

**EINE QUANTITATIVE ANALYSE ANHAND VON BIG DATA AUS DEM  
KATALOG:  
WOHIN ENTWICKELT SICH DIE ROMANISTIK?**

***FACHSPEZIFISCHE ANGEBOTE DER SUB GÖTTINGEN - VON ANGLISTIK BIS  
ZENTRALASIENKUNDE, 5.5.2022***

---

**DR. JOSÉ CALVO TELLO**



NIEDERSÄCHSISCHE STAATS- UND  
UNIVERSITÄTSBIBLIOTHEK GÖTTINGEN | SUB

- *Digital Humanities?* Ok
- *Big Data* für die Geisteswissenschaften?

**ME WHEN SOMEBODY TALKS ABOUT**



**BIG DATA APPLIED TO HUMANITIES**

# **KATALOGDATEN ALS FORSCHUNGSOBJEKT**



**Titel:** [Horses and riding equipment in Indian art](#) / [Jean Deloche](#)  
**VerfasserIn:** [Deloche, Jean](#)  
**Ausgabe:** English ed  
**Sprache/n:** Englisch  
**Sprache des Originals:** Franzoesisch|  
**Veröffentlichungsangabe:** Madras : Indian Heritage Trust, 1990  
**Umfang:** 82 S : zahlr. Ill ; 30 cm  
**Einheitssachtitel:** [Le cheval et son harnachement dans l'art indien <engl.>](#)  
**Anmerkung:** Includes bibliographical references  
**ISBN:** : Rs150.00  
**Schlagwörter:** \*[Indien](#) / [Kunst](#) / [Pferd](#)  
\*[Horses in art](#) ; [Horses -- Equipment and supplies -- Pictorial works](#) ; [Art, Indic](#)  
**Sachgebiete:** [20.41 Asiatische Kunst](#)  
**Mehr zum Thema:** Klassifikation der Library of Congress: [N7668.H6](#)  
Dewey Dezimal-Klassifikation: [704.9432](#) ; [704.9/432](#)  
Regensburger Verbund-Klassifikation: [LH 84680: Kunstgeschichte / Allgemeines. Allgemeine Kunstgeschichte / Kunstgeschichte einzelner Gattungen der Kunst / Ikonographie, Ikonologie / Weltliche Ikonographie / Ikonographie der Tiere \(Bestiarien\) / Einzelne Tiere \(soweit sie nicht ausschließlich als biblische oder christliche Symbole aufgefaßt sind\) / Sonstige](#) ; [LO 88380: Kunstgeschichte / Kunst nach Ländern bzw. Kontinenten / Asien und Vorderer Orient / Kunstgeschichte Südasiens / Kunstgeschichte Indiens / Indische Kunstgeschichte nach Gattungen / Sonstige Gattungen \(u.a. Ornament\)](#)  
**Sachgebiete:** [GED 130 Buddhistische Kunst. Hinduistische Kunst {Geschichte der bildenden Kunst}](#);  
[GEC 700 Profane Ikonographie {Kunst}](#)  
**Standort:** [SUB Zentralbibliothek](#)  
**Signatur:** **A 96 B 35008**  
**Ausleihstatus:** Ausleihbestand  
Derzeit verfuegbar ➔ [Bestellen](#)

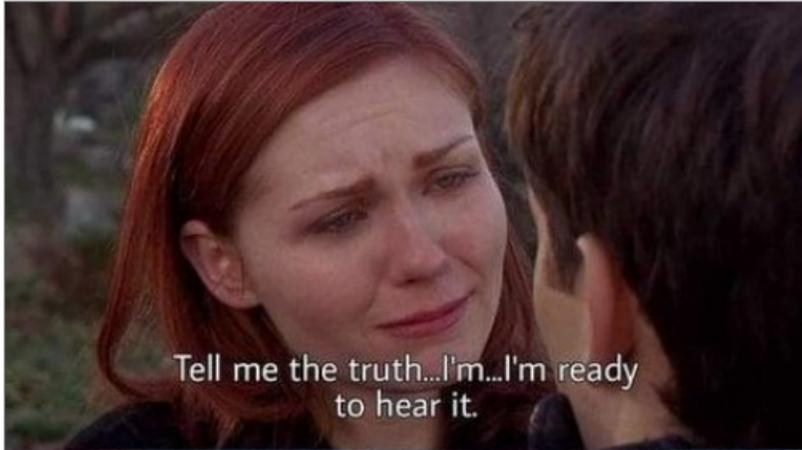
## KATALOGDATEN ALS FORSCHUNGSOBJEKT

- Wachsendes Interesse an Daten aus Bibliographien und Katalogen
- *Bibliographic Data Science*
- Vorteile
  - Große Datenmengen (*Big Data*)
  - Kuratierte Daten (in vielen Fällen normiert)

**BIG CURATED DATA!**

- Allerdings...





Tell me the truth...I'm...I'm ready  
to hear it.



**EVEN WHEN LIBRARY CATALOGUES ARE  
BIG AND CURATED, THEY ARE ALSO  
INCOMPLETE, HETEROGENEOUS AND CONTAIN BIASES**





## ZIELE

- Entwicklung der Romanistik in den letzten 40 Jahren?
- Erwartung für die Romanistik für die nächsten Jahre?

**WOFÜR BRAUCHE ICH DIESE ANALYSE ALS FACHREFERENT?**

## WOFÜR BRAUCHE ICH DAS ALS FACHREFERENT?

- Um die Angebote der Bibliothek an die aktuellen Entwicklungen der Romanistik anzupassen
  - Bestand
  - Lesesaal
  - Digital vs. Print
  - Beratung

**DATEN**

# QUELLEN: VERBUNDKATALOGE

- K10plus (GBV und SWB)
- HeBIS (Hessisches BibliotheksInformationsSystem)
- Art der Veröffentlichung: Monographien, Zeitschriften, Serien, E-Books (auch Open Access)
- Fachinformationsdienste (FIDs):
  - Fachinformationsdienst Romanistik
  - Fachinformationsdienst Lateinamerika, Karibik und Latino Studies

## AUSWAHL DER VERÖFFENTLICHUNGEN DER ROMANISTIK

- Göttinger Online-Klassifikation (GOK).

# ROMANISTIK IN KLASSIFIKATIONSSYSTEMEN:

- Basisklassifikation (BK)
  - regex = "18.[23]\d"
- Sachgruppe der DNB
  - regex = "^ (4[456]\d|8[456])"
- Dewey Decimal Classification (DDC)
  - regex = "^ (4[456]\d|8[456])"
- Regensburger Verbundklassifikation (RVK)
  - regex = "^|. \*?"
- Library of Congress Classification (LCC)
  - regex = "^p[ cq]. \*?"
- Göttinger Online Klassifikation (GOK)
  - regex = "^|H-Z. \*?"

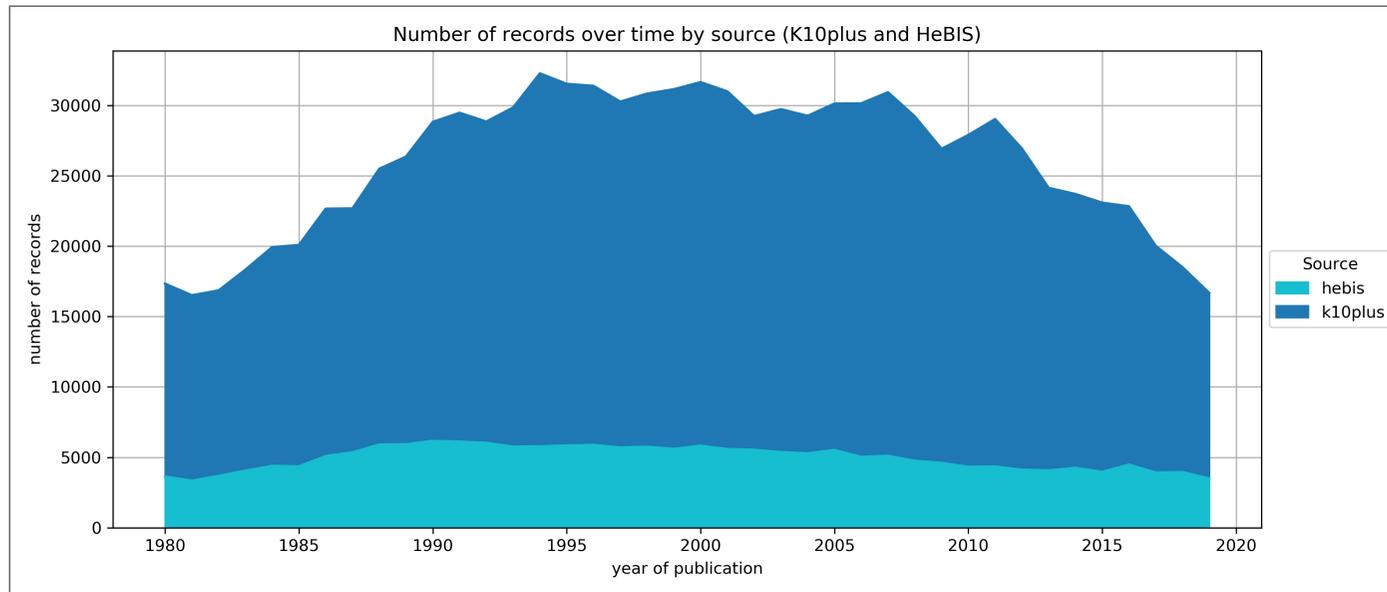
# SEKUNDÄRLITERATUR

- Folgende Klassen wurden ausgeschlossen:
  - Sachgruppen: B
  - BK: 17.97 und 17.98
  - RVK: Klassen mit "Gesammelte Werke" oder "Einzelwerke"
  - Art des Inhalts: Fiktionale Darstellung, Anthologie, Briefsammlung

# ZUSAMMENSTELLUNG UND EXTRAKTION

1. Zeitliche Grenzen: 1980-2019
2. Daten über Schnittstellen (API) herunterladen
3. Daten extrahieren (xPaths und RegEx) und in eine Tabelle konvertieren
4. Berechnung der relativen Häufigkeit anhand der Anzahl pro Jahr

# DATENSATZ



- 1.041.404 Titel (unter Berücksichtigung der Anzahl von Bibliotheken, die den Titel im Bestand haben)
- 81% K10plus
- 19% HeBIS
- Daten und Skripte bereits in GitLab veröffentlicht:  
<https://gitlab.gwdg.de/jose.calvotello/romance-languages-studies-2021>

# METHODE

# LINEARE REGRESSION

- Bereich des überwachten maschinellen Lernens
- ~~Klassifikation für kategoriale Werte~~
- Ziel ist die Vorhersage eines numerischen Werts

# SEITEN: BEISPIEL

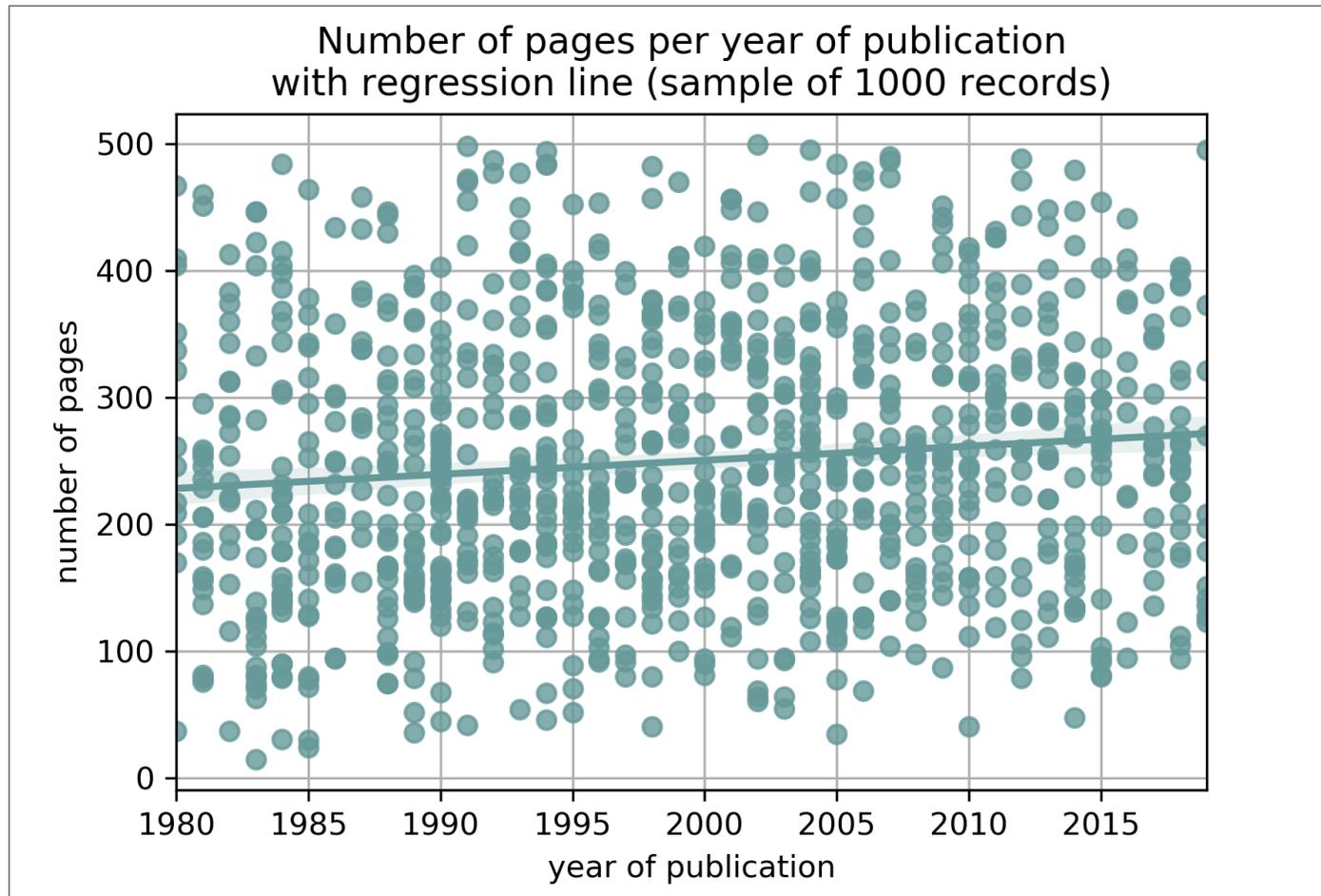
- **Forschungsfrage:** Werden die Veröffentlichungen länger? Nimmt die Anzahl der Seiten im Laufe der Zeit zu?

	<b>Titel:</b> <a href="#">Horses and riding equipment in Indian art</a> / <a href="#">Jean Deloche</a>
<b>VerfasserIn:</b>	<a href="#">Deloche, Jean</a>
<b>Ausgabe:</b>	English ed
<b>Sprache/n:</b>	Englisch
<b>Sprache des Originals:</b>	Franzoesisch
<b>Veröffentlichungsangabe:</b>	Madras : Indian Heritage Trust, 1990
<b>Umfang:</b>	82 S : zahlr. Ill ; 30 cm
<b>Einheitsachtitel:</b>	<a href="#">Le cheval et son harnachement dans l'art indien &lt;engl.&gt;</a>
<b>Anmerkung:</b>	Includes bibliographical references
<b>ISBN:</b>	: Rs150.00
<b>Schlagwörter:</b>	* <a href="#">Indien</a> / <a href="#">Kunst</a> / <a href="#">Pferd</a> * <a href="#">Horses in art</a> ; <a href="#">Horses -- Equipment and supplies -- Pictorial works</a> ; <a href="#">Art, Indic</a>
<b>Sachgebiete:</b>	<a href="#">20.41 Asiatische Kunst</a>
<b>Mehr zum Thema:</b>	Klassifikation der Library of Congress: <a href="#">N7668.H6</a> Dewey Dezimal-Klassifikation: <a href="#">704.9432</a> ; <a href="#">704.9/432</a> Regensburger Verbund-Klassifikation: <a href="#">LH 84680: Kunstgeschichte / Allgemeines. Allgemeine Kunstgeschichte / Kunstgeschichte einzelner Gattungen der Kunst / Ikonographie, Ikonologie / Weltliche Ikonographie / Ikonographie der Tiere (Bestiarien) / Einzelne Tiere (soweit sie nicht ausschließlich als biblische oder christliche Symbole aufgefaßt sind) / Sonstige</a> ; <a href="#">LO 88380: Kunstgeschichte / Kunst nach Ländern bzw. Kontinenten / Asien und Vorderer Orient / Kunstgeschichte Südasiens / Kunstgeschichte Indiens / Indische Kunstgeschichte nach Gattungen / Sonstige Gattungen (u.a. Ornament)</a>
<b>Sachgebiete:</b>	<a href="#">GED 130 Buddhistische Kunst. Hinduistische Kunst {Geschichte der bildenden Kunst}</a> . <a href="#">GEC 700 Profane Ikonographie {Kunst}</a> .
<b>Standort:</b>	<a href="#">SUB Zentralbibliothek</a>
<b>Signatur:</b>	<b>A 96 B 35008</b>
<b>Ausleihstatus:</b>	Ausleihbestand Derzeit verfuegbar ➔ <a href="#">Bestellen</a>

## SEITEN: FIKTIONALES BEISPIEL

- 1980: 80 Seiten
- 1981: 81 Seiten
- 1982: 82 Seiten
- 1990: 90 Seiten
- 2020: 120 Seiten
- 2030: ?

