# Performance of linking graduates to researchers

Flavio & Christoph

05 December, 2022

## Contents

## Overview

### SQL example for sourcing number of authors with same name

```sql
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?

```r
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"
         ) %>%
```

```r
  # 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)

  )
```
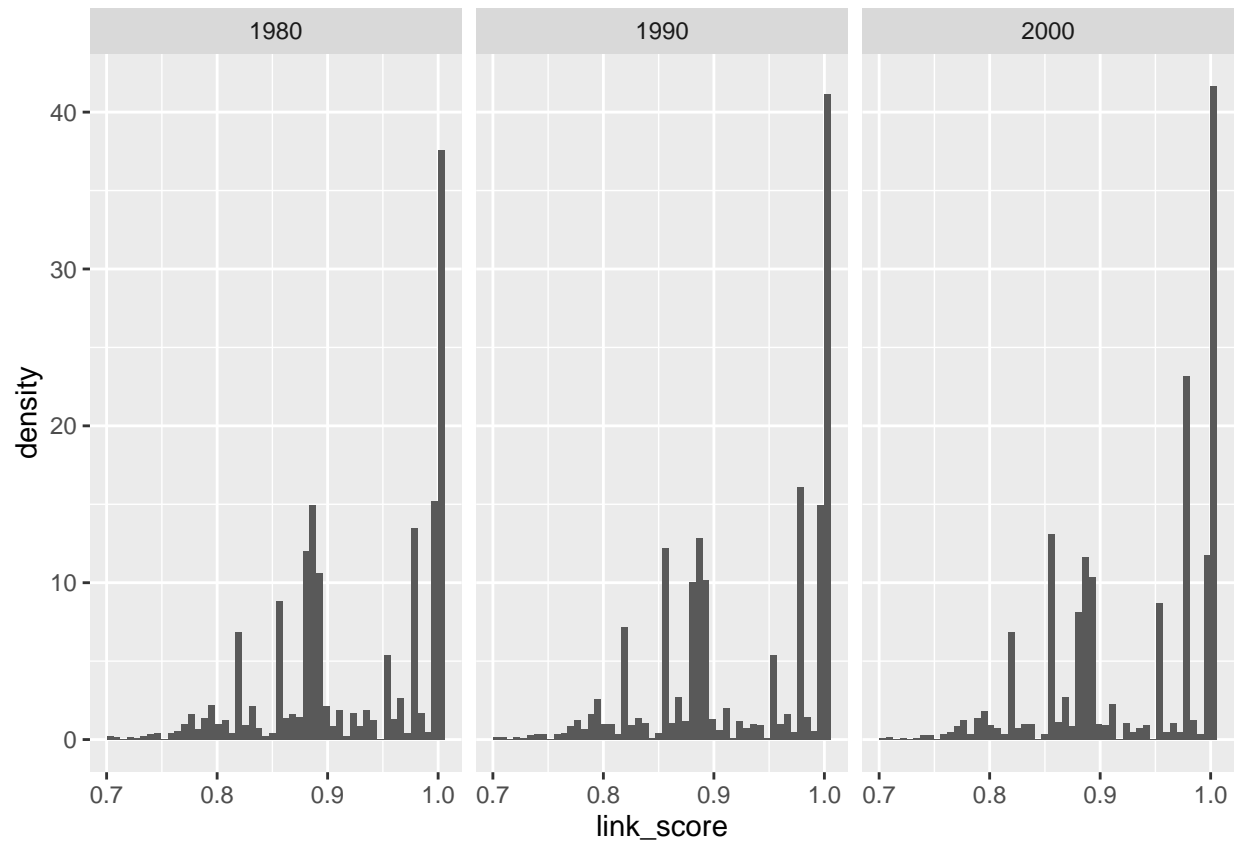
```
d_links$base <- d_links$base %>% filter(year_pq <= 2005)
```
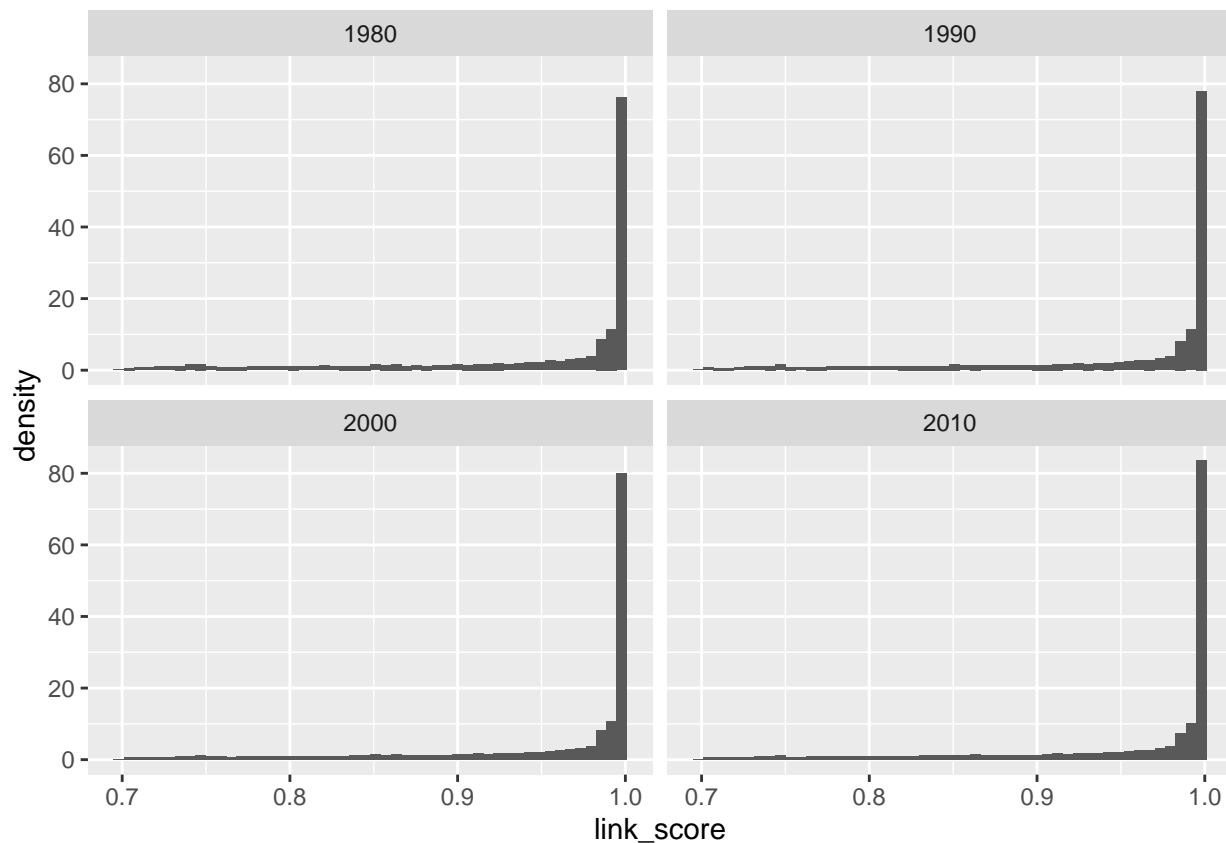
## Some histograms

### link score by field

