

Partitioning specialisation in multilevel textual corpora

A hierarchical model of pairwise similarities

Carlo Debernardi

CS2Italy2025

16/01/2025

Background

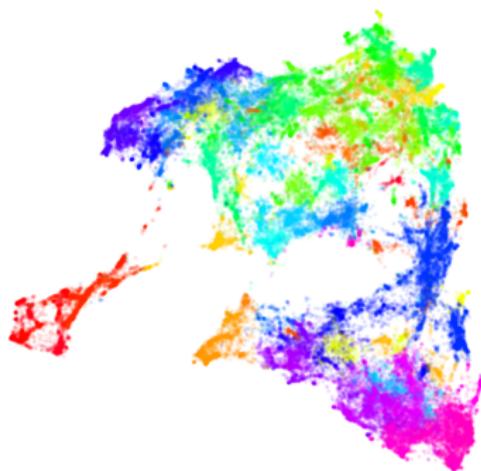
Recent advances in text analysis
(Mikolov et al. 2013, Vaswani et al. 2017, Reimers & Gurevych 2019)

	Text Embedding	Generative LLMs
Architecture	Encoder	Decoder
Output	Vectors	Text
Example	BERT	GPT

Close relatives, but one is vastly more popular than the other.
Not all potential applications have been explored yet.

Background

Typical application of text-embedding:
dimensionality reduction and clustering
(e.g: BERTopic - Grootendorst 2022)



Great for exploration but:

- Approximate lower dimensionality
- Dichotomization (technical and theoretical issue)
- Difficult to integrate into inferential analysis

Background

How to use the output of text embedding algorithms in inferential analyses?

Output: high dimensional vectors, usually optimized to encode relations as cosine similarities

- As a covariate: to be explored (Gaussian processes?)
- As the outcome: focus of this work

I.e. are texts, on average, more (or less) similar as a function of x ?

A case study

Objective

Study the specialisation in the scientific literature,
disentangling field and journal-level effects.

Specialisation operationalised as the average similarity among pairs of abstracts (i.e. higher average similarity = higher specialisation).

Ideally we want to be able to introduce covariates (e.g. time trend).

Note: the modeling exercise forces us to take care of details that turn out to be interesting in and of themselves.

The data

Data collection:

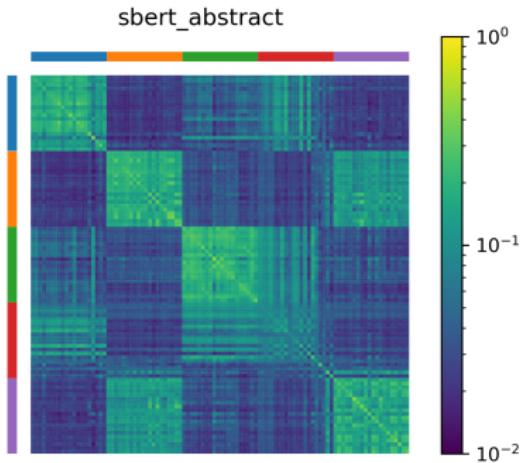
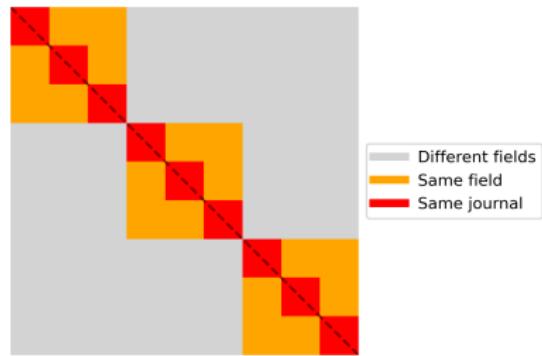
- Data collected from OpenAlex (CC0 license!)
- 5 fields (ASJC classification, Scopus)
- 20 journals per field (high rank, at least 100 pubs)
- 100 random papers per journal (timeframe 2010-2020)
- Results in 10k abstracts

Note: random subsampling in testing phase, 1.9M obs per field!

