

A Service of

ZBW

Leibniz-Informationszentrum Wirtschaft Leibniz Information Centre for Economics

Hertweck, Friederike; Jonas, Lukas; Thome, Boris; Yasar, Serife

Research Report RWI-UNI-SUBJECTS: Complete records of all subjects across German HEIs (1971-1996)

RWI Datenbeschreibung

Provided in Cooperation with: RWI – Leibniz-Institut für Wirtschaftsforschung, Essen

Suggested Citation: Hertweck, Friederike; Jonas, Lukas; Thome, Boris; Yasar, Serife (2024) : RWI-UNI-SUBJECTS: Complete records of all subjects across German HEIs (1971-1996), RWI Datenbeschreibung, RWI - Leibniz-Institut für Wirtschaftsforschung, Essen

This Version is available at: https://hdl.handle.net/10419/307990

Standard-Nutzungsbedingungen:

Die Dokumente auf EconStor dürfen zu eigenen wissenschaftlichen Zwecken und zum Privatgebrauch gespeichert und kopiert werden.

Sie dürfen die Dokumente nicht für öffentliche oder kommerzielle Zwecke vervielfältigen, öffentlich ausstellen, öffentlich zugänglich machen, vertreiben oder anderweitig nutzen.

Sofern die Verfasser die Dokumente unter Open-Content-Lizenzen (insbesondere CC-Lizenzen) zur Verfügung gestellt haben sollten, gelten abweichend von diesen Nutzungsbedingungen die in der dort genannten Lizenz gewährten Nutzungsrechte.

Terms of use:

Documents in EconStor may be saved and copied for your personal and scholarly purposes.

You are not to copy documents for public or commercial purposes, to exhibit the documents publicly, to make them publicly available on the internet, or to distribute or otherwise use the documents in public.

If the documents have been made available under an Open Content Licence (especially Creative Commons Licences), you may exercise further usage rights as specified in the indicated licence.



WWW.ECONSTOR.EU



Datenbeschreibung

RWI - Leibniz-Institut für Wirtschaftsforschung

RWI-UNI-SUBJECTS: Complete Records of All Subjects Across German HEIs (1971–1996)

November 2024

Friederike Hertweck Lukas Jonas Boris Thome Serife Yasar

Impressum

Herausgeber:

RWI – Leibniz-Institut für Wirtschaftsforschung Hohenzollernstraße 1–3 | 45128 Essen, Germany

Postanschrift: Postfach 10 30 54 | 45030 Essen, Germany

Fon: +49 201-81 49-0 | E-Mail: rwi@rwi-essen.de www.rwi-essen.de

Vorstand Prof. Dr. Dr. h. c. Christoph M. Schmidt (Präsident) Prof. Dr. Thomas K. Bauer (Vizepräsident) Dr. Stefan Rumpf (Administrativer Vorstand) Prof. Dr. Kerstin Schneider (Mitglied des erweiterten Vorstands)

© RWI 2024

Der Nachdruck, auch auszugsweise, ist nur mit Genehmigung des RWI gestattet.

RWI Datenbeschreibung

Schriftleitung: Prof. Dr. Dr. h. c. Christoph M. Schmidt Gestaltung: Magdalena Franke, Claudia Lohkamp

RWI-UNI-SUBJECTS: Complete Records of All Subjects Across German HEIs (1971–1996)

November 2024 Friederike Hertweck Lukas Jonas Boris Thome Serife Yasar

Datenbeschreibung

RWI – Leibniz-Institut für Wirtschaftsforschung

RWI-UNI-SUBJECTS: Complete Records of All Subjects Across German HEIs (1971–1996)

November 2024

Friederike Hertweck Lukas Jonas Boris Thome Serife Yasar



Contents

1	Overview and Analytic Options	2
	1.1 Introduction	2
	1.2 Analytic Options	3
	1.3 Descriptives	3
	1.4 Data access	5
2	Data preparation	5
	2.1 Data source	5
	2.2 Automated extraction of tables	7
	2.3 Additional data preparation	12
	2.4 Data quality and limitations	14
3	References	15
A	Academic Coding System	17
в	Details on temporal evolution of institutions	19
	B.1 Number of subjects and institutions over time	19
	B.1.1 Number of subjects over time	19
	B.1.2 Number of institutions over time	20
С	Error rates per year	21
D	Codebook	22

1 Overview and Analytic Options

1.1 Introduction

The dataset RWI-UNI-SUBJECTS is based on the guide on "Study and Career Choice" (in German: "Studien- und Berufswahl" and recently "Studienwahl"). The guides provide high school graduates and prospective students with comprehensive information on choosing a higher education program in Germany. The book is considered as being the official guide to studying in Germany and is provided and distributed on behalf of the Federal Employment Agency. The book contains detailed information on higher education in Germany, including information on institutions, subjects, study opportunities, financial aid, student housing, and related topics.

The first guide was published in 1971 and has been updated yearly. Based on digitized copies of the guides covering the years 1971 to 1996,¹ the dataset documents all higher education institutions across Germany during this period, along with the array of subjects available for study at these institutions on a yearly basis. Notably, the dataset RWI-UNI-SUBJECTS encompasses universities and universities of applied sciences (UAS, in German: *"Fachhochschulen"*) as well as additional higher education institutions such as colleges of arts and music, theological colleges, and the universities of the German Armed Forces (in German: *"Universität der Bundeswehr"*).

Figure 1 illustrates the geographic distribution of higher education institutions for the years 1971 and 1996. It is important to note that data pertaining to the higher education landscape of the German Democratic Republic (GDR) is not available. As a result, the eastern states of Germany were integrated into the dataset in 1991, the year subsequent to reunification. Any information regarding locations in Eastern Germany before 1991 is thus absent from the dataset. Moreover, alongside the establishment of higher education institutions in the newly formed federal states following reunification, Figure 1 also highlights a notable surge in the number of higher education institutions within the old federal states, particularly evident in regions such as Lower Saxony and Bavaria.



Figure 1: Higher education locations in Germany in 1971 and 1996

Note: Data on the higher education landscape of the German Democratic Republic is not available. Thus, data pertaining to institutions in the Eastern states is only available from the year 1991 onwards.

 $^{^{1}}$ Our dataset covers the year from 1971 to 1996, as there exists a structural break in the design of the study guidance books. Extracting the same information from the books from 1997 onwards needs another methodology. We are currently working on this.

Depending on the year, it is possible to distinguish 87 to 164 different subjects (in German: "Studienfach"). These are grouped into up to 58 different subject areas (in German: "Studienbereiche") and eight broad subject groups (in German: "Fächergruppe"). These subject groups are: Humanities (01), Sports (02), Law, Economics and Social Sciences (03), Mathematics and Natural Sciences (04), Human Medicine and Health Sciences (05), Agricultural, Forestry, and Nutritional Sciences, as well as Veterinary Medicine (07), Engineering (08), and Arts and Art Sciences (09). The aggregation of subjects follows the academic coding system provided by Destatis (2023). Please also refer to Appendix A for a more detailed description of the fields.

1.2 Analytic Options

The RWI-UNI-SUBJECTS exhibits significant potential for addressing various research questions due to its comprehensive coverage of the expansion of higher education in Germany over 25 years. The dataset offers insights into the evolution of the German higher education landscape and tracks fluctuations in the number of higher education institutions across both time and regions, capturing the dynamic progression of academic subjects. Thus, it enables nuanced analyses of the expansion of specific subjects, subject areas, and subject groups, thereby facilitating examinations of institutional dynamics, subject trends, and regional shifts starting from the 1970s onwards.

Furthermore, the dataset provides opportunities for linkages with other datasets through its temporal and regional attributes. In order to facilitate regional analyses, the official municipality code (in German: "Amtlicher Gemeindeschlüssel, AGS") has been integrated into the dataset. The municipality code enables seamless connection with educational panel datasets, such as the National Educational Panel Study (NEPS-Netzwerk, 2023), the Sample of Integrated Employment Biographies (SIAB) provided by the Federal Employment Agency (IAB) (Schmucker, Seth, & vom Berge, 2023), or the DZHW Graduate Panel 1989 (German Centre for Higher Education Research and Science Studies (DZHW), n.d.). These linkages enhance the dataset's utility and extend its potential for comprehensive research endeavors across various domains.

To the best of our knowledge, only two similar datasets exist to RWI-UNI-SUBJECTS. The first is the College Scorecard dataset. It contains comprehensive institution-level data from 1996 to 2023 and subject-level data from 2014 to 2020 on enrollment, financial aid, costs, debt, repayment, and post-graduation earnings. It also includes crosswalk files that link colleges' identification codes (OPEID) with identifiers from the Integrated Postsecondary Education Data System (IPEDS UnitID). The second dataset, the Catalogue of First and Second Cycle Degree Programmes of the University of Bologna provides detailed program data by academic year, with resources covering the academic year from 2004/2005 to 2020/2021, but exclusively from the University of Bologna. Neither dataset fully covers or substitutes for the RWI-UNI-SUBJECTS dataset, as both provide information on only a subsample, cover different periods, and do not address the years of higher education (2024) and University of Bologna (2024).

1.3 Descriptives

Between 1971 and 1996, both the number of higher education institutions and the number of subjects increased. The upper panel of Figure 2 illustrates the temporal evolution of the total number of subjects available at universities and UAS. In 1971, students had a choice among 1,349 subjects at universities and 586 subjects at UAS. By 1990, these figures had risen to 3,389 subjects at universities and 959 subjects at UAS. By 1996, the total number of subjects had further increased to 6,056 at universities (including locations in the former GDR) and 1,512 subjects at UAS.

The lower panel of Figure 2 illustrates the temporal evolution of the total number of higher education institutions. In 1971, students could enroll in one of 44 universities, a figure that had increased to already 61 by 1990. By 1996, the number of universities in the data had further risen to 103. Any information on UAS was aggregated at the city level in the guides until 1985, after which they were separately listed. By 1996, the number of UAS increased to 198. Notably, even excluding the spike in 1991 attributable to the integration of the former GDR, the data shows a noteworthy expansion both in terms of institutions and subjects. The substantial increase in the number of subjects in 1996 can be attributed to a structural shift in teacher training and the reclassification of institutions of type "Other HEI" as universities in 1996, also affecting the number of universities in 1996.

Overall, the dataset encompasses all study options documented in the guides on "Study and Career Choice". It is, however, possible that the guides' editors inadvertently omitted certain study options. Thus, any mistakes in the official guides may have been brought forward into the dataset and are discussed in section 2.4.



Figure 2: Temporal evolution of higher education landscape

Type 🗝 Municipalities with UAS 🗝 UAS 🛥 Uni

Note: The upper panel presents the number of subjects separately for UAS and universities. The lower panel provides the number of higher education institutions over time. In the lower panel, universities are represented in yellow, UAS in dark

blue, and aggregated UAS information at the city level is depicted in light blue until 1985. Prior to 1986, if two or more UAS were situated in the same city, the guides aggregated the information at the city level without distinguishing between individual UAS. However, a clear differentiation between separate UAS is provided from 1986 onward. Underlying numbers are available in Tables 1 and 2 in Appendix B.1.

1.4 Data access

The dataset is available as a public-use file via the Research Data Center Ruhr (FDZ Ruhr). Forwarding the data without the consent of the FDZ Ruhr is prohibited. The FDZ Ruhr expects the dataset to be used responsibly. The dataset must be cited as follows:

Hertweck, F., Jonas, L., Thome, B., and Yasar, S., 2024. RWI-UNI-SUBJECTS: Complete records of all subjects across German HEIs (1971 - 1996). RWI-Micro. Version: 1. RWI – Leibniz Institute for Economic Research. Dataset. doi:10.7807/studi:buch:suf:v1.

2 Data preparation

The creation of the RWI-UNI-SUBJECTS dataset involves several key steps. It begins with converting PDF tables into images, followed by the use of OpenCV to recognize table grids and symbols. To ensure the integrity of the extracted information, extensive accuracy checks are conducted, including visual reviews and comparisons with original sources. Furthermore, additional data preparation steps involve correcting the names of higher education institutions and subjects, aligning them with unique identifiers from the German Federal Statistical Office, and incorporating official municipality codes to facilitate seamless connection with other datasets.

2.1 Data source

The dataset RWI-UNI-SUBJECTS is based on tables from the guides on "Study and Career Choice" (in German: "Studien- und Berufswahl") from 1971 to 1996, provided by the Federal Employment Agency (Hirschfeld, 1971, 1972; Bock, 1973, 1974, 1975, 1976, 1977, 1978, 1979, 1980, 1981, 1982, 1983, 1984, 1985, 1986, 1987, 1988, 1989, 1990, 1991, 1992, 1993, 1994, 1995, 1996). These guides have been digitised and stored in Portable Document Format (PDF for short).

These guides inform high school students towards the end of upper secondary school about higher education and occupational career paths in Germany. Every year, high schools distribute these guides to their students free of charge. At the time of publication, the Federal Employment Agency (in German: *"Bundesagentur für Arbeit"*) describes these guides as the "official study guide for Germany" (in German: *"offiziellen Studienführer für Deutschland"*). Due to its color, it is also known as "the green book" (see Figure 3).



Figure 3: Cover pages of the guides from 1971, 1986, 1994 and 1996

Note: A sample of cover pages. These are based on Hirschfeld (1971) and Bock (1986, 1994, 1996).

The guides have been distributed to high school students from 1971 onwards. The state monopoly mandate for guidance on career choice (in German: "staatlicher Monopolauftrag einer öffentlichen Berufsberatung") was confirmed in the Labor Promotion Law (in German: "Arbeitsförderungsgesetz") on 25 June 1969 and transferred to the Federal Labor Agency (formerly known as "Bundesanstalt für Arbeit", now "Bundesagentur für Arbeit"). This law created the institutional framework for uniform organization and development of career guidance in Germany (Meisel, 1978). Since the first edition in 1971, the guides to "Study and Career Choice" are annually updated, published, and distributed to high school students. Over the past 50 years, several publishers were involved in the publishing process.

Each of these guides provides comprehensive details on all subjects available across Germany and the higher education institutions offering them. Specifically, these guides comprise tables that summarize the availability of subjects across various locations (see Figure 4). In some instances, even second campuses of higher education institutions are listed separately within these tables.

These tables consistently adhere to a standardized structure across all books spanning the years 1971 to 1996. Columns correspond to locations, universities, or UAS (from books starting in 1986 onwards), while rows represent different subjects. The second column denotes the chapter in the book where additional information related to each topic can be further accessed. Information regarding the availability of a subject, program type (e.g., full-time or part-time), the starting term (e.g., winter or summer term), and entry requirements are succinctly summarized using different symbols within the tables. Figure 4 illustrates an extract from one of these tables from the year 1980 (Bock, 1980). Each table spans multiple pages.

Studienmöglichkeiten		8	den				6		÷.				şl	5 5		He	ien (Net	ersad	heer	ċ.	h	orde	÷	me	etale				1	1	t:
an Universitäten, Technischen Universitäten/Hochschulen und Gesamthochschulen		Ĩ		11	1				112-04				1	Hambu	11 June 10	H		2		2.									3	3		E I	Soldier
	3	12	21	dal.	1.1		bl-	48	1.2	.12	12	ы			31			. 17	2 -	32	Ы	2)2		Ы	3	23	2 1	3	3	a B	16	116	
	8	616	1	12	151	1	t i	10		11	- 2	10	212	215	18	15	23	11	34			5.5	65	12	.18	11	e 2.	8	218	3.3	181-	LR	1
Zutenungsbeschränkungen siehe 3.2.2 und zwiinflo	Same X	Pietonia Pietonia	-	Constant		1000	A que		1		11	Br.D.	11	1 miles	Hant	1.111	Configuration of the local section of the local sec	11	0		Harris		Autor	2	1	Dane of	1	1	100	1	Variation of		0.00
lefe in geningt in	6,40							10		0				6				0						П				0			10		11
leven-ingre	6.41	0				0			1.1	0		101	0	0		1.1		9	LK	1.1				1.2	0			10	0		0	10	LT
Amerikanustik.	6,40				21.1	0	10		Q8	0	0	0	0			108	0.1	00	1.3	1.1						C			O C	100	0		
logistik (Mapister/%	6.40	00		0	20	0	00	Sec.	108	0	90	0	010				0.0	00	C.R		101	1.1	010	201	2	0	0	0	00	00	1.0		101
Inglanik (Symnasium Sek Al)	6.40	010		101	10	đ. 1	0.0		100	01.1	00	101	010	00	L K	101	010	00	L K	<u> </u>	10	0.0	010	101	200	010	101	101	010	1007	1.10	1.5	101
Anglistik (Raulachule/Sek.()	6.40				10		140	6.5	100	0	00	10	0.0	00			20		1.15	1.L.	0	230	010	101	200	lo k	101	10	00	1010	1.10	1010	12
k-th-opologie	6.7					0				0		1				10			C.X												1.0		10
vbeitslehre Technik (Sek /~/i) 12	6.74		1.1					11		1			10	CL.		102	010			11		oI. Io			10					1.12	11	11	
rchielogie	6.41	00		11		01			100	0	_0	10	210	0		103		21.	1.10	1.1.	1.1			191	21.1			10	0	1.1	1.10	2010	101
vchitektur (Digilan)	6.1		1.1	<u>a 1</u>	-17				1.1	10			-12		1.16	a	10		1	1.1.	101		Q.,	1.1	10		11		- 0	191	K.,		11
rchitektur (Beruft 5./Sek./I)	6.1		1.1	11	11				11				-			11	11		1.1.		1.1		11	1.1			11				КЛ.:	1	11
stronomia, Astrophysik	6.43			11	-1-1	91			1.1	4	-	101	-	_ P	-		-	-	LK	5. J.	1.1	1.1		101	à		11	1	0			4	D.
any clopie	6,41	120	1.1	11	-1-2	Q.,			193	0			-	- 12	-	1.1	1.1	1	1.15	1.1.	1.4	-		1.1			1.1	-	-				11
aungenieurweiten (Diption)	6.2		1.15	0.1	-52	-			1.1	-15		-	-191		- 6	<u>n i</u>	-121	-12	1.	1.1.	151	-	12	1.4	-12		-10	-	-	1.12	M.,	4	11
aungenieurweisen Stleruft 5./Set.M	6.2		1.15	1	12			14	-	19			-19		1.15	4.1	-		1.1.	1.1.	19		21	11			121			1.12	1.	1	11
egtes (Opion)	6.3		1.1	11	-			14	-		-	- 1	-12			-	-		101	14	1.1	-	21.	1.1			1.1					1	11
erghau /Beruft 5./Sek.//J	6.3			11				1.1	1.1			1.1	1.1		-	1.1						1.1	21.	11			1.1	1.1				ц.	11
etratisavita/heftslehra (Dyslam)	6.4			1.1	20		0.0	221.	192	2	010	10	212	12	-	100	2	9	1.45	1.1.	1.1	-12	_10	1.1	-10	-	101	-101	0.0	2		100	12.
bliothekswissenschaft	6.5		1.1	11				11	191				21		1	1.1			11	11	1.1	1		11				10		1.1.			11
ochemie	6.6	120	1.1	191	1.1	21.		1.1	12	91.1		193	210	19	- 8	391	1	10	1.5	1	19		1.19	1.1	1		11	1.1			1.19	11	11
Lologie (Digitan)	6.7	120	19	10	15	25		121	12	99	19	2.3	21.1	12	1. K	22	2.1	-	6. K	1. J.	10	883	120	1.1	S.,	12	14	12	91		122		11
Lologie (Gymnasium/Sell.11)	6.7	100		101	12	257	1	12	10	90	10	1.1	210	120	5. K	20	20	540	6.19	1	12	23.7	120	10	20	10%	20	19	0		1010	1.15	11
Lologie (Realschule/Sel./)	6.7			11	11			12	15.7	91	12	127	010	510	1	572	20	- 10	4.19	1	10	235	100		0.0	100	22	12	0.0	2010	1.12	1.15	14
iotechnik (Beruff S. Sek II)	6.64		1.1	11	11			11				-			14	1.1			11	1.1.	191						152						11
autoren	6.29			11	11	1		11		10					14	1.1			11	11					1					11		1	11
enne-eitechnik	6.29			11				1.1							14	1.1					11												11
ytant motik	6.41								1.1	21			21.1			1.1					1.1				01			-101	2		1.12	1	
hema (Oyton)	6.8	120	1.1		12	25.	1	101	102	0.5	C	101	0.0	20		1	Q., .	100	1.4	1	10	QL	2.0	10		23	22	- 101	20	1929	1.5	1.15	12
Nama (Gymnasum/Sak.11)	6.8	120		201	12	1		191	157	20	19	191	212	10	1.15	10	919	100	1.4		150	21	229	10	22.	101	22	12	919	1929	1.10	1.19	121
hamia (Raulschule/Sek //	6.8			11	10			191.	10	0	19	101	010	90	1	101	0.0	10	196	1.1.	191	0	020	101	010	228	10	10	970	1010	10	1.19	121
anisch (Cymnesium/Sell II)	6.17			11									-			1.1	-	-			1.1	-											121
anisch /Realschule/Sek.//	6.17			11	1.1			11	1.1			- 1			1	11		_	11		1.1		1	1.4			11				1	1	101
lak trotachinik	6.10			<u>6 I</u>	157			11	19	19			19	1	1.15	1.1	10	10	1	11	19		1	191	12	1.5		0	10	107	15	1.15	11
Nektronechnik (Beruff 5./Sek./1) ¹⁰	6.10		1.1	ā.,	10				10	10			-12	Δ	1.16	2.1	.10	-			32	-	Ø.,	101	37	1.5	1.1		10	200	11.	1.10	1.
Instrungswissenschaft (Diplom)	6.20		101	11				1.1		10		- 1	1	10	1	1.1	2	-	1.1	1.1	1.1			1.2	g.,		1.1					1	12
mahrungswissenschaft (Beruff 5. (Set 11)	6.20							1.1				ΕT			- T	11					101			1.0								1.10	1.1

Figure 4: Example page of a subject table at universities (1980)

Note: Extract of the table of subjects at universities of the year 1980, based on Bock (1980). The structure of these tables is standardized throughout the years.

2.2 Automated extraction of tables

Pre-processing of tables

In the first step, the relevant pages of the PDF files were converted into image files. These image files were saved as *Portable Network Graphics* (PNG for short) with a pixel density of 300. Subsequently, Python-based image processing functions from OpenCV were utilized for further processing (Bradski, 2000).

The interior area of the tables within the image files was isolated and, if necessary, rotated to ensure that only the symbols and the table grid (without axis labels) were visible. This approach offers the advantage of preventing axis labels from being erroneously interpreted as table contents. The axis labels were separately recorded for each table and subsequently integrated into the grids through automated procedures.

Recognition of table grid and symbols

The extraction of the tables can be divided into two tasks: recognition of the table grid and extraction of the symbols. To recognize the **table grid**, OpenCV functions were employed, leveraging horizontal and vertical line recognition facilitated by a kernel. To circumvent thicker intermediate lines being identified as multiple lines, a minimum distance between lines was specified, a value that could vary depending on the respective guide. The identified lines were then arranged in ascending order based on their pixel coordinates and subsequently filtered using the minimum distance. This ensured that the area between

two lines corresponded to a row or column. The prior rotation of the tables (if necessary) ensured that horizontal lines commence and conclude at the same y coordinates. Similarly, vertical lines exhibit consistent x coordinates at their top and bottom, eliminating shifts due to rotation. This streamlines the storage of recognized line coordinates and simplifies the subsequent insertion of symbol positions into the table.

Figure 5: Illustration of the table structure



Note: The figure illustrates a schematic example of the structure of the tables. The higher education institutions are listed in columns, while subjects are listed in rows.

Figure 5 illustrates how the coordinates are stored efficiently as a tuple. With this data structure, each row tuple uses two y coordinates and each column tuple uses two x coordinates. For example, the tuple (x_1, x_2) would describe a column that starts at position x_1 and ends at position x_2 . The subsequent column (x_2, x_3) starts at the position x_2 , i.e., where the preceding column concluded, and terminates at position x_3 . The storage of row coordinates follows a similar structure, with y coordinates representing the height position. These tuples, describing the rows or columns, construct a representation of the table grid and are stored in the first row or column of a data frame.

The second step in the process of extracting the information from the tables involves recognizing the various symbols for which OpenCV's $template matching^2$ implementation was used.

Template matching entails identifying a small pattern or image within a larger image. To implement this procedure, sample images, known as *templates*, are necessary. The template matching algorithm compares a portion of the larger image with the template by systematically moving the template across the entire image. At each position, a comparison is conducted between the template and the corresponding section of the larger image, yielding a similarity value based on color values. If this similarity value surpasses a predefined *threshold*, the template is recognized at the corresponding position. Thus, for the template matching algorithm to operate effectively, at least one sample image is required for each symbol intended for recognition within the table. This sample image should ideally serve as the most accurate representation (e.g., same size, rotation, etc.) of the respective symbol.

²https://docs.opencv.org/4.x/d4/dc6/tutorial_py_template_matching.html

Figure 6: Sample of symbol templates



Note: The figure presents several symbols that appear in the tables. The shape of the symbol and the presence and position of the black box within each symbol signify specific attributes such as program type and starting semester.

The templates for the symbols were extracted from the original images (refer to Figure 6) and passed to the template matching algorithm. Given the diversity of symbols in the guides and variations in scan quality, templates were generated separately for university and UAS tables within each guide. In certain instances, employing multiple templates for symbols that posed recognition challenges proved advantageous. Due to the inherent similarity among some symbols, the algorithm selects the symbol with the highest similarity value.

The template matching algorithm returns the x and y coordinates of each recognized symbol along with the similarity value of the evaluated template. If the similarity value exceeds the specified threshold, an associated symbol code and similarity value are stored in the data frame. Data is organized using pre-determined row and column tuples, along with the coordinates of recognized symbols: if symbol coordinates (x, y) fall between column values (x_i, x_{i+1}) and row values (y_j, y_{j+1}) , the symbol is placed in column *i* and row *j*. This process iterates for each symbol, ensuring only the symbol with the highest similarity value is retained in the data frame.

A complication arose from smaller additional symbols, depicted as small dots adjacent to the main symbol (see Figure 6), that convey secondary meanings, such as *studies begin in the winter term* or *recommended to start in winter term*. Thus, a template icon for each symbol combination was supplied to the algorithm. However, these supplementary symbols, which accompanied a main symbol, were not consistently and precisely identified due to the variability in the spacing between the main symbol and its accompanying dot. Yet, even if the smaller symbol was not recognized, the recognition of the larger symbol remained correct.

In the final step, the algorithm compares the similarity scores for each icon in each cell. Only the numerical code with the highest similarity score is kept in the cells. If no icon template is recognized for a cell, the cell remains empty. This results in a table that holds up to one numerical code in each cell. Axis labels containing subject and location names are then reintegrated into the tables, replacing coordinate tuples for rows and columns. Given that overview tables span multiple pages, tables extracted from individual pages are concatenated into a single comprehensive table for universities and UAS, respectively.

Evaluation of the extraction process

The image-based extraction of the tables from the guides on "Study and Career Choice" proved to be a major challenge. The unique format of symbol-based tables necessitated a tailored technical solution. Given the dynamic presentation of tables in the guides, the algorithm had to be highly adaptable. Moreover, the quality of the PDF scans varied considerably, with common defects including blurring, rotation, overexposure, or contamination in the scan.

To ensure data quality, the outcomes of the table extraction underwent thorough examination throughout the extraction process. Multiple interim checks were conducted as follows in iterative processes:

- Visual Inspection of intermediate results: Visualizations of the intermediate results were generated, aiding in the identification of any inconsistencies
- Manual comparison between extractions and original: Extracted tables were compared against the original images to verify accuracy and completeness.

• Error Identification: Any errors or discrepancies were promptly identified and flagged for further investigation.

By incorporating these steps into the extraction process, potential issues were promptly addressed, thereby enhancing the overall reliability and accuracy of the extracted data.

Figure 7 provides an example of **visual inspection**. It illustrates the table grid of a particular table, highlighting in blue the horizontal lines detected by the algorithm. Similar visualizations were also utilized for vertical line detection. This visualization method facilitated rapid error analysis and was employed for every table recognition instance. In cases where lines were not detected, or an excessive number of lines were identified, algorithm parameters such as the minimum length of a line or the minimum distance between two lines were adjusted accordingly. The flexibility to adjust these parameters was crucial because the display format of tables was not entirely consistent over the years.

		1.16.20	101	1.11	0.01					0	0
0	101		101-1	0.0	 Q 		0		1.195	10(0)	00
	0 0	10.01	- DOI:	TO ICH OH	- E - E - E - E - E - E - E - E - E - E	 HOROWING 	0 101 -		1 101 107	1-101c201	0 0
00 00	000	ololo		1000	0.0	lololoto	0 0	0 100		lalolololo	0000
121121	10000	100.00	0.00406	A 50101	0.0710	a second second	101 101	- IOHOHO O	Protocol Activation	- IORHOHORO	00000
	105	CLOV -		000	0.010	lolold.	0 0	00000	ala la la la la la	lolololold	0000
	401	1.1.1	1111		101	101	0 0				0 KB
					0.0	Idiol/d		101 101	101.	1.10	
00	0 0		001	10102	0.0	12101 0	1.1.62		10101	10101	0000
1 1 107	of I		1.1.1.1		0	(A) 1 KA	105	01 07	0.0	I ROLOR	
10.0	53		1 1 1 1	1.0	1 1		0		0303	121	0
1.38/31	10		10101-1		1.0		101				
	0.1		1 1 163		0	KA 1072	10	0.10	NY 105- 10	10	
101	101		1.1.101		1.1	C C		101 101	1 1 1 1 1 1		
					0		101	1.12			
								0			
	00-1-1	LICA.II	COD R	1000	0.0	- OKE -IL	0	1 1 10	0 - D	101010101	000
	1	1.11	100	1.1.1.0						01111	
102.110	0		1000	1.10 01	11 11	KACK J. L.	80F 10F	. icf. 1	Charles and the second		0 0
shinks to be	A- ICTOR	1.00	101050	1.102101	1.9621	REAL REAL REAL	0 0	a front for	COLORED COL	1000	10 10 10
of the loss	 Idokt 	1.104	1.5.5.5	0.00	100	Michael Contraction	ALC: NOT	Cloid	CINE ACCOUNT	1050	0.01 0.01
	1.6	1.105	100.0	101071014	0.000	etches.	0.0	a figure of	CONTRACTOR OF	10.000	0 0 0
								100	10		
			1 1 1 1								
			1.201.1	1.1.101					100	10101 111	IOL III
CO 010	5 Infold	0	inichel.	10.00	0.0101	00000	0000	0.00	and a state to the t	Distance.	0.01 0.01
a final statement	8-1-6-6-4	1.101		1250500	0000	NOICE NO.	0000	10501 12	and the second second	the second se	0.00
	0.1	1.0	10105 1	100000	0.0501	OLOUG .	CUCIO	0.100	all design for the first	101200-00	0 00
						121212	-				
1.	1.6		1.5. 1.5		1	0.0	1.1	0.0	10. 10. 10.	0	6 6
1.1.6	1.4		1.6 1.6		0.0	0.00		10 10	18 18 18 1	100.00	1 1 1
1.0			1.1.1.04		100	1 0 1					0
								1.1	10		0
_		_									

Figure 7: Visual example of grid recognition

Furthermore, the subjects offered by each university were plotted over the years, and all resulting images underwent manual review. If any data exhibited implausible fluctuations, original and generated datasets were cross-checked manually, and the extraction algorithm was adjusted accordingly. An example of such a plot is shown in Figure 8.



Figure 8: Appearance of subjects over time

Note: The figure illustrates the temporal evolution of a subject within a particular institution, based on the name of the subject in the guides. The upper panel exemplifies data from the University of Stuttgart, while the lower panel presents data from TU Darmstadt. Notably, variations or discontinuities in the naming of subjects may occur over time.

A manual comparison between extractions and the original sources was conducted to ensure data accuracy. These comparisons encompassed not only the axis labels but also an examination of the entire tables. As the **axis labels** were recorded independently from the table grid recognition process, a comparison was made between the number of manually recorded axis labels and the number of recognized rows or columns. If discrepancies arose where the number of manually recorded axis labels did not align with the number of recognized rows or columns, both sides were thoroughly examined. Depending on the source of the error, adjustments were made either to the axis labels or to the table grid recognition for the affected tables. This iterative approach ensured alignment between the axis labels and the table grid, thus maintaining data accuracy and integrity.

As an additional measure to assess the quality of the data, the resulting tables underwent **manual comparison with the original tables** in the guides. During each evaluation, two to six rows and two to six columns of each table were randomly selected and cross-checked against the corresponding rows and columns on different pages of the original tables. On average, these randomly selected cells constituted 5.8 % of all cells. To assess whether random selection of cells is an appropriate method for evaluating the quality of the data extraction algorithm, the table for the year 1992 was chosen due to its lower scan quality, and a larger number of cells were inspected. Despite the increased number of checked cells, the

error rates remained within the previously stated range, even considering the lower quality of the guide scans. This suggests that random selection of cells is sufficient for determining error rates and allows for extrapolation of error rates across the entire dataset.

2.3 Additional data preparation

Further action was needed to create the final version of RWI-UNI-SUBJECTS. These steps comprise the harmonization of names of higher education institutions as well as subjects. All subjects were furthermore classified in the German academic coding system as provided by Destatis (2023). Finally, locations and municipality codes were added to the data.

Names and identifiers of higher education institutions

In the first step, all axis labels of universities and UAS were manually adjusted so that any reading errors, duplicates, and inconsistent names over time were harmonized. In a second step, higher education institution codes from the German Federal Statistical Office were matched to the data as they provide the unique identifiers used in German higher education statistics (Destatis, 2022).

In Destatis (2022), the Federal Statistical Office provides names and unique identifiers for all German higher education institutions. In total, 401 different names from the guides had to be matched to the identifiers from Destatis (2022). In roughly 70% of cases, an assignment was easy to implement due to the high similarity in names between the guides and Destatis (2022). For instance, "Düsseldorf U" could be assigned to "University of Düsseldorf" without further ado. For the remaining institutions, further investigation was required based on their appearance and the subjects offered. For example, "Berlin Sozial FH", could be correctly assigned to "Alice-Salomon-Hochschule Berlin (FH)" based on the fields offered. In a few exceptional cases, the guides were manually checked for additional information regarding the relevant higher education institutions. For example, each guide includes contact details for each university's disability officer, allowing for the recognition of the abbreviated name of the university. This meticulous approach ensured accurate identification and assignment of institutions where automated methods were insufficient.

During the process of adding unique identifiers based on Destatis (2022) to the dataset, several special cases occurred:

- 1. Change of names: Name changes may occur. In such cases, the current name and, therefore, the current identifier are assigned, and the variable "HE_Change" is set to 1. Additionally, the last previous name is included in the dataset.
- 2. Mergers and integration: A merger of two institutions or integration of one institution into another may occur. In such cases, the current name is assigned, and if the separate location still maintains its own identifier, that identifier is utilized. The variable "HE_Change" is set to 2 in this scenario. If an integration takes place and the separate campus no longer retains its own identifier, the current name is assigned, and variable "HE_Change" is set to 3. In both instances, the last previously valid name is separately included.
- 3. Separate campus: It may occur that a higher education institution offers a small and very limited range of subjects at one location, resulting in this campus never being assigned its own identifier. In such cases, the indentation in the guides, as depicted in Figure 9, was utilized to allocate the campus to its associated institution. For instance, the example in Figure 9 provides that "Gummersbach" is a campus of UAS Cologne. In addition, the websites of higher education institutions were manually checked for additional information. If the location could be identified as a campus, the name and

thus the identifier of the associated higher education institution are assigned, and "HE_Change" is set to 4, indicating the existence of a campus without its own identifier.

Figure 9: Identification of a separate campus by indentation in the guides

derborn FH der Wirtschaft Rheinfeld (FH Niederrhein chengladbach Augustin (FH erborn U-GH kath. aderborn' eschede H Æ 뎡 Į unster ē Ŧ Ŧ

The campus "Gummersbach" of Cologne UAS is in the second column of the table. It is just a campus, not an independent higher education institution, as indicated by the intended position. Further to the right are other campuses, such as "Höxter", "Meschede", and "Soest", all of these are part of the University of Paderborn.

4. Comprehensive University: Comprehensive Universities (in German: "Gesamthochschulen (GH)") was a form of higher education institution that combined features of universities and UAS. These existed in Germany at individual university locations between 1971 and 2003. Because comprehensive universities combined features of universities and UAS, they appear in the guides as both a university and a UAS, each offering respective subjects. To denote this exceptional circumstance, the variable "HE_Change" is set to 5. Additionally, the institution is assigned the most recent name.

In five cases, the name of the university or UAS as written in the guides could not be matched against a name from Destatis (2022), resulting in no unique identifier being assigned. Similarly, eight institutions in the category "Other HEI" could not be traced back.

Locations

The names of the institutions were then used to extract the location, specifically the name of the city which is often part of the institution's name. Subsequently, the official municipality code (in German: "Amtlicher Gemeindeschlüssel, AGS") in accordance with GeoBasis-DE / BKG 2013 (as of the territorial status of 31 December 2013) was added in an additional column. The municipality code enables, for instance, direct linkage to all data in the National Education Panel (NEPS), given the identical territorial status in both datasets.

Subjects and the academic coding system

Similarly, the names of the subjects were manually harmonized to eliminate reading errors and spelling inconsistencies. Then, the official codes from the German academic coding system were appended to the subject names as provided by Destatis (2023).

In Germany, the Federal Statistical Office provides an academic coding system (in German: "*Fächer-systematik*") that distinguishes 276 subjects, summarised into 64 subject areas and further aggregated into 9 subject groups. The complete assignment of the subject names as provided in the guides to those of the academic coding system is shown in Table A.1 in the Online Appendix.

2.4 Data quality and limitations

Data quality

At the end of the intense data preparation exercises, a final check was conducted to understand the final quality of the dataset. Similar to the interim checks conducted, rows and columns of the final dataset were randomly chosen and compared to the original guides. The resulting error rate was determined as the proportion of incorrect entries overall randomly chosen cells in a final manual comparison. The extrapolated average error rate of the final dataset is 0.2 %, meaning that one in every five hundred entries may contain an error. Please refer to Table 3 for more details.

Limitations

Structural breaks within the guides introduce complexities to the comparative analysis of subject offerings over time. These breaks encompass significant events such as German reunification, the aggregation of UAS at the city level until 1985, the introduction of tables dedicated to highly specific higher education institutions, and the reorganized presentation of teacher training programs.

- German reunification: The German reunification united Germany and the German Democratic Republic (GDR) in 1990. Data on higher education institutions in the former GDR is not available (see also Figure 2).
- Display of UAS before 1985: Until 1985, the guides only provided information on UAS aggregated at the city level. In cases where two or more UAS were situated within the same city, the guides consolidated the information at the city level, without distinguishing between individual UAS. From 1986 onwards, it is possible to distinguish two or more UAS located in the same city
- Introduction of new tables: From 1971 to 1977, the guides displayed the subjects offered at *universities* and *UAS*, constituting the two primary tables featured in each edition. Between 1978 and 1996, several additional tables were introduced, and the institutions listed here are classified as "Other HEI". Some of these additional tables are not included in the final dataset due to their peculiarities or very specific contexts. These are the additional tables on distance learning, colleges of teacher education and educational sciences, and UAS for public administration. For instance, UAS for public administration primarily displays subjects that are scarcely related to conventional study subjects. The introduction of additional tables was as follows:
 - 1978: Additional tables detailing the subjects offered at *church colleges, theological colleges, colleges and academies of fine arts, colleges of music, and other colleges were subsequently incorporated.* Moreover, a table on teacher education was introduced, specifying the type of school for which the training was intended, but not the detailed subjects.
 - 1986: The scope of subjects offered at UAS for public administration was also integrated into the dataset.
 - 1991: Universities of the German armed forces were introduced as an additional table in the dataset.
 - 1995: A table detailing subjects available through distance learning was introduced, but already removed from 1997 onwards.
 - 1996: Information on the universities of the German armed forces, church colleges, philosophicaltheological colleges, and those listed in the table of other higher education institutions were incorporated into the main table of universities. This means that they are removed from the

type Other higher education institutions and assigned the type university, which explains the significant increase observed in the latter in 1996, as depicted in Figure 2. Moreover, a detailed table on teacher education was introduced.

• University of Kassel in 1972: In the guide of 1972, the University of Kassel is missing. Most likely, it has been inadvertently omitted. However, information regarding available subjects for this university is included in the guides for the years 1971 and 1973 and all subsequent years after 1973.

3 References

- Bock, K. H. (1973). Studien- und Berufswahl Entscheidungshilfen für Abiturienten und Absolventen der Fachoberschulen. (B.-L.-K. für Bildungsplanung und Bundesanstalt für Arbeit, Ed.) (No. 1). Verlag Karl Heinrich Bock.
- Bock, K. H. (1974). Studien- und Berufswahl Entscheidungshilfen für Abiturienten und Absolventen der Fachoberschulen. (B.-L.-K. für Bildungsplanung und Bundesanstalt für Arbeit, Ed.) (No. 1). Verlag Karl Heinrich Bock.
- Bock, K. H. (1975). Studien- und Berufswahl Entscheidungshilfen für Abiturienten und Absolventen der Fachoberschulen. (B.-L.-K. für Bildungsplanung und Bundesanstalt für Arbeit, Ed.) (No. 1). Verlag Karl Heinrich Bock.
- Bock, K. H. (1976). Studien- und Berufswahl Entscheidungshilfen f
 ür Abiturienten und Absolventen der Fachoberschulen. (B.-L.-K. f
 ür Bildungsplanung und Forschungsf
 örderung und Bundesanstalt f
 ür Arbeite, Ed.) (No. 1). Verlag Karl Heinrich Bock.
- Bock, K. H. (1977). Studien- und Berufswahl Entscheidungshilfen f
 ür Abiturienten und Absolventen der Fachoberschulen. (B.-L.-K. f
 ür Bildungsplanung und Forschungsf
 örderung und Bundesanstalt f
 ür Arbeite, Ed.) (No. 1). Verlag Karl Heinrich Bock.
- Bock, K. H. (1978). Studien- und Berufswahl Entscheidungshilfen f
 ür Abiturienten und Absolventen der Fachoberschulen. (B.-L.-K. f
 ür Bildungsplanung und Forschungsf
 örderung und Bundesanstalt f
 ür Arbeite, Ed.) (No. 1). Verlag Karl Heinrich Bock.
- Bock, K. H. (1979). Studien- und Berufswahl Entscheidungshilfen f
 ür Abiturienten und Absolventen der Fachoberschulen. (B.-L.-K. f
 ür Bildungsplanung und Forschungsf
 örderung und Bundesanstalt f
 ür Arbeite, Ed.) (No. 1). Verlag Karl Heinrich Bock.
- Bock, K. H. (1980). Studien- und Berufswahl Entscheidungshilfen f
 ür Abiturienten und Absolventen der Fachoberschulen. (B.-L.-K. f
 ür Bildungsplanung und Forschungsf
 örderung und Bundesanstalt f
 ür Arbeite, Ed.) (No. 1). Verlag Karl Heinrich Bock.
- Bock, K. H. (1981). Studien- und Berufswahl Entscheidungshilfen f
 ür Abiturienten und Absolventen der Fachoberschulen. (B.-L.-K. f
 ür Bildungsplanung und Forschungsf
 örderung und Bundesanstalt f
 ür Arbeite, Ed.) (No. 1). Verlag Karl Heinrich Bock.
- Bock, K. H. (1982). Studien- und Berufswahl Entscheidungshilfen f
 ür Abiturienten und Absolventen der Fachoberschulen. (B.-L.-K. f
 ür Bildungsplanung und Forschungsf
 örderung und Bundesanstalt f
 ür Arbeite, Ed.) (No. 1). Verlag Karl Heinrich Bock.
- Bock, K. H. (1983). Studien- und Berufswahl Entscheidungshilfen f
 ür Abiturienten und Absolventen der Fachoberschulen. (B.-L.-K. f
 ür Bildungsplanung und Forschungsf
 örderung und Bundesanstalt f
 ür Arbeite, Ed.) (No. 1). Verlag Karl Heinrich Bock.
- Bock, K. H. (1984). Studien- und Berufswahl Entscheidungshilfen f
 ür Abiturienten und Absolventen der Fachoberschulen. (B.-L.-K. f
 ür Bildungsplanung und Forschungsf
 örderung und Bundesanstalt f
 ür Arbeite, Ed.) (No. 1). Verlag Karl Heinrich Bock.

- Bock, K. H. (1985). Studien- und Berufswahl Entscheidungshilfen f
 ür Abiturienten und Absolventen der Fachoberschulen. (B.-L.-K. f
 ür Bildungsplanung und Forschungsf
 örderung und Bundesanstalt f
 ür Arbeite, Ed.) (No. 1). Verlag Karl Heinrich Bock.
- Bock, K. H. (1986). Studien- und Berufswahl Entscheidungshilfen f
 ür Abiturienten und Absolventen der Fachoberschulen. (B.-L.-K. f
 ür Bildungsplanung und Forschungsf
 örderung und Bundesanstalt f
 ür Arbeite, Ed.) (No. 1). Verlag Karl Heinrich Bock.
- Bock, K. H. (1987). Studien- und Berufswahl Entscheidungshilfen f
 ür Abiturienten und Absolventen der Fachoberschulen. (B.-L.-K. f
 ür Bildungsplanung und Forschungsf
 örderung und Bundesanstalt f
 ür Arbeite, Ed.) (No. 1). Verlag Karl Heinrich Bock.
- Bock, K. H. (1988). Studien- und Berufswahl Entscheidungshilfen f
 ür Abiturienten und Absolventen der Fachoberschulen. (B.-L.-K. f
 ür Bildungsplanung und Forschungsf
 örderung und Bundesanstalt f
 ür Arbeite, Ed.) (No. 1). Verlag Karl Heinrich Bock.
- Bock, K. H. (1989). Studien- und Berufswahl Entscheidungshilfen f
 ür Abiturienten und Absolventen der Fachoberschulen. (B.-L.-K. f
 ür Bildungsplanung und Forschungsf
 örderung und Bundesanstalt f
 ür Arbeite, Ed.) (No. 1). Verlag Karl Heinrich Bock.
- Bock, K. H. (1990). Studien- und Berufswahl Entscheidungshilfen f
 ür Abiturienten und Absolventen der Fachoberschulen. (B.-L.-K. f
 ür Bildungsplanung und Forschungsf
 örderung und Bundesanstalt f
 ür Arbeite, Ed.) (No. 1). Verlag Karl Heinrich Bock.
- Bock, K. H. (1991). Studien- und Berufswahl Entscheidungshilfen f
 ür Abiturienten und Absolventen der Fachoberschulen. (B.-L.-K. f
 ür Bildungsplanung und Forschungsf
 örderung und Bundesanstalt f
 ür Arbeite, Ed.) (No. 1). Verlag Karl Heinrich Bock.
- Bock, K. H. (1992). Studien- und Berufswahl Entscheidungshilfen f
 ür Abiturienten und Absolventen der Fachoberschulen. (B.-L.-K. f
 ür Bildungsplanung und Forschungsf
 örderung und Bundesanstalt f
 ür Arbeite, Ed.) (No. 1). Verlag Karl Heinrich Bock.
- Bock, K. H. (1993). Studien- und Berufswahl Entscheidungshilfen f
 ür Abiturienten und Absolventen der Fachoberschulen. (B.-L.-K. f
 ür Bildungsplanung und Forschungsf
 örderung und Bundesanstalt f
 ür Arbeite, Ed.) (No. 1). Verlag Karl Heinrich Bock.
- Bock, K. H. (1994). Studien- und Berufswahl Entscheidungshilfen für Abiturienten und Absolventen der Fachoberschulen. (B.-L.-K. für Bildungsplanung und Forschungsförderung und Bundesanstalt für Arbeite, Ed.) (No. 1). Verlag Karl Heinrich Bock.
- Bock, K. H. (1995). Studien- und Berufswahl Entscheidungshilfen f
 ür Abiturienten und Absolventen der Fachoberschulen. (B.-L.-K. f
 ür Bildungsplanung und Forschungsf
 örderung und Bundesanstalt f
 ür Arbeite, Ed.) (No. 1). Verlag Karl Heinrich Bock.
- Bock, K. H. (1996). Studien- und Berufswahl Entscheidungshilfen f
 ür Abiturienten und Absolventen der Fachoberschulen. (B.-L.-K. f
 ür Bildungsplanung und Forschungsf
 örderung und Bundesanstalt f
 ür Arbeite, Ed.) (No. 1). Verlag Karl Heinrich Bock.
- Bradski, G. (2000). The OpenCV Library. Dr. Dobb's Journal of Software Tools.
- Destatis. (2022). Schlüsselverzeichnisse für die Studenten-, Prüfungsstatistik, Promovierendenstatistik. Stand: 2021/2022 & SS 2022.
- Destatis. (2023). Systematik der Fächergruppen, Studienbereiche und Studienfächer.
- German Centre for Higher Education Research and Science Studies (DZHW). (n.d.). DZHW Graduate Survey Series. Study Series: DZHW Graduate Survey Series. (https://doi.org/10.21249/DZHW: gra1989:2.0.0)
- Hirschfeld, G. (1971). Studium und Beruf Informationen f
 ür Abiturienten und Absolventen der Fachoberschulen (B.-L.-K. f
 ür Bildungsplanung und Bundesanstalt f
 ür Arbeit, Ed.) (No. 1). aspekte verlag gmbh.

- Hirschfeld, G. (1972). Studium und Beruf Informationen f
 ür Abiturienten und Absolventen der Fachoberschulen (B.-L.-K. f
 ür Bildungsplanung und Bundesanstalt f
 ür Arbeit, Ed.) (No. 1). aspekte verlag gmbh.
- Meisel, H. (1978). Die Deutsche Berufsberatung Gesamtüberblick. Stuttgart, Berlin, Köln, Mainz: Kohlhammer.
- NEPS-Netzwerk. (2023). Nationales Bildungspanel, Scientific Use File der Startkohorte Erwachsene. Bamberg. (https://doi.org/10.5157/NEPS:SC6:14.0.0)
- Schmucker, A., Seth, S., & vom Berge, P. (2023). Stichprobe der Integrierten Arbeitsmarktbiografien (SIAB) 1975 - 2021 (Tech. Rep. No. 02/2023). Nürnberg: FDZ Datenreport. (https://doi.org/ 10.5164/IAB.FDZD.2302.de.v1, language = de)
- University of Bologna. (2024). Degree Programmes Data University of Bologna. Retrieved from https://data.europa.eu/data/datasets/8cbba6bd-1686-4828-b223-6f4dab477434?locale= en (Catalog of first and second cycle degree programs by academic year, from 2004/2005 to 2020/2021)
- U.S. Department of Education. (2024). College Scorecard Data. Retrieved from https:// collegescorecard.ed.gov/data/ (Data on student enrollment, completion, debt, repayment, and earnings for U.S. institutions)

A Academic Coding System

The guides to "Study and Career Choice" from 1971 to 1996 differentiate between subjects. These are assigned to subjects, subject areas, and subject groups based on the academic coding system from Destatis (2023). Notably, no subject was assigned to the ninth subject group from the academic coding system, labeled "Outside the subject area classification". A comprehensive list of subjects from the guides, along with their assigned subjects from the subject classification system, can be found in Table A.1 of the Online Appendix. Below is an example of the assignment for each subject group. The hierarchical structure delineates the subject group, subject area, and corresponding subjects from the academic coding system, along with their respective assigned subject names from the guides.

```
1 Geisteswissenschaften
```

- (a) 02 Evang. Theologie, -Religionslehre
 - i. 053 Evang. Theologie, -Religionslehre Name as in the guide: *Theologie ev.*

...

2 Sport

- (a) 22 Sport, Sportwissenschaft
 - i. 029 Sportwissenschaft:

Names as in the guide: Sport, Sportwissenschaft

•••

3 Rechts-, Wirtschafts- und Sozialwissenschaften

- (a) 25 Politikwissenschaft
 - i. 129 Politikwissenschaft/Politologie:

Names as in the guide: Politik/Sozialkunde, Politikwissenschaft, Politikwissenschaft/Politologie, Politologie

...

4 Mathematik, Naturwissenschaften

- (a) 37 Mathematik
 - i. 105 Mathematik:

Names as in the guide: Mathematik, Mathematik (Diplom), Mathematik (Diplom/Magister), Mathematik (Gymnasium), Mathematik (Realschule)

...

5 Humanmedizin/Gesundheitswissenschaften

- (a) 50 Zahnmedizin
 - i. 185 Zahnmedizin:

Names as in the guide: Zahnmedizin, Zahnmedizin (Stomatologie)

...

7 Agrar-, Forst- und Ernährungswissenschaften, Veterinärmedizin

- (a) 58 Agrarwissenschaften, Lebensmittel- und Getränketechnologie
 - i. 003 Agrarwissenschaft/Landwirtschaft:

Names as in the guide: Agrarwissenschaft, Landwirtschaft, Landwirtschaft/Agrarwirtschaft (Magrarwirtschaft) and Magrarwirtschaft (Magrarwirtschaft) and Magrarwirtsch

•••

8 Ingenieurwissenschaften

- (a) 63 Machinenbau/Verfahrenstechnik
 - i. 104 Maschinenbau/-wesen:

Names as in the guide: Maschinenbau, Maschinenbau/Stahlbau, Maschinenbau allgemein, Kraftfahrzeugbau

•••

9 Kunst, Kunstwissenschaft

- (a) 74 Kunst, Kunstwissenschaften allgemein
 - i. 091 Kunsterziehung:

Names as in the guide: Kunst/Kunsterziehung, Kunst/freie Kunstpädagogik, Kunst/-pädagogik, Kunst (Gymnasium)

•••

B Details on temporal evolution of institutions

B.1 Number of subjects and institutions over time

B.1.1 Number of subjects over time

Table 1: Number	of subjects by	year and type of higher	education institution

Year	UAS	Uni	UAS & Uni
1971	586	1349	1935
1972	476	1202	1678
1973	666	2510	3176
1974	692	2492	3184
1975	696	2652	3348
1976	691	2650	3341
1977	725	2768	3493
1978	697	2974	3671
1979	711	3027	3738
1980	702	3039	3741
1981	695	3060	3755
1982	754	3191	3945
1983	763	3110	3873
1984	776	3161	3937
1985	788	3300	4088
1986	836	3206	4042
1987	841	3235	4076
1988	884	3270	4154
1989	906	3273	4179
1990	959	3389	4348
1991	1060	3849	4909
1992	1176	4145	5321
1993	1329	4330	5659
1994	1371	4449	5820
1995	1420	4432	5852
1996	1512	6056	7568

Note: The table summarizes the number of subjects by year and type of institution.

B.1.2 Number of institutions over time

Year	Number of Unis	Number of UAS	Number of municipalities with UAS
1971	44	-	95
1972	44	-	102
1973	49	-	100
1974	52	-	103
1975	52	-	104
1976	55	-	104
1977	55	-	104
1978	56	-	106
1979	57	-	109
1980	56	-	106
1981	58	-	104
1982	59	-	104
1983	61	-	104
1984	61	-	105
1985	61	-	104
1986	61	135	106
1987	61	136	106
1988	61	137	106
1989	61	135	107
1990	61	136	107
1991	74	152	121
1992	80	169	132
1993	79	177	137
1994	79	184	143
1995	79	193	151
1996	103	198	155

Table 2: Number of higher education institutions by institutional type over time

Note: The table shows the number of higher education institutions based on the count of unique higher education institutions. Since UAS prior to 1986 are only listed by place of study, the number of municipalities with UAS is presented separately.

C Error rates per year

Table 3 provides the number of randomly selected cells, the number of incorrect entries, and the resulting error rates for each guide. Depending on the quality of the tables in the guides as well as on the quality of scans, the error rate varies over the years and ranges from 0 % to 1.27 %. The extrapolated average error rate is 0.2 %, meaning that one in every five hundred entries may contain an error.

Year	Universities			UAS		
	Number of checked cells	Number of incorrect en- tries	Error rate	Number of checked cells	Number of incorrect en- tries	Error rate
1971	396	3	0.75%	390	0	0.00%
1972	299	1	0.33%	395	0	0.00%
1973	392	5	1.27%	400	0	0.00%
1974	591	7	1.18%	398	2	0.50%
1975	597	1	0.16%	399	3	0.75%
1976	612	0	0.00%	421	1	0.24%
1977	410	3	0.73%	528	0	0.00%
1978	678	1	0.14%	561	0	0.00%
1979	684	0	0.00%	564	0	0.00%
1980	513	2	0.38%	540	1	0.19%
1981	690	4	0.58%	429	1	0.23%
1982	582	2	0.34%	444	1	0.23%
1983	352	0	0.00%	549	1	0.18%
1984	470	0	0.00%	444	1	0.23%
1985	474	0	0.00%	442	0	0.00%
1986	530	3	0.57%	512	0	0.00%
1987	590	0	0.00%	588	0	0.00%
1988	594	3	0.51%	669	0	0.00%
1989	536	0	0.00%	672	0	0.00%
1990	735	0	0.00%	558	0	0.00%
1991	558	0	0.00%	864	0	0.00%
1992	1973	17	0.86%	634	0	0.00%
1993	738	0	0.00%	646	0	0.00%
1994	819	1	0.12%	662	0	0.00%
1995	548	0	0.00%	873	3	0.34%
1996	903	0	0.00%	598	0	0.00%

Table 3: Error rate of checked cells per year and institution

Note: The table provides a summary of the results from the random manual review of item identification for UAS and universities. The average error rate across all years from 1971 to 1996 is 0.20 %. Notably, the guide from 1992 exhibited low scan quality, necessitating a broader selection of cells for quality control. Results for the higher education institutions of type Other are comparable, with the average error rate hovering around 0.1 %.

D Codebook

Variable	Definition	Characteristics					
Year	Year of guide and therefore observation	Every year from 1971 until 1996.					
		Numeric variable.					
Type	Institutional Type	University, UAS or of type Other higher education institution.					
		String variable.					
HE_name_orig	Institution name as found in the guide	This name can change over time.					
	of the respective year	String variable.					
HE_number	Institution number according to	Unique, time constant four-digit identifier which may start with a 0.					
	Destatis (2022)	Missings are assigned a 0.					
		String variable.					
HE_name_destat	Institution name according to Destatis	Time constant current institution name. Missings are assigned a 0.					
	(2022)	String variable.					
Subject_orig	Original subject name as found in the	Name for specific subject may change slightly over time.					
	guide of the respective year	String variable.					
Subject_code	Systematic number of subject	Three-digit identifier for subject which may start with a 0. Subjects					
	according to Destatis (2023)	that cannot be assigned a specific $Subject_code$ are assigned a 0.					
		String variable.					
Subject_area_code	Systematic number of subject area	Two-digit identifier for subject area which may start with a 0. Subjects					
	according to Destatis (2023)	that cannot be assigned a specific $Subject_area_code$ are assigned a 0.					
		String variable.					
Subject_group_code	Systematic number of subject group	One-digit identifier for subject group. Subjects that cannot be assigned					
	according to Destatis (2023)	a specific <i>Subject_group_code</i> are assigned a 0.					
		Numeric variable.					

Table 4: Codebook

Continued on next page

Variable	Definition	Characteristics
Subject	Designation of subject according to Destatis (2023)	Designation of subject corresponding to <i>Subject_code</i> . Subjects that cannot be assigned a specific <i>Subject_code</i> are assigned a 0. String variable.
Subject_area	Designation of subject area according to Destatis (2023)	Designation of subject area corresponding to Subject_area_code. Subjects that cannot be assigned a specific Subject_area_code are assigned a 0. String variable.
Subject_group	Designation of subject group according to Destatis (2023)	Designation of subject group corresponding to Subject_group_code. Subjects that cannot be assigned a specific Subject_group_code are assigned a 0. String variable.
AGS	Municipality code on district level as of 31.12.2013	Four-digit identifier (five-digit for former East German states) for district (in German: "Landkreis" or "kreisfreie Stadt"). Numeric variable.
Location_name	Location of the respective institution	Location, usually a city, of higher education institution. String variable.
HE_name_destat_last	Last previous institution name according to Destatis (2023) in case it changed over time	When changes occurred this variable displays the last previous institution name, otherwise it is assigned a 0. String variable.

fable i commuted from providuo page

Continued on next page

Variable Definition	Characteristics
HE_change Categorical Variable indicating changes in the institution characteristics	 0: No observed change. 1: Change of institution name. 2: Merger or integration after which the separate location still has its own <i>HE_number</i>. 3: Merger or integration after which the location does not have a separate <i>HE_number</i> anymore. 4: Separate campus which never had its own <i>HE_number</i>. 5: Former comprehensive university. Numeric variable.

Table 4 – Continued from previous page

Continued on next page

Variable	Definition	Characteristics
Study_Type	Variable indicating the type of the	1: Full study
	study modes	11: Full study Winter term (WS) required
		12: Full study Winter term (WS) recommended
		13: Full study Summer term (SS) required
		2: Full study admission-restricted
		21: Full study admission-restricted Winter term (WS) required
		22: Full study admission-restricted Winter term (WS) recommended
		3: Specialization
		31: Specialization Winter term (WS) required
		32: Specialization Winter term (WS) recommended
		4: Specialization admission-restricted
		5: Advanced study
		51: Advanced study Winter term (WS) required
		52: Advanced study Winter term (WS) recommended
		6: Advanced study admission-restricted
		7: Partial study
		71: Partial study Winter term (WS) required
		72: Partial study Winter term (WS) recommended
		7a: Partial study starting from
		7b: Partial study until
		8: Minor subject
		81: Minor subject Winter term (WS) required
		82: Minor subject Winter term (WS) recommended
		X: No new students
		Xa: No new students soon

Table 4 – Continued from previous page

Note: The table shows the variables in the dataset, their definitions and some characteristics.

Datenbeschreibung

Datenbeschreibung



Leibniz-Institut für Wirtschaftsforschung

Das RWI wird vom Bund und vom Land Nordrhein-Westfalen gefördert.