```
## $base
```
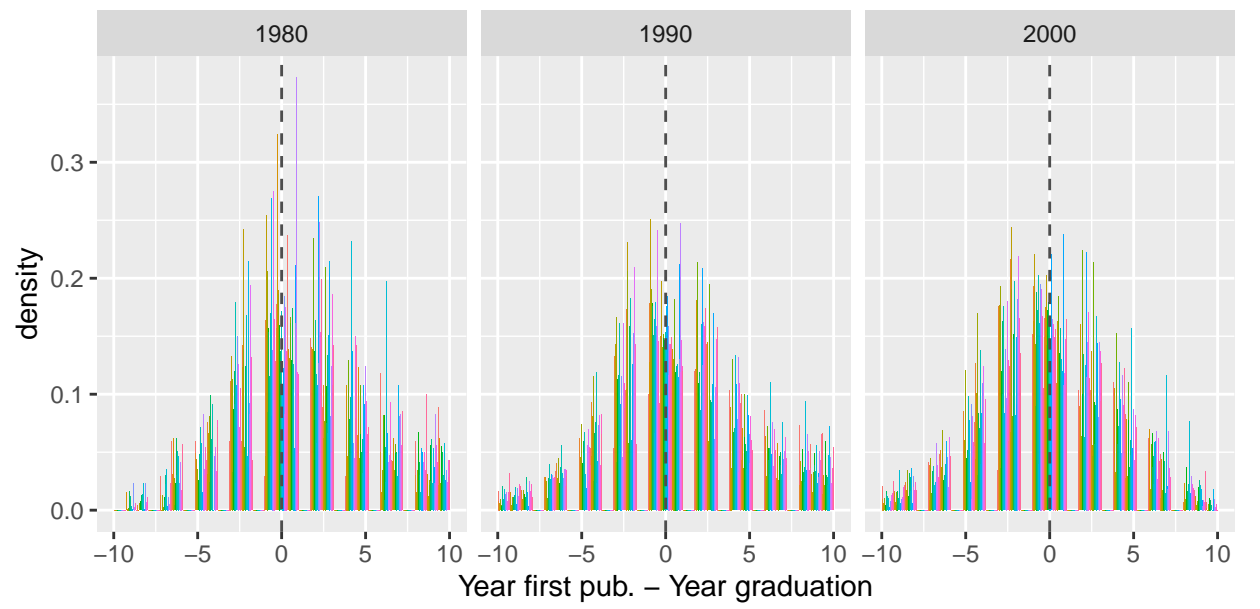


```
##
## $revise
```

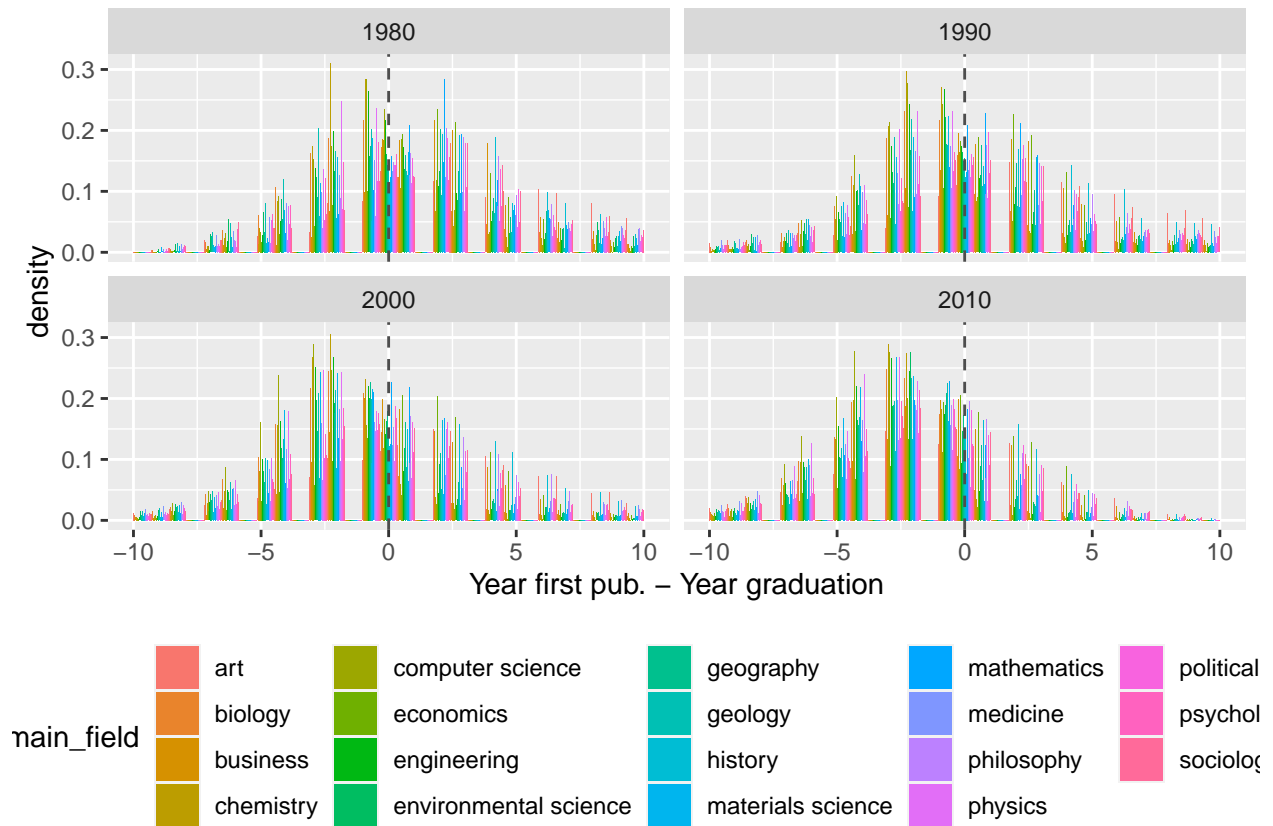**Year between first pub and graduation**

- 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`.
```
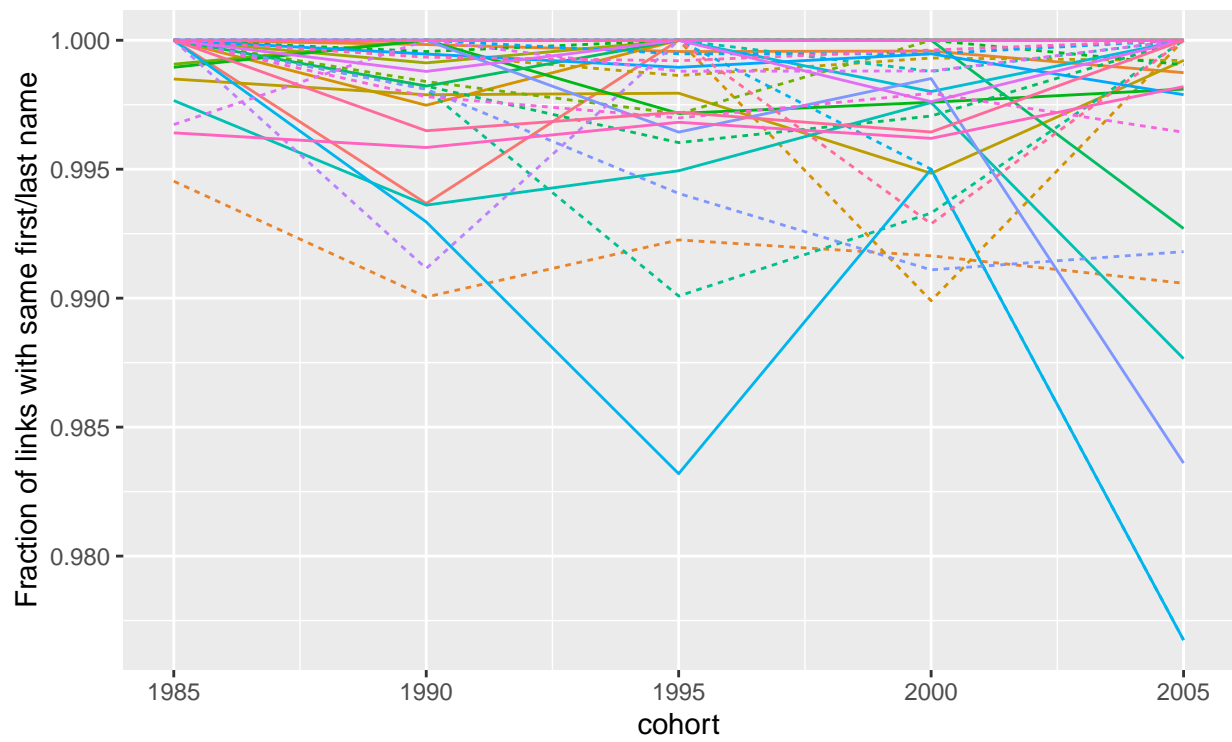
## 
## $revise