## SEITEN: VISUALISIERUNG



- Durchschnitt von Seiten in **1980: 219**
- Durchschnitt von Seiten in **2019: 256**

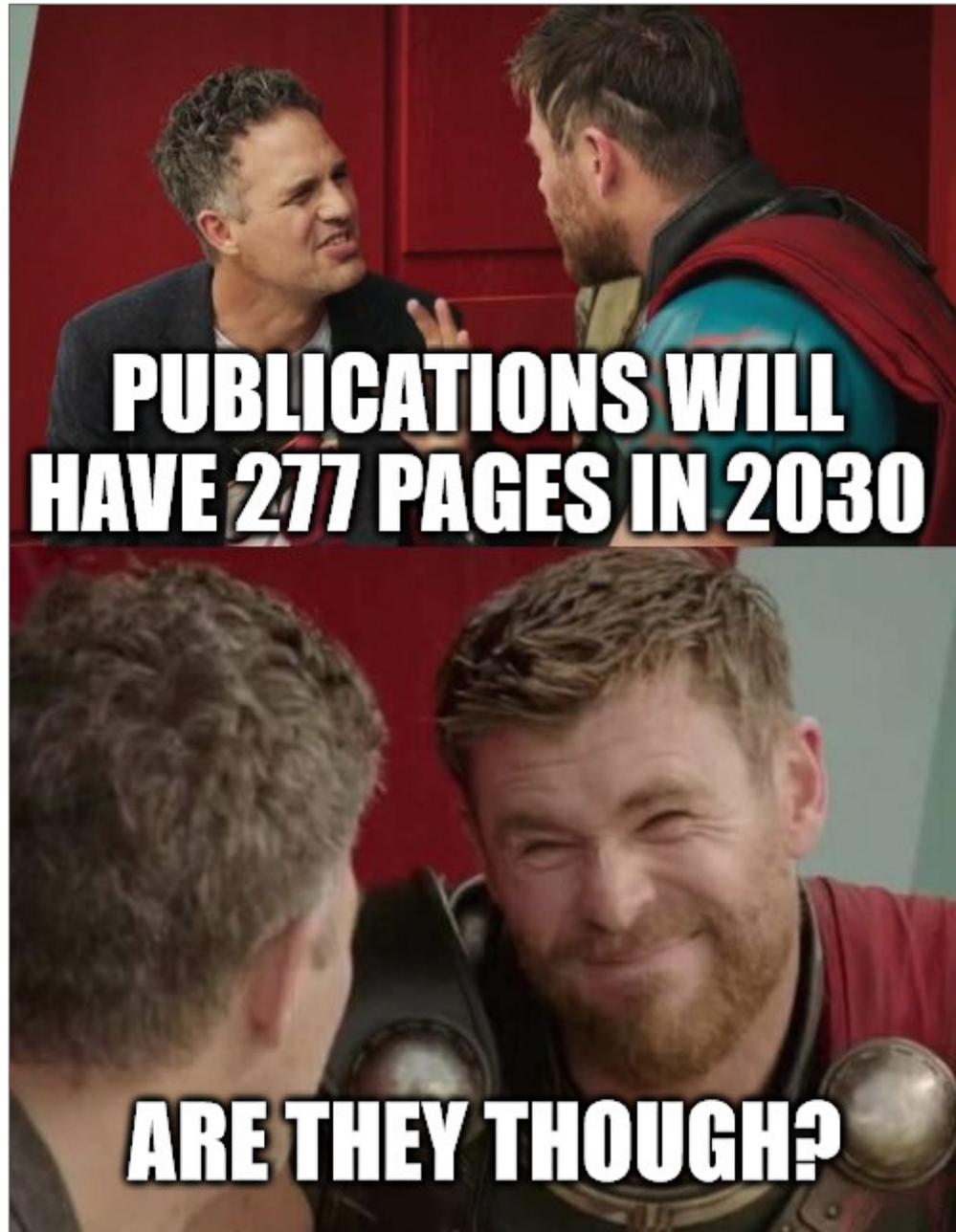


## SEITEN: REGRESSION

- slope = 0,9, p-value < 0,001
- Erwartung:
  - Publikationen werden in 10 Jahren im Durchschnitt 9 Seiten länger sein
  - 2030 Publikationen von 277
  - 2040 Publikationen von 286

**ABER...**





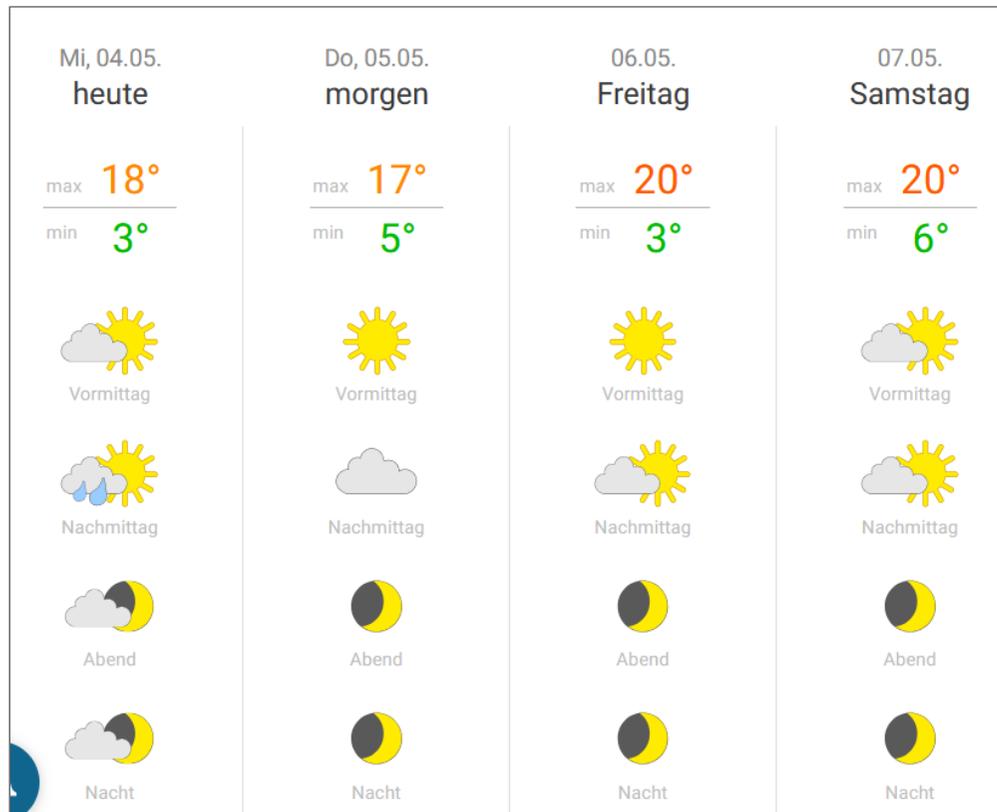
**PUBLICATIONS WILL  
HAVE 277 PAGES IN 2030**

**ARE THEY THOUGH?**



**NEIN, DIESE VORHERSAGEN WERDEN NICHT GANZ  
ZUTREFFEND**

# WETTERVORHERSAGEN SIND AUCH NIE PERFECT, TROTZDEM NÜTZLICH



# UM ENTSCHEIDUNGEN ZU TREFFEN

- Entweder für den Trend
- Oder dagegen!

# **GENERELLES PROBLEM DER METHODEN DES AKTUELLEN (ÜBERWACHTEN) MASCHINELLEN LERNENS**

## **KONSERVATIV**

- Mit Daten, die wir bereits gesammelt haben, neue Fälle vorhersagen
- Mit Daten der Vergangenheit zukünftige Entwicklungen vorhersagen