Fields:

- Condensed Matter Physics
- Economics and Econometrics
- Genetics
- Multidisciplinary
- Sociology and Political Science

Structure in the data



Highly structured data! → Hierarchical model
Non-independence & in/out journal

The model

$$Y_{ij} \sim Beta(\mu_{ij}, \phi)$$

Beta model to respect rescaled outcome domain [0-1]

Varying intercept as a baseline version

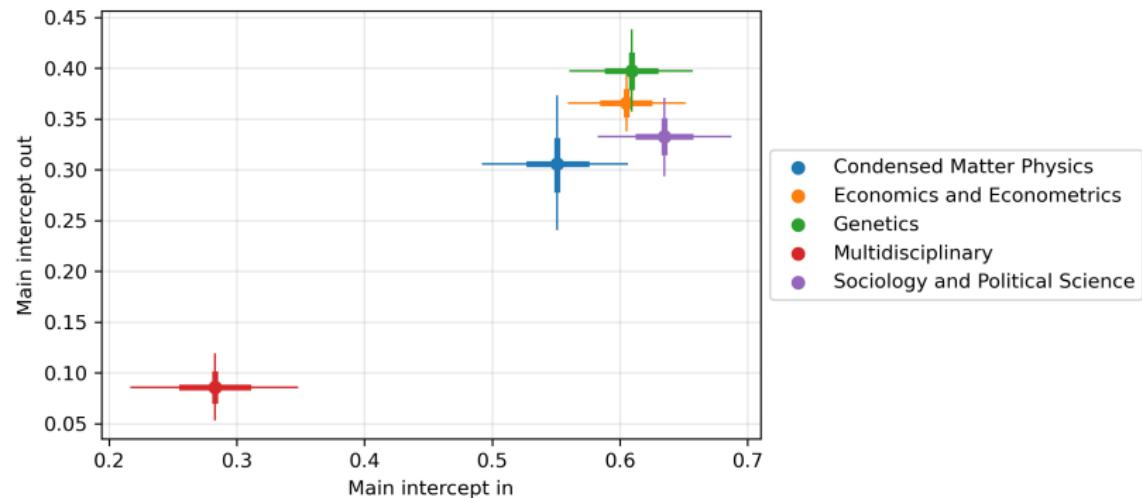
$$\mu_{ij} = logit^{-1}(F + J + P_i + P_j), \text{ if } J_i = J_j$$

Avg journal spec Journal spec Paper centrality (in journal)

$$\mu_{ij} = logit^{-1}(F + J_i + J_j + P_i + P_j), \text{ if } J_i \neq J_j$$

Field spec Journal centrality (in field) Paper centrality (in field)

