

Performance of linking graduates to researchers

Flavio & Christoph

11 September, 2023

Contents

Overview	1
SQL example for sourcing number of authors with same name	1
Which linking iterations to keep?	1
Some histograms	3
link score by field	3
Year between first pub and graduation	4
First and last name matches by cohort and field	6
How do fields of ProQuest map into fields in MAG?	8
Fraction matched by year and field	28
Checking non-linked entities that should be a link	29

This document compares the links we obtain for all fields in the latest iteration. But it does not consider the further processing done in `prep_linked_data.py`. For better information about the final linked sample, see `quality_linking_graduates_chemistry.Rmd`.

Overview

SQL example for sourcing number of authors with same name

```
select *
from author_sample
inner join (
  select authorid, normalizedname, papercount, citationcount
  from authors
  where normalizedname = "lawrence b slobodkin"
) using (authorid)
inner join (
  select authorid, fieldofstudyid
  from author_fields
  where fieldclass = "first"
) using (authorid)
```

Which linking iterations to keep?

```
keep_iter_ids_base <- linking_info %>%
  filter(date <= date_method_change
         & keywords == "False"
  )
```

```

keep_iter_ids_revise <- linking_info %>%
  filter(date > date_method_change
        & keywords == "True"
        ) %>%
  # keep only the latest iteration here
  group_by(field) %>%
  filter(iteration_id == max(iteration_id)) %>%
  ungroup()
stopifnot(nrow(keep_iter_ids_revise) == n_distinct(keep_iter_ids_revise$field))

```

```

keep_iter_ids <- list(
  base = keep_iter_ids_base,
  revise = keep_iter_ids_revise
)

```

```

keep_iter_ids <- map(
  .x = keep_iter_ids,
  .f = ~.x %>%
    filter(field %in% select_fields) %>%
    pull(iteration_id)
)

```

```

linked_ids <- map(
  .x = keep_iter_ids,
  .f = ~linked_ids %>%
    filter(iteration_id %in% .x)
)

```

```

d_links <- map(
  .x = linked_ids,
  .f = ~.x %>%
    left_join(mag_authors %>%
      select(AuthorId,
             year_mag = year,
             firstname_mag = firstname,
             lastname_mag = lastname,
             field_mag = fieldofstudy,
             field0_mag = mag_field0),
      by = "AuthorId") %>%
    left_join(pq_authors %>%
      select(goid,
             year_pq = year,
             firstname_pq = firstname,
             lastname_pq = lastname,
             field_pq = fieldofstudy,
             field0_pq = mag_field0),
      by = "goid") %>%
    mutate(year_diff = year_mag - year_pq,
           same_firstname = ifelse(firstname_mag == firstname_pq, 1, 0),
           same_lastname = ifelse(lastname_mag == lastname_pq, 1, 0)) %>%
    left_join(field_names_id %>%
      rename(main_field = NormalizedName),
      by = c("field0_pq" = "FieldOfStudyId")) %>%
    filter(goid != 305107842) %>% # this is some author which was linked but should not have been in

```

```

    filter(link_score > min_link_score
           & abs(year_diff) <= max_year_diff)

)

d_links$base <- d_links$base %>% filter(year_pq <= 2005)

```

Some histograms

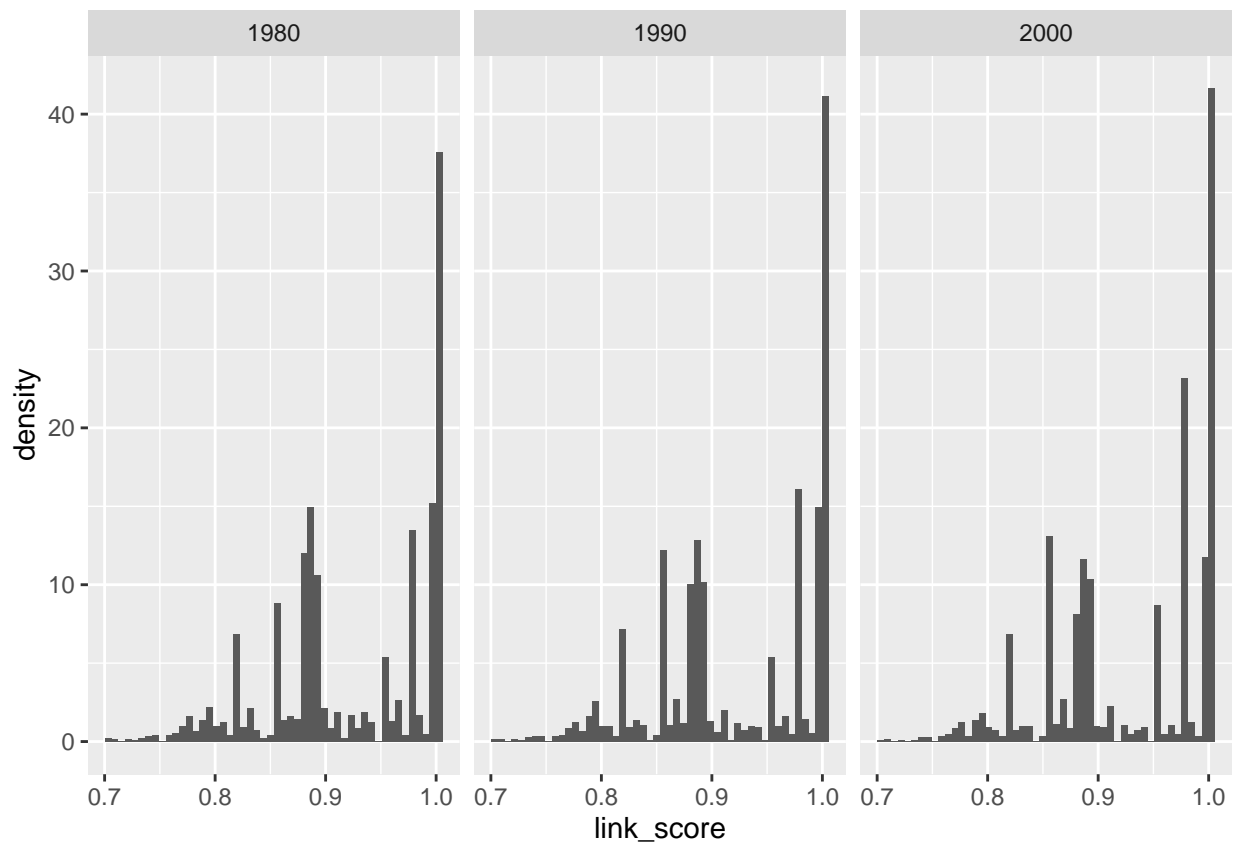
link score by field

```
## $base
```

```

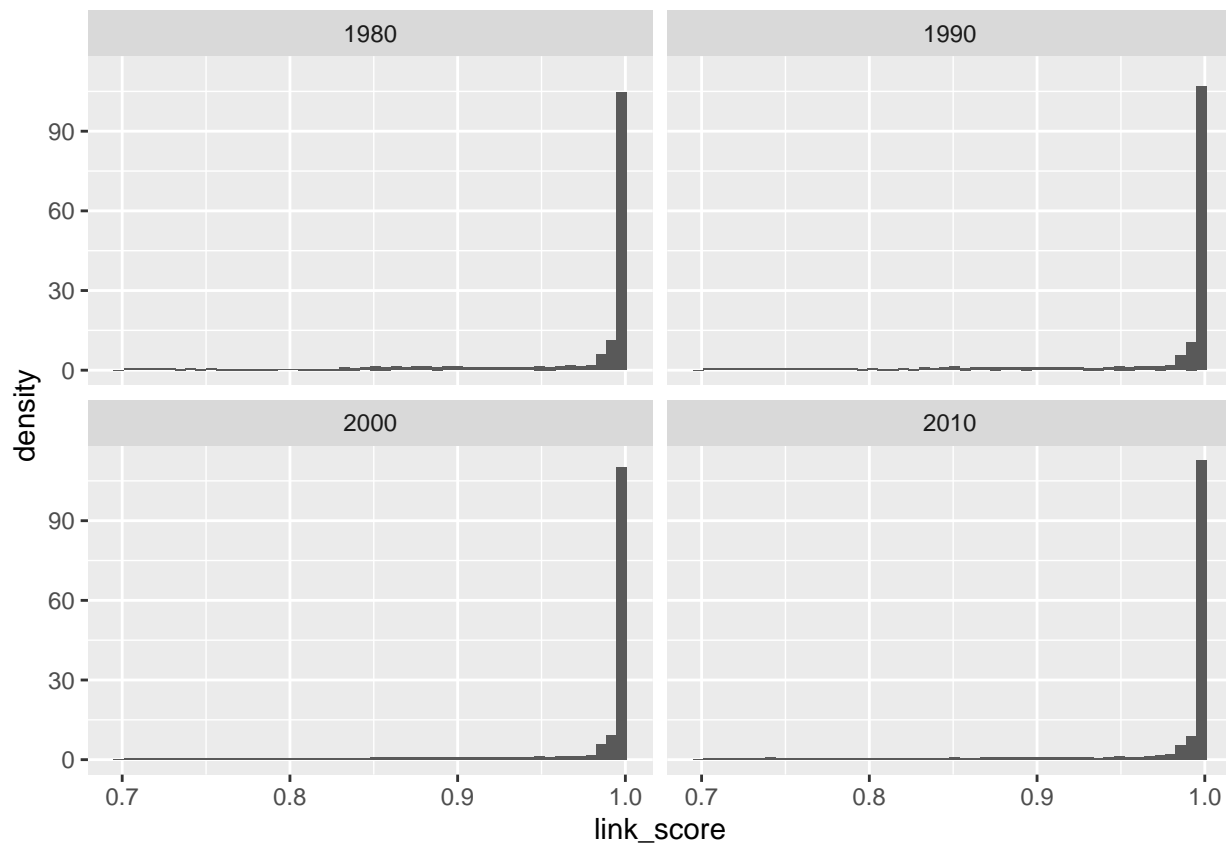
## Warning: The dot-dot notation (`..density..`) was deprecated in ggplot2 3.4.0.
## i Please use `after_stat(density)` instead.
## This warning is displayed once every 8 hours.
## Call `lifecycle::last_lifecycle_warnings()` to see where this warning was
## generated.