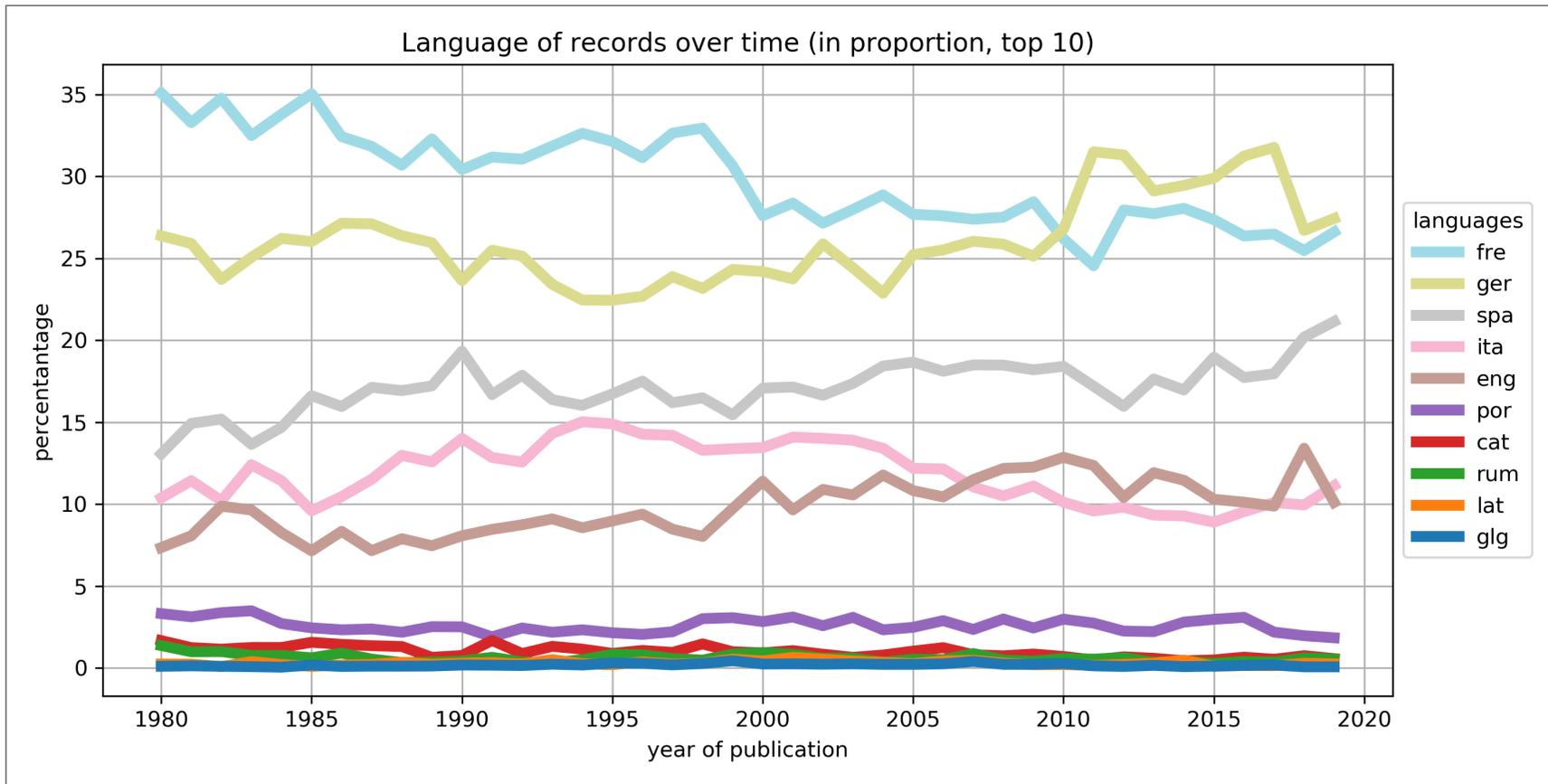
**WICHTIG: IMMER NEUE KURATIERTE DATEN EINTRAGEN**

# KATEGORIEN

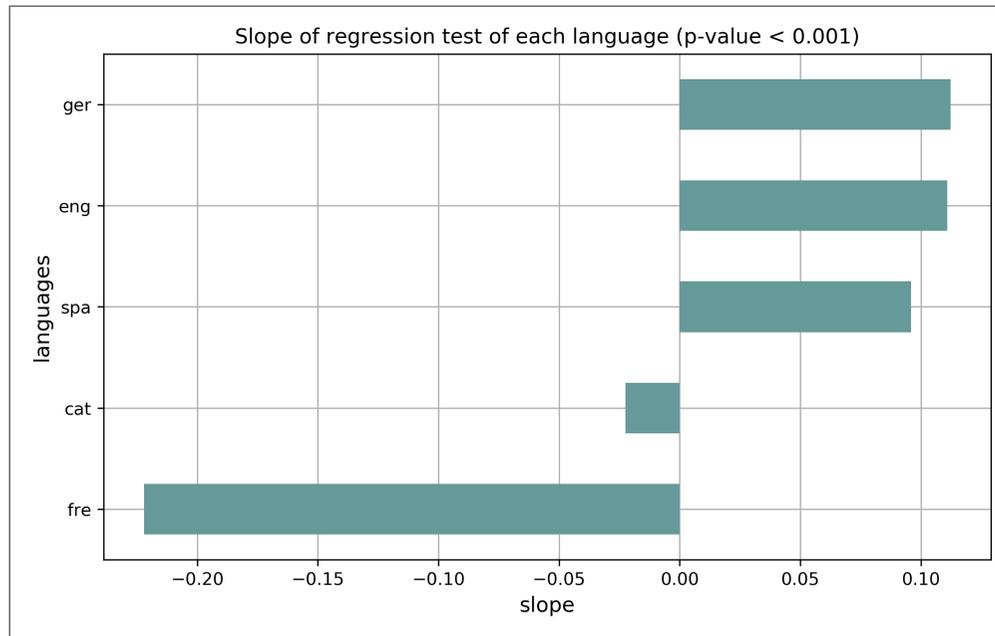
# SPRACHEN

- Veröffentlichungssprache
- ~~Sprache als Forschungsobjekt (später!)~~

# SPRACHEN: ENTWICKLUNG



# SPRACHEN: SLOPE

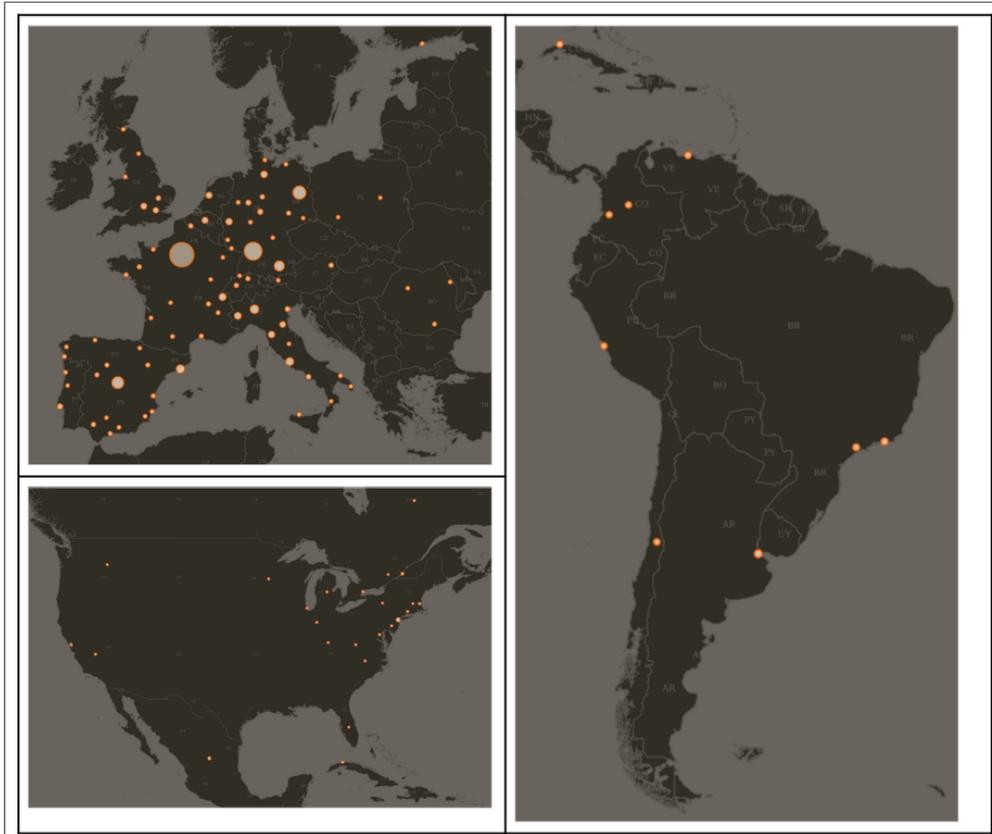


- Erwartung für 2030
  - 29 % Deutsch
  - 24 % Französisch
  - 20 % Spanisch (2. Sprache ab 2040)
  - 13 % Englisch