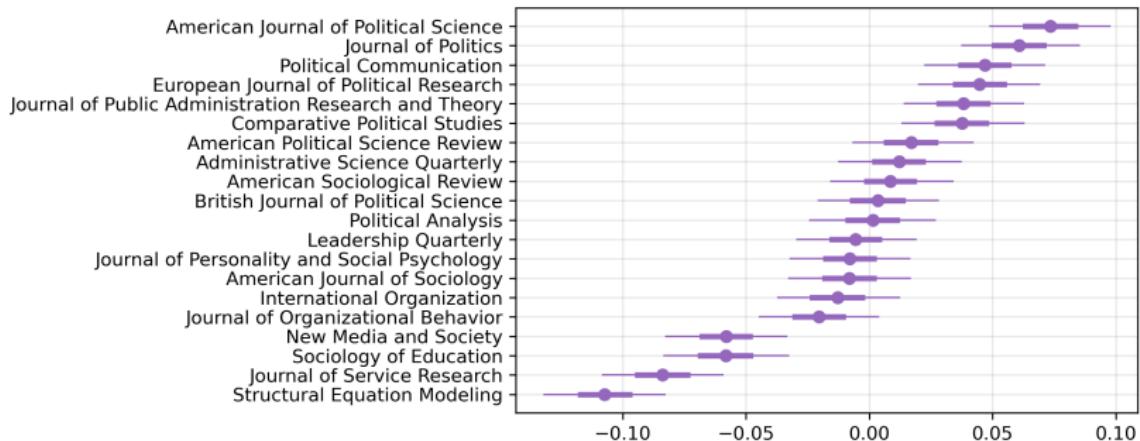
Field intercepts



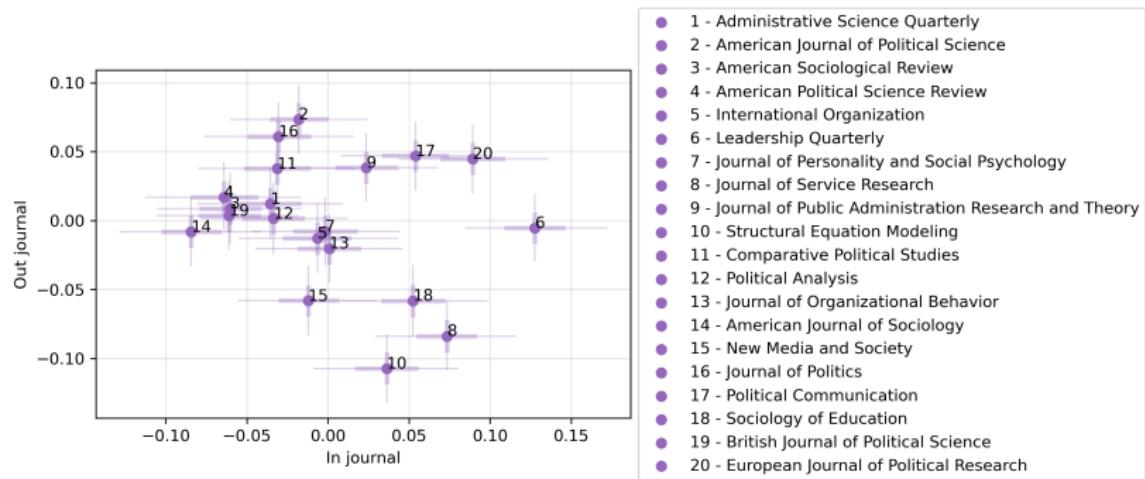
Journal intercepts - specialisation



Journal intercepts - centrality



Journal intercepts



Summary

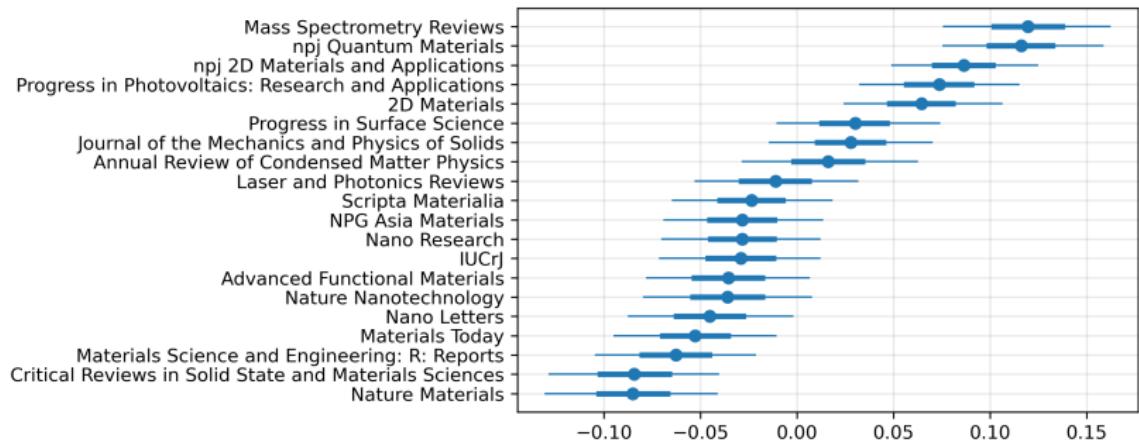
To sum up:

- There are still unexplored applications of text embedding
- It can be employed together with inferential analysis
- We gain the ability to estimate marginal quantities
- We are forced to deal with the structure of the data
- It is possible to disentangle structural levels that contribute to specialisation at varying degree