```



```
##
```

```
## $revise
```

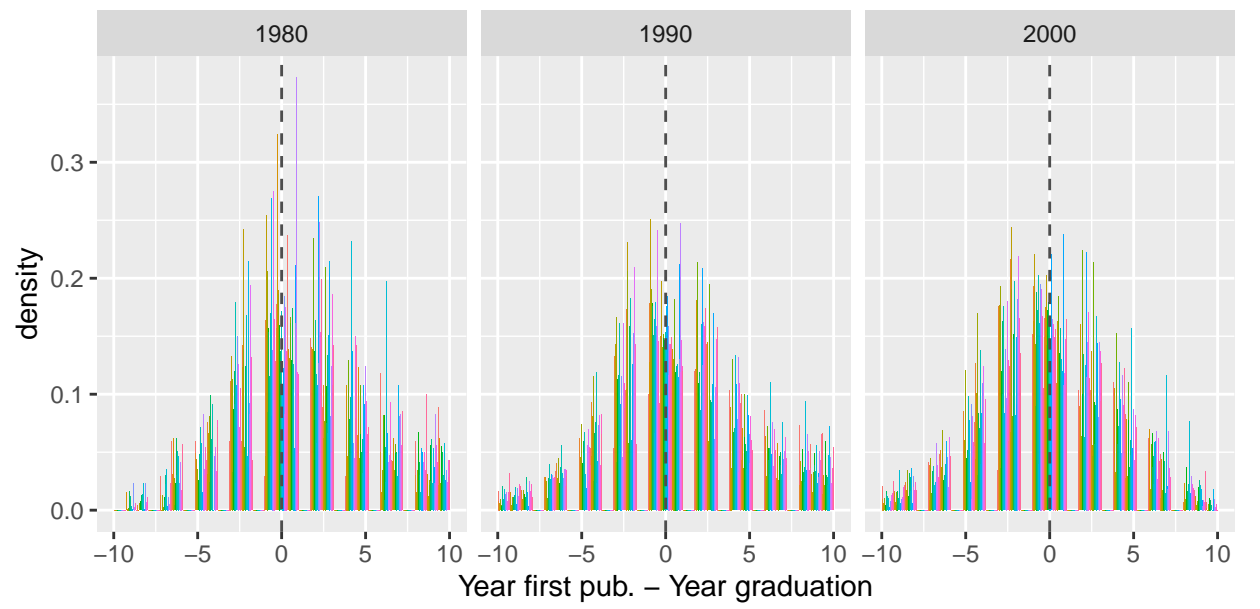


Year between first pub and graduation

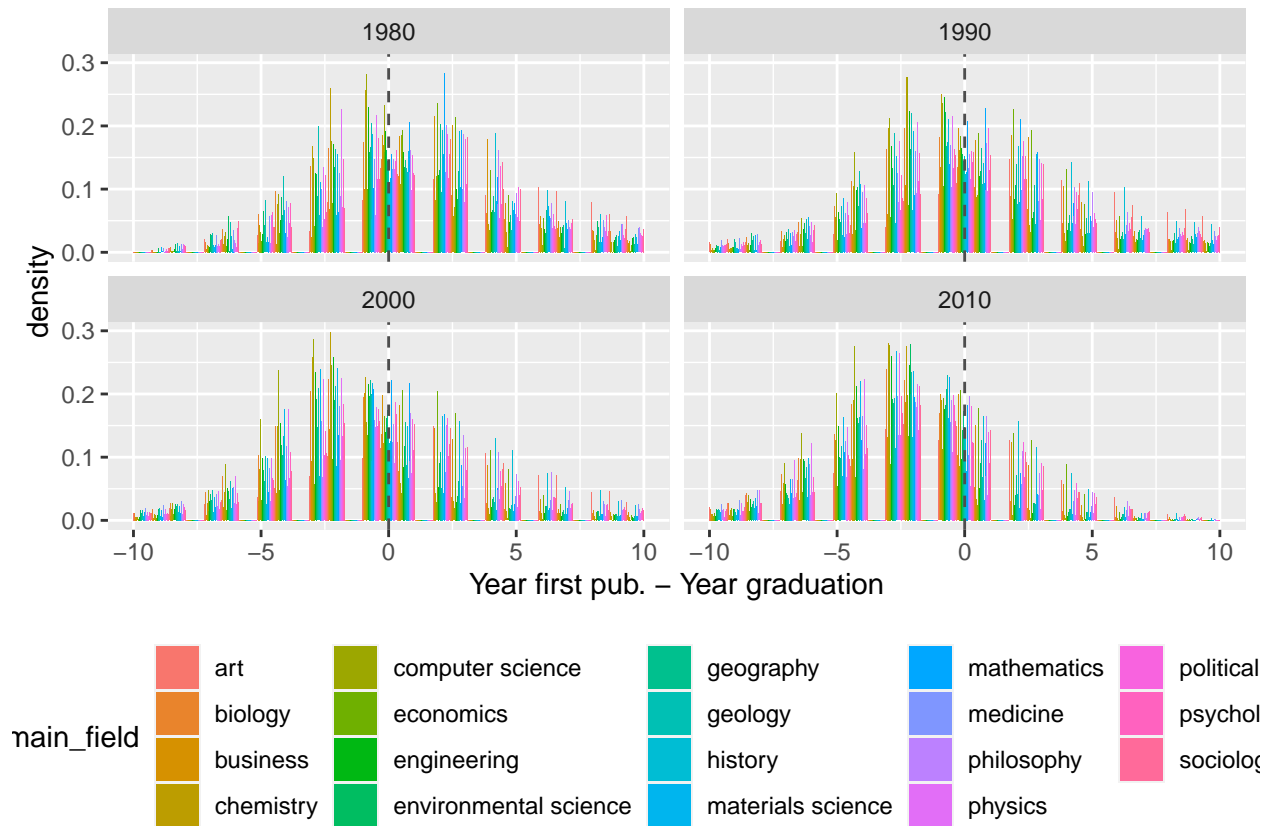
- why are there other fields than maths/biology for the following two figures?
- this is because we sample persons whenever they are in any of the linking fields
 - thus, a graduate can be linked in a biology iteration if her first field is chemistry
 - compare this with the advisor links!
 - this also means the join above should take care of this, and indicate the multiplicity of the graduates!

```
## $base
```

```
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```

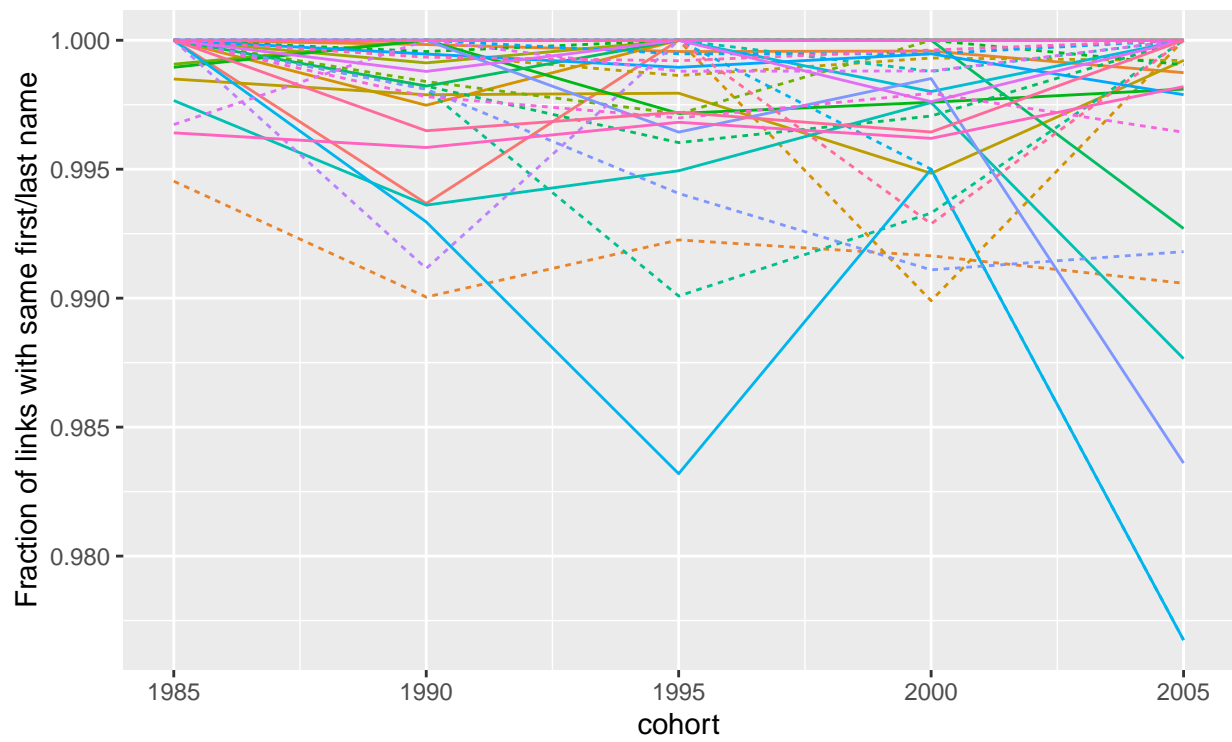


```
##
## $revise
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```



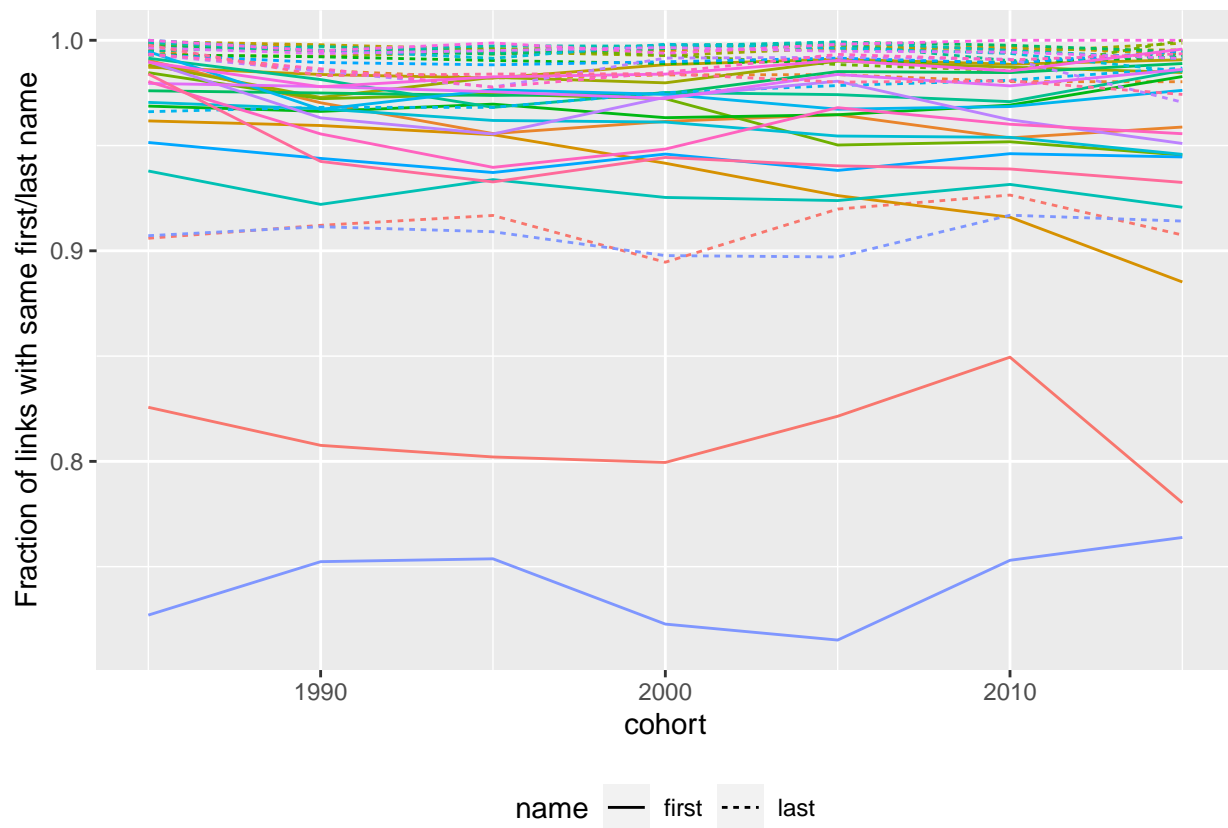
First and last name matches by cohort and field