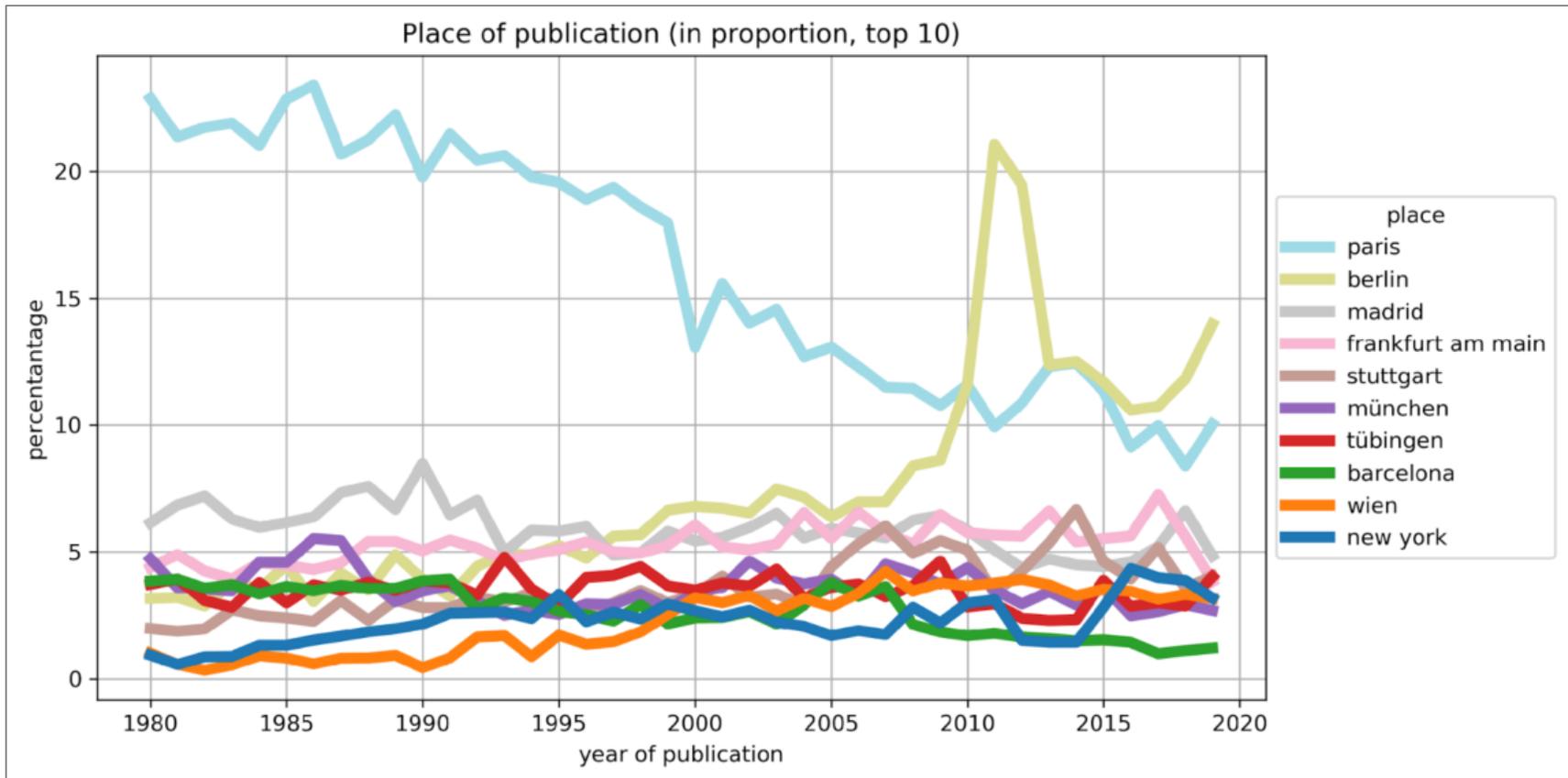
# **ORT DER VERÖFFENTLICHUNG**

## ORT DER VERÖFFENTLICHUNG: VISUALISIERUNG

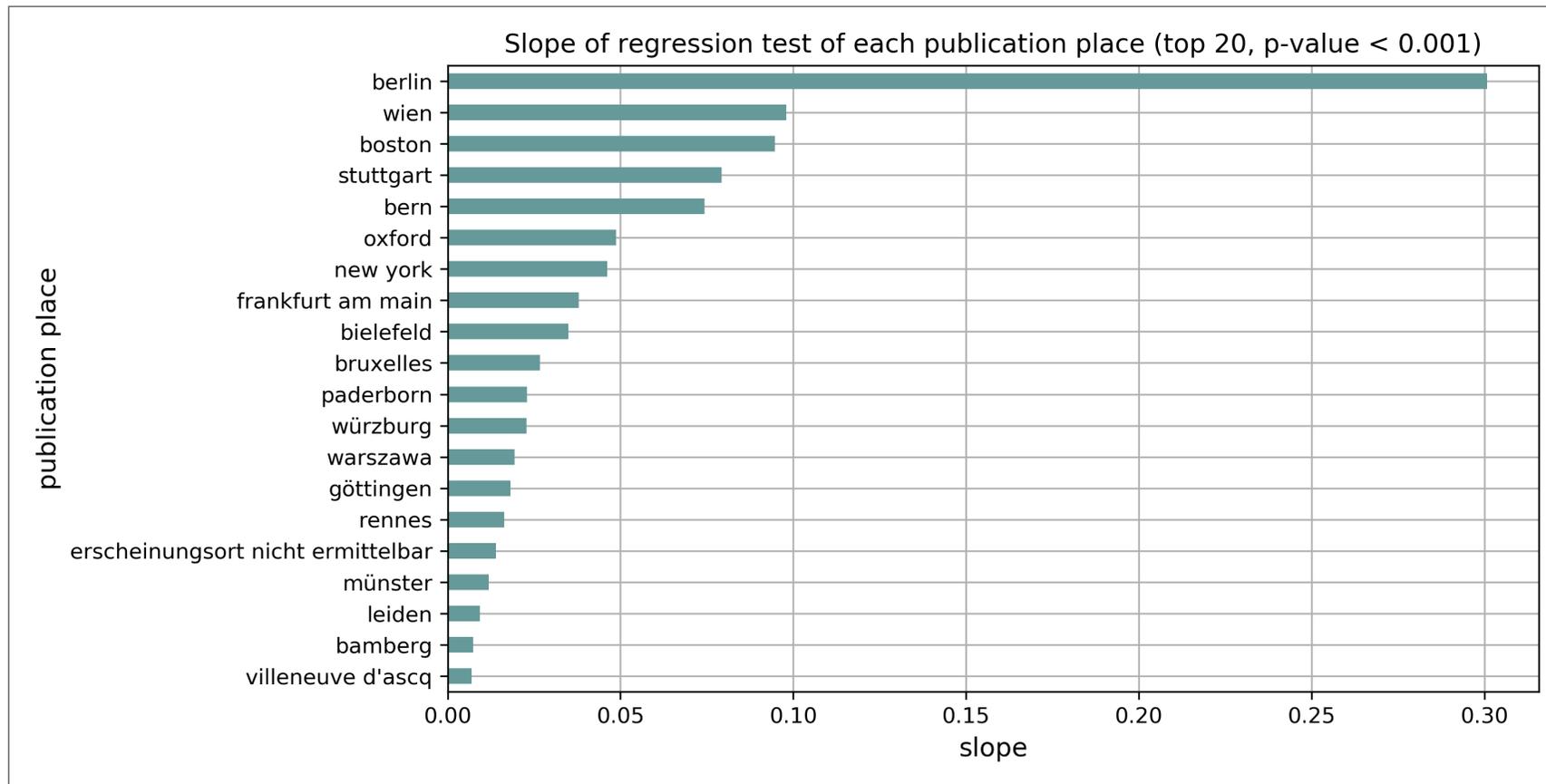
- Sample mit 5.000 Publikationen
- In DARIAH-DE Geo-Browser visualisiert und veröffentlicht