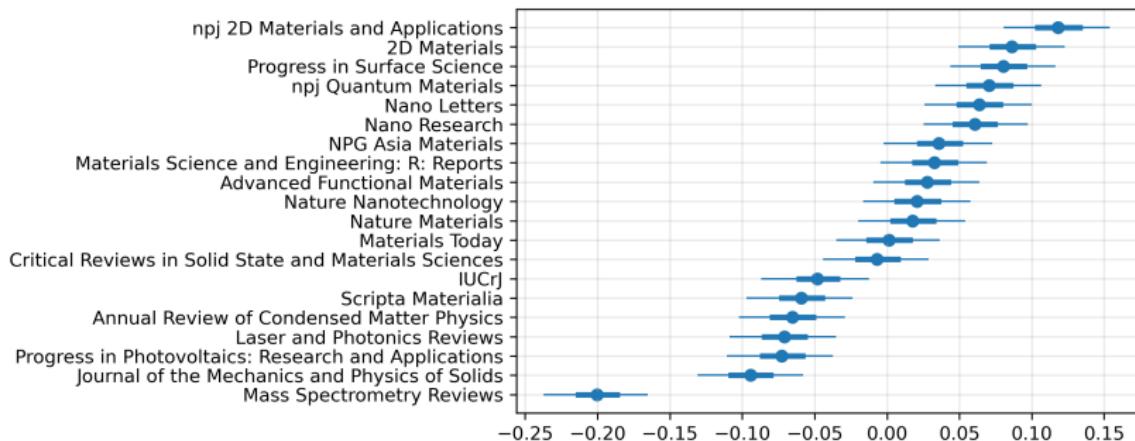
Thanks for the attention

carlo.debernardi@unisi.it

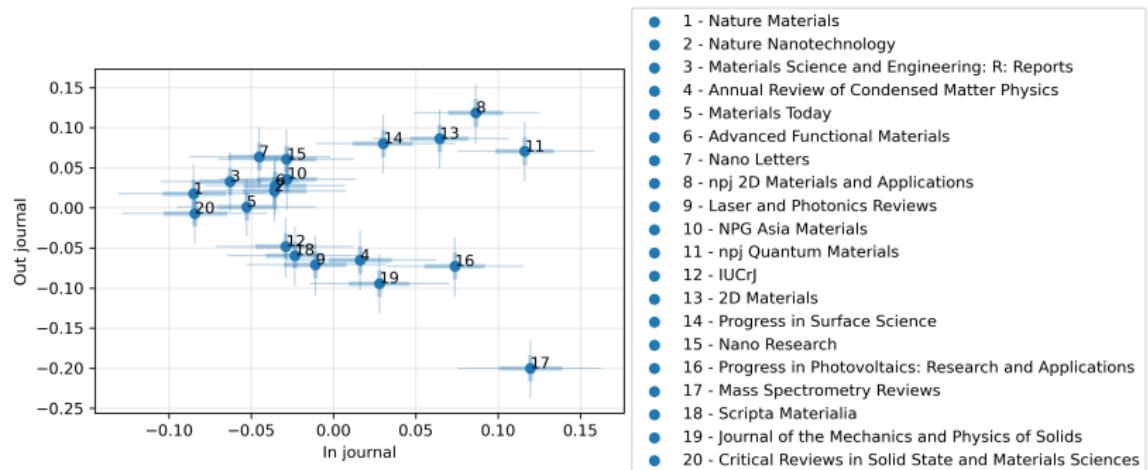
Journal intercepts - specialisation



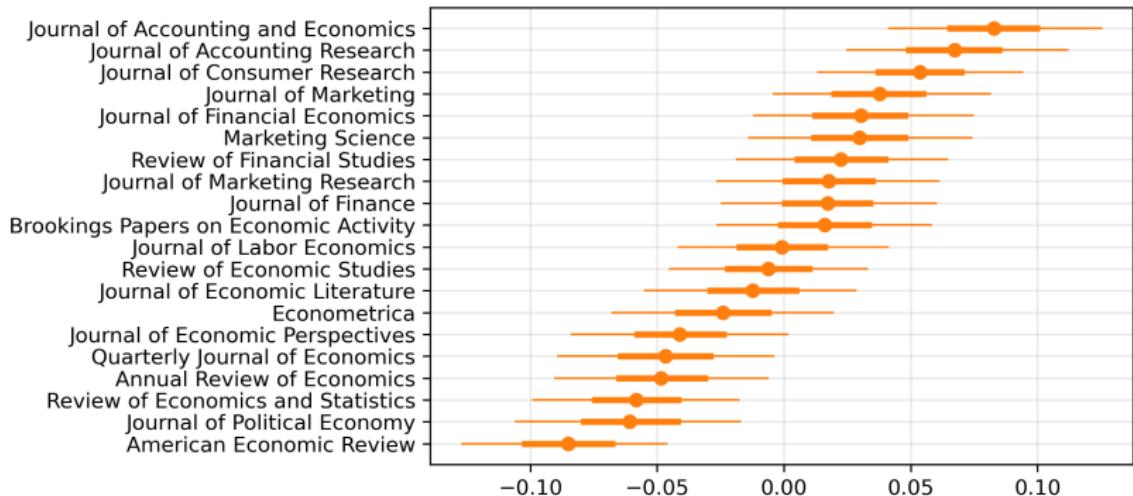
