

SIMPSON'S SEARCH – ZWISCHENBERICHT

INFORMATION RETRIEVAL SS 2019

DER DATENSATZ

- 4 CSV Dateien
 - Script_Lines
 - Locations
 - Episodes
 - Characters
- Datenbankstruktur

	Zeilen/ Einträge	Spalten	Inhalt
Script_Lines	158.315	13	u.a. Spoken Words inkl. Normalisierung, Character Referenz, Location Referenz, Word Counts
Locations	4.459	3	Liste aller Locations
Episodes	600	12	u.a. Episodenname, Views, Ratings
Characters	6.722	4	Liste aller Charaktere

DER DATENSATZ

- **Script Lines** 158.315 Zeilen (bis Gesamt-Episode 568; Season 26; Episode 16)
 - Spalten: id, episode_id, number, raw_text, timestamp_in_ms,
 - speaking_line, character_id, location_id, raw_character_text,
 - raw_location_text, spoken_words, normalized_text, word_count)
- **Locations** 4.459 Zeilen
 - Spalten: id, name, normalized_name
- **Episodes**: 600 Zeilen
 - Spalten: id
 - title, original_air_date, production_code, season,
 - number_in_season, number_in_series, us_viewers_in_millions, views,
 - imdb_rating, imdb_votes, image_url, video_url)
- **Characters**: 6.722 Zeilen
 - Spalten: id, name, normalized_name, gender

STAND DER SUCHMASCHINE

- .net Core Webanwendung
- Indexierung und Suche ermöglicht durch lucene.net
- Lucene.net: c# clone der Apache.Lucene Search Library
- Indexierung des Datensatzes und der Suche funktionieren bereits

DERZEITIGES VERHALTEN DER SUCHMASCHINE

- Stichwortsuche gesprochener Wörter/Sätze von Charakteren
- Die Suchanfrage wird mit dem string in dem Feld „raw_text“ verglichen
- Rating der Ergebnisse anhand exakter Übereinstimmung mit Suchanfrage

TOPICS UND TOPICKLASSEN - DATENBANKBASIERT

- [character name]
 - Out: Folgen, in denen der Charakter vorkommt
- [character name] [character name]
 - Out: Szenen, in denen beide Charaktere vorkommen / interagieren
- [character name] catchphrase
 - Out: Sätze, die der Charakter am häufigsten benutzt
- [character name] [location]
 - Out: Szenen mit dem Charakter an angegebenem Ort *oder* direkt Rede des Charakters an dem Ort

TOPICS UND TOPICKLASSEN – INTERPRETATIVE SUCHE

- [technology]
internet, computing
AI, artificial intelligence
- [politics]
abortion
gun control
Trump
Bush Jr
politics (in general)
Elections
immigrants
- [violence]
crime
military
Iraq War
terrorism
- [culture]
travel
super heroes
sports
musicals, singing
fantasy (genre)

TOPICS UND TOPICKLASSEN – INTERPRETATIVE SUCHE

- [social issues]

feminism

alcohol, alcoholism

vegetarians, vegans

patriotism

racism

LGBT

marriage

mobbing, bullies

(illegal) drugs

- [Holidays]

Halloween

Christmas

Valentine's day

St. Patricks Day

Thanks Giving

TO DO:

- verbinden der Datensätze in einem Index, für komplexere Suchanfragen
- Topic abhängige Suche implementieren
- Beeinflussen des Ratings der Suchergebnisse