



The Incredibly Differentiated Labor Market Evidence from Job Offers

Nenad Pantelic, Vera Maria Charvat, Raven Adams, Saranya
Balasubramanian, Jakob Sonnberger, Hanna Misera,
Jörn Kleinert, Manfred Füllsack, Georg Vogeler

“What’s Past is Prologue”

NewsEye International Conference, 17th of March, 2021

Why Job Ads?

1) Job advertisements are an indicator for the role of labor in society

2) Job advertisements are sources for the change of perception of labor and job market organisation



Digital methods

- Data generation
- Data analysis

Goals:

- Creating a unique, reusable data set
- Illustrating the strong change in labor relationships in many dimensions
- Develop digital methods for text mining and text analyses

Methodology

1) Job advertisements are an indicator for the role of labor in society

2) Job advertisements are sources for the change of perception of labor and job market organisation

Digital Methods



Data generation:

- finding relevant pages
- region detection
- OCR

Data Analysis (Machine Learning):

- create ground truth
- evaluate existing methods

Problems with data generation

Goals:

- Creating a unique, reusable data set
- Illustrating the strong change in labor relationships in many dimensions
- Develop digital methods for text mining and text analyses

ANNO Repository

Österreichische
Nationalbibliothek

ANNO
Historische Zeitungen
und Zeitschriften



Über ANNO

[Was ist ANNO?](#)

[10 Jahre ANNO](#)

[FAQs](#)

[Suchen in ANNO](#)

[Drucken aus ANNO](#)

[Kooperationspartner](#)

[Kontakt](#)

[Impressum](#)

ANNO - AustriaN Newspapers Online

Historische österreichische Zeitungen und Zeitschriften online.



Listen

- [Alphabetische Liste der Zeitungen und Zeitschriften](#)
- [Jahresübersicht der Zeitungen und Zeitschriften](#)
- [Thematischer Einstieg](#)



Suchen

- [ANNO-Suche: Volltextsuche in historischen Zeitungen und Zeitschriften](#)

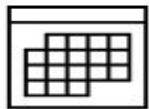
1689 - 1947

Anno Corpus



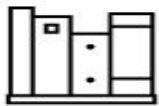
22 M. pages

=



1,5 M. editions

=



1000 titles



93 % are in German



71 % between 1821 - 1918



42 % printed in Vienna



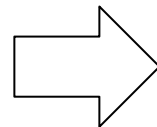
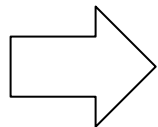
95 % are newspapers

Our Corpus

Only **40** journals have more than **10.000** editions.

Title	years	editions
Wiener Zeitung	146	46.223
Innsbrucker Nachrichten	92	27.536
Neue Freie Presse	76	26.650
(Linzer) Tages-Post	80	23.291
Linzer Volksblatt	73	21.364
Neues Wiener Tagblatt	63	21.074
(Neuigkeits) Welt Blatt	70	20.773
Salzburger Volksblatt	72	19.844

Finding pages with job ads



1. Step - Manual Selection

- Browsing pages
- Keyword search
- Publication patterns

2. Step - Algorithmic Detection

- Keyword search -> 70 % - 90 % accuracy in identifying pages with job ads
- Problems:
 - No identification of job ads on the page: bad OCR, missing keywords
 - Findings due to coincidence

3. Step - Improve Detection

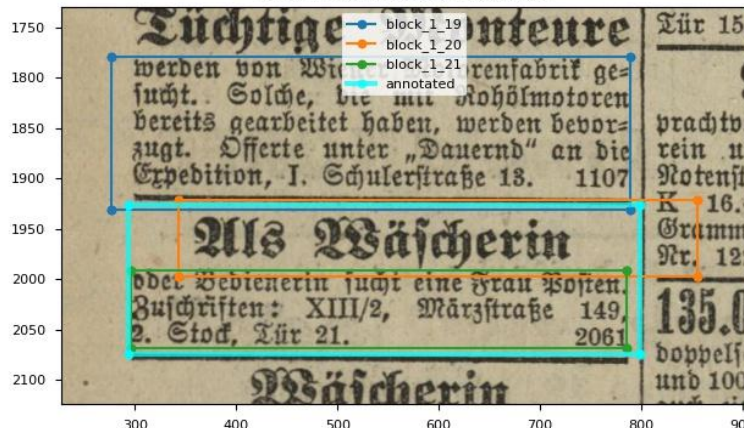
- Train classifiers combining dirty OCR and images
- To be done...