Journal intercepts - centrality



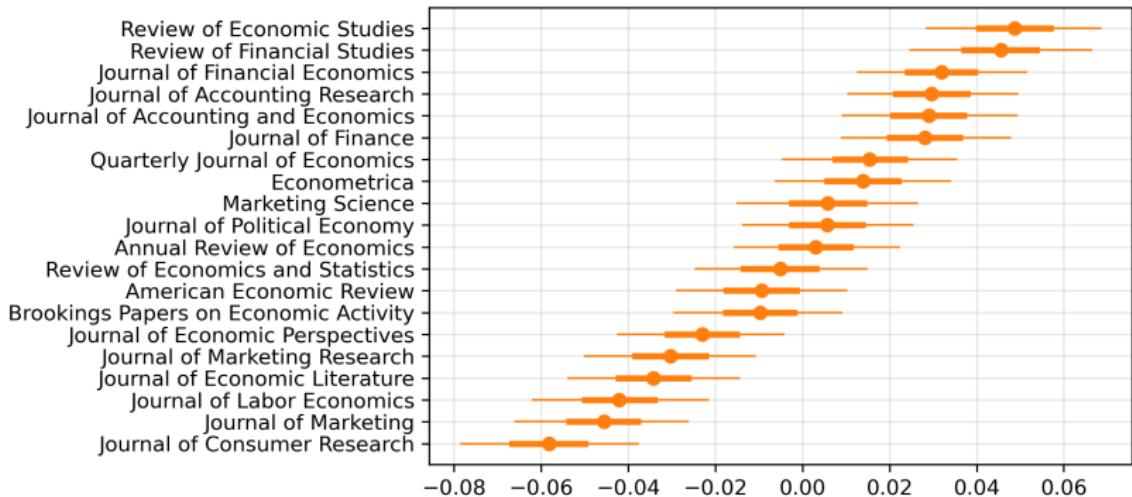
Journal intercepts



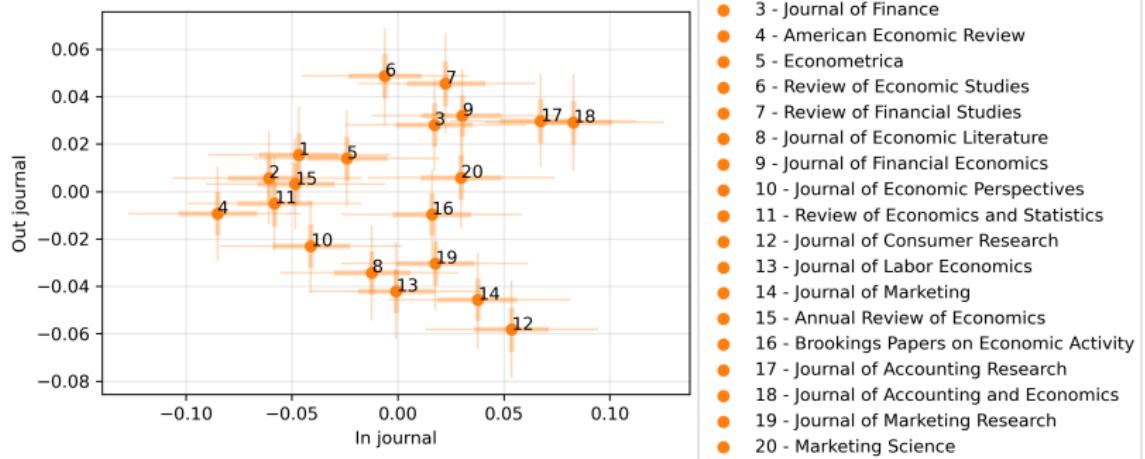
Journal intercepts - specialisation



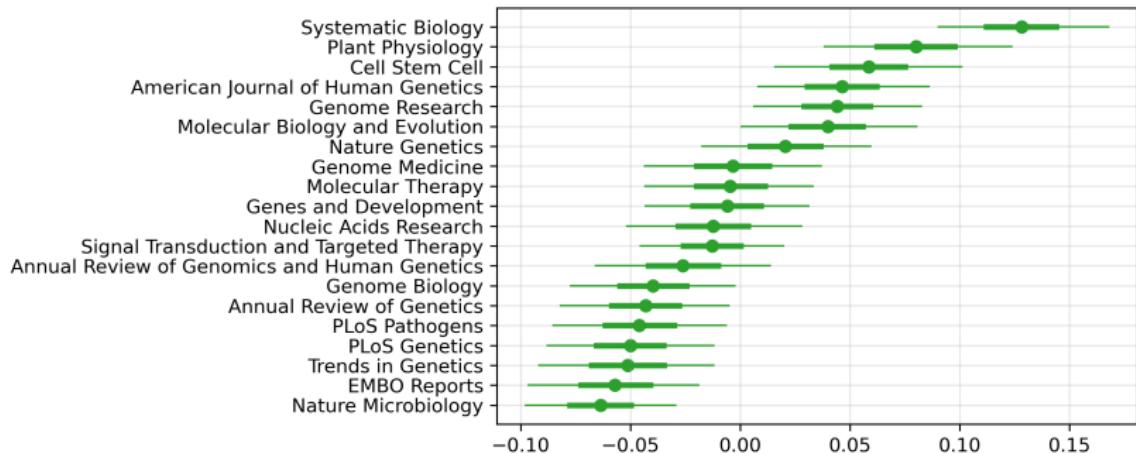