\$base



name — first ---- last

\$revise

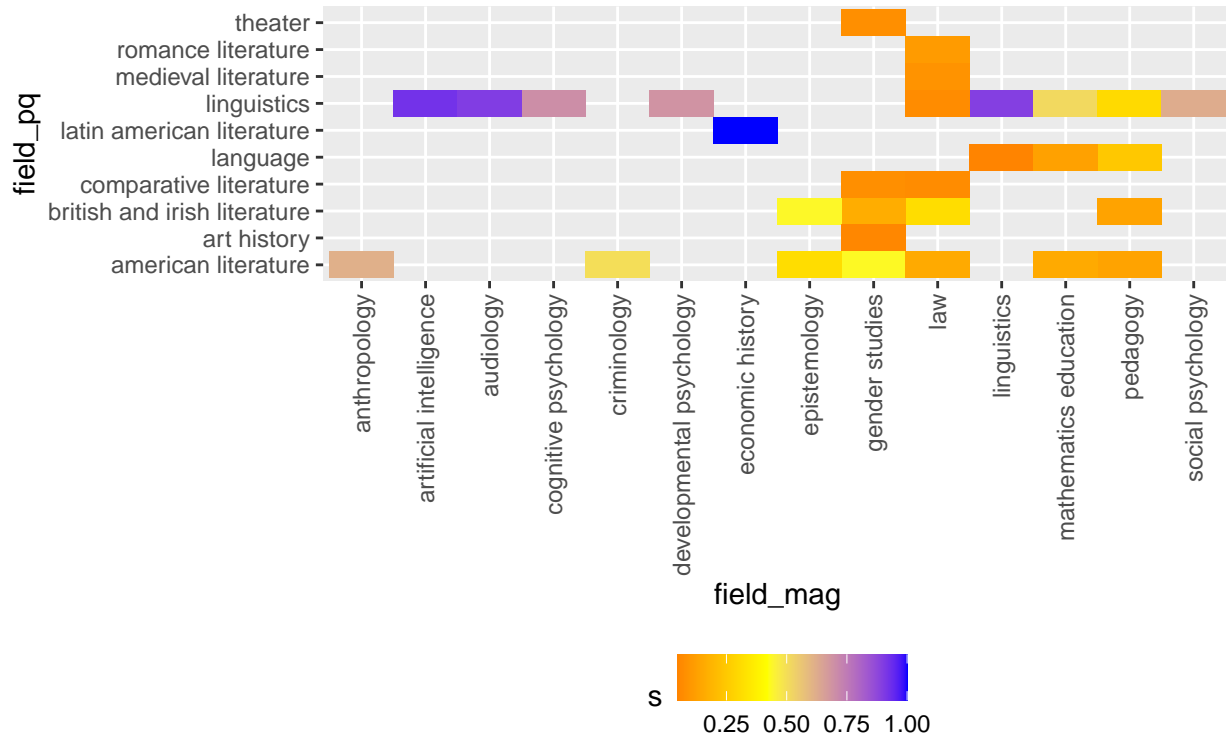


How do fields of ProQuest map into fields in MAG?

[[1]]

Fraction of field ProQuest into field MAG

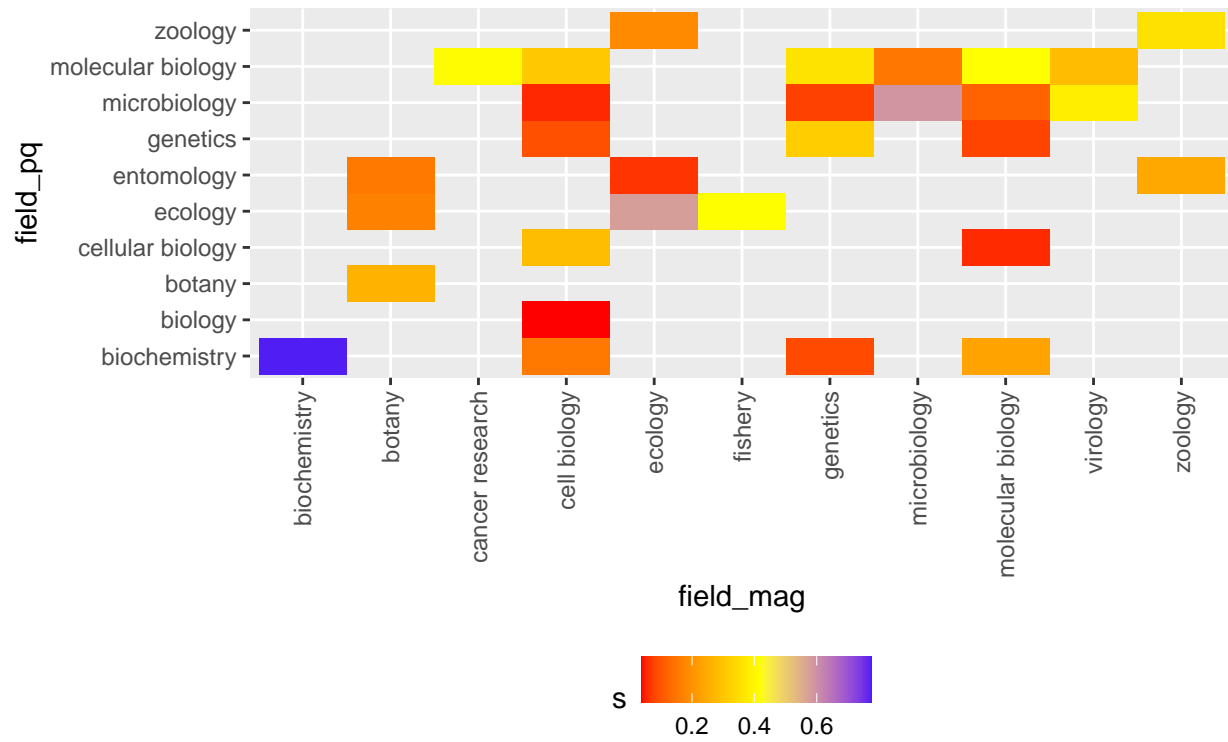
Field: art



[[2]]

Fraction of field ProQuest into field MAG

Field: biology

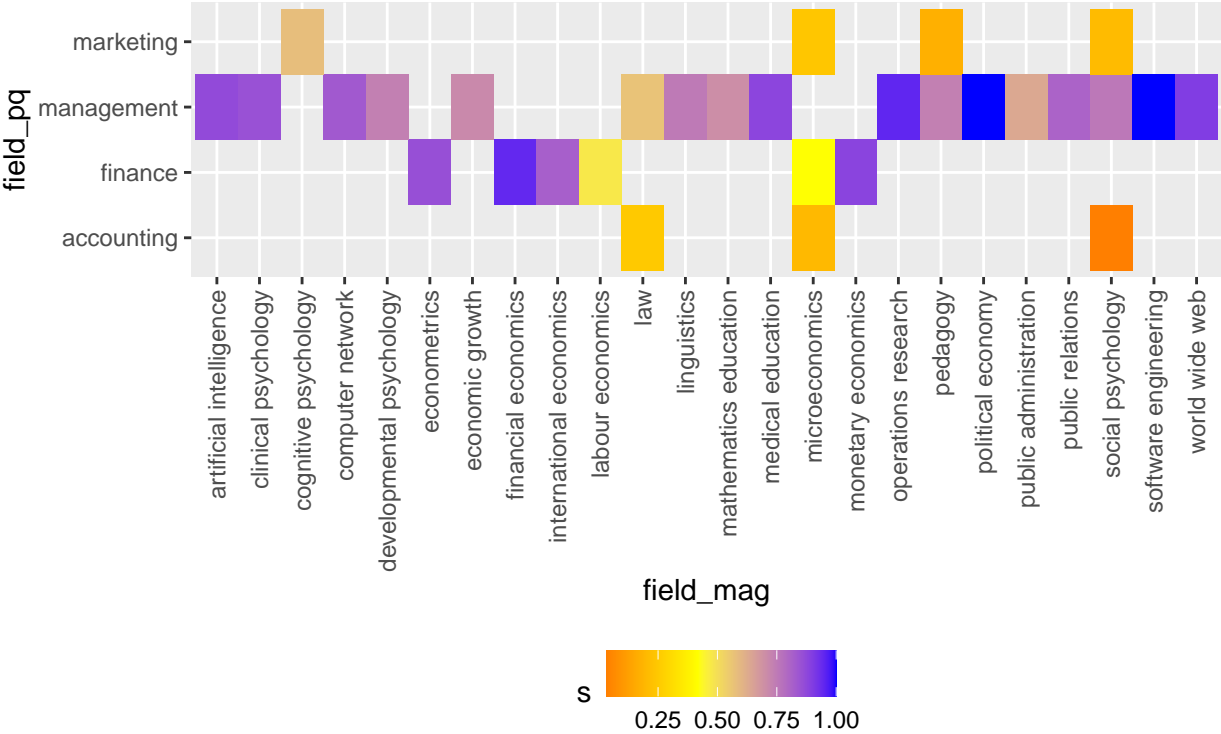


##

[[3]]

Fraction of field ProQuest into field MAG

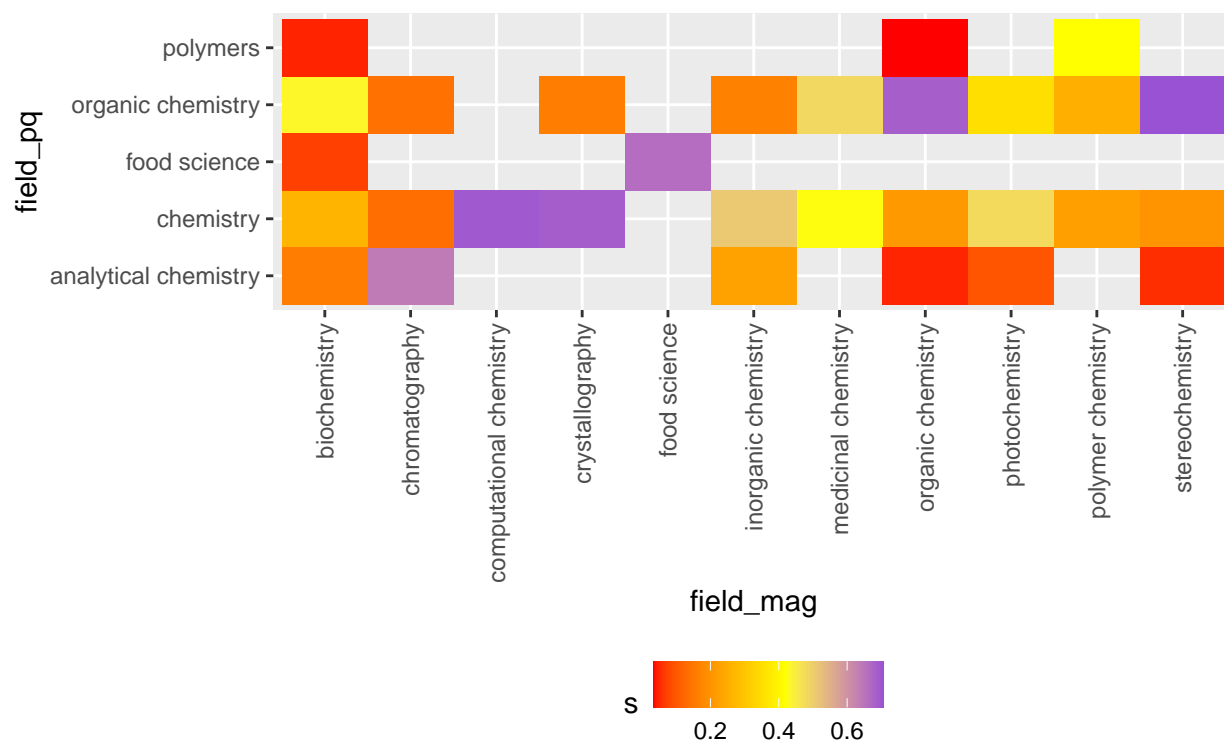
Field: business



[[4]]

Fraction of field ProQuest into field MAG

Field: chemistry

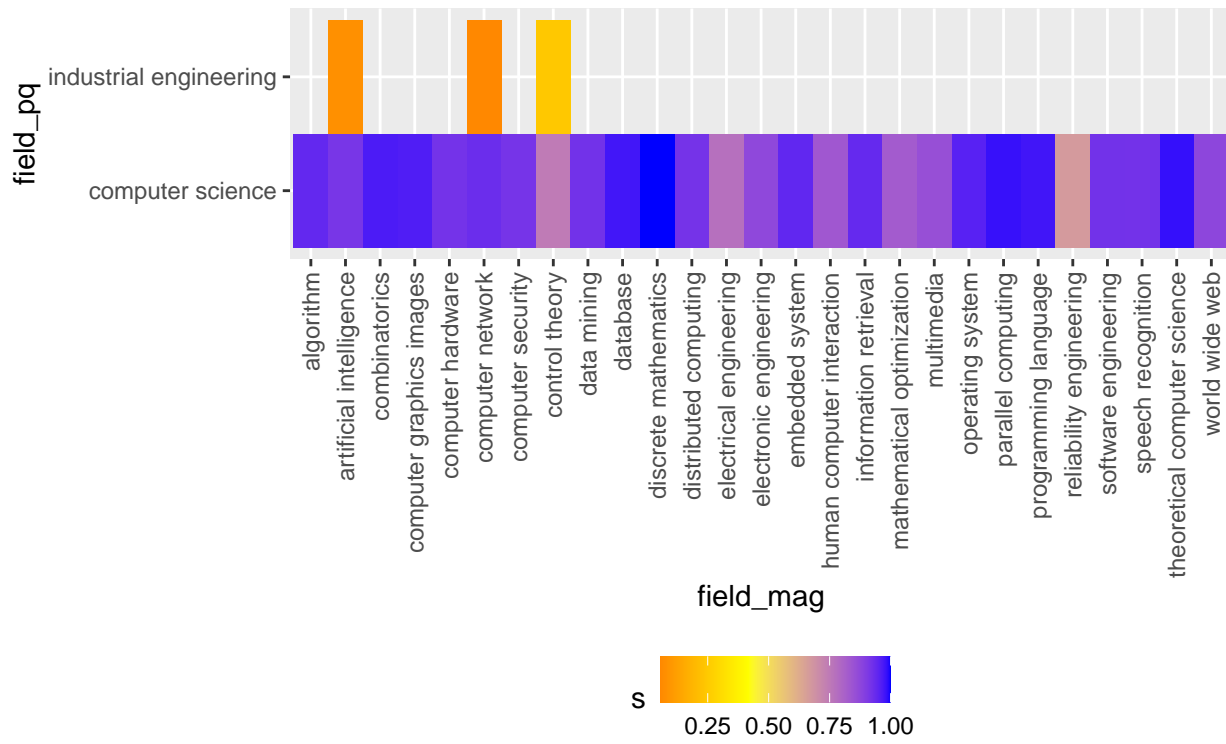


##

[[5]]

Fraction of field ProQuest into field MAG

Field: computer science

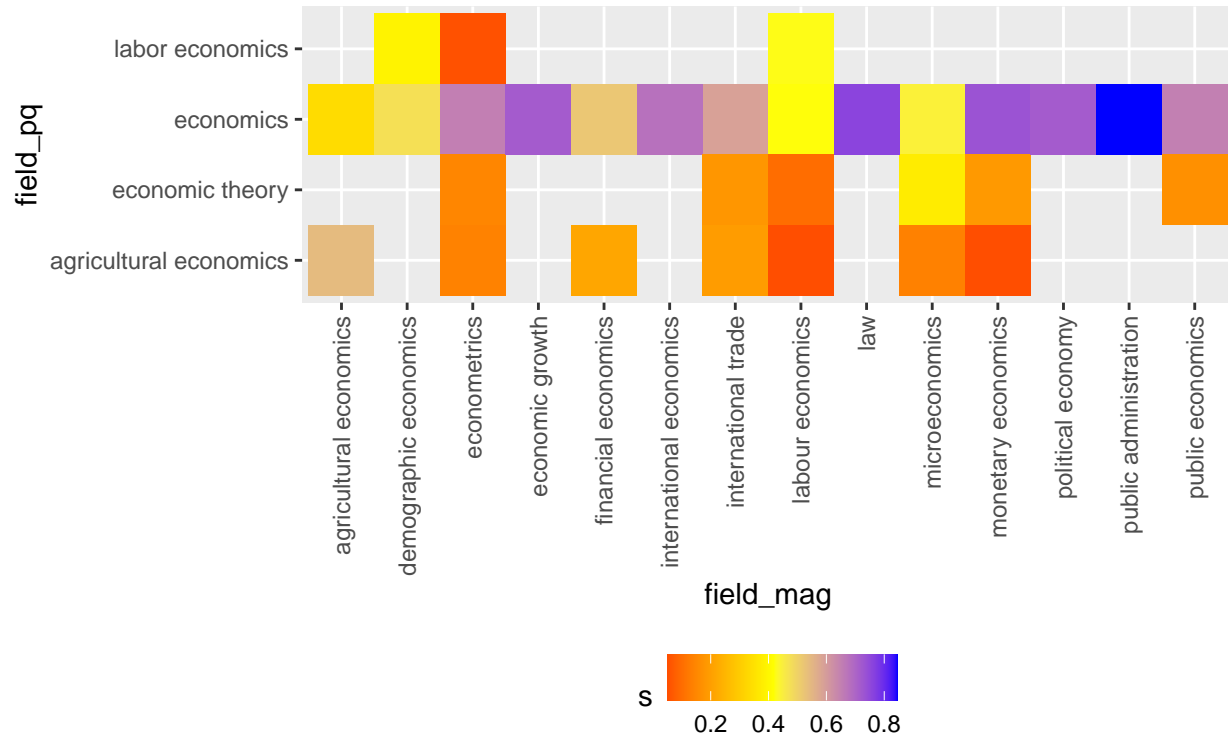


##

[[6]]

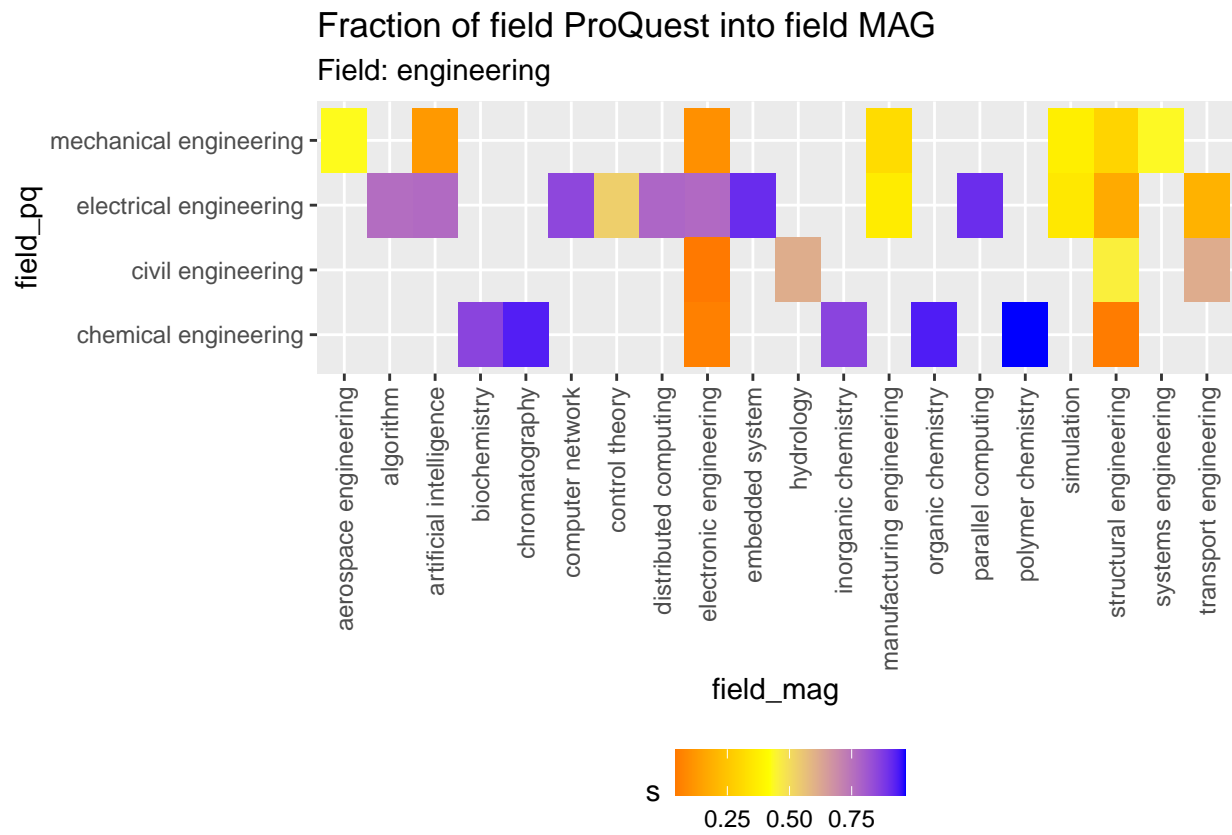
Fraction of field ProQuest into field MAG

Field: economics



##

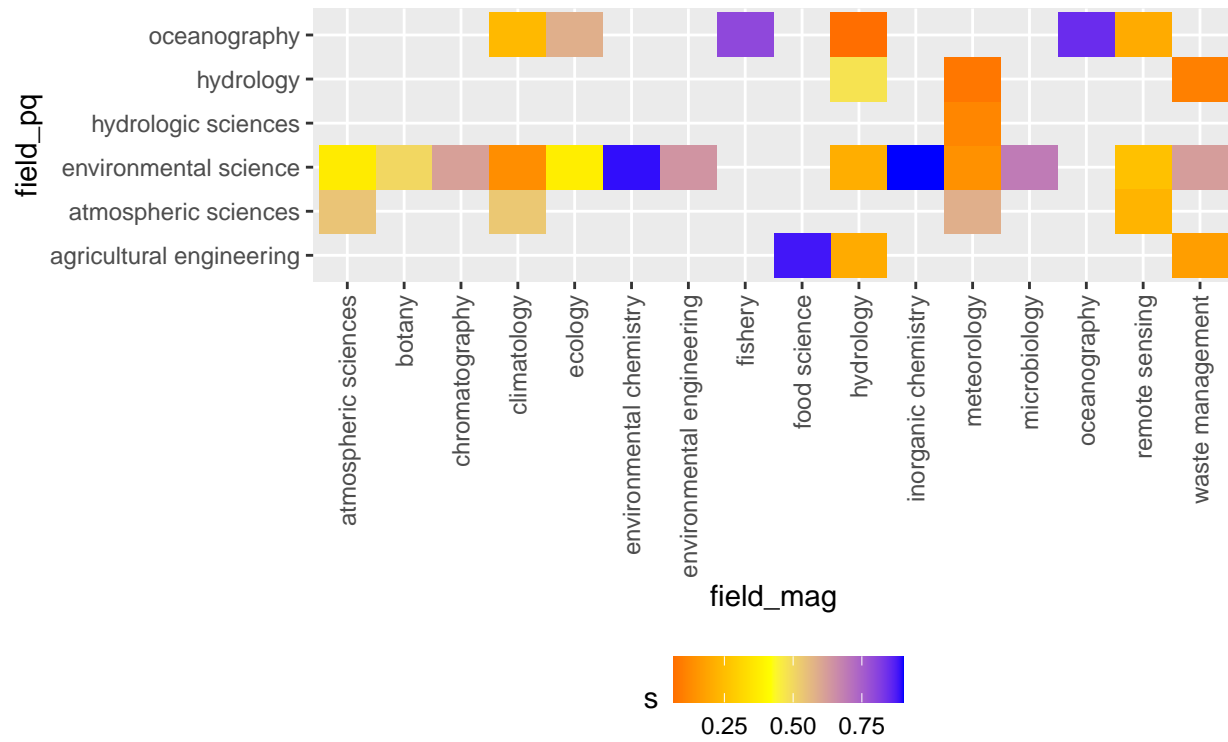
[[7]]



[[8]]

Fraction of field ProQuest into field MAG

Field: environmental science

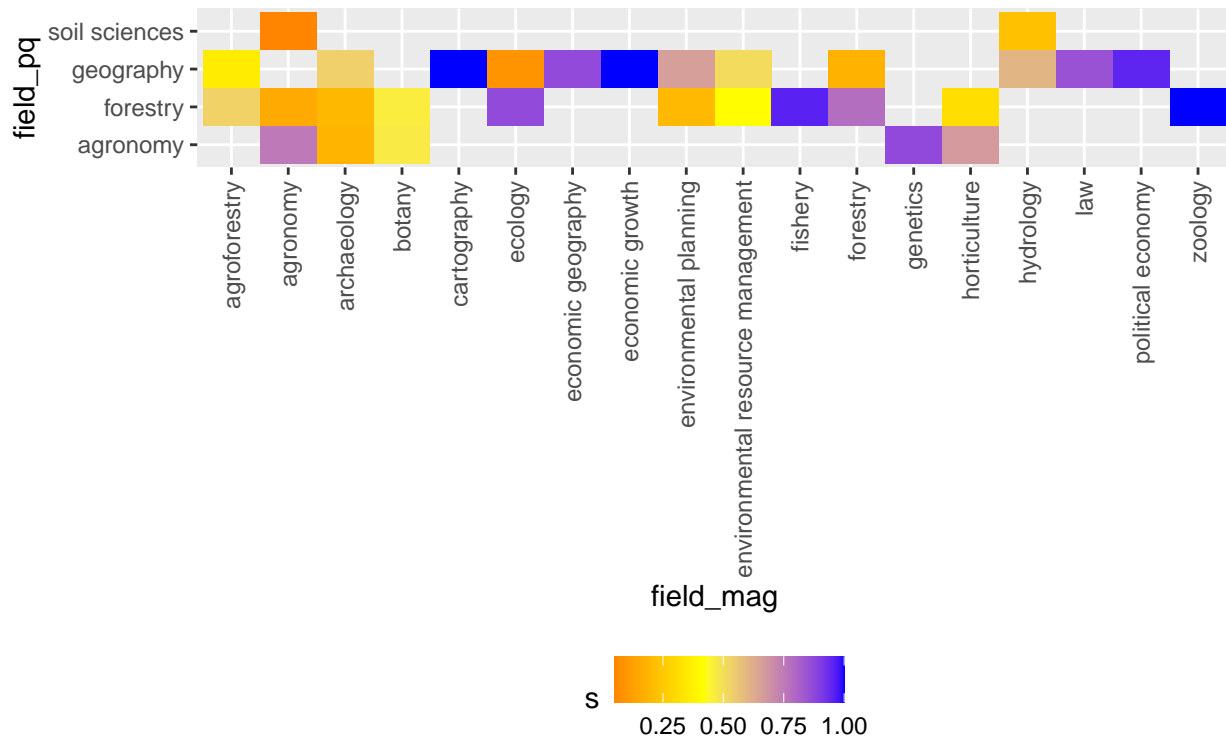


##

[[9]]

Fraction of field ProQuest into field MAG

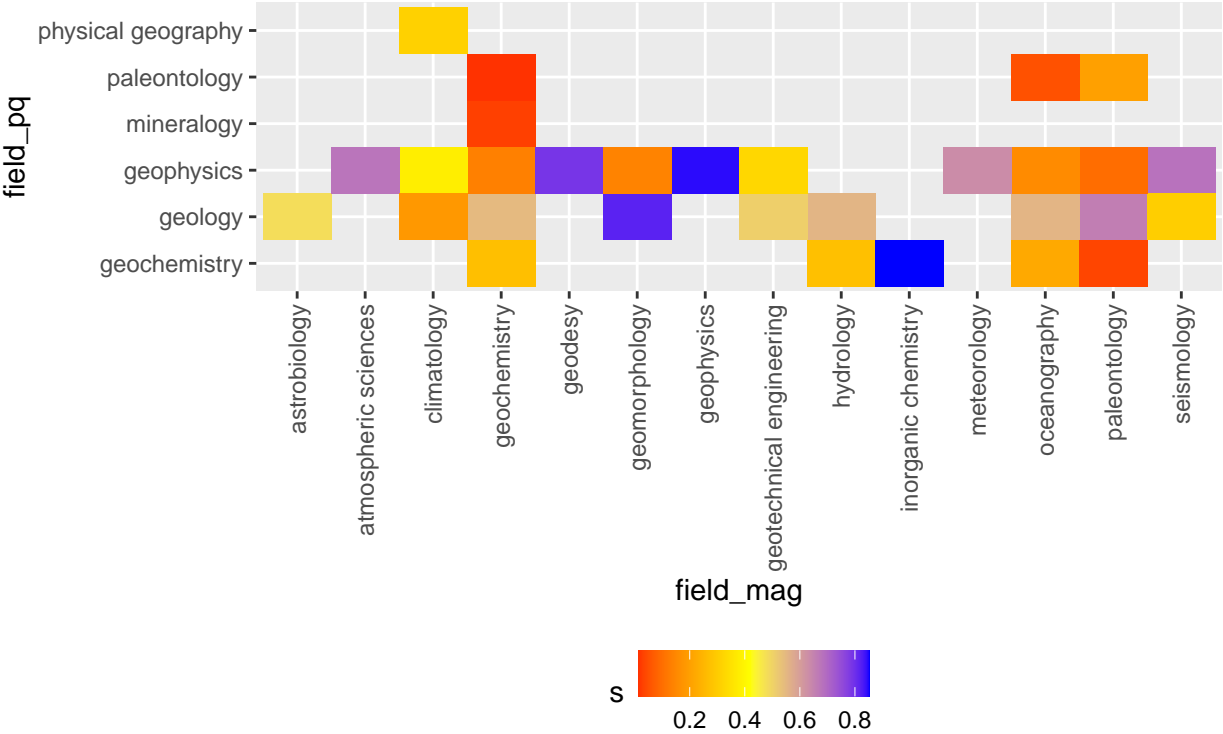
Field: geography



[[10]]

Fraction of field ProQuest into field MAG

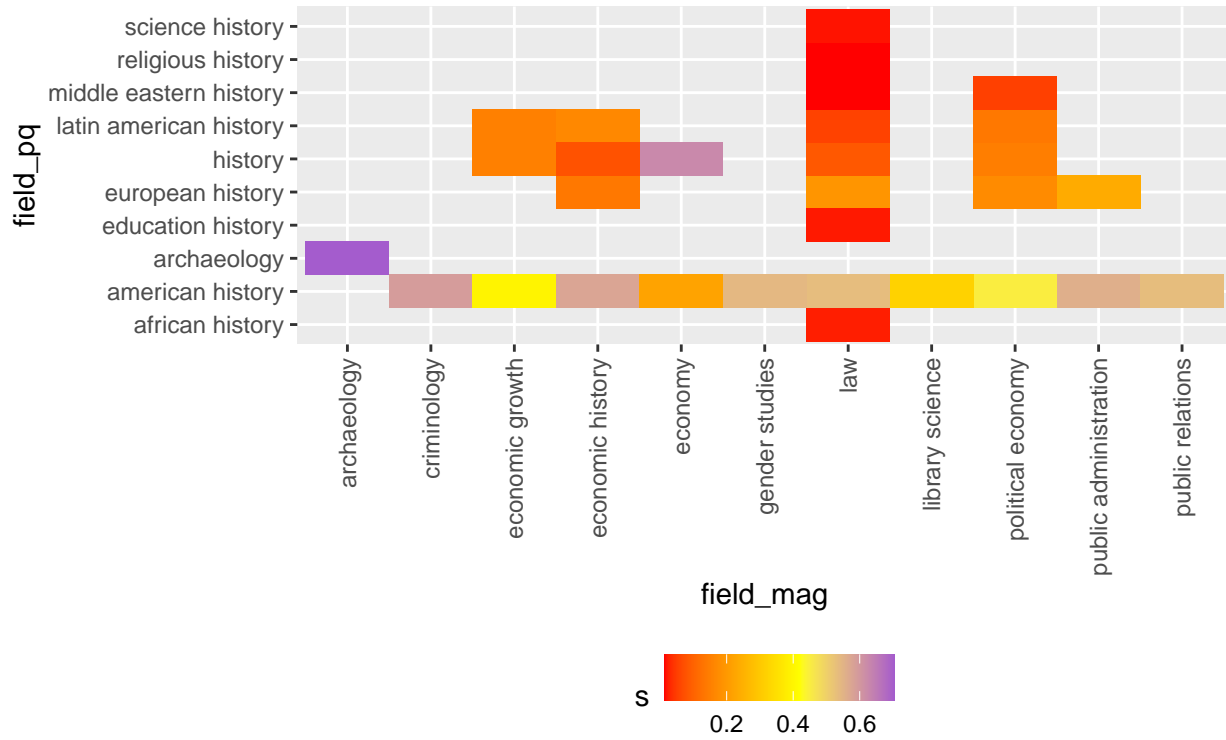
Field: geology



[[11]]

