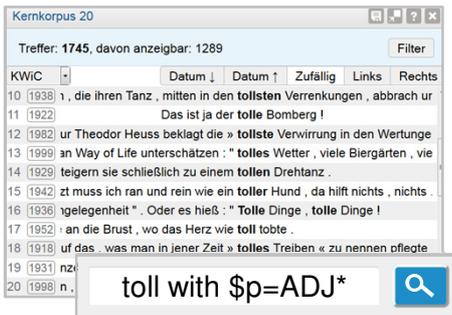


Korpus-basierte linguistische Recherche und Analyse mit Hilfe von Data-Mining („KobRA“)

Erprobung innovativer Data-Mining-Verfahren für die empirische Arbeit mit strukturierten Sprachressourcen



Empirische Sprachforschung mithilfe großer Textkorpora



Ausgangslage: Notwendigkeit aufwändiger manueller Analyse langer KWIC-Ergebnislisten

- Aussondern falsch positiver Treffer
- Disambiguieren nach Lesarten
- Nachklassifikation und Annotation gesuchter Belege
- Auffinden seltener, aber interessanter Treffer

Ziel: Beschleunigung und Verbesserung der Ergebnisse, neue Erkenntnisse in Bezug auf Voraussetzungen, Zuschnitt und Nutzen geeigneter Data-Mining-Verfahren

Fallstudien

Diachrone Linguistik

Entwicklung und Ausdifferenzierung von Stützverbgefügen (*ins Rollen bringen*)

Korpusbasierte Lexikographie

Rekonstruktion von Bedeutungswandel

Varietätenlinguistik

Sprachmerkmale und Variation in Genres computervermittelter Kommunikation



Sprachressourcen

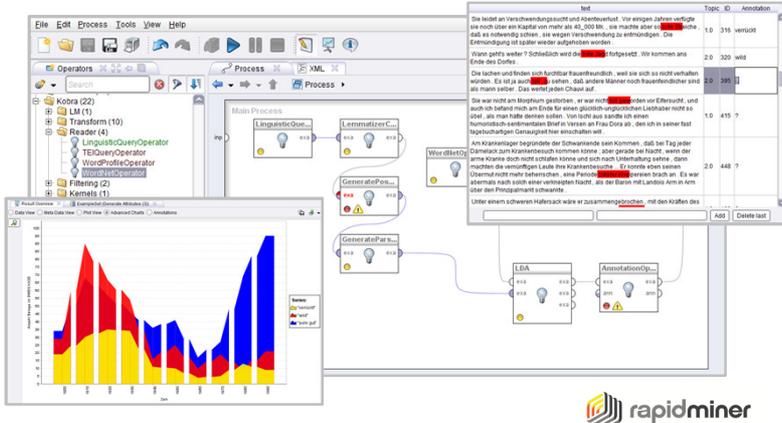
Sprachressourcen im Rahmen des DWDS (Digitales Wörterbuch der deutschen Sprache) Deutsches Textarchiv (DTA)

Tübinger Baubanken WebLicht

Deutsches Referenzkorpus (DeReKo) Wikipedia



Data-Mining-Umgebung und KobRA-Erweiterung



Basis: RapidMiner

- eine der meistgenutzten Data-Mining-Umgebungen
- Vielzahl an schachtelbaren Operatoren für Vorverarbeitung, Analyse und Visualisierung

Erweiterung: KobRA

- direkter Zugriff auf Sprachressourcen
- Methoden zur Extraktion linguistischer Merkmale und Metadaten
- Methoden zur Klassifikation und Disambiguierung
- integrierte Annotationsumgebung