Journal intercepts - centrality



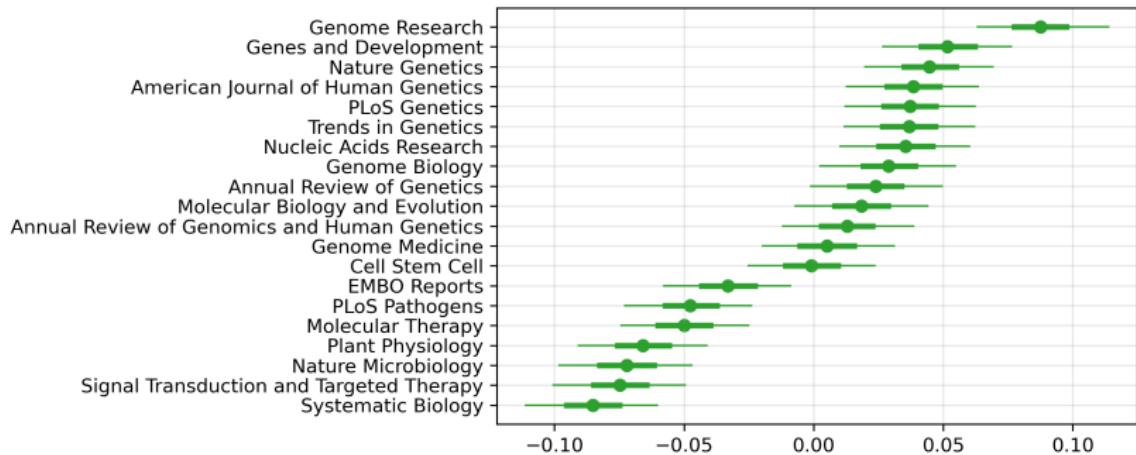
Journal intercepts



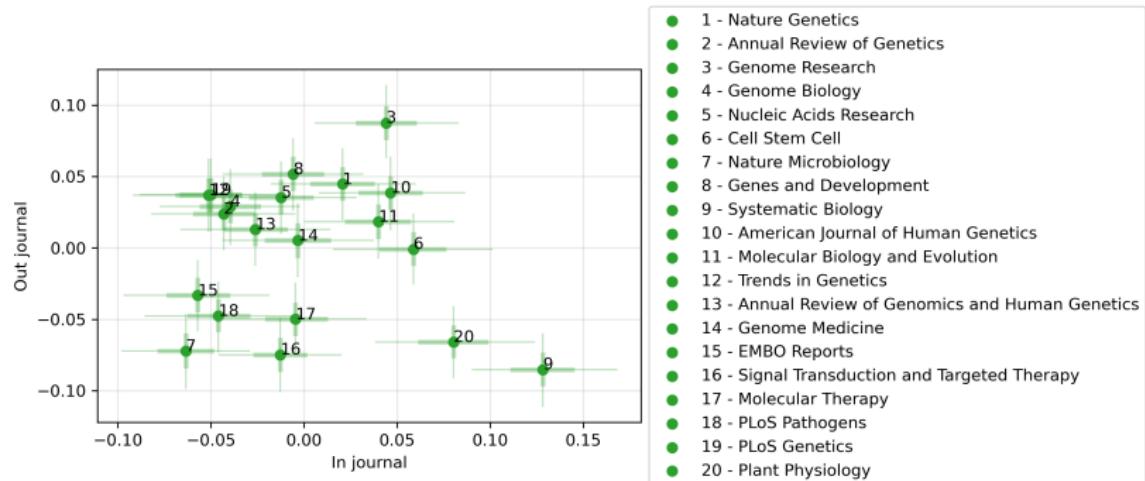
Journal intercepts - specialisation



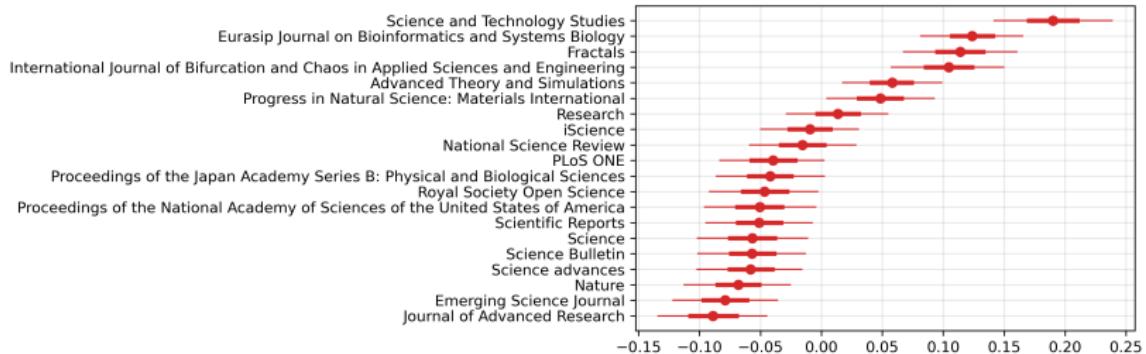