Fraction of field ProQuest into field MAG

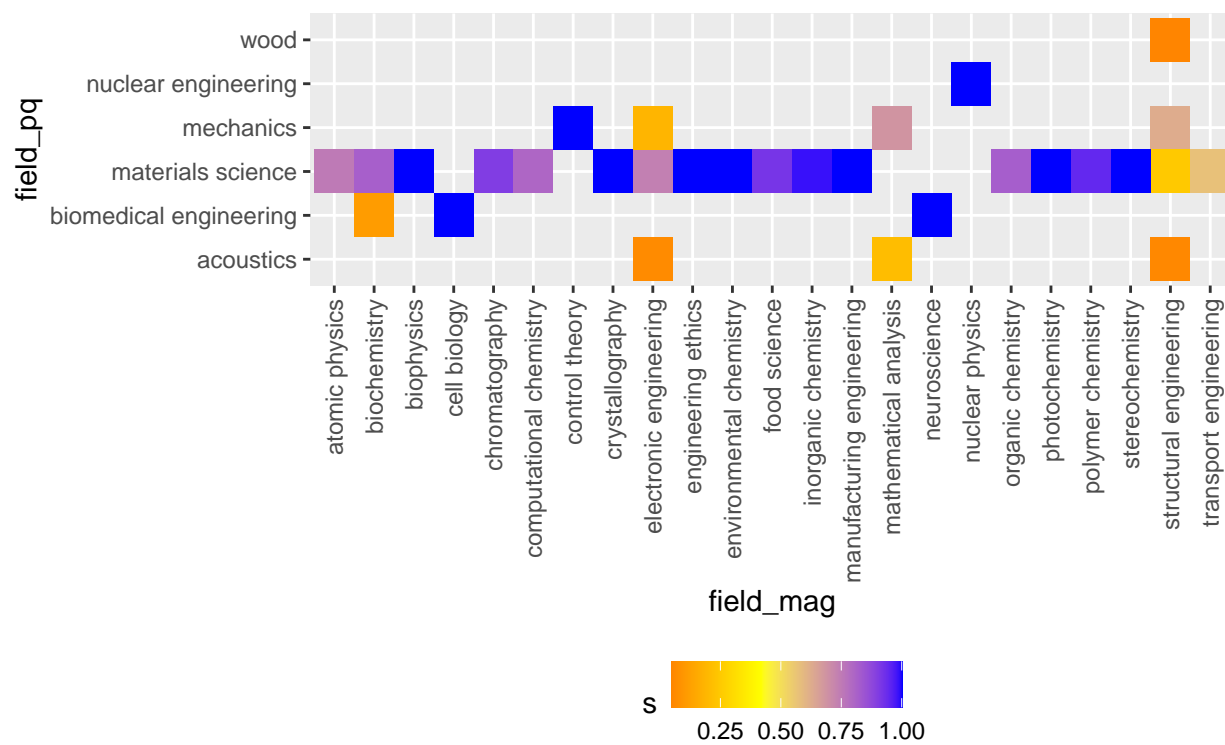
Field: history



[[12]]

Fraction of field ProQuest into field MAG

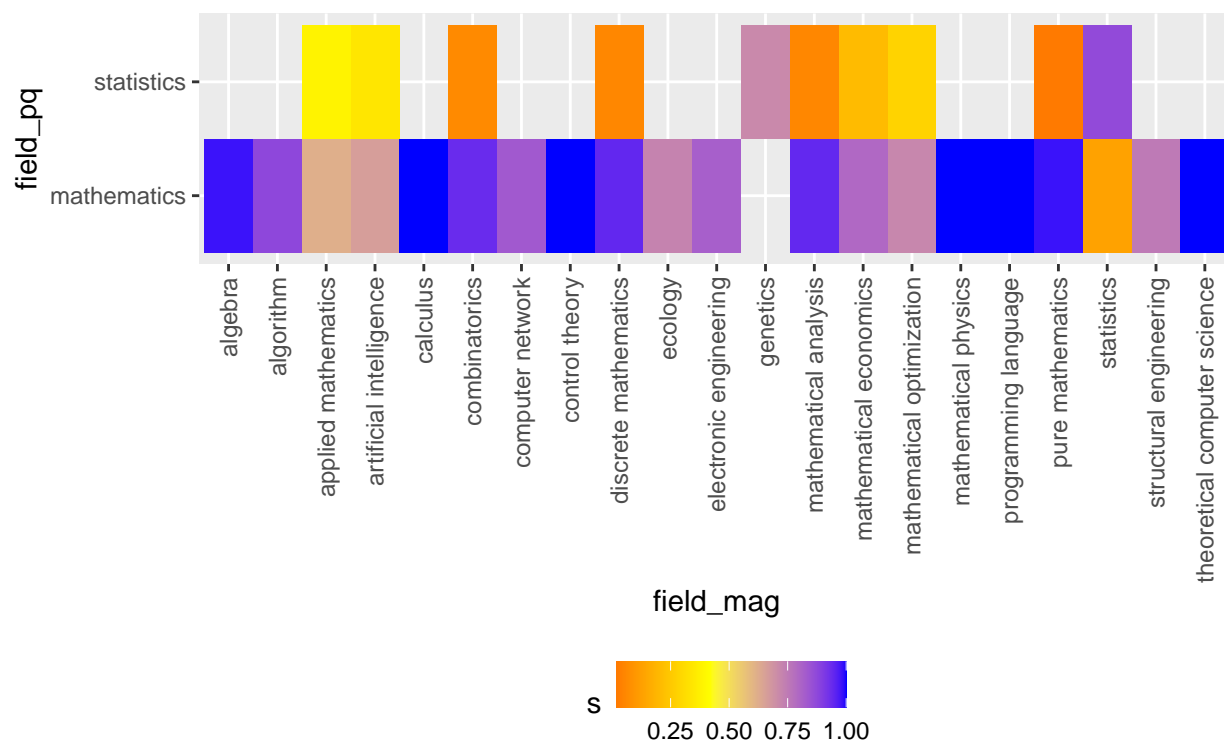
Field: materials science



[[13]]

Fraction of field ProQuest into field MAG

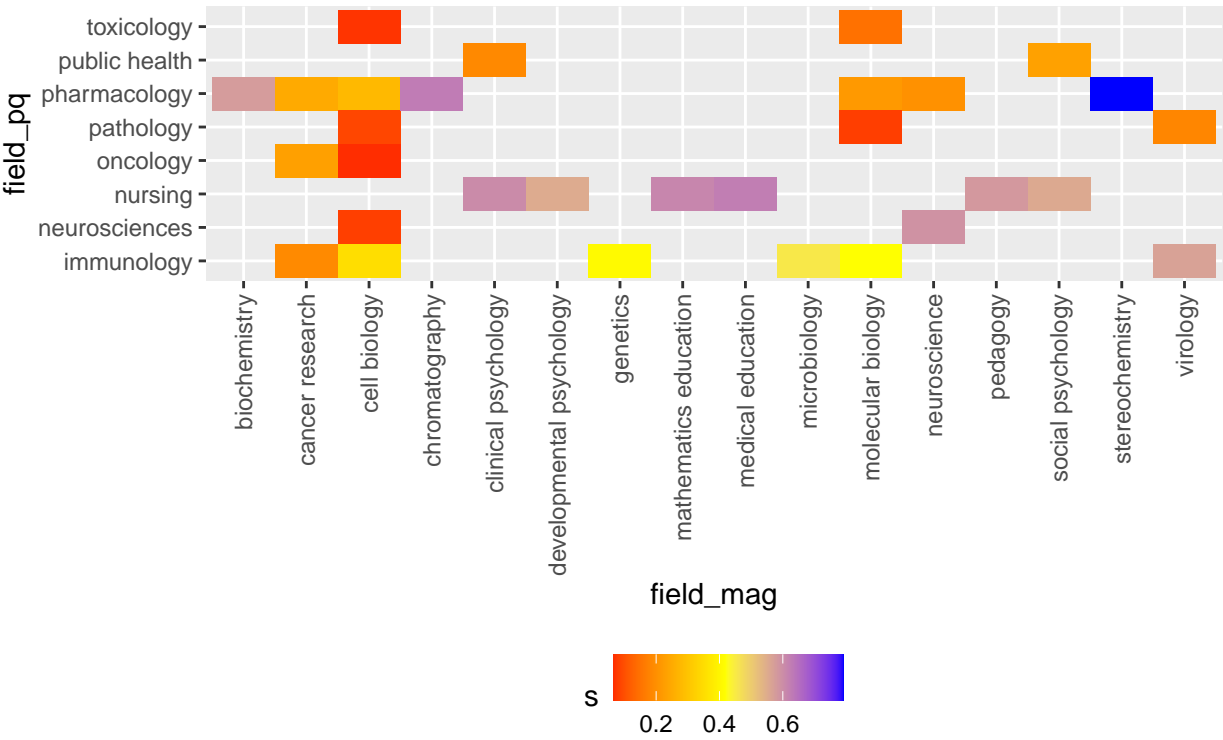
Field: mathematics



[[14]]

Fraction of field ProQuest into field MAG

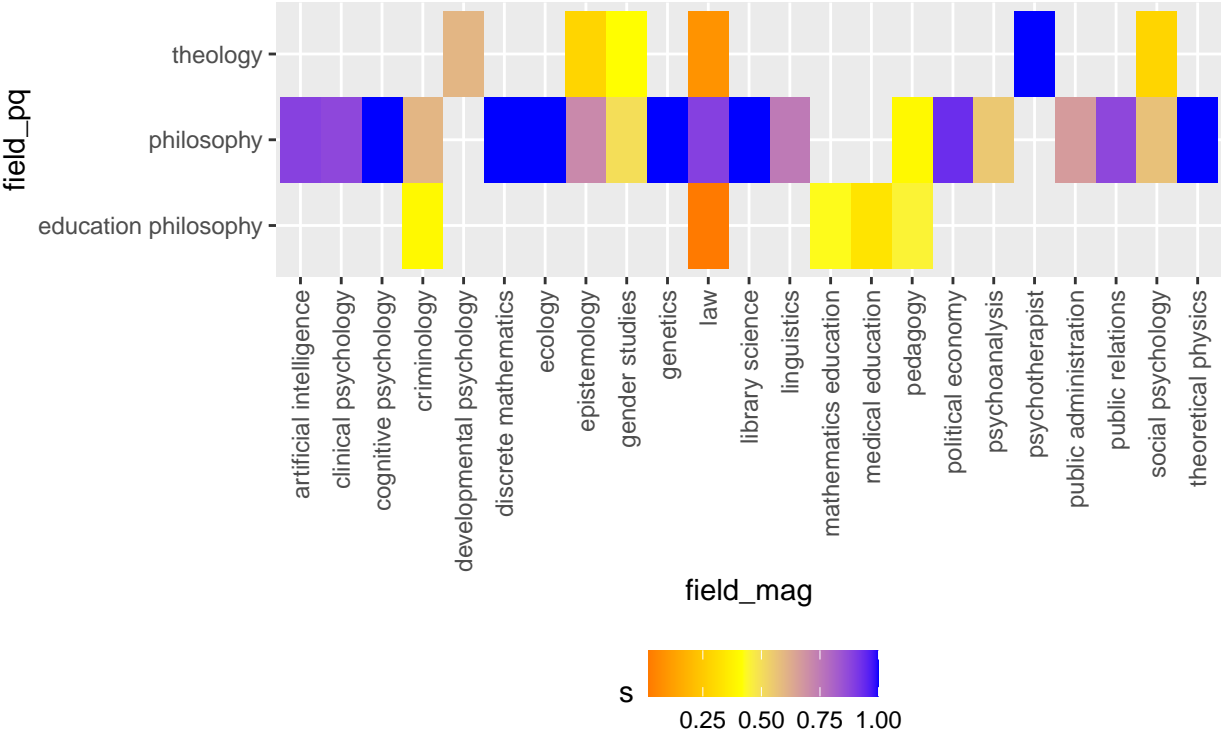
Field: medicine



[[15]]

Fraction of field ProQuest into field MAG

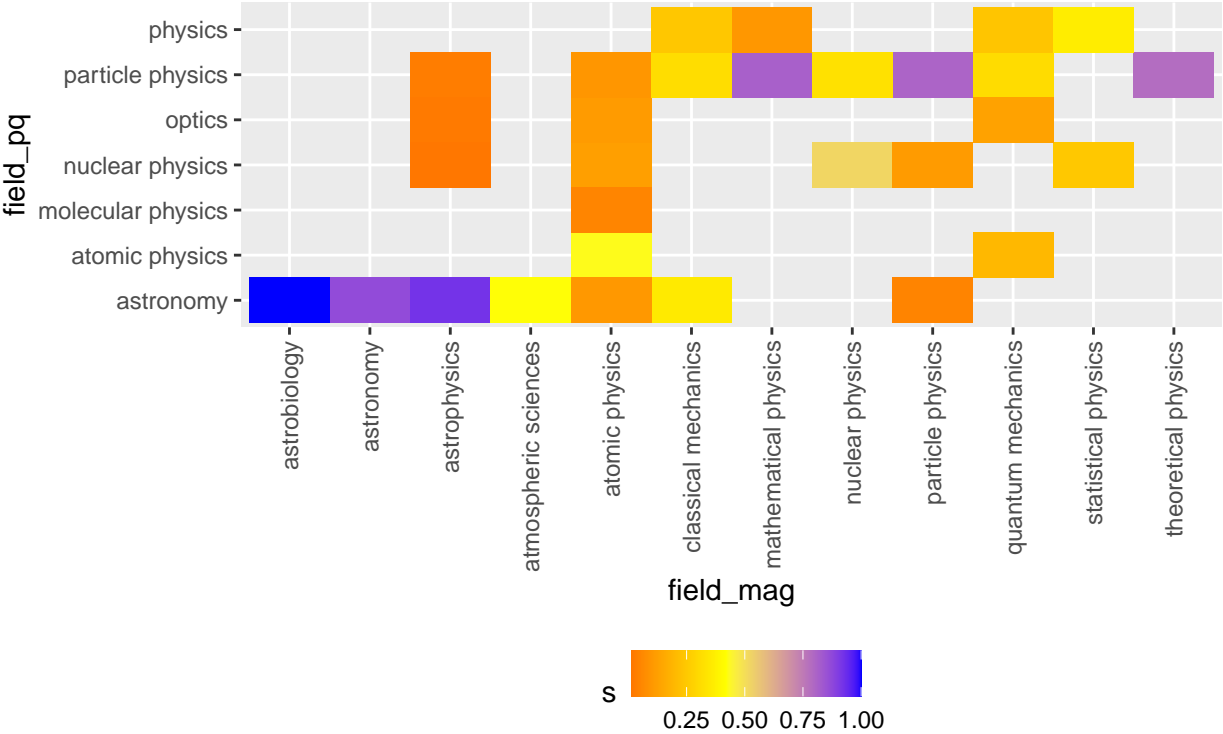
Field: philosophy



[[16]]

Fraction of field ProQuest into field MAG

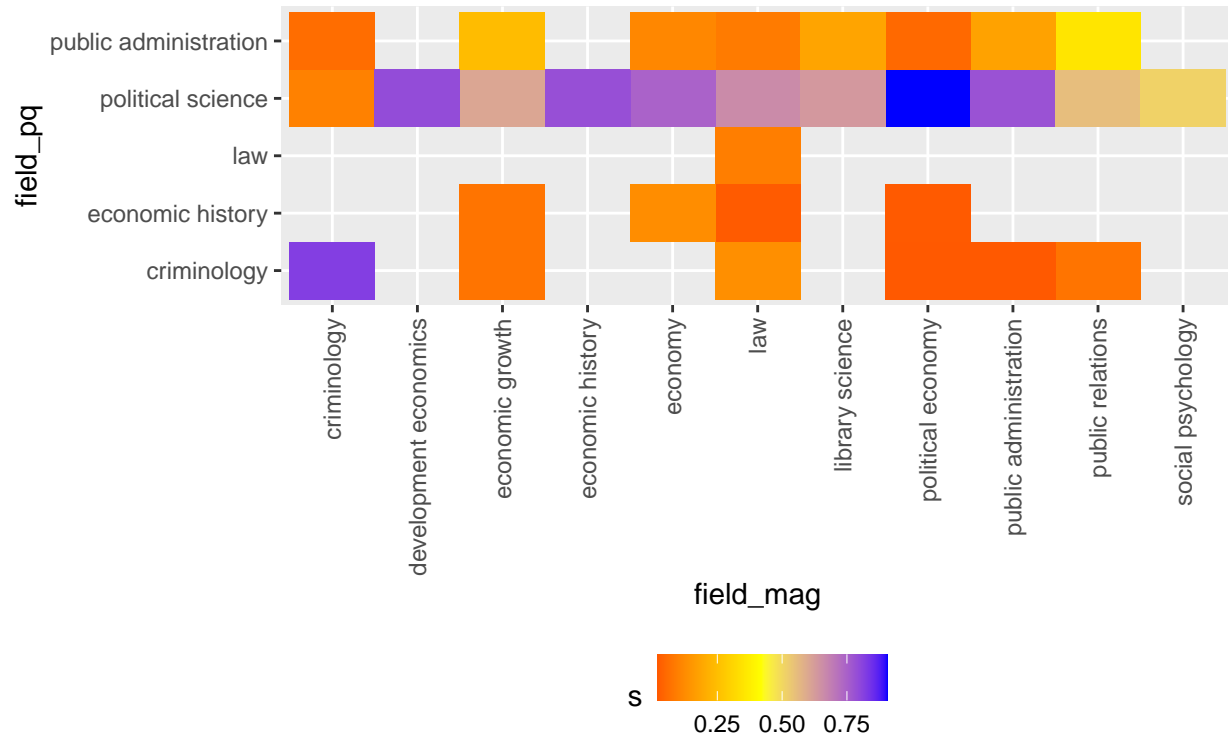
Field: physics



[[17]]

Fraction of field ProQuest into field MAG

Field: political science

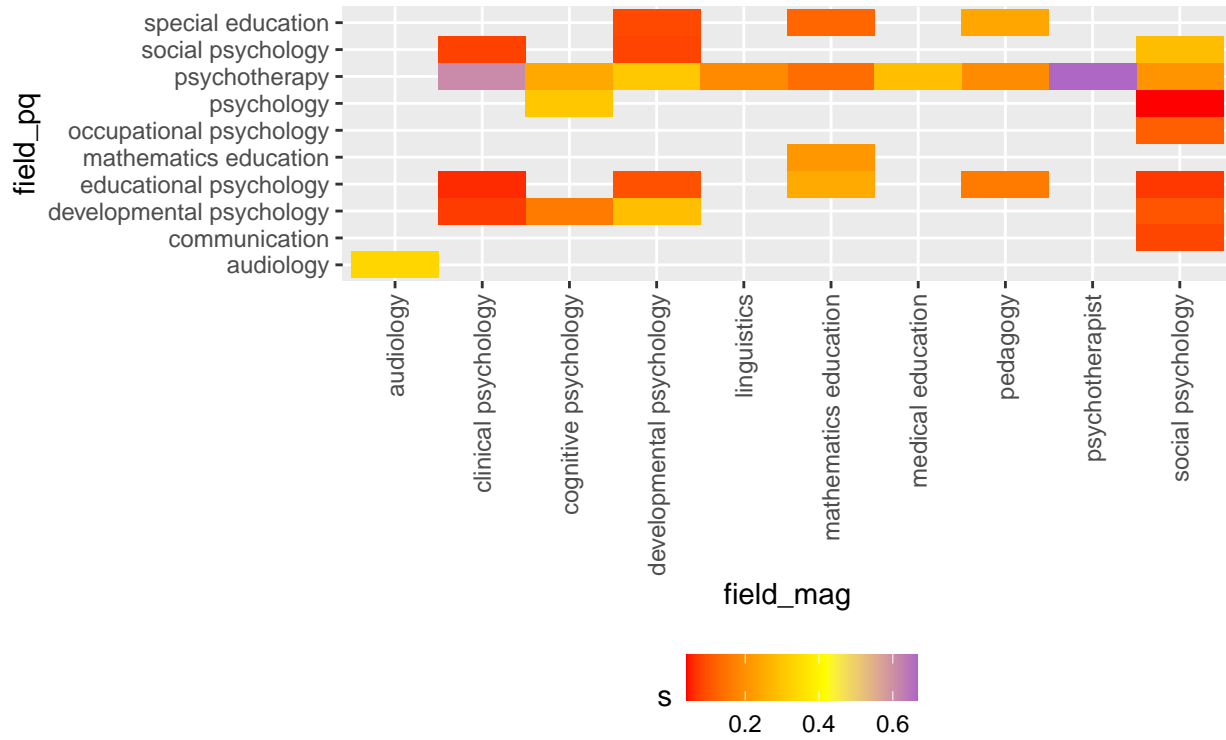


##

[[18]]

Fraction of field ProQuest into field MAG

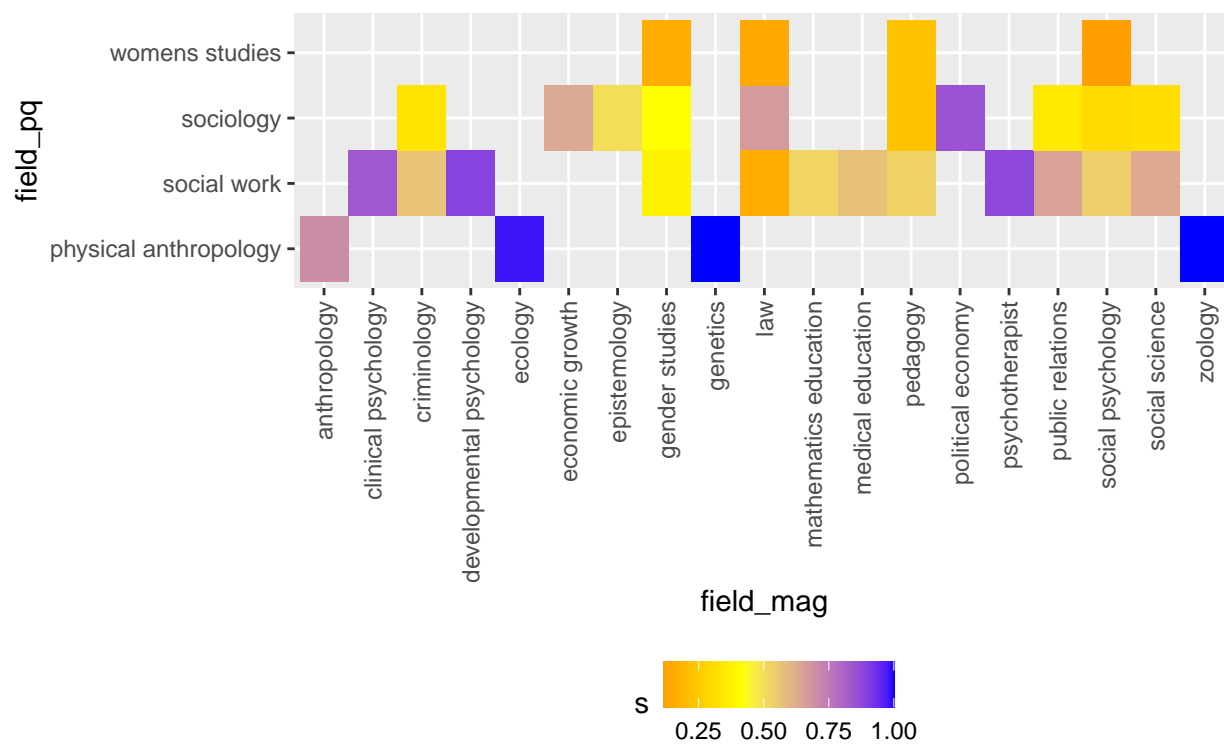
Field: psychology



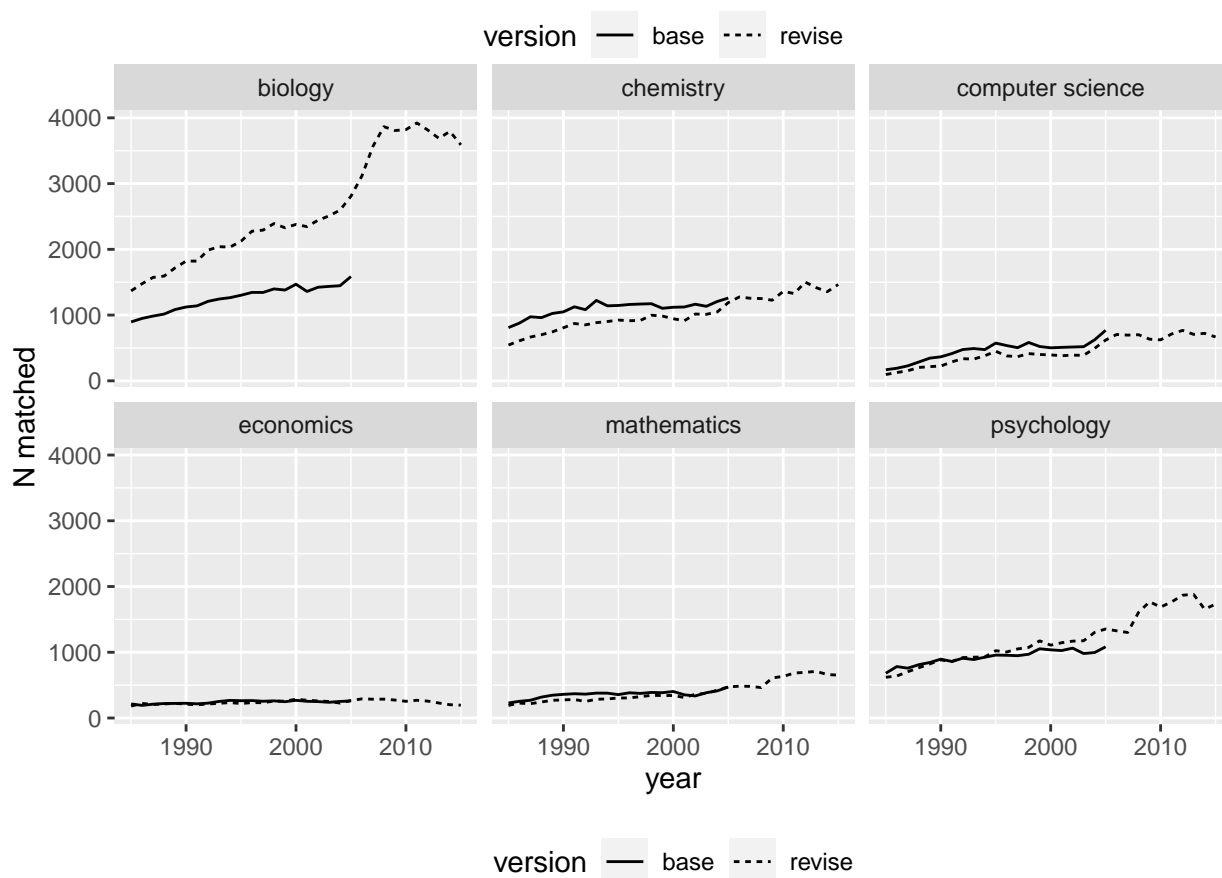
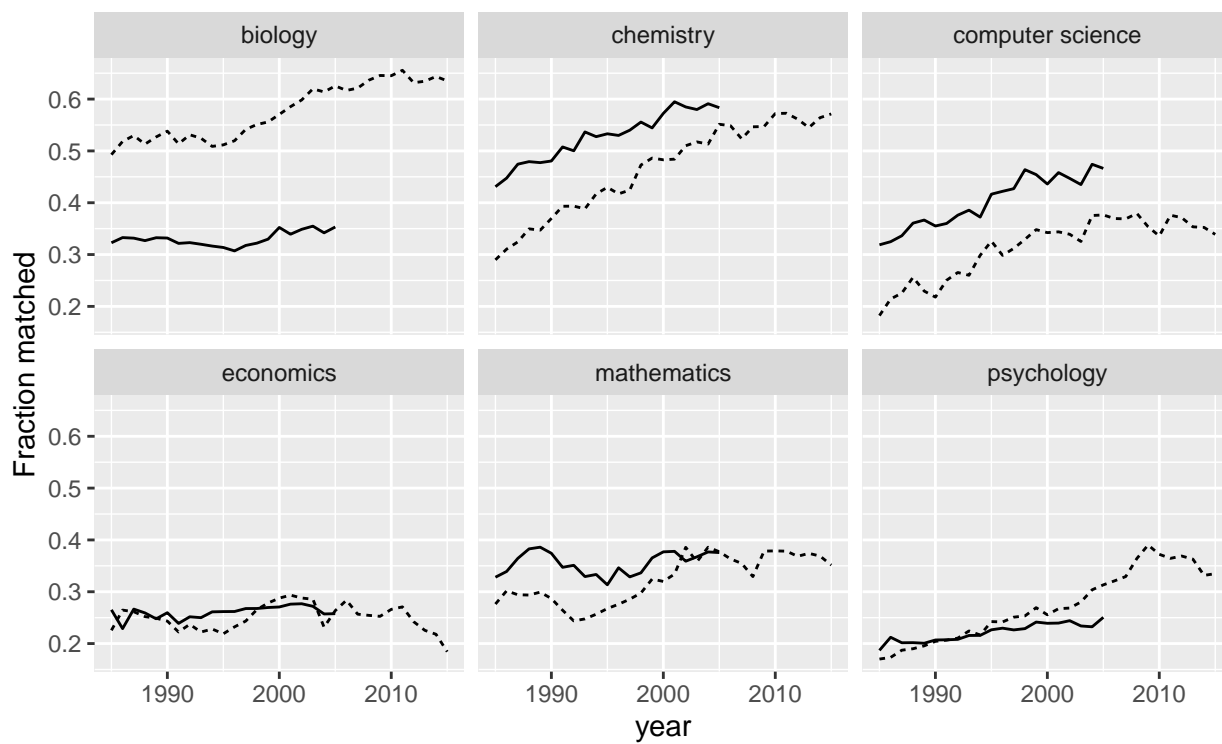
[[19]]

Fraction of field ProQuest into field MAG

Field: sociology



Fraction matched by year and field



Checking non-linked entities that should be a link

```
d_chem <- pq_authors %>%
  left_join(field_names_id %>%
    rename(main_field = NormalizedName),
    by = c("mag_field0" = "FieldOfStudyId")) %>%
  mutate(link = ifelse(goid %in% d_links$revise$goid, "linked", "not linked")) %>%
  filter(main_field == "chemistry")

pq_unis <- tbl(con, "pq_authors") %>%
  left_join(tbl(con, "pq_unis") %>%
    select(university_id, normalizedname),
    by = "university_id") %>%
  select(goid, uni_name = "normalizedname") %>%
  collect()

d_chem <- d_chem %>%
  left_join(pq_unis, by = "goid")

d_chem %>%
  filter(year == 1995 & uni_name == "stanford university" & link == "not linked") %>% head(10)

## # A tibble: 10 x 11
##       goid  year firstname lastname  middlename fieldofstudy mag_field0
##   <int64> <int> <chr>      <chr>      <chr>      <chr>          <int>
## 1 304229925 1995 nancy      hansen      fisher      chemistry      185592680
## 2 304229722 1995 mark      pavlosky    alan        chemistry      185592680
## 3 304228620 1995 kristin   sannes      ann          chemistry      185592680
## 4 304218381 1995 glenn      jones       clark        chemistry      185592680
## 5 304201950 1995 david      offord      alan         chemistry      185592680
## 6 304238172 1995 robert     guettler    david        chemistry      185592680
## 7 304202002 1995 eric       remy        david        chemistry      185592680
## 8 304229882 1995 thomas     schoch      k            chemistry      185592680
## 9 304229838 1995 philip     merrill     bradley      chemistry      185592680
## 10 304218488 1995 claude     maechling   ricketts     chemistry      185592680
## # i 4 more variables: university_id <int>, main_field <chr>, link <chr>,
## #   uni_name <chr>

#unique(d_chem$fieldofstudy)
## comparing to candidates:
# harvard:
# weldon in materials science
# beltrame in chemistry
# mit:
# lapointe is chemistry
# duff is chemistry
# stanford:
# shear in chemistry
# marcus is in biology
# hansen is in biology
# tokmakoff is in materials science

# update, chemistry check 8/11/22
# - tokmakoff still not linked; b/c of year first pub? -- yes, the linking score is 0.66...
```

- nancy fisher hansen (2649181519) is not linked (unclear if she should be linked)
- hopefully the keywords from topic models would help us here?
- maybe david h offord (304201950) would also be linked with the keywords?