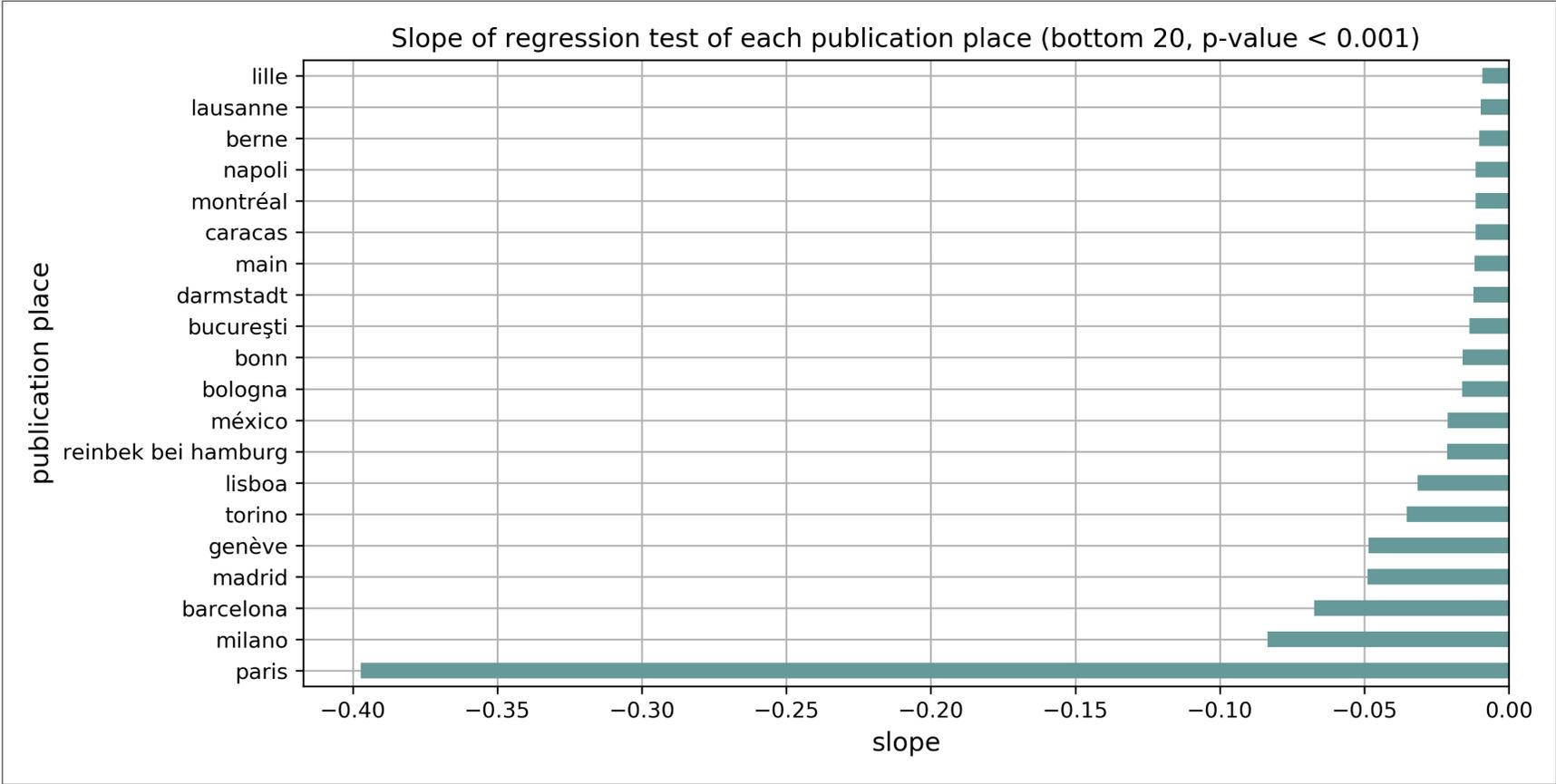
# ORT DER VERÖFFENTLICHUNG: ENTWICKLUNG



# ORT DER VERÖFFENTLICHUNG: SLOPE (TOP)

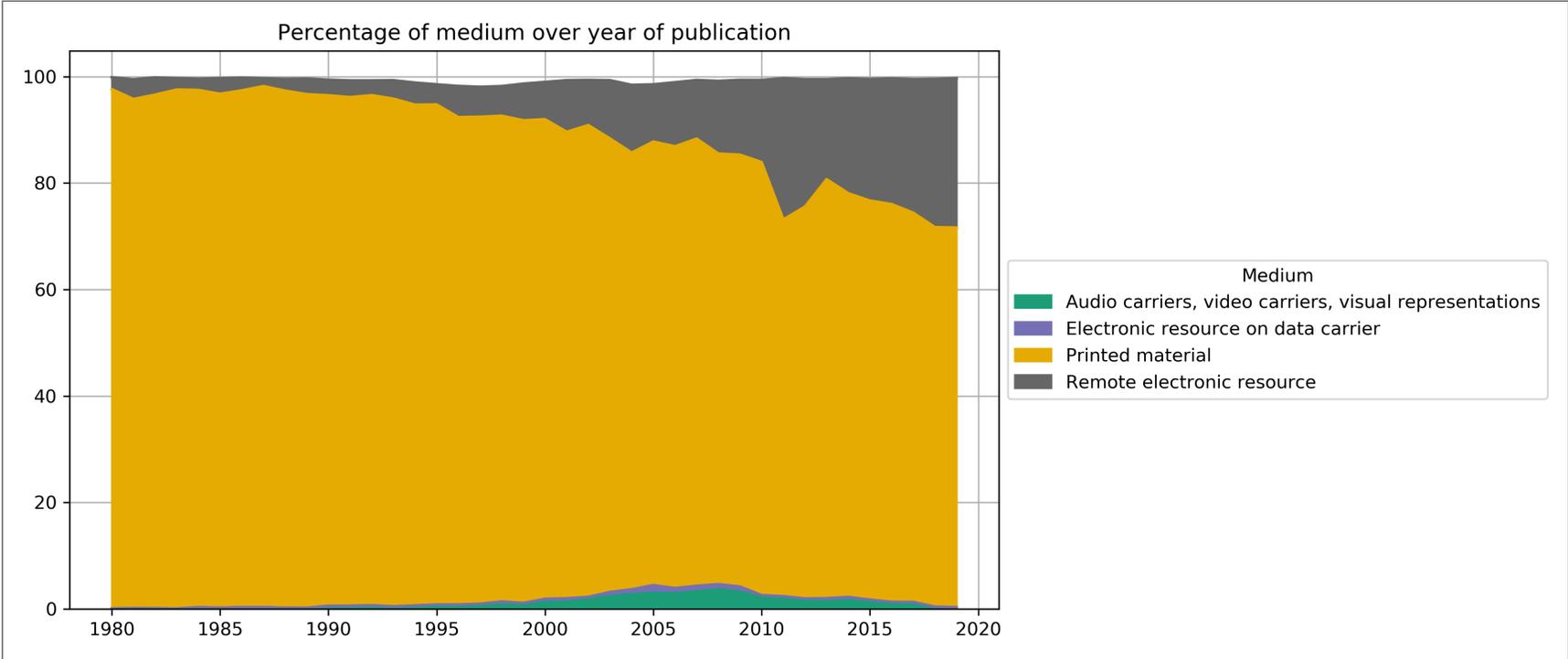


# ORT DER VERÖFFENTLICHUNG: SLOPE (BOTTOM)



## **MEDIUM (PRINT VS. E-BOOK)**

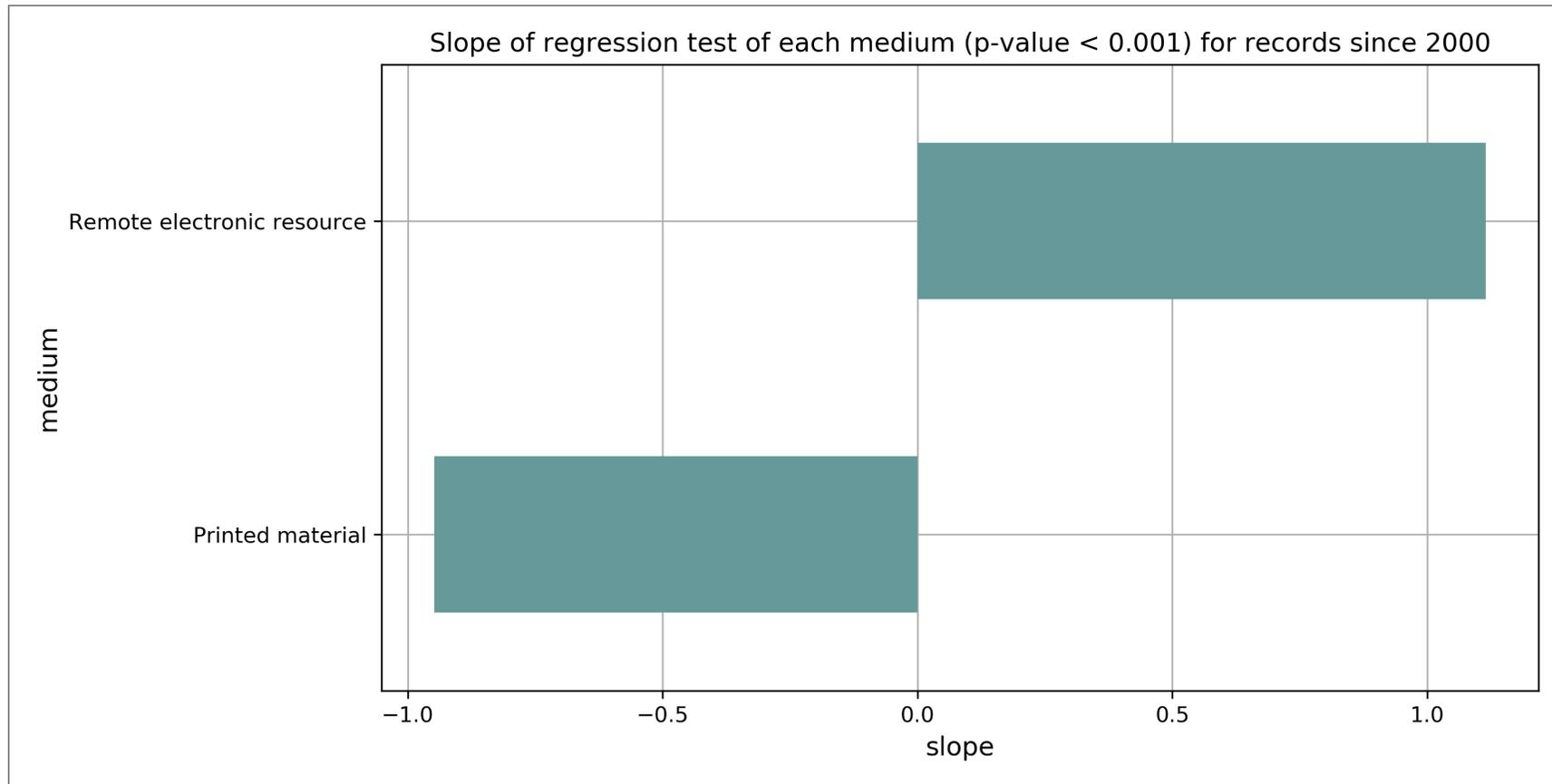
# MEDIUM: ENTWICKLUNG



## **MEDIUM: ENTWICKLUNG**

- Gesamte Periode
  - Print: ca. 90 %
  - E-Books: ca. 10 %
- Nach 2010
  - E-Books > 20 %

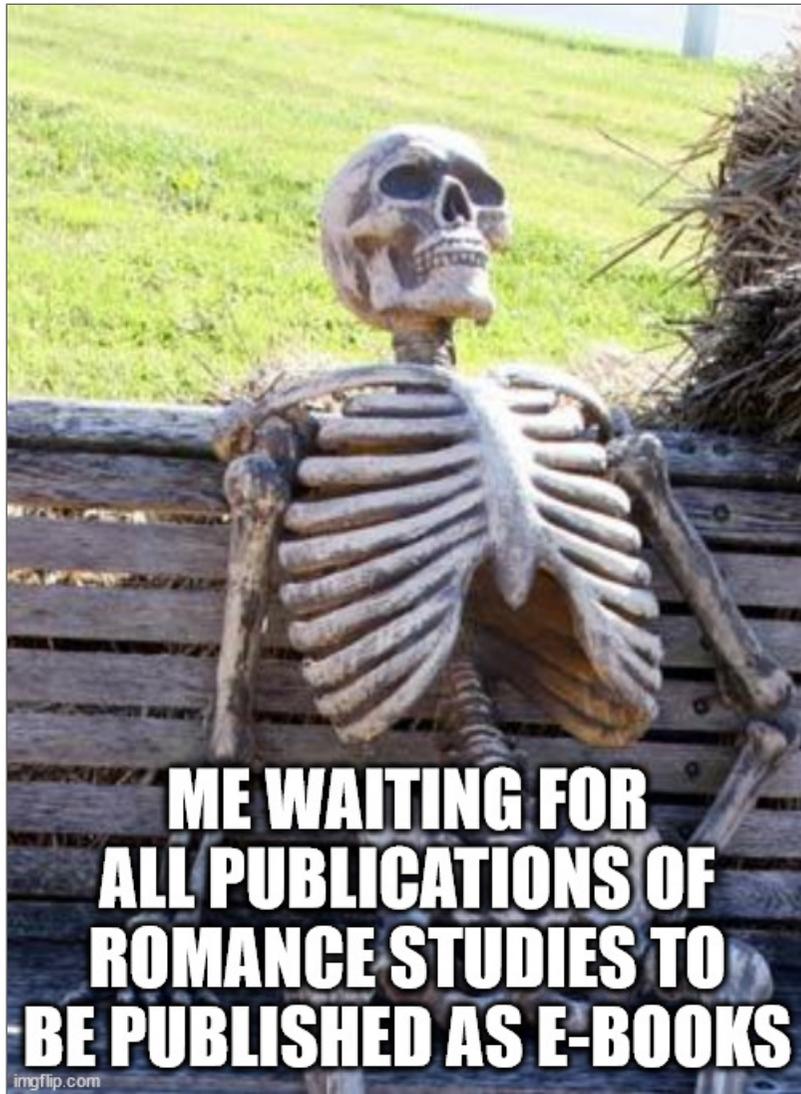
## MEDIUM: SLOPE



- Erwartung: 1,1 % mehr E-Books pro Jahr

**WANN WERDEN ALLE PUBLIKATIONEN DER ROMANISTIK ALS E-BOOK  
VORLIEGEN?**

**100% E-BOOKS IN 2080!**



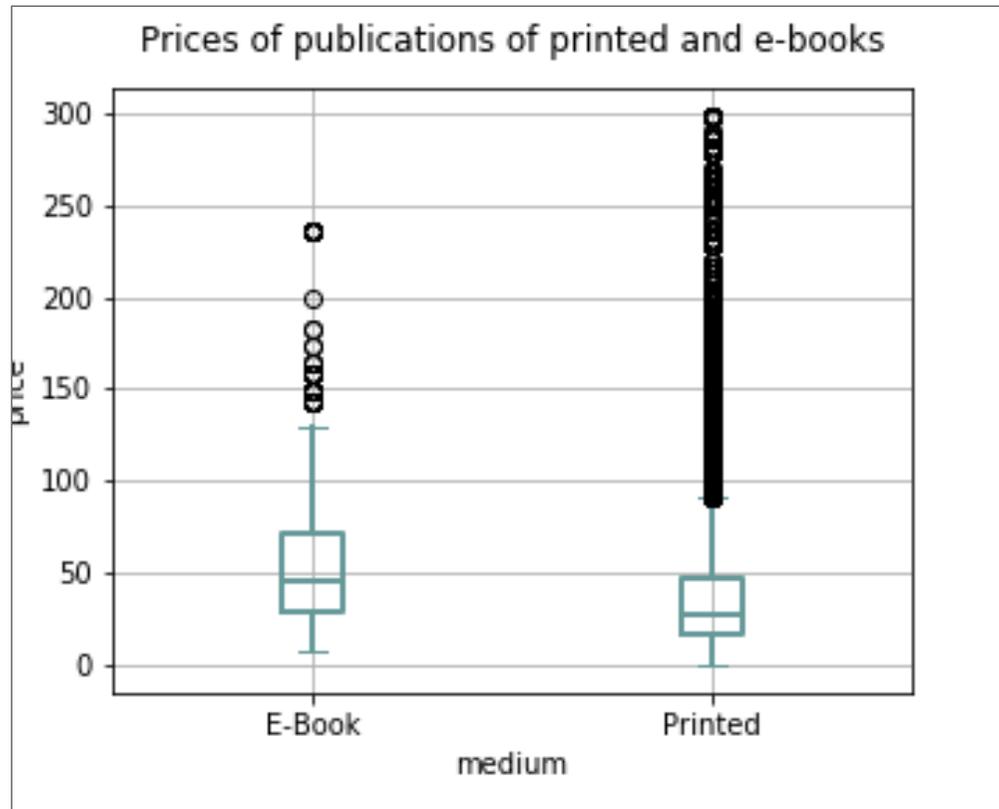


## **PREIS: PRINT VS. E-BOOK**

- Was ist günstiger? E-Books oder Print?

**BIBLIOTHEKEN MÜSSEN BESONDERE LIZENZEN FÜR E-BOOKS  
ERWERBEN!**

# PREIS: PRINT VS. E-BOOK



- Print = 28 € Median Preis
- E-Book = 46 € Median Preis
- P-value < 0,001 in Welch's t-test
- Die Ergebnisse sind nicht endgültig!



# THEMEN

# **THEMEN: SCHLAGWÖRTER**

# SCHLAGWÖRTER IM KATALOG

 PPN: 1735146498 

Titel: [Phraséologie et stylistique de la langue littéraire : approches interdisciplinaires / Ludwig Fesenmeier/Iva Novakova \(eds.\)](#)  
[Phraseology and stylistics of literary language : interdisciplinary approaches](#)

Person/en: [Fesenmeier, Ludwig \\*1970-\\*](#) [HerausgeberIn]  ; [Novakova, Iva](#) [HerausgeberIn] 

Körperschaft/en: [Phraséologie et stylistique de la langue littéraire <Veranstaltung> <2019, Erlangen>](#) [VerfasserIn]   
[Peter Lang GmbH](#) [Verlag] 

Konferenz: [Phraséologie et stylistique de la langue littéraire ; \(Erlangen\) : 2019](#)

Ort/Jahr: Berlin ; Bern ; Bruxelles ; New York ; Oxford ; Warszawa ; Wien : Peter Lang, [2020]

Sprache/n: Englisch, Französisch

Art des Inhalts: [Konferenzschrift \(2019, Erlangen\)](#)

Umfang: 376 Seiten : Illustrationen

Anmerkung: "Le présent ouvrage est issu du colloque international "Phraséologie et Stylistique de la langue littéraire/Phraseology and S  
à la Friedrich-Alexander-Universität Erlangen-Nürnberg" (Introduction)  
Beiträge teilweise französisch, teilweise englisch

ISBN: 978-3-631-81137-5  
Weitere Ausgaben: 978-3-631-83632-3, 978-3-631-83633-0, 978-3-631-83634-7

Identifizier: DOI: [10.3726/b17628](#)

