# How to manage big data (well)

## A digitized process-chain

*Inga Brentel*
*Heinrich Heine -*
*Universität Düsseldorf*

*CESSDA Training Day*
*27-28 November 2019, Cologne*

cessda.eu  @CESSDA_Data

cessda

# How to manage big data (well) - a digitized process-chain

Inga Brentel

HEINRICH HEINE
UNIVERSITÄT DÜSSELDORF

FORSCHUNGSVERBUND NRW
DIGITALE GESELLSCHAFT

## Definition & Relevance

- Digitaliztion brings data tracks (of human behavior) = *datafication*

    → **Big-Data** as unstructured, heterogeneous data bulk

- Exploitation and Management of data

    → A huge chance for (social) science

    → Nessecarity of high quality standards

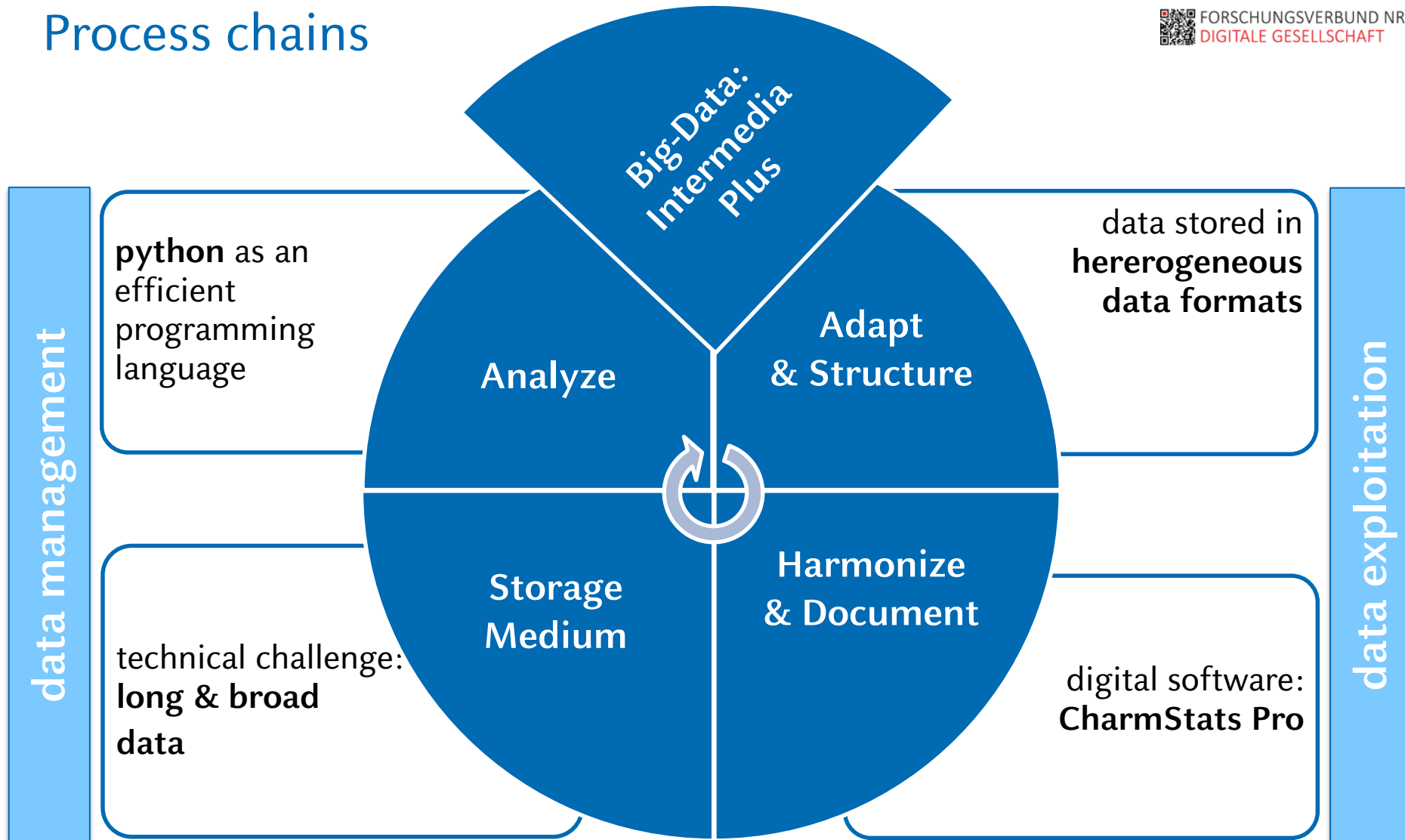# Big-Data & (Social-)Science

## Big-Data as a challenge & a chance

- *Exploitation of Big-Data* sources for CSS
  → setup of a new data infrastructure

- *Data management* for a sustainable access to research data of high quality

- *Technical interfaces* to combine process chains


- **Solution:** *digitized process chains for Big-Data exploitation and analysis*

## Big-Data & Social Science

- Yearly gathered, cross sectional study for German media reach; conducted by commercial interests

- **Data typ:** survey & technical measures/tracking

- **Big-Data:** over 18.000 variables & 1,6 mio cases (for 2014 to 2016)

- **Content:**

  - Media reach for & media use of ~12.000 online offerings, 150 magazines, 100 newspapers, 100 radio channels, 10 tv-channels

  - Further info on respondents like socio demographics, daily routine, free-time activities, habits and household

- **Target:** *a structured, longitudinal dataset for academic research*

# Solution Model I.

## Process chains

FORSCHUNGSVERBUND NRW
DIGITALE GESELLSCHAFT

**data management**

**python** as an efficient programming language

technical challenge: **long & broad data**

Big-Data: Intermedia Plus

**Analyze**

**Adapt & Structure**

**Storage Medium**

**Harmonize & Document**

data stored in **hererogeneous data formats**

digital software: **CharmStats Pro**

**data exploitation**

# Solution Model II.

## Technical Interfaces

Half-automized matching with SPSS-data

Transfer of the relevant Info
(target & source variable,
coding of values)

**Charmstats Pro**

CharmStatsPro

- Adapt & structure the data formats for CharmStats

**SPSS & Excel**

- harmonization
- documentation

- MySQL as storage medium
- python as efficient language for statistical analysis

POWERED BY
MySQL

python

**CharmStana**

Import

# Interface I.: Data documentation



Big-Data & (Social-)Science | IntermediaPlus | solution models | **interfaces** | deep dive | learnings     www.hhu.de

HEINRICH HEINE
UNIVERSITÄT DÜSSELDORF

FORSCHUNGSVERBUND NRW
DIGITALE GESELLSCHAFT

Export Syntax

xTab

xTab

SPSS V

n Pay TV: H
012
n Pay TV: V

## Gesamtangebote Online

### Variablenname:

O_GA_airl_CT  -  Full website: airliners.de
O_GA_yh_CT  -  Full website: Yahoo! Deutschland
O_GA_micro_CT  -  Full website: Microsoft
O_GA_mtv_CT -  Full website: MTV
O_GA_vge_CT  -  Full website: VOGUE.DE
O_GA_bike_CT  -  Full website: bike-magazin.de
O_GA_bld_CT  -  Full website: BILD.de
O_GA_ebay_CT  -  Full website: eBay.de

### Kodierung:

-7, nicht ermittelt

### Übersichtstabelle der Variable und Harmonisierung für die Jahre 2010 bis 2015

|  | 2014 | 2015 |
|---|---|---|
| Variablenname | G11<br>G111, G401, G406, G901, G2166, G2191, G3586, G19046 | G11<br>G111, G401, G406, G901, G2166, G2<br>G3586, G19046 |

HEINRICH HEINE
UNIVERSITÄT DÜSSELDORF

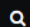FORSCHUNGSVERBUND NRW
DIGITALE GESELLSCHAFT

Export Syntax

**Gesamtangebote Online**

*Variablenname:*

Pay TV: H
012

*Übersichtstabelle der Variable und Harmonisierung für die Jahre 2010 bis 2015*

| | 2014 | 2015 | 2016 |
|---|---|---|---|
| Variablenname | G11<br>G111, G401, G406, G901, G2166, G2191, G3586, G19046 | G11<br>G111, G401, G406, G901, G2166, G2191, G3586, G19046 | G11<br>G111, G401, G406, G901, G2166, G2191, G3586, G19046 |
| Frage | | | |
| Interviewer-anweisung | | | |
| Kodierung | -7, nicht ermittelt | -7, nicht ermittelt | -7, nicht ermittelt |
| Harmonisierungs-Syntax | RECODE G11 (-7=-7) (ELSE=COPY) INTO O_GA_rtlb_CT.<br>MISSING VALUES O_GA_rtlb_CT (-7).<br>VARIABLE LABELS O_GA_rtlb_CT 'full website: RTL 104.6'.<br>VALUE LABELS O_GA_rtlb_CT -7 'nicht ermittelt'.<br>EXECUTE. | RECODE G11 (-7=-7) (ELSE=COPY) INTO O_GA_rtlb_CT.<br>MISSING VALUES O_GA_rtlb_CT (-7).<br>VARIABLE LABELS O_GA_rtlb_CT 'full website: RTL 104.6'.<br>VALUE LABELS O_GA_rtlb_CT -7 'nicht ermittelt'.<br>EXECUTE. | RECODE G11 (-7=-7) (ELSE=COPY) INTO O_GA_rtlb_CT.<br>MISSING VALUES O_GA_rtlb_CT (-7).<br>VARIABLE LABELS O_GA_rtlb_CT 'full website: RTL 104.6'.<br>VALUE LABELS O_GA_rtlb_CT -7 'nicht ermittelt'.<br>EXECUTE. |
| Vermerk | | | |

*Übersicht Virtuelle Variablen siehe Anhang 1*

# Interface II.: gesis Harmonization-HUB

Big-Data & (Social-)Science | IntermediaPlus | solution models | **interfaces** | deep dive | learnings          www.hhu.de

# Interface II.: gesis Harmonization-HUB

Big-Data & (Social-)Science | IntermediaPlus | solution models | **interfaces** | deep dive | learnings          www.hhu.de

# Interface III.: Charmstana

Big-Data & (Social-)Science | IntermediaPlus | solution models | **interfaces** | deep dive | learnings

Entwickelt mit
Tanja Roeder (BA)

HEINRICH HEINE
UNIVERSITÄT DÜSSELDORF

FORSCHUNGSVERBUND NRW
DIGITALE GESELLSCHAFT

| | Struktur | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|

| BefragtenID | Jahr | HVarName | HVarLabel | HVarWert | HVarWerteLabel | is_missing | HVarLevel | quell_id | ziel_id |
|---|---|---|---|---|---|---|---|---|---|
| 100_2014 | 2014 | A235 | Befragte/r: Lehre ohne Abschluss | 0 | HVarWert ist nicht in Tabelle categories gepflegt | | 2 | 160 | 77 |
| 100_2014 | 2014 | A236 | Befragte/r: Lehre mit Abschluss (Gehilfen, Geselle... | 0 | HVarWert ist nicht in Tabelle categories gepflegt | | 2 | 160 | 77 |
| 100_2014 | 2014 | A237 | Befragte/r: Gewerbeschule, Fachschule mit Abschlus... | 0 | HVarWert ist nicht in Tabelle categories gepflegt | | 2 | 160 | 77 |
| 100_2014 | 2014 | A238 | Befragte/r: Fachhochschulabschluss (auch Ingenieur... | 0 | HVarWert ist nicht in Tabelle categories gepflegt | | 2 | 160 | 77 |
| 100_2014 | 2014 | A239 | Befragte/r: Hochschulabschluss | 0 | HVarWert ist nicht in Tabelle categories gepflegt | | 2 | 160 | 77 |
| 100_2014 | 2014 | A240 | Befragte/r: andere Art der Berufsausbildung | 0 | HVarWert ist nicht in Tabelle categories gepflegt | | 2 | 160 | 77 |
| 100_2014 | 2014 | A241 | Befragte/r: nichts davon | 0 | HVarWert ist nicht in Tabelle categories gepflegt | | 2 | 160 | 77 |
| 100_2014 | 2014 | alterg | Alter des/der Befragten | 5 | 35-39 Jahre | 0 | 1 | 133 | 62 |
| 100_2014 | 2014 | altergk | Alter des/der Befragten (gruppiert) | 5 | 60 Jahre und älter | 0 | 1 | 133 | 61 |
| 100_2014 | 2014 | BAnole | Befragte/r: noch in der Lehre | 0 | HVarWert ist nicht in Tabelle opera_prescriptions ... | | 1 | 159 | 77 |
| 100_2014 | 2014 | bund | bund | 5 | Nordrhein-Westfalen | 0 | 1 | 94 | 44 |
| 100_2014 | 2014 | eink | Einkommen des/der Befragten | 2 | 500€ bis unter 1.000€ | 0 | 1 | 139 | 64 |
| 100_2014 | 2014 | einkbez | Personen mit eigenem Einkommen | 2 | 2 Personen | 0 | 1 | 127 | 57 |
| 100_2014 | 2014 | einkbezk | Personen mit eigenem Einkommen | 2 | 2 Personen | 0 | 1 | 130 | 58 |
| 100_2014 | 2014 | erwerb | Berufstaetigkeit der/des Befragten | 0 | HVarWert ist nicht in Tabelle opera_prescriptions ... | | 1 | 115 | 75 |
| 100_2014 | 2014 | erwerbk | Berufstätigkeit der/des Befragten | 0 | HVarWert ist nicht in Tabelle | | 1 | 117 | 54 |

Tabelle
☐ analyse
☐ quellvar
☐ zielvari
3 Tabell

Entwickelt mit
Tanja Roeder (BA)

# DEEP-DIVE

## How to structure (Big-)Data?

# Solution Model I.

## Process chains



**data management**

**python** als as an efficient programming language

technical challenge: **long & broad data**

**Big-Data: Intermedia Plus**

Analyze

Storage Medium

Adapt & **Structure**

Harmonize & Document

data stored in **hererogeneous data formats**

digital software: **CharmStats Pro**

**data exploitation**

## Variables for media use

- 100 radio broadcasts, usage per hour

- 150 magazins & 100 newspapers (incl. local newspapers)

- 10 TV-channels, usage per 30 minutes

- ~ 3000 Online sites (net coverage on a daily, weekly, monthly basis; cross coverage on a monthly and quarter basis)

  - ~ 733 overall online media

  - 335 single online sites

  - 52 homepages

➔ **Goal:** harmonization of Intermedia Plus 2014 to 2017+

HEINRICH HEINE
UNIVERSITÄT DÜSSELDORF

## The case: MA Intermedia Plus, online-tranche

FORSCHUNGSVERBUND NRW
DIGITALE GESELLSCHAFT

| Felder | | | | | Ausweisungs-einschränkung | |
|---|---|---|---|---|---|---|
| PI-Summen | | Tagesbasis | | | | |
| 3 Monate | | Monat | Woche | Tag | T = Nur Woche und Monat | |
| von | bis | | | | F = Fallzahlgrenze nicht erreicht (351 ung. WNKs) | |
| | | | | | L = Leerfelder | Angebotsnamen |
| 321 | 322 | 323 | 324 | 325 | | babyclub.de Gesamt |
| 326 | 327 | 328 | 329 | 330 | | Baby-Vornamen.de Gesamt |
| 331 | 332 | 333 | 334 | 335 | | Baden Online Gesamt |
| 336 | 337 | 338 | 339 | 340 | | OMS Badische Zeitung Online Gesamt |
| 341 | 342 | 343 | 344 | 345 | | Basketball Bund Gesamt |
| 346 | 347 | 348 | 349 | 350 | F | OMS Kreiszeitung Böblinger Bote Gesa |
| 351 | 352 | 353 | 354 | 355 | F | OMS bbv-net Gesamt |
| 356 | 357 | 358 | 359 | 360 | | Bergfex.de Gesamt |
| 361 | 362 | 363 | 364 | 365 | | OMS Berliner Kurier Online Gesamt |
| 366 | 367 | 368 | 369 | 370 | | OMS Berlin.de Gesamt |
| 371 | 372 | 373 | 374 | 375 | | OMS BerlinOnline Gesamt |
| 376 | 377 | 378 | 379 | 380 | F | OMS Berliner Rundfunk.de Gesamt |
| 381 | 382 | 383 | 384 | 385 | | OMS Berliner Zeitung Gesamt |
| 386 | 387 | 388 | 389 | 390 | | Best of Home Gesamt |
| 391 | 392 | 393 | 394 | 395 | | OMS bigFM.de Gesamt |
| 396 | 397 | 398 | 399 | 400 | | bigpoint.com Gesamt |
| 401 | 402 | 403 | 404 | 405 | | bike-magazin.de Gesamt |
| 406 | 407 | 408 | 409 | 410 | | BILD.de Gesamt |
| 411 | 412 | 413 | 414 | 415 | | Bild der Frau Gesamt |
| 416 | 417 | 418 | 419 | 420 | | billiger.de Gesamt |
| 421 | 422 | 423 | 424 | 425 | L | |
| 426 | 427 | 428 | 429 | 430 | | Bisafans.de Gesamt |
| 431 | 432 | 433 | 434 | 435 | F | OMS Backnanger Kreiszeitung Gesamt |
| 436 | 437 | 438 | 439 | 440 | L | |
| 441 | 442 | 443 | 444 | 445 | | boerse-frankfurt.de Gesamt |
| 446 | 447 | 448 | 449 | 450 | | BOERSE.de Gesamt |
| 451 | 452 | 453 | 454 | 455 | | BOERSE-ONLINE.de Gesamt |
| 456 | 457 | 458 | 459 | 460 | | boersennews.de Gesamt |
| 461 | 462 | 463 | 464 | 465 | | Boerse Stuttgart Gesamt |
| 466 | 467 | 468 | 469 | 470 | F | BOOTE-Magazin online Gesamt |
| 471 | 472 | 473 | 474 | 475 | F | OMS Borkener Zeitung Gesamt |
| 476 | 477 | 478 | 479 | 480 | F | brainguide.de Gesamt |
| 481 | 482 | 483 | 484 | 485 | L | |
| 486 | 487 | 488 | 489 | 490 | | BRAVO Online Gesamt |
| 491 | 492 | 493 | 494 | 495 | | bremen.de Gesamt |

## The case: MA Intermedia Plus, online-tranche

| PI-Summ 3 Monat von | Felder | | | | | Ausweisungs- einschränkung | |
|---|---|---|---|---|---|---|---|
| | PI-Summen | | Tagesbasis | | | | |
| | 3 Monate | | Monat | Woche | Tag | T = Nur Woche und Monat | |
| 321 | von | bis | | | | F = Fallzahlgrenze nicht erreicht (351 ung. WNKs) | |
| 326 | | | | | | L = Leerfelder | **Angebotsnamen** |
| 331 336 | 4551 | 4552 | 4553 | 4554 | 4555 | F | Netzwelt.de Home |
| 341 | 4556 | 4557 | 4558 | 4559 | 4560 | T | Netzwelt.de Internet |
| 346 | 4561 | 4562 | 4563 | 4564 | 4565 | T | Netzwelt.de Mobile |
| 351 356 | 4566 | 4567 | 4568 | 4569 | 4570 | T | Netzwelt.de Software |
| 361 | 4571 | 4572 | 4573 | 4574 | 4575 | T | Netzwelt.de Video |
| 366 371 | 4576 | 4577 | 4578 | 4579 | 4580 | L | |
| 376 | 4581 | 4582 | 4583 | 4584 | 4585 | L | |
| 381 | 4586 | 4587 | 4588 | 4589 | 4590 | L | |
| 386 391 | 4591 | 4592 | 4593 | 4594 | 4595 | L | |
| 396 | 4596 | 4597 | 4598 | 4599 | 4600 | L | |
| 401 | 4601 | 4602 | 4603 | 4604 | 4605 | L | |
| 406 411 | 4606 | 4607 | 4608 | 4609 | 4610 | F | vorname.com - Startseite |
| 416 | 4611 | 4612 | 4613 | 4614 | 4615 | | vorname.com - Suche |
| 421 426 | 4616 | 4617 | 4618 | 4619 | 4620 | L | |
| 431 | 4621 | 4622 | 4623 | 4624 | 4625 | L | |
| 436 441 | 4626 | 4627 | 4628 | 4629 | 4630 | F | BILD.de Digital - Downloads |
| 446 | 4631 | 4632 | 4633 | 4634 | 4635 | T | BILD.de Digital - Handy |
| 451 | 4636 | 4637 | 4638 | 4639 | 4640 | T | BILD.de Digital - Internet |
| 456 461 | 4641 | 4642 | 4643 | 4644 | 4645 | T | BILD.de Digital - Computer |
| 466 | 4646 | 4647 | 4648 | 4649 | 4650 | F | BILD.de Auto - Tuning & Zubehör |
| 471 476 | 4651 | 4652 | 4653 | 4654 | 4655 | T | BILD.de Auto - Service |
| 481 | 4656 | 4657 | 4658 | 4659 | 4660 | T | BILD.de Auto - Tests |
| 486 | 4661 | 4662 | 4663 | 4664 | 4665 | T | BILD.de Auto - Gebrauchtwagen |
| 491 | 4666 | 4667 | 4668 | 4669 | 4670 | F | BILD.de Auto - Neuwagen |
| | 4671 | 4672 | 4673 | 4674 | 4675 | | BILD.de Auto |
| | 4676 | 4677 | 4678 | 4679 | 4680 | T | BILD.de Auto - News |
| | 4681 | 4682 | 4683 | 4684 | 4685 | | BILD.de Digital |
| | 4686 | 4687 | 4688 | 4689 | 4690 | | BILD.de Geld |

HEINRICH HEINE
UNIVERSITÄT DÜSSELDORF

FORSCHUNGSVERBUND NRW
DIGITALE GESELLSCHAFT

## Structure by online format

| full entity | single entity |
|---|---|

## Structure **FULL** entity by business model

| e-Commerce | Context | **Content** | Connection | Games |
|---|---|---|---|---|

## Structure **SINGLE** entity with **CONTENT** business model by genre

## Use structure for CharmStats and data documentation

| full entity | single entity | business model (content) | genre |
|---|---|---|---|

## The case: MA Intermedia Plus, online-tranche

FORSCHUNGSVERBUND NRW
DIGITALE GESELLSCHAFT

**2016**

| offering name | variable stem | online format | business model | genre | terminal device | measurement name | measurement label |
|---|---|---|---|---|---|---|---|
| | Spalte22 | Spalte3 | Spalte4 | Spalte5 | Spalte 73 | Spalte6 | Spalte7 |
| BILD (df Gesamt) | bild | GA | CT | | df | GA_bild_CT_pidf | Page Impressions BILD, Gesamtangebot: Content (df) |
| BILD (if Website Angebot) | bild | GA | CT | | if | GA_bild_CT_piif | Page Impressions BILD, Gesamtangebot: Content (if) |
| BILD (mf Android Phone App) | bild | GA | CT | | mf Android Phone App | GA_bild_CT_pimfa | Page Impressions BILD, Gesamtangebot: Content (mf Andr |
| BILD (mf Gesamt) | bild | GA | CT | | mf | GA_bild_CT_pimf | Page Impressions BILD, Gesamtangebot: Content (mf) |
| BILD (mf iPhone App) | bild | GA | CT | | mf iPhone App | GA_bild_CT_pimfi | Page Impressions BILD, Gesamtangebot: Content (mf iPhor |
| BILD (mf MEW Angebot) | bild | GA | CT | | mf MEW | GA_bild_CT_pimfm | Page Impressions BILD, Gesamtangebot: Content (mf MEW |
| BILD (mf Phone App Angebot) | bild | GA | CT | | mf Phone App | GA_bild_CT_pimfp | Page Impressions BILD, Gesamtangebot: Content (mf Phor |
| BILD (mf Windows Phone App) | bild | GA | CT | | mf | GA_bild_CT_pimf | Page Impressions BILD, Gesamtangebot: Content (mf) |
| BILD Auto (df BE) | bild | EA | CT | Auto | df | EA_bild_CT_A_pidf | Page Impressions BILD - Auto #Auto (df) |
| BILD Auto (if BE) | bild | EA | CT | Auto | if | EA_bild_CT_A_piif | Page Impressions BILD - Auto #Auto (if) |
| BILD Auto (mf Android Phone App BE) | bild | EA | CT | Auto | mf Android Phone App | EA_bild_CT_A_pimfa | Page Impressions BILD - Auto #Auto (mf Android Phone App |
| BILD Auto (mf iPhone App BE) | bild | EA | CT | Auto | mf iPhone App | EA_bild_CT_A_pimfi | Page Impressions BILD - Auto #Auto (mf iPhone App) |
| BILD Auto (mf MEW BE) | bild | EA | CT | Auto | mf MEW | EA_bild_CT_A_pimfm | Page Impressions BILD - Auto #Auto (mf MEW) |
| BILD Bundesliga (df BE) | bild | EA | CT | Auto | df | EA_bild_CT_A_pidf | Page Impressions BILD - Auto #Auto (df) |
| BILD Bundesliga (if BE) | bild | EA | CT | Fussball | if | EA_bild_CT_Fu_piif | Page Impressions BILD - Bundesliga #Fussball (if) |
| BILD Bundesliga (mf MEW BE) | bild | EA | CT | Fussball | mf MEW | EA_bild_CT_Fu_pimfm | Page Impressions BILD - Bundesliga #Fussball (mf MEW) |
| BILD BYou (df BE) | bild | EA | CT | Markenname | df | EA_bild_CT_M_pidf | Page Impressions BILD - Byou #Markenname (df) |
| BILD BYou (if BE) | bild | EA | CT | Markenname | if | EA_bild_CT_M_piif | Page Impressions BILD - Byou #Markenname (if) |
| BILD BYou (mf MEW BE) | bild | EA | CT | Markenname | mf MEW | EA_bild_CT_M_pimfm | Page Impressions BILD - Byou #Markenname (mf MEW) |
| BILD Community (df BE) | bild | EA | CT | Forum | df | EA_bild_CT_Fo_pidf | Page Impressions BILD - Community #Forum (df) |
| BILD Digital - Computer (if BE) | bild | EA | CT | Digital | if | EA_bild_CT_DCo_piif | Page Impressions BILD - Digital #Digital (if) |
| BILD Digital - Handy (if BE) | bild | EA | CT | Digital | if | EA_bild_CT_DH_piif | Page Impressions BILD - Digital #Digital (if) |
| BILD Digital (df BE) | bild | EA | CT | Digital | df | EA_bild_CT_D_pidf | Page Impressions BILD - Digital #Digital (df) |

## Structure your (Big-) Data *smart*

- What is my research interest?

- What info in my data can I use as structure?

- What info do I need/want to add my data anyway?

- How should the final result look like?

  → CharmStats helps providing a formal structure

# How to manage Big-Data (well)

- Data exploitation is relevant for social science

  - New data sources through digitalization

  - Need for a set up of digital-data infrastucture to support secondary data analysis and re-analysis

- Data management as key to Big-Data in social science, follwoing the quality standards of academia:

  - transparency,

  - intersubjectivity,

  - replicability, sustainability, etc.

  → **Solution:** *digitized process chains*

# BACKUP-SLIDES

## CharmStats

## The case: add the information you want

# Step 3: data-processing with CharmStats Pro

## Categorical & nominal Variables

## metric Variables

## metric Variables

# Step 3: data-processing with CharmStats Pro

## The case of virtual Variables (text based)

## The case of virtual Variables (text based)

## Report as data documentation

- The **report** feature is used to output project content. Templates allow **user-defined** extraction of content, using an html formated document inter-spersed with keywords.

- Keywords and their interactions with the template are provided by plugins.

## Further Features

- CS Libary as an archive and exchange plattform

- Literature

- Work in Teams

  - In Mail

  - User-specified „tracking" and reports

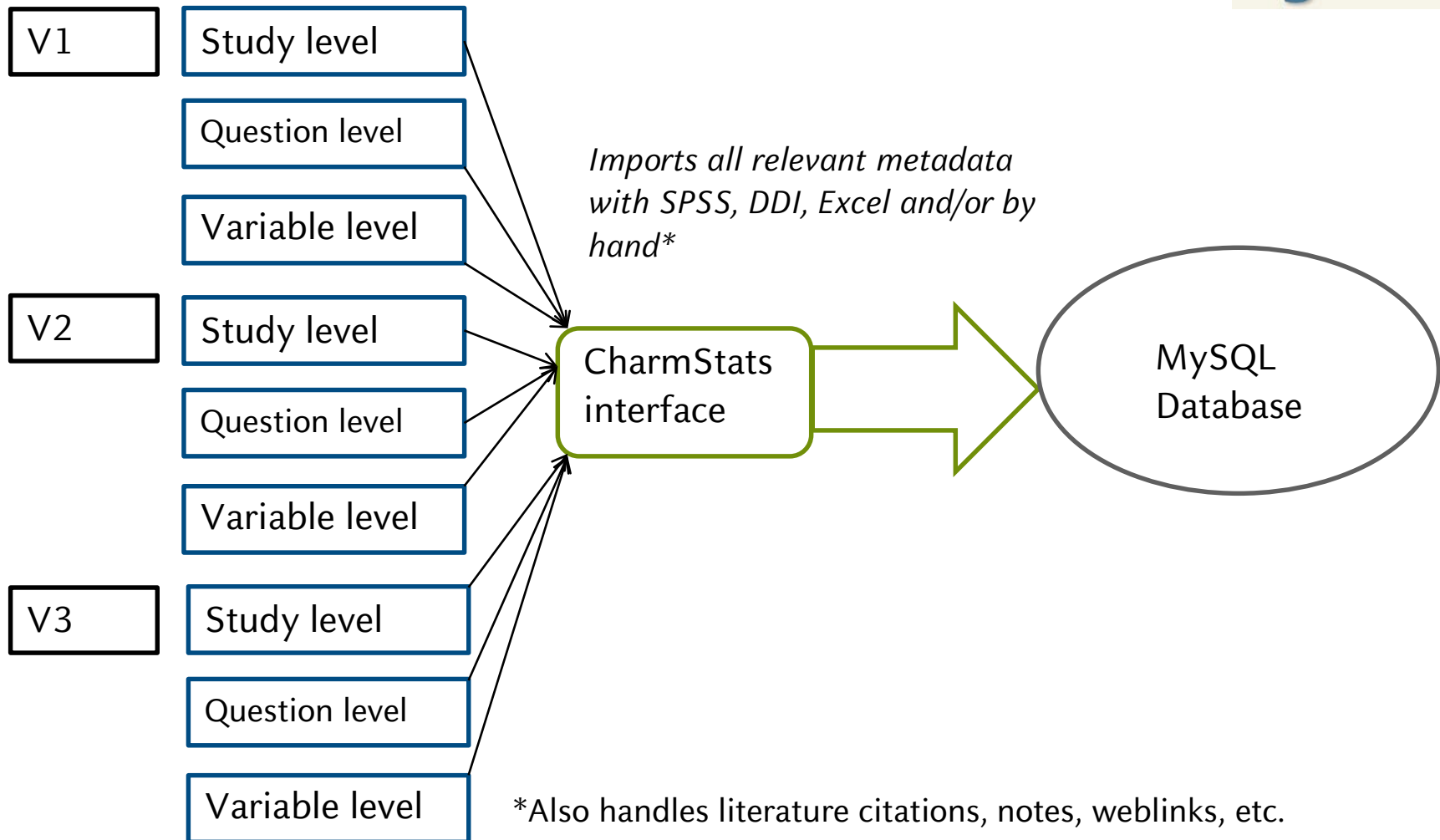- Usage for experiments: Make Intensions, concepts and tretments visible

## What is CharmStats for?

<u>Situation in (social) science</u>: **lack of tracability** through insufficient data-documentation, esp. in case of big data

<u>Advantages of CharmStats:</u>

1) Documentation of measurements *(shown today)*

2) Documentation of stimulus, treatments like pictures, questionnaire *(not shown today)*

3) Documentation of data processing and syntax *(shown today)*
   →esp. through CharmStats-Library *(launched for 2019)*

4) Output for different statistical programs *(shown today)*

## A digital solution



V1
- Study level
- Question level
- Variable level

V2
- Study level
- Question level
- Variable level

V3
- Study level
- Question level
- Variable level

*Imports all relevant metadata with SPSS, DDI, Excel and/or by hand**

CharmStats interface

MySQL Database

*Also handles literature citations, notes, weblinks, etc.

- yearly gathered cross-sectionally by ag.ma measuring media use for commercial purpose since 1954

- two main datasets: pressmedia (1954-2015) & radio (1977-2015)

  - Pressmedia: ~7.600 Variables, < 1,2 Mio cases

  - Radio: ~25.000 Variables, < 1,6 Mio cases

- One „new" dataset: online

- One combined dataset: Intermedia (Plus) (1999-2014-2017)

  - since 2014: ~18.600 Variables (downsized), < 1,2 Mio cases

  - Intermedia Plus is a result of a joint venture of ag.ma, agof & AGF/GfK

## Variables

- Socio-demographic variables

- variables regarding the daily routine (per hour and more detailed), e.g. sleeping, eating, driving to work, taking the bus, housekeeping, etc.

- free time activities, e.g. read books, do sports, go out, go to cinema, etc.

- habits of household and respondent, e.g. have a Laptop, a car, mobile, telephone, TV-Abo (e.g. Sky), a flat, pets, etc.

- media use, e.g. reading newspapers (SZ, FAZ, tz, …), magazines (Automobil, AutoBild, Lisa, Bravo, …), listening to radio per hour and more detailed (Antenne Bayern, Hessischer Rundfunk 1-3, WDR, 1 Live, KissFM)