

Data support for Translation Imbalances in Wikipedia language editions

Notebook dependencies and setup

```
1 Mix.install([
2   {:csv, "~> 3.0"},
3   {:mediawiki_client, "~> 0.5.0"},
4   {:kino_vega_lite, "~> 0.1.0"},
5   {:vega_lite, "~> 0.1.0"}
6 ])
```

Evaluated*



:ok

Introduction

The data here is support for an extended abstract presented at Wiki Workshop 2025, "[Translation Imbalances Between Wikipedia Language Editions](#)". Although we originally did this work using private data sources, the limited results published in the paper can be recalculated from public sources and we attempt to do so below. Note that there may be slight differences due to the time ranges included.

This file is an Elixir [livebook](#) and it should be possible to run locally if you wish.

Fetch site statistics for all Wikipedia language editions

```
1 {:ok, all_sites} = Wiki.SiteMatrix.get_all()
2
3 wikipedias =
4   Enum.filter(all_sites, fn site ->
5     site.project == "wiki" && !site.closed && !site.private
6   end)
```

Evaluated*



```
1 site_stats =
2   wikipedias
3   |> Enum.map(fn site ->
4     {:ok, %{result: %{"query" => %{"statistics" => result}}}} =
5     site
6     |> Wiki.Action.new()
7     |> Wiki.Action.get(action: :query, meta: :siteinfo, siprop: :statistics)
8
9     {site.dbname, result}
10   end)
11   |> Enum.into(%{})
```

Evaluated*



Fetch Content Translation statistics

There used to be an API `contenttranslationstats` but this was disabled due to performance concerns (see task [T392839](#)). We're forced to use a stored response here, dated 2023-06-12.

In the future we should be able to refresh the data using sources being developed for a new [public dashboard for Content Translation](#).

```
1 content_translation_pairs =
2   File.stream!(Path.join(__DIR__, "content_translation_stats.csv"))
3   |> CSV.decode!(headers: true)
4   |> Enum.filter(&(&1["status"] == "published"))
```

Evaluated*



Aggregate by language

Here we collapse the raw data into a summary for each language, with the following aggregate calculations:

- Number of article pages in language edition
- Number of active editors for language
- Number of outgoing translations from language
- Number of incoming translations to language
- Translation hegemony (outgoing : incoming)

```

1 lang_stats =
2   wikipedias
3   |> Enum.reject(&(&1.dbname == "simplewiki"))
4   |> Enum.map(fn %{lang: lang, dbname: dbname} ->
5     stats = %{
6       article_count: site_stats[dbname]["articles"],
7       active_editor_count: site_stats[dbname]["activeusers"],
8       outgoing_translations:
9         content_translation_pairs
10         |> Enum.filter(&(&1["sourceLanguage"] == lang))
11         |> Enum.sum_by(&String.to_integer(&1["count"])),
12       incoming_translations:
13         content_translation_pairs
14         |> Enum.filter(&(&1["targetLanguage"] == lang))
15         |> Enum.sum_by(&String.to_integer(&1["count"])))
16     }
17
18     stats =
19       Map.put(
20         stats,
21         :translation_hegemony,
22         case stats.incoming_translations do
23           0 -> 0
24           _ -> stats.outgoing_translations / stats.incoming_translations
25         end
26       )
27
28     {lang, stats}
29   end)
30   |> Enum.into(%{})

```

Evaluated*

```

1 flat_lang_stats =
2   lang_stats
3   |> Enum.map(fn {lang, stats} ->
4     Map.put(stats, :lang, lang)
5   end)

```

Evaluated*

Save the aggregated data back to disk so it can be used by other packages.

```

1 flat_lang_stats
2 |> CSV.encode(
3   headers: ~w(lang article_count active_editor_count outgoing_translations incoming_translations translation_hegemony)
4 )
5 |> Enum.into(File.stream!(Path.join(__DIR__, "aggregate_s

```

Evaluated*

Scatterplots of translation hegemony

```
1 alias Vegalite, as: VL
```

Evaluated*

Vegalite

```

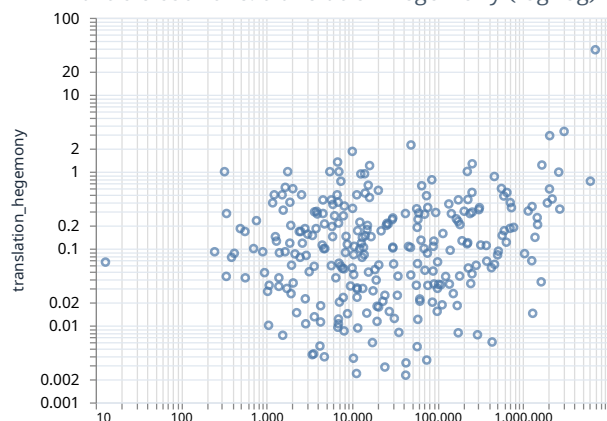
1 zero_safe =
2   flat_lang_stats
3   |> Enum.reject(& &1.translation_hegemony == 0)

1 VL.new(width: 400, height: 300, title: "Wiki article count vs. translation hegemony (log-log)")
2 |> VL.data_from_values(zero_safe, only: [:article_count, :translation_hegemony])
3 |> VL.mark(:point)
4 |> VL.encode_field(:x, "article_count", type: :quantitative, scale: [type: :log])
5 |> VL.encode_field(:y, "translation_hegemony", type: :quantitative, scale: [type: :log])

```

Evaluated*

Wiki article count vs. translation hegemony (log-log)



article_count

```
1 Vl.new(width: 400, height: 300, title: "Wiki active editor count vs. translation hegemony (log-log)"
2 )
3 |> Vl.data_from_values(zero_safe, only: [:active_editor_count, :translation_hegemony])
4 |> Vl.mark(:point)
5 |> Vl.encode_field(:x, "active_editor_count", type: :quantitative, scale: [type: :log])
6 |> Vl.encode_field(:y, "translation_hegemony", type: :quantitative, scale: [type: :log])
```

Wiki active editor count vs. translation hegemony (log-log)