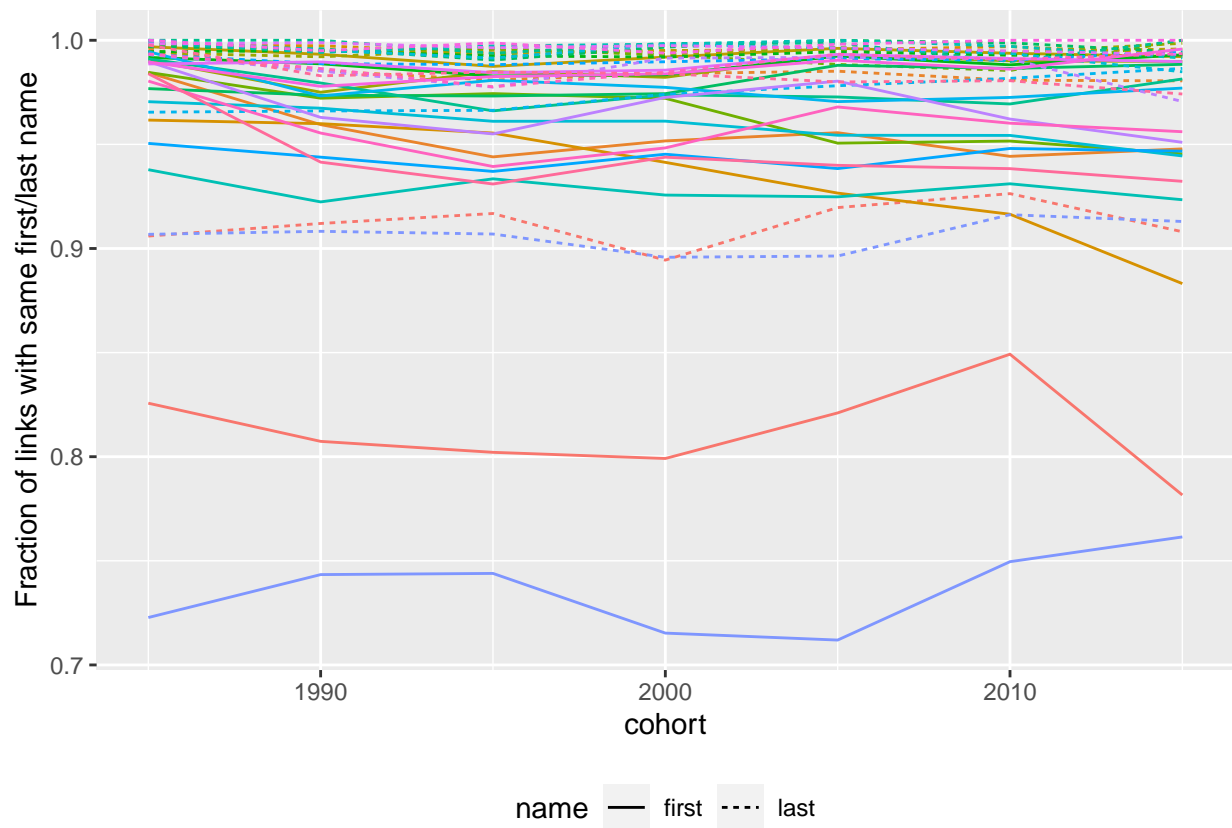## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.

First and last name matches by cohort and field
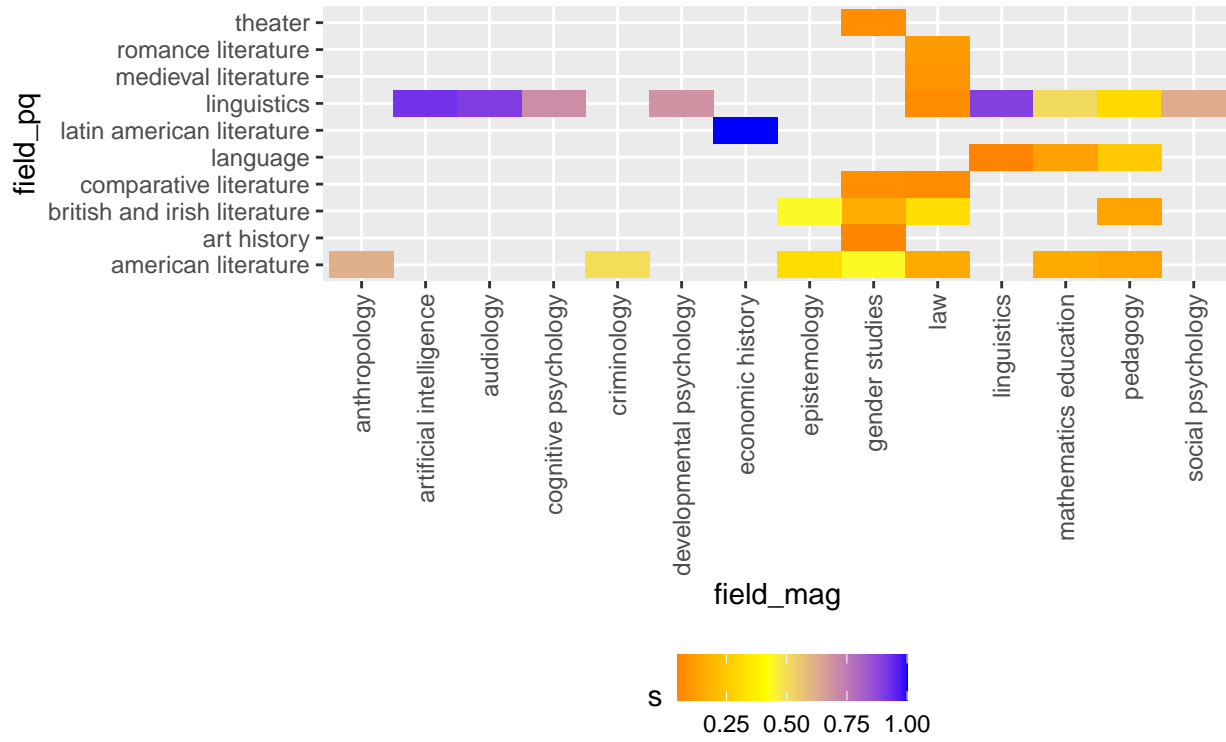
```
## $base
```

```
##
## $revise
```