Layout detection: Tesseract standard model

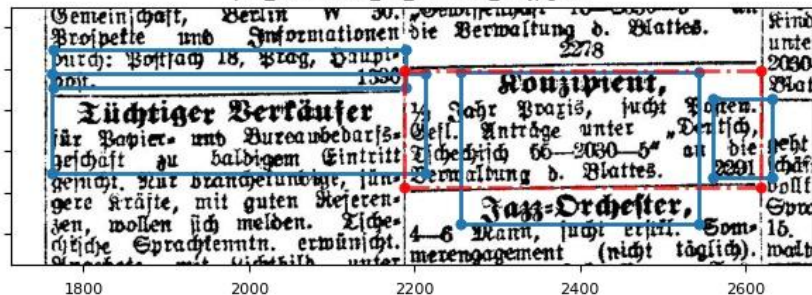
Segmentation - Too Big

Segmentation - Broken regions

aze_19130917_12_OFFER_4.jpg (3)

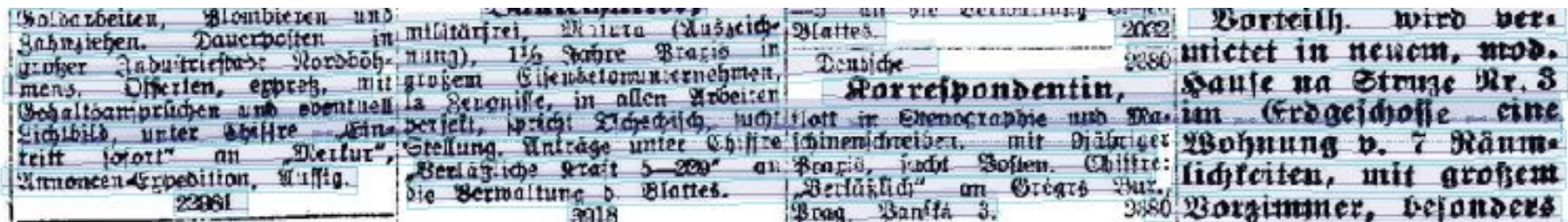


ptb_19300531_15_SEARCH_28.jpg (4)



Layout detection: Transkribus

Transkribus experiment looks better, but ...



ptb_19300531_15.jpg: text line crossing four columns

... or no headings at all
e.g. Innsbrucker Nachrichten (1854 - 1945)

Mineral-Öel-Lampen
(ihres sehr hellen und billigen Lichtes wegen Jedermann empfehlenswerth)
sind in verschiedener Form vorrätig bei
12 Conrad Stocker, Museumsstraße.

Kapital-Gesuch. Es werden 300 bis 400 fl. Kapital gegen
Georgi-Ziel ~~anzunehmen~~ gesucht. Das Nähere bei der Expedition dieses Blattes. 13

Es wird ein braver, ordentlicher Hausknecht gesucht. Wo? sagt die Expedition. 23

K. K. National-Theater in Innsbruck.

9. Vorstellung im VII. Abonnement.

Donnerstag, den 24. Jänner 1861:

Erste Gastrolle des Herrn Schulz

ehemaliges Mitglied des k. k. Hofburgtheaters in Wien.

Der Erbsöster,

oder

Der Mord im heimlichen Grund.

Schauspiel in 5 Akten von Otto Ludwig.

Der Fettel hiezu als Beilage.

Freitag, den 25. Jan. Zum Benefiz des Herrn Kapellmeister Czerny
die Oper: „Don Juan.“

Geist der:		Telegraphische Börsen-Course in Wien.	
Geld-Sorten.		Am 23. Jänner 1861.	
Am 21. Jan. 1861.		In öherr. W. zu 5% für 100 fl.	
R. Kronen	20.65	Metaliques	62.90
R. Ming-Ducaten	7.14	Nationalanleihen	75.
R. Rand-Ducaten	7.13	Actien der Nationalbank per Stück	724.
Napoleon'sche	12.06	Actien der Credit-Anstalt für Handel und	
Souveräin'sche	20.90	Gewerbe	157.50
Russische Imperiale	12.38	Silber ?	150.50
Preussische	12.80	Ponden 10 Pfund Sterling	150.70
Englische Sovereigns	16.16	Kaiserliche Münzkassaten	7.13
Preuss. Kassen-Anweisungen	2.28		

Frankfurt a. M. 22. Jan., Wien: 76%.

Verantwortlicher Redakteur: G. Hauschild.

Verlag der Wagner'schen Buchhandlung. — Druck der Wagner'schen Buchdruckerei.