DOI in Druckwerken: 10.3726/b17628

**Schlagwörter:** [\\*Dichtersprache](#)  / [Phraseologie](#)  / [Literarischer Stil](#) 

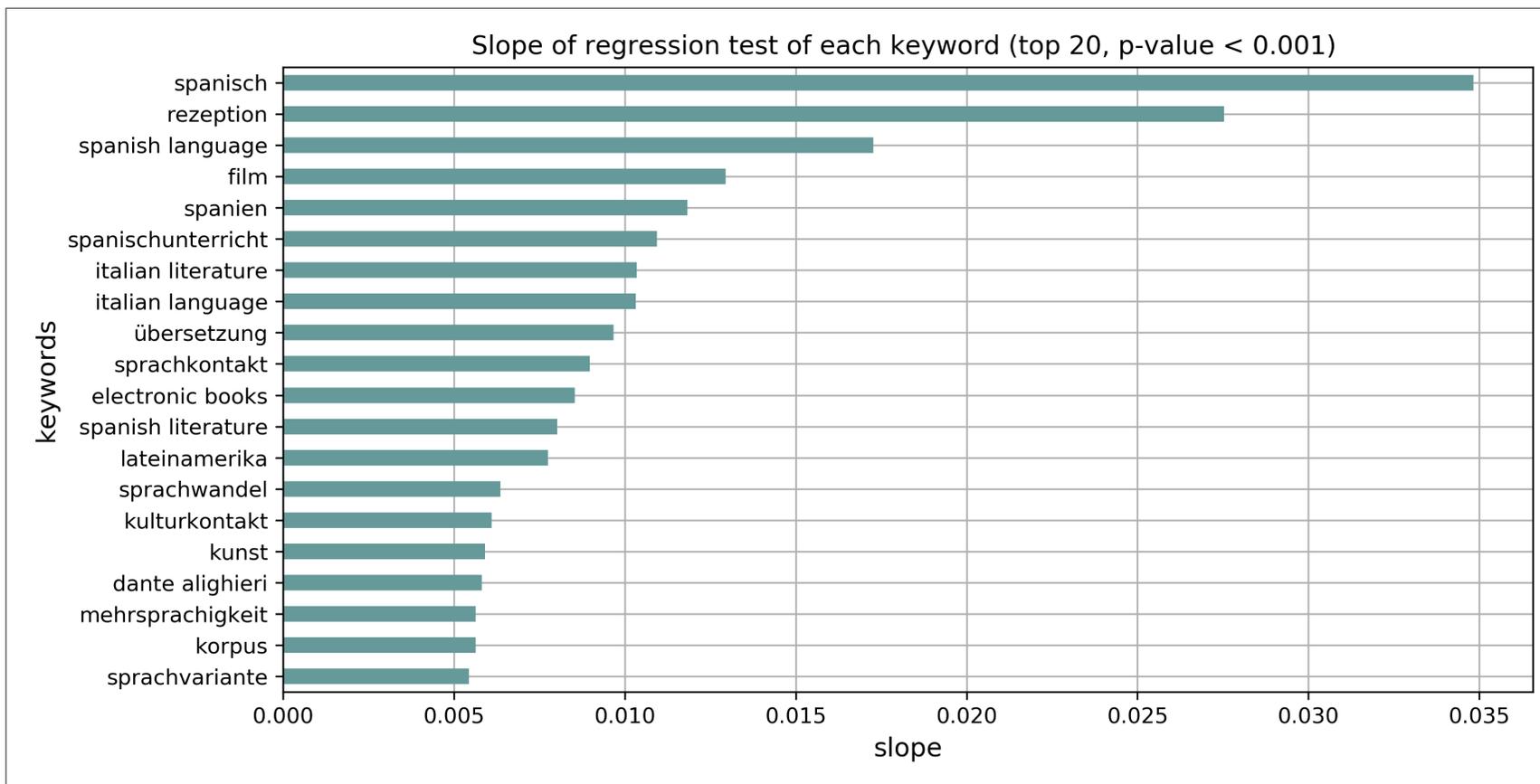
Klassifikation: Dewey Decimal Classification: 808

 Links zum Titel: [Inhaltstext](#)



- Intellektuelle Annotation durch Fachleute

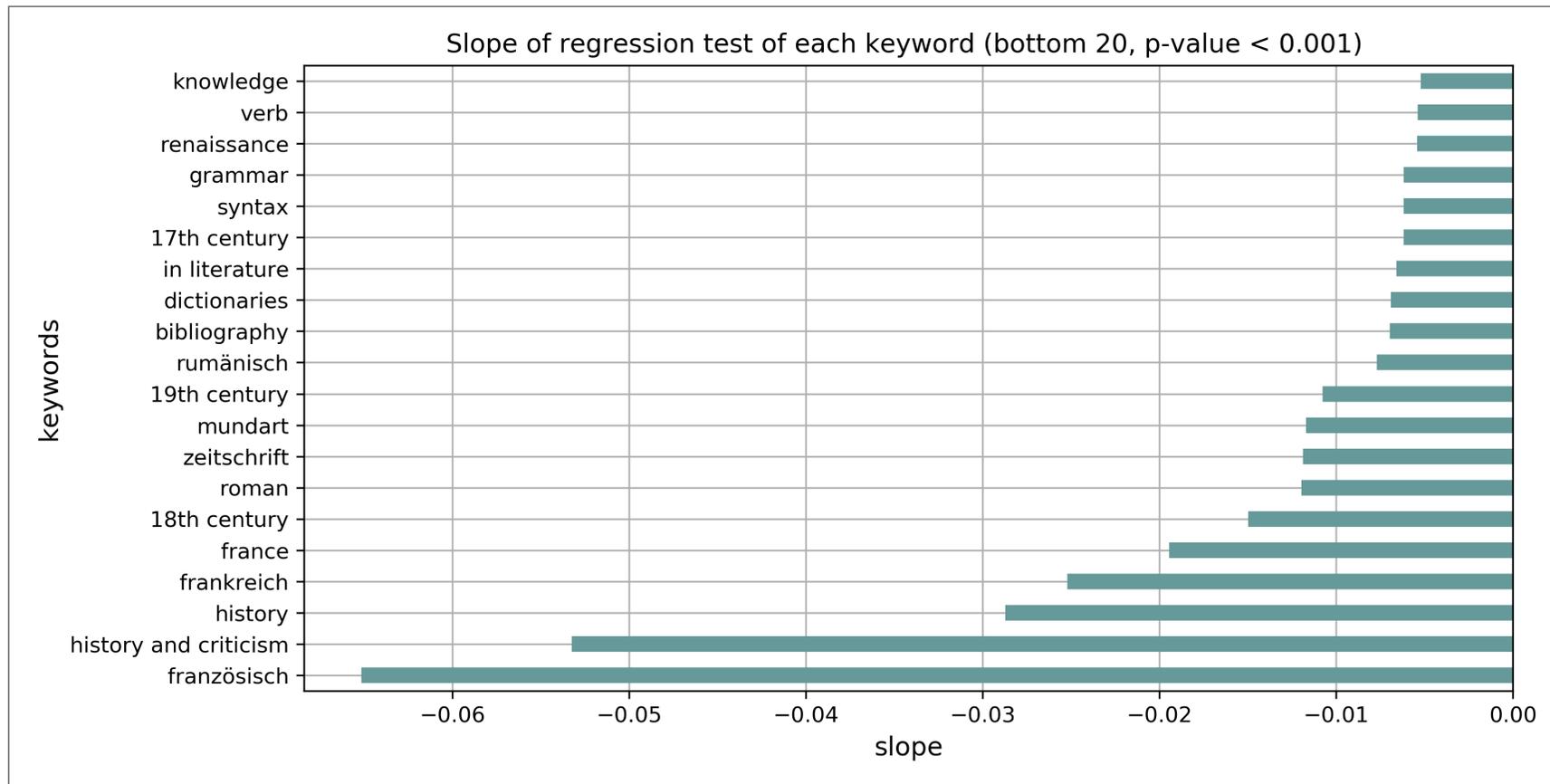
# SCHLAGWÖRTER: SLOPE (TOP)



- Spanisch, Italienisch
- Rezeption, Film, E-Books, Korpus
- Sprachwandel, Sprachkontakt, Mehrsprachigkeit



# SCHLAGWÖRTER: SLOPE (BOTTOM)



- Auf und über Französisch und Rumänisch
- Klassische Gattungen und Medien
- 15. bis 19. Jahrhundert



# SCHLUSSFOLGERUNGEN

## SCHLUSSFOLGERUNGEN: TENDENZEN

- Entwicklung der Themen
- Längere Publikationen
- Weniger Veröffentlichungen in französischer Sprache
- Mehr auf **Deutsch**, Spanisch und Englisch
- Ort der Veröffentlichung im deutschsprachigen Raum
- Gedrucktes immer noch sehr dominant
- E-Books nehmen langsam zu
- E-Books sind deutlich teurer

## **SCHLUSSFOLGERUNGEN: DATENSATZ**

- Enormes Potenzial von Daten aus Katalogen
- Kombination
  - Nutzung von kuratierten Daten
  - Anwendung von quantitativen Methoden
- Mit realistischen Erwartungen
- Fachkenntnisse von mehreren Bereichen



Library

Qualitative  
Data

Quantitative  
Approaches

## **DANKE AN ALLE, DIE DIE DATEN IM KATALOG**

- eintragen
- bearbeiten
- annotieren
- überprüfen
- (sach)erschließen
- verwalten
- konvertieren
- ...

