extensions are copyright of their respective owners.<br>    Information and contribution at https://typo3.org/<br>--><br><br><title>Digitale Frachtpapiere und Quantenbiologie // Universität Oldenburg</title><br><meta http-equiv="x-ua-compatible" content="IE=edge"><br><meta name="generator" content="TYPO3 CMS"><br><meta name="description" content="Exzellente Forschung und herausragende Promotion: Der Physiker Ilia Solov&#039;yov und der Jurist David Saive erhalten in diesem Jahr die Preise der Universitätsgesellschaft Oldenburg."><br><meta name="robots" content="index,follow"><br><meta name="viewport" content="width=device-width, initial-scale=1"><br><meta property="og:title" content="Digitale Frachtpapiere und Quantenbiologie"><br><meta property="og:type" content="article"><br><meta property="og:url" content="https://uol.de/aktuelles/artikel/digitale-frachtpapiere-und-quantenbiologie-4637"><br><meta property="og:image" content="https://uol.de/fileadmin/_processed_/8/1/csm_hoersaal-a14-remmers-martin_kl_8af04d7c56.jpg"><br><meta property="og:image:width" content="1200"><br><meta property="og:image:height" content="675"><br><meta property="og:description" content="Exzellente Forschung und herausragende Promotion: Der Physiker Ilia Solov&#039;yov und der Jurist David Saive erhalten in diesem Jahr die Preise der Universitätsgesellschaft Oldenburg."><br><meta name="twitter:card" content="summary"><br><meta name="apple-mobile-web-app-title" content="Carl von Ossietzky Universität Oldenburg"><br><br><link rel="stylesheet" href="/_assets/d433799a7157dc4ba6bb3ee09f89cb2e/Assets/Css/fonts.css?1766074465" media="all"><br><link rel="stylesheet" href="/typo3temp/assets/compressed/media2click-e253fe1d4d64ad94fe26027e560e51de.css.gz?1766074555" media="all"><br><link rel="stylesheet" href="/typo3temp/assets/compressed/ecb3e01623a42af68084ccb4b40d310f-769b6db91a577963d2a8b3576b41a6f2.css.gz?1766074555" media="all"><br><link rel="stylesheet" href="/typo3temp/assets/compressed/ods_osm-7b045edb1eef455ae1b04208db5a681a.css.gz?1766074555" media="all"><br><link rel="stylesheet" href="/typo3temp/assets/compressed/bootstrap.min-dca42613438050e67b18fa259abc6f96.css.gz?1766074555" media="all"><br><link rel="stylesheet" href="/typo3temp/assets/compressed/uol-fontawesome-45793179f75260ea5767d545f31ef6ba.css.gz?1766074555" media="all"> | 20 |

4.For the each "item" below, create a separate column based on the downloaded content for each URL:

–page title

**Add column based on column Page_content**

New column name    Page_titie

On error    ◉ set to blank ○ store error ○ copy value from original column

Expression    Language General Refine Expression Language (GREL) ✓

```
value.parseHtml().select("title")[0].htmlText()
```
No syntax error.

**Preview**    History    Starred    Help

| row | value | value.parseHtml().select("titl ... |
|---|---|---|
| 1. | <!DOCTYPE html><br><html class="uol" lang="de"><br><head><br><br><meta charset="utf-8"><br><!--<br>    This website is powered by TYPO3 - inspiring people to share!<br>    TYPO3 is a free open source Content Management Framework initially created by Kasper | Digitale Frachtpapiere und Quantenbiologie // Universität Oldenburg |

OK    Cancel

–publication date

## Add column based on column Page_content

New column name  publication_date

On error  ⦿ set to blank  ○ store error  ○ copy value from original column

Expression  Language  General Refine Expression Language (GREL) ▾

```
value.parseHtml().select("time")[0].htmlText()
```
No syntax error.

**Preview**  History  Starred  Help

| row | value | value.parseHtml().select("time ... |
|-----|-------|-----|
| 1. | <!DOCTYPE html><br><html class="uol" lang="de"><br><head><br><br><meta charset="utf-8"><br><!--<br>     This website is powered by TYPO3 - inspiring people to share!<br>     TYPO3 is a free open source Content Management Framework initially created by Kasper Skaarhoi and | 09.12.2020 |

**OK**  Cancel

–content of the main article for each given URL

## Add column based on column Page_content

New column name  article_content

On error  ⦿ set to blank  ○ store error  ○ copy value from original column

Expression  Language  General Refine Expression Language (GREL) ▾

```
value.parseHtml().select(".news-text-wrap")[0].htmlText()
```
No syntax error.

**Preview**  History  Starred  Help

| row | value | value.parseHtml().select(".new ... |
|-----|-------|-----|
| 1. | <!DOCTYPE html><br><html class="uol" lang="de"><br><head><br><br><meta charset="utf-8"><br><!--<br>     This website is powered by TYPO3 - inspiring people to share!<br>     TYPO3 is a free open source Content Management Framework initially created by Kasper | Exzellente Forschung und herausragende Promotion: Der Physiker Ilia Solov'yov und der Jurist David Saive erhalten in diesem Jahr die Preise der Universitätsgesellschaft Oldenburg. Mit dem „Preis für exzellente Forschung" den Physiker Prof. Dr. Ilia Solov'yov und den „Preis für herausragende Promotion" für den Juristen Dr. David Saive würdigt die Universitätsgesellschaft Oldenburg e. V. (UGO) die wichtige Pionierarbeit der beiden Wissenschaftler. Der Forschungspreis |

**OK**  Cancel

–*(optional -remove HTML tags)

**Custom text transform on column Page_content**

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

```
value.replace(/<[^>]+>/, "")
```
No syntax error.

**Preview**  History  Starred  Help

```
<strong>Navigation</strong></a>
[<code>Access-Key 1</code>]</p>
        <ul>
            <li><a id="sprung_sprache"
accesskey="8" href="#language_menu"              Sprache wechseln
lang="en">Choose other language</a>              Change Language
[<code>Access-Key 8</code>]</li>
            <li><a id="sprung_inhalt"
accesskey="2" href="#content">Zum Inhalt
springen</a> [<code>Access-Key 2</code>]
</li>
            <li><a id="sprung_suche"
```

On error    ◉ keep original        ☐ Re-transform up to [10] times until no change
            ○ set to blank
            ○ store error

                                                    [OK]  [Cancel]

# Tasks –Part IV

1.Import data from the "phm-collection-4000.tsv" file into your OpenRefineproject

2.In the column "Categories" change values into uppercase

**Custom text transform on column Categories**

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

```
value.toUppercase()
```
No syntax error.

**Preview**  History  Starred  Help

| row | value | value.toUppercase() |
| --- | --- | --- |
| 1. | null | Error: toUppercase expects a string |
| 2. | Mineral Samples-Geological | MINERAL SAMPLES-GEOLOGICAL |
| 3. | null | Error: toUppercase expects a string |
| 4. | Spacesuits | SPACESUITS |
| 5. | Spacecraft\|Models\|Space Technology | SPACECRAFT\|MODELS\|SPACE TECHNOLOGY |

On error    ◉ keep original        ☐ Re-transform up to [10] times until no change
            ○ set to blank
            ○ store error

                                                    [OK]  [Cancel]

3.What is the oldest item in the dataset according to the provided information/data?

```
✕ – start_year                          change reset



        1,040 — 2,020

☑ Numeric   ☐ Non-numeric   ☑ Blank   ☐ Error
  1941           0             2058        0
```

4.Create a new column, which contains the length of the "Description" column

## Add column based on column Description.

New column name: len_desc

On error: ● set to blank ○ store error ○ copy value from original column

Expression    Language: General Refine Expression Language (GREL) ⌄

```
value.length()
```
No syntax error.

**Preview**   History   Starred   Help

| row | value | value.length() |
|-----|-------|---------------|
| 2625. | Diazo print, poly chrome roof plan and sections Acklands Hotel Woollahra, paper / ink, Cyril Ruwald, 9 Martin Place, Sydney, New South Wales, Australia, printed by Commercial Copying Company Pty Ltd, 28 Martin Place, Sydney, New South Wales, Australia, 1st November 1937 A diazo print depicting the roof plan ('ROOF PLAN') and 2 sections ('SECTION "A"-"A" ' and 'SECTION "B"-"B" ') of Acklands Hotel, on the corner of Queen and Moncur Streets, Woollahra. The outline of the plan and each section is in brown and coloured in orange, yellow, red and blue. Notations handwritten in black ink can be seen in the top right corner of the print and additional notations handwritten in red ink can be seen down the right side | 957 |

OK    Cancel

5.Create multiple new columns and keep only numbers inside newly created columns. Do it for the following columns "Height", "Width", "Depth"

## Add column based on column Height

New column name: height_num

On error: ● set to blank ○ store error ○ copy value from original column

Expression    Language: General Refine Expression Language (GREL) ⌄

```
value.replace(/[^0-9.]/, "")
```
No syntax error.

**Preview**   History   Starred   Help

| row | value | value.replace(/[^0-9.]/, "") |
|-----|-------|------------------------------|
| 2625. | 335 mm | 335 |
| 2819. | 365 mm | 365 |
| 3156. | 337 mm | 337 |
| 2166. | 455 mm | 455 |
| 2272. | 455 mm | 455 |
| 2366. | 455 mm | 455 |

OK    Cancel

## Add column based on column Width

New column name    width_num

On error    ● set to blank    ○ store error    ○ copy value from original column

Expression    Language    General Refine Expression Language (GREL) ▼

`value.replace(/[^0-9.]/, "")`    No syntax error.

**Preview**    History    Starred    Help

| row | value | value.replace(/[^0-9.]/, "") |
|-----|-------|------------------------------|
| 2625. | 450 mm | 450 |
| 2819. | 480 mm | 480 |
| 3156. | 357 mm | 357 |
| 2166. | 690 mm | 690 |
| 2272. | 690 mm | 690 |
| 2366. | 690 mm | 690 |

OK    Cancel

## Add column based on column Depth

New column name    depth_num

On error    ● set to blank    ○ store error    ○ copy value from original column

Expression    Language    General Refine Expression Language (GREL) ▼

`value.replace(/[^0-9.]/, "")`    No syntax error.

**Preview**    History    Starred    Help

| row | value | |
|-----|-------|---|
| 411. | null | Error: replace expects three strings, or one string, one regex, and one string |
| 430. | 10 mm | 10 |
| 657. | 58 mm | 58 |
| 669. | 36 mm | 36 |
| 230. | null | Error: replace expects three strings, or one string, one regex, and one string |
| 566. | null | Error: replace expects three strings, or one string, one regex, and one string |
| 570. | null | Error: replace expects three strings, or one string, one regex, and one string |
| 618 | null | Error: replace expects three strings, or one string, one regex, and one string |

OK    Cancel

6.Create a new column that will calculate volume of each item (if measurements are provided). Which measurement unit should be used here?

## Add column based on column height_num

New column name    vol_mm3

On error              ◉ set to blank   ○ store error   ○ copy value from original column

Expression             Language  General Refine Expression Language (GREL) ▾

```
cells["height_num"].value * cells["width_num"].value *
cells["depth_num"].value
```
No syntax error.

**Preview**     History     Starred     Help

| | | |
|---|---|---|
| 618. | 26 | null |
| 827. | null | null |
| 199. | 60 | 24000 |
| 437. | null | null |
| 346. | 10 | null |
| 399. | null | null |
| 568. | null | null |

**OK**   Cancel