Journal intercepts - centrality



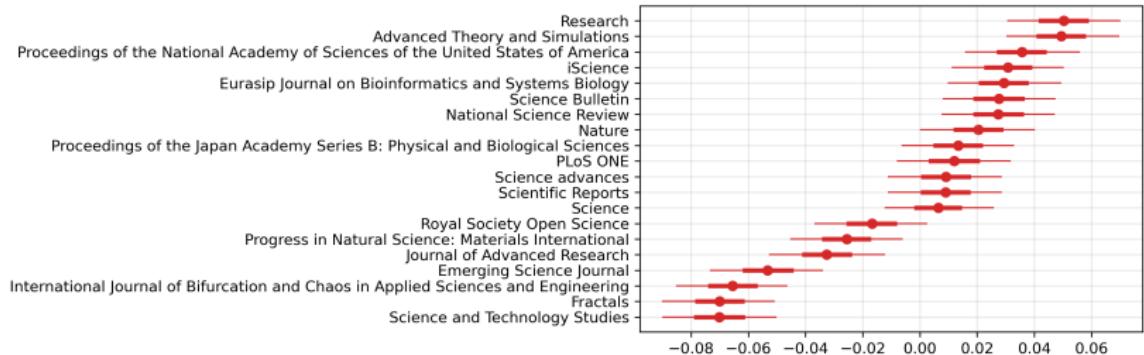
Journal intercepts



Journal intercepts - specialisation



Journal intercepts - centrality



Journal intercepts