Es wird ein braver, ordentlicher Hausknecht gesucht. Wo? sagt die Expedition. 23

1861-01-24, p. 8

Banern-Theater in Pradl.

(Lobronhofer Hof.)

Samstag den 31. Mai wird in dem besagten von berühmten Theater aus der 1. u. 2. Klasse. Der Kaiserin Elisabeth (Kaiserin Elisabeth) in der 1. u. 2. Klasse. Der Kaiserin Elisabeth (Kaiserin Elisabeth) in der 1. u. 2. Klasse.

Kaiserin Elisabeths Tod

Der Tod der Jüdin.

Die Kaiserin Elisabeths Tod. Der Kaiserin Elisabeths Tod. Der Kaiserin Elisabeths Tod. Der Kaiserin Elisabeths Tod. Der Kaiserin Elisabeths Tod.

Die Kaiserin Elisabeths Tod. Der Kaiserin Elisabeths Tod. Der Kaiserin Elisabeths Tod. Der Kaiserin Elisabeths Tod. Der Kaiserin Elisabeths Tod.

Morgen Sonntag nach dem Theater

Concert

der Kaiserin Elisabeths Tod. Der Kaiserin Elisabeths Tod. Der Kaiserin Elisabeths Tod. Der Kaiserin Elisabeths Tod. Der Kaiserin Elisabeths Tod.

Braves Mädchen

Die Kaiserin Elisabeths Tod. Der Kaiserin Elisabeths Tod. Der Kaiserin Elisabeths Tod. Der Kaiserin Elisabeths Tod. Der Kaiserin Elisabeths Tod.

Zu verkaufen

Die Kaiserin Elisabeths Tod. Der Kaiserin Elisabeths Tod. Der Kaiserin Elisabeths Tod. Der Kaiserin Elisabeths Tod. Der Kaiserin Elisabeths Tod.

Bediende, geübte Helferinnen

Die Kaiserin Elisabeths Tod. Der Kaiserin Elisabeths Tod. Der Kaiserin Elisabeths Tod. Der Kaiserin Elisabeths Tod. Der Kaiserin Elisabeths Tod.

Offene Stellen

Die Kaiserin Elisabeths Tod. Der Kaiserin Elisabeths Tod. Der Kaiserin Elisabeths Tod. Der Kaiserin Elisabeths Tod. Der Kaiserin Elisabeths Tod.

Zu Vermietten

Die Kaiserin Elisabeths Tod. Der Kaiserin Elisabeths Tod. Der Kaiserin Elisabeths Tod. Der Kaiserin Elisabeths Tod. Der Kaiserin Elisabeths Tod.

Fremden-Concerte des Orchesters der Innsbrucker Musikkapelle.

Wolfgang R. Sperr.

Concerte des Herrn R. Sperr. Concerte des Herrn R. Sperr. Concerte des Herrn R. Sperr. Concerte des Herrn R. Sperr. Concerte des Herrn R. Sperr.

Concert-Anzeige.

Concert-Anzeige. Concert-Anzeige. Concert-Anzeige. Concert-Anzeige. Concert-Anzeige.

Agnes Banger,

Agnes Banger. Agnes Banger. Agnes Banger. Agnes Banger. Agnes Banger.

Wilten. Vis-à-vis dem Gasthof „Templ“ Wilten.

englisch-amerikanische

Schiff-Schaukel

Schiff-Schaukel. Schiff-Schaukel. Schiff-Schaukel. Schiff-Schaukel. Schiff-Schaukel.

Gasthaus-Empfehlung.

Gasthaus-Empfehlung. Gasthaus-Empfehlung. Gasthaus-Empfehlung. Gasthaus-Empfehlung. Gasthaus-Empfehlung.

„Goldenen Bräut“

„Goldenen Bräut“. „Goldenen Bräut“. „Goldenen Bräut“. „Goldenen Bräut“. „Goldenen Bräut“.

George Kappold,

George Kappold. George Kappold. George Kappold. George Kappold. George Kappold.

Sonntag den 31. Mai

Garteneröffnung beim Neuwirt in Hötting,

Garteneröffnung. Garteneröffnung. Garteneröffnung. Garteneröffnung. Garteneröffnung.

Offene Stellen

Offene Stellen. Offene Stellen. Offene Stellen. Offene Stellen. Offene Stellen.

Working with



SuperAnnotate

Default Class	
h	heading
o	offer
s	search

Braves Mädchen
mit guten Zeugnissen sucht als Küchenmädchen oder zu Kinder auf 1. oder 15. Juni Stelle. Näh. im Ann.-Bur. Birchner.

Tüchtige, geübte Helferinnen
im Kleidermachen werden aufgenommen Museumstraße 10, dritten Stock.

Offene Stellen
1 oder 2 Mann, die zu Wasserleitungsarbeiten praktisch sind, finden Beschäftigung. Näh in Winklers Ann.-Bur. unter Nr. 1154.
Selbständige Köchin wird zu einer kleinen Familie ohne Kinder neben Stubenmädchen auf 15. Juni gesucht. Näh Maria Theresienstraße 8, dritten Stock. 1155—311

```

"instances": [
  {
    "type": "bbox",
    "classId": 3,
    "probability": 100,
    "points": {
      "x1": 192.6,
      "x2": 1069.7,
      "y1": 1891.7,
      "y2": 2152.3
    },
    "groupId": 0,
    "pointLabels": {},
    "locked": false,
    "visible": true,
    "attributes": []
  }
]
    
```

e.g. Innsbrucker Nachrichten

1896-05-30, p. 8

Further OCR attempts: Tesseract

Basic German and ONB trained models in Tesseract → works best

Fraktur!

Tischlergehilfen
werden aufgenommen in der Möbel-
fabrik Siegm. Oppenheim, XVII. Bez.,
Comeniussgasse 8. 1832

Long s - f

Basic German model

Tischlergehilfen

werben. aufgenommen im der Möbel:
fabrik Siegm, Oppenheim, XVII. Bez,
Comeninsnasse 2. . 1832 |

ONB trained model

Tischlergehilfen

werden aufgenommen in der Möbel-
fabrik Siegm. Oppenheim, XVII. Bez.,
Comeniugasse 3. 1832

Text Mining as OCR support

Cleaning messy OCR results with FASTtext + Jaro-Winkler similarity (JW):

Assumptions:

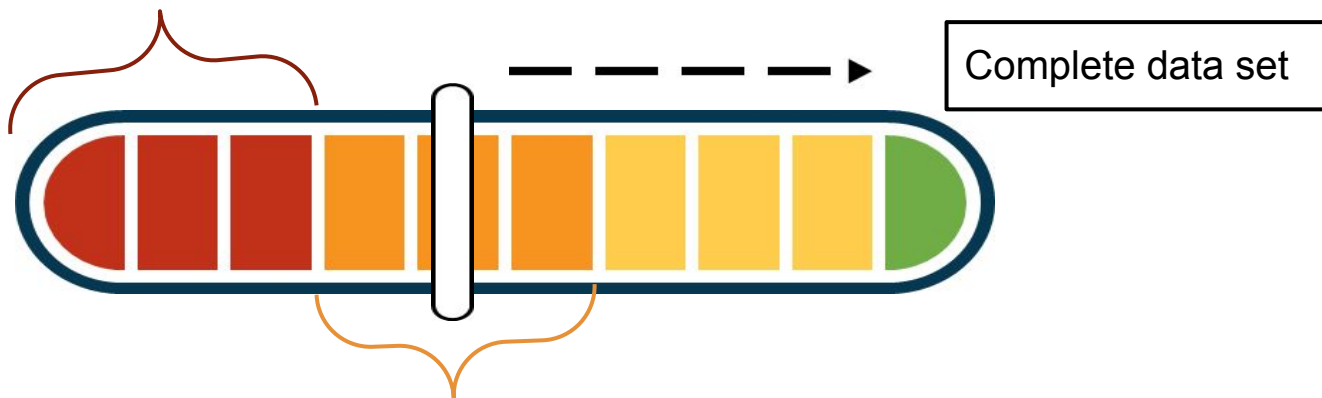
- a. mis-recognized words appear in the same context as the correct words
- b. correctly spelled words are dominant in the word embeddings
- c. FASTtexts character based embeddings further aid in finding correct words

Process:

1. crosscheck OCR results with german dictionary → get misspelled words
2. correct single-letter mistakes with high JW threshold (0.98)
3. generate word embeddings with FASTtext
4. calculate vector similarity for misspelled words → apply JW to suggestions
5. insert corrections into corpus → use cleaned corpus for classification task

Creation of a new, unique text corpus that is useful for digital research on the labor market and the perception of work *is* possible:
e.g. from “Mädchen mit besten Zeugnissen” to
“praktischer Ökonom / Praktikant mit schöner Schrift”

Conclusion



- Go beyond the standard Tesseract OCR/layout detection
- Word embeddings + edit distance is an interesting quality enhancement method



Thank you for your attention!



ÖAW

AUSTRIAN
ACADEMY OF
SCIENCES