How do fields of ProQuest map into fields in MAG?

```
## [[1]]
```

# Fraction of field ProQuest into field MAG

## Field: art



```
## 
## [[2]]
```
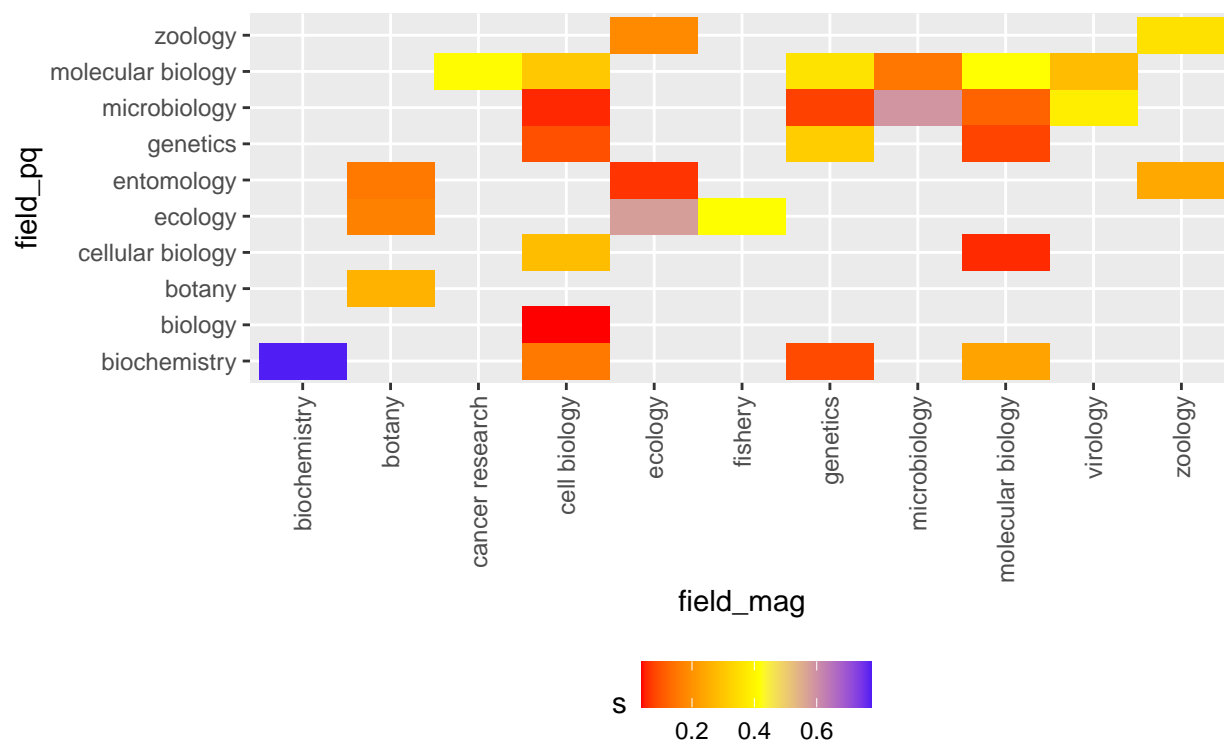
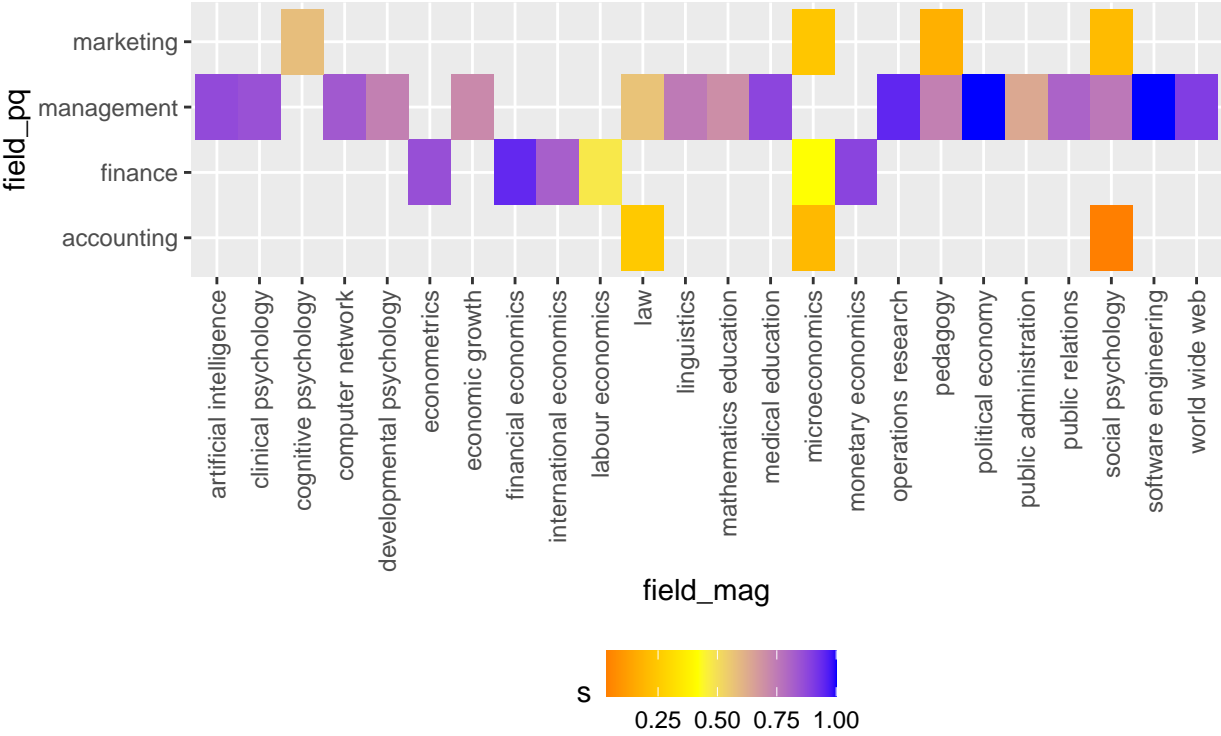Fraction of field ProQuest into field MAG

Field: biology

## 
## [[3]]

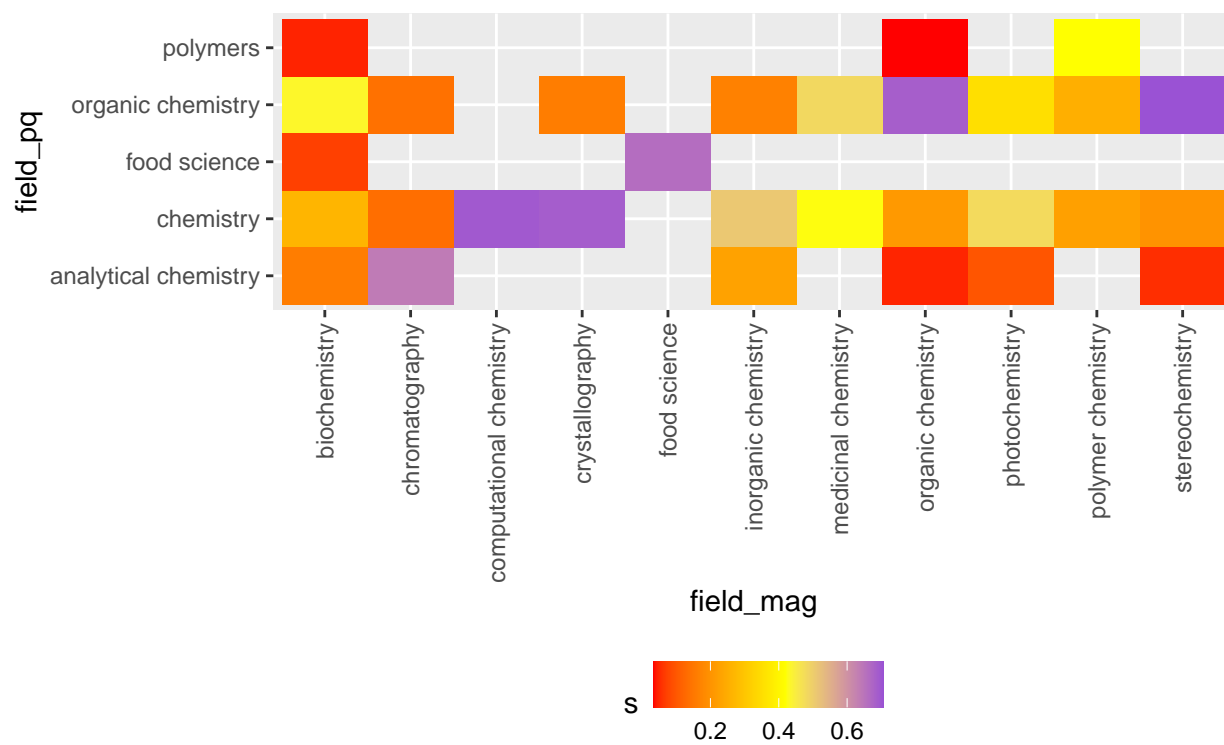Fraction of field ProQuest into field MAG

Field: business

```
##
## [[4]]
```

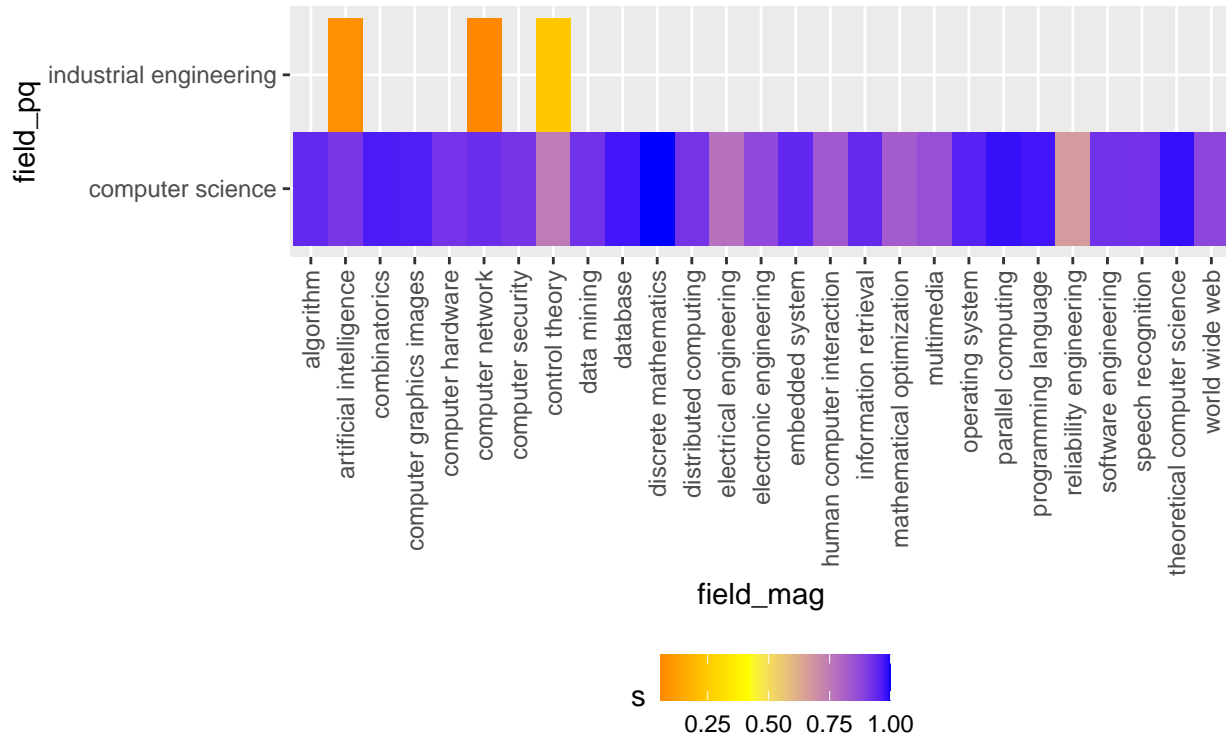Fraction of field ProQuest into field MAG

Field: chemistry

```
##
## [[5]]
```

# Fraction of field ProQuest into field MAG

## Field: computer science



```
##
## [[6]]
```

## Fraction of field ProQuest into field MAG
### Field: economics



```
## 
## [[7]]
```

Fraction of field ProQuest into field MAG
Field: engineering

```
##
## [[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]]
```
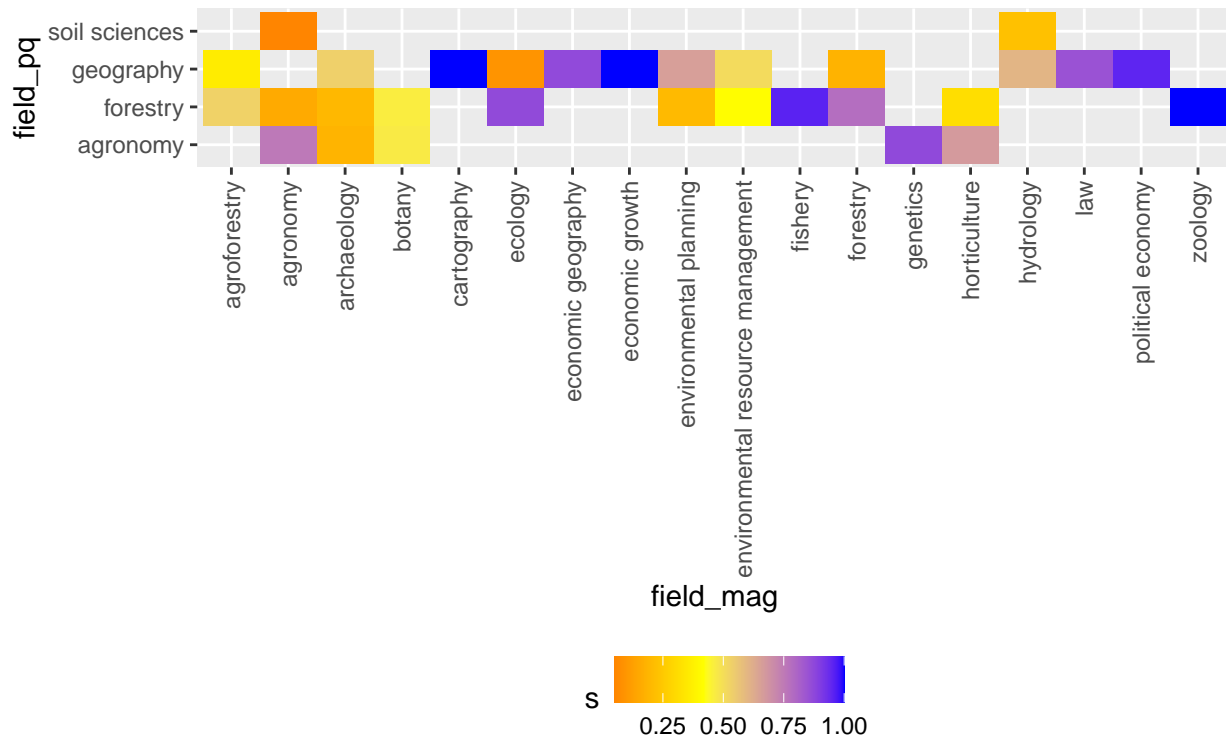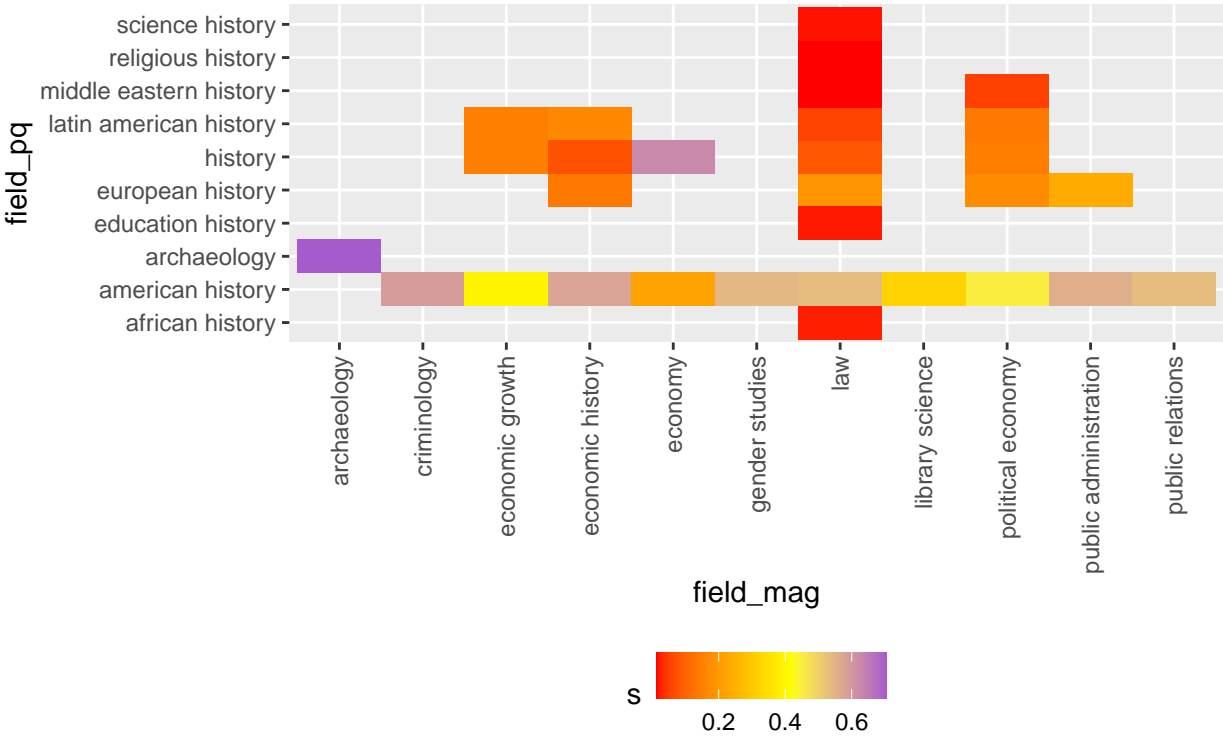
Fraction of field ProQuest into field MAG

Field: mathematics

field_pq

statistics

mathematics

field_mag

algebra
algorithm
applied mathematics
artificial intelligence
calculus
combinatorics
computer network
control theory
discrete mathematics
ecology
electronic engineering
genetics
mathematical analysis
mathematical economics
mathematical optimization
mathematical physics
programming language
pure mathematics
statistics
structural engineering
theoretical computer science

s

0.25  0.50  0.75  1.00

```
##
## [[14]]
```

Fraction of field ProQuest into field MAG

Field: medicine

```
##
## [[15]]
```

# Fraction of field ProQuest into field MAG
## Field: philosophy



field_pq

theology

philosophy

education philosophy

field_mag: artificial intelligence, clinical psychology, cognitive psychology, criminology, developmental psychology, discrete mathematics, ecology, epistemology, gender studies, genetics, law, library science, linguistics, mathematics education, medical education, pedagogy, political economy, psychoanalysis, psychotherapist, public administration, public relations, social psychology, theoretical physics
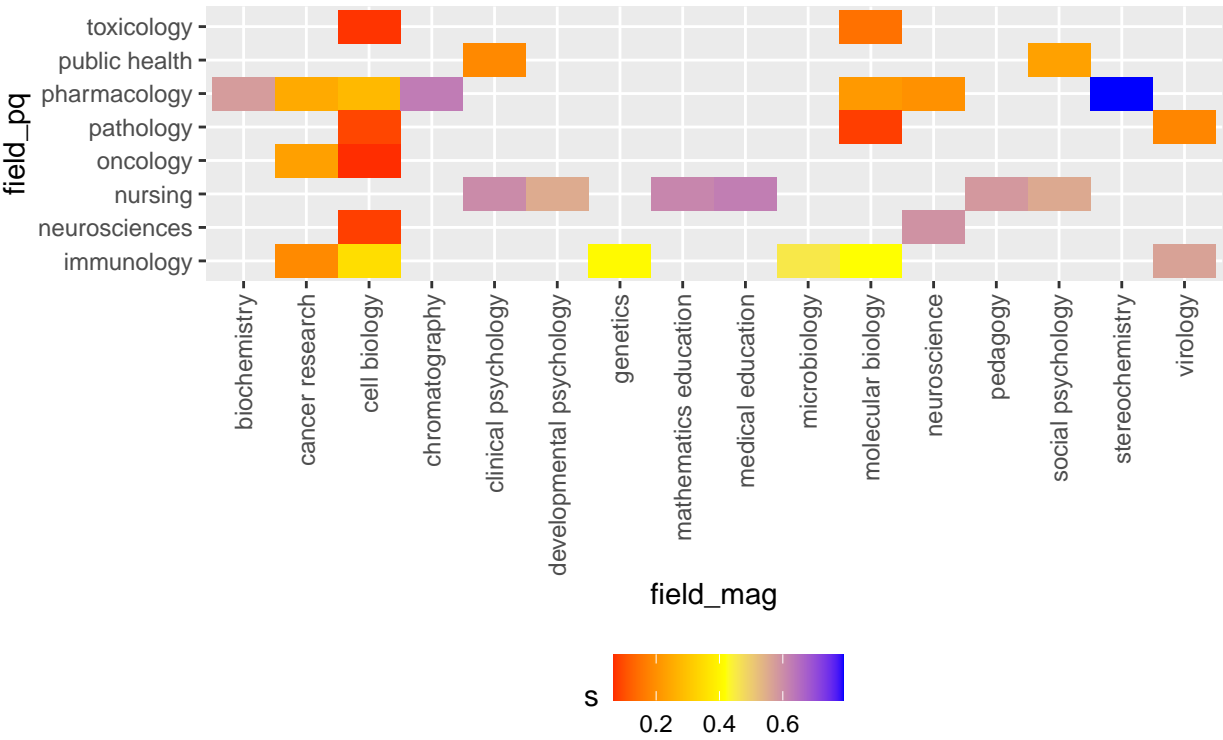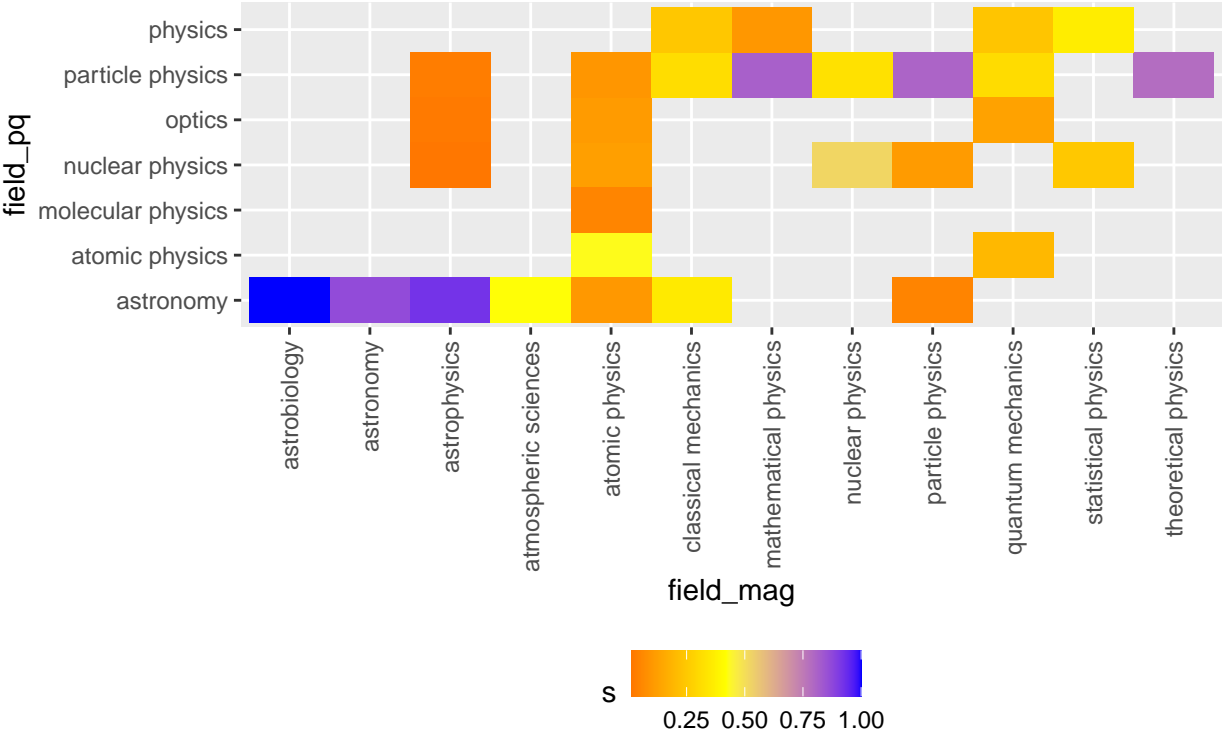
field_mag

s  0.25 0.50 0.75 1.00

```
##
## [[16]]
```

Fraction of field ProQuest into field MAG

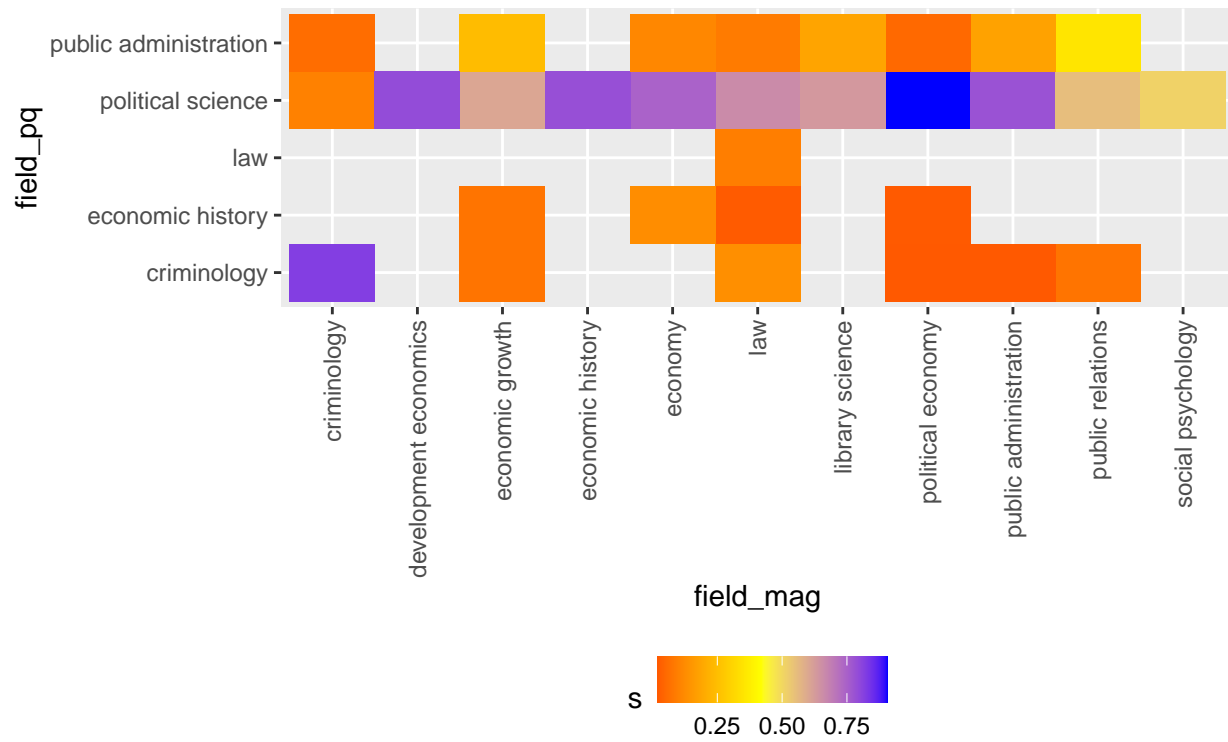Field: physics

##
## [[17]]

# Fraction of field ProQuest into field MAG
## Field: political science



```
## 
## [[18]]
```

# Fraction of field ProQuest into field MAG
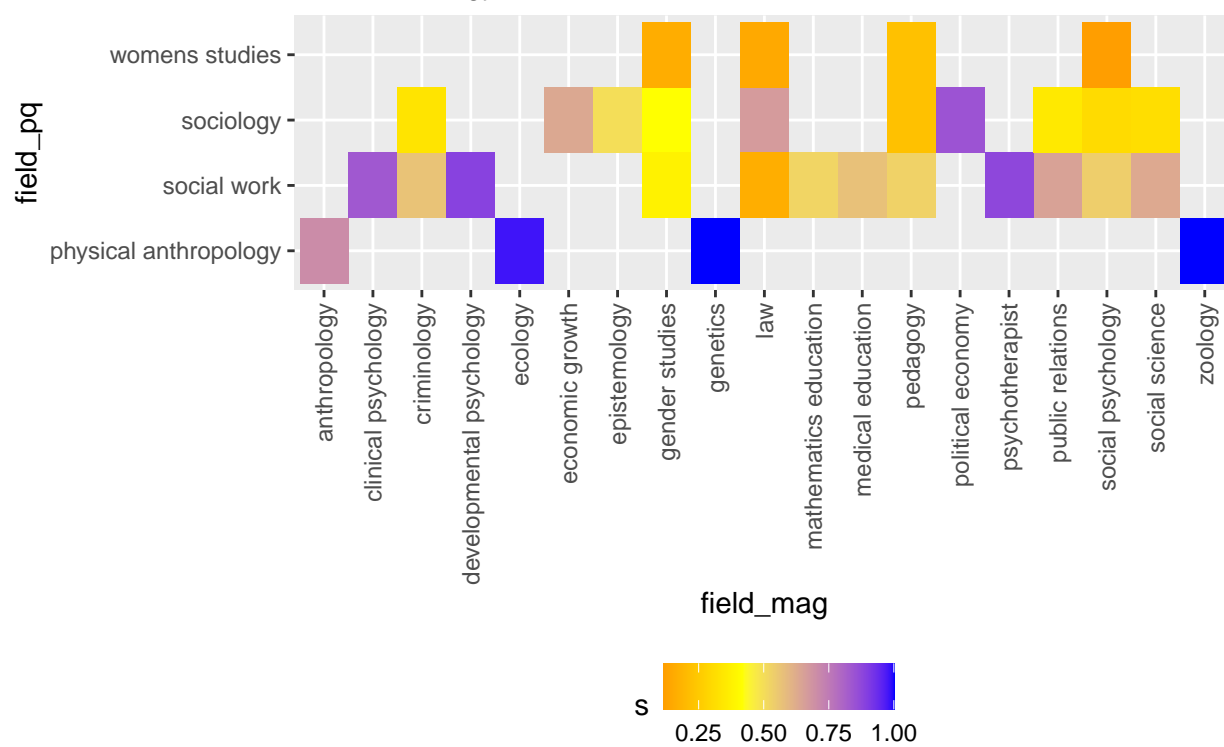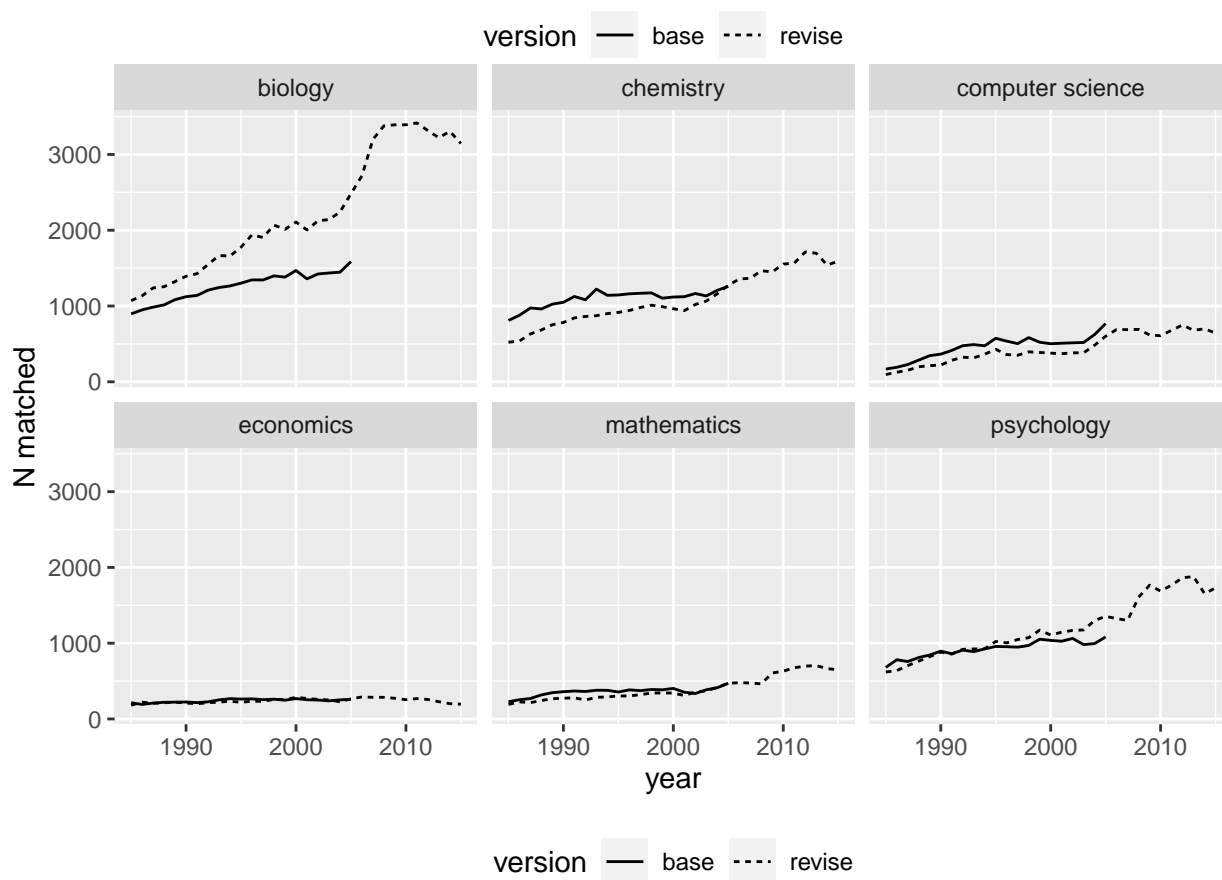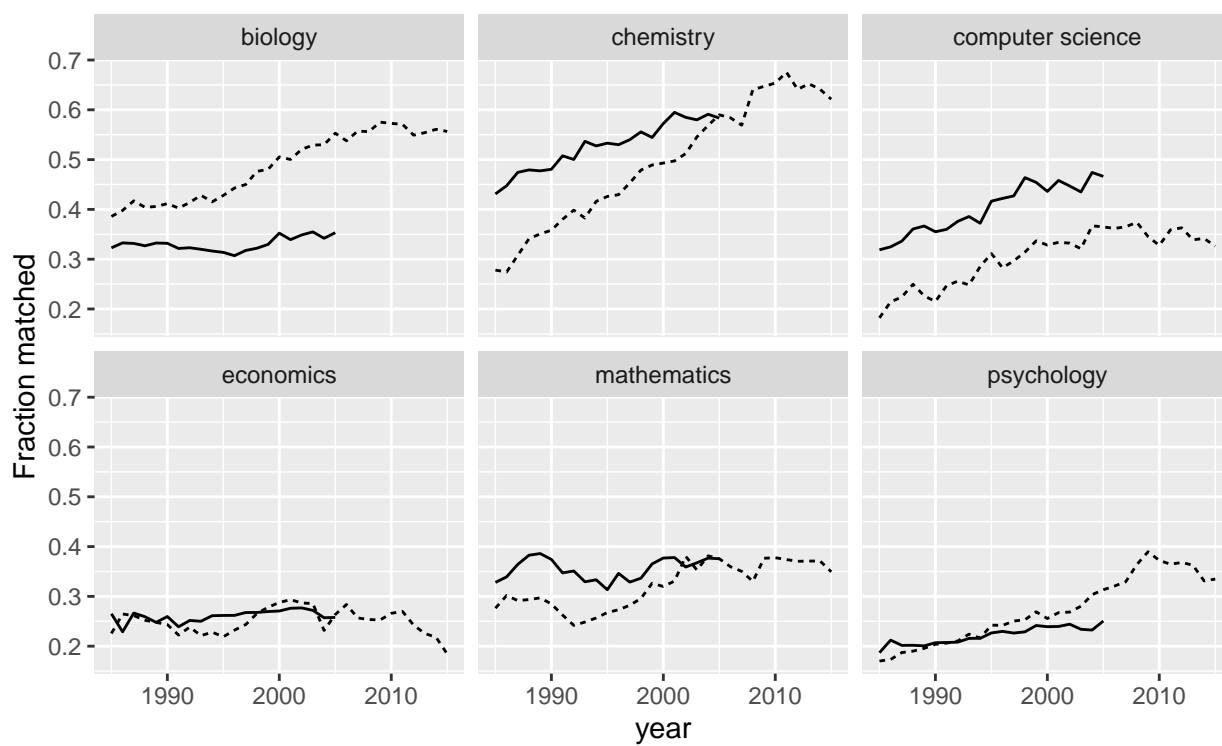
## Field: psychology



```
##
## [[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 10
##         goid  year first~1 lastn~2 middl~3 field~4 mag_f~5 main_~6 link  uni_n~7
##        <int64> <int> <chr>   <chr>   <chr>   <chr>     <int> <chr>   <chr> <chr>
##  1 304229925  1995 nancy   hansen  fisher  chemis~ 1.86e8 chemis~ not ~ stanfo~
##  2 304229722  1995 mark    pavlos~ alan    chemis~ 1.86e8 chemis~ not ~ stanfo~
##  3 304228620  1995 kristin sannes  ann     chemis~ 1.86e8 chemis~ not ~ stanfo~
##  4 304238241  1995 andrei  tokmak~ <NA>    chemis~ 1.86e8 chemis~ not ~ stanfo~
##  5 304218381  1995 glenn   jones   clark   chemis~ 1.86e8 chemis~ not ~ stanfo~
##  6 304218443  1995 david   brown   earl    chemis~ 1.86e8 chemis~ not ~ stanfo~
##  7 304201950  1995 david   offord  alan    chemis~ 1.86e8 chemis~ not ~ stanfo~
##  8 304238172  1995 robert  guettl~ david   chemis~ 1.86e8 chemis~ not ~ stanfo~
##  9 304202002  1995 eric    remy    david   chemis~ 1.86e8 chemis~ not ~ stanfo~
## 10 304238397  1995 james   brown   william chemis~ 1.86e8 chemis~ not ~ stanfo~
## # ... with abbreviated variable names 1: firstname, 2: lastname, 3: middlename,
## #   4: fieldofstudy, 5: mag_field0, 6: main_field, 7: uni_name
```

```
#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...
```

29

```
# - 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?
```