

Datensatzakquisition (engl. dataset acquisition)				Datensatzverarbeitung (engