

Hands-on-Übung zur Datenbereinigung mit



OpenRefine

Agnes Brauer
a.brauer@ub.uni-frankfurt.de

Heute ...

- Zellen teilen und (wieder) vereinigen mit Hilfe von Separatoren
- Facetieren, Filtern + Clustern von Daten
- Anreichern eigener Daten aus externen Quellen (Beispiel: Wikidata)

Einfaches Beispiel zum „Aufräumen“ von Daten

Beispiel: Schlagworte zum Thema „Digital Humanities“

- Erstellen Sie ein OpenRefine-Projekt mit Daten der Übungsdatei „DH.txt“
- Legen Sie eine einspaltige Tabelle an, in der in jeder Zeile ein Schlagwort steht („Edit Cells“ → „Split multi-valued cells...“ → als Separator ein Leerzeichen wählen)
- Verschaffen Sie sich einen ersten Überblick über die Daten mit Hilfe der Funktion „Text facet“
- Welche Schlagworte kommen im Datensatz mehrmals vor?
- Probieren Sie auf der Grundlage der Facette die Funktion „Cluster“ aus und beseitigen Sie Tippfehler
- Vergleichen Sie die Ergebnisse beim Einsatz unterschiedlicher Clustering-Methoden und -Algorithmen

Übung 2: „Aufräumen“

- Erstellen Sie das OpenRefine-Projekt mit Daten aus dem Dokument „AuszugBib.txt“
- Explorieren Sie die Daten und betreiben Sie etwas Datenbereinigung u.a. mithilfe der Funktionen „Edit Columns“ → „Split into several columns...“ Was ist ein geeigneter Separator?
- Transponieren Sie die Daten in mehrere Spalten: [Transpose](#) → [Columnize by Key/Value Columns](#)
- Überlegen Sie sich weitere sinnvolle „Aufräumarbeiten“ und probieren Sie sie einfach aus

Anwendungsmöglichkeiten – fortgeschrittene Funktionen

Reconcile & Match

- Vergleichen / Angleichen der eigenen Daten anhand von Datenbanken (z.B. Wikidata)
- Anreicherung von Daten (z.B. mit Identifiern oder geographischen Koordinaten)
- Verlinkung von Daten

Reconciliation is the process of matching name strings to identifiers of entities in a database like an authority file, Wikidata etc. This is useful whenever you want to merge differing name strings for the same person in your data or when you want to fetch additional data from the target database you are reconciling against.

Reconcile & Match

- Legen Sie ein neues Projekt anhand der Datei „FriedenspreisAb2000.txt“ an
- Splitten Sie die Daten ggf. mithilfe eines geeigneten Separators
- Führen Sie eine Reconciliation mit Wikidata durch
- Fügen Sie neue Spalten auf Basis der Reconciliation hinzu, z.B.:
 - Geburtsort
 - Koordinaten des Geburtsortes
 - Geburtsdatum
 - Nationalität

Tipps & Tricks / Links / Literatur

Blogs

<https://histhub.ch/cat/net/blog/openrefine/> (Blog-Serie zur Arbeit an historischen Daten mit OpenRefine)

<http://blog.lobid.org/2018/08/27/openrefine.html> (einfache Anreicherung von Daten in OpenRefine mit [Personen-] Daten aus der GND via lobid.org)

Literatur

Ruben Verborgh/Max de Wilde, Using OpenRefine. The essential OpenRefine guide that takes you from data analysis and error fixing to linking your dataset to the Web (Community experience distilled), Birmingham, Mumbai 2013. [[Online-Ressource über UB FFM](#)]

Danke für Ihre Aufmerksamkeit!

Workshop konzipiert in Anlehnung an "[Library Carpentry: OpenRefine Lessons for Librarians.](#)" (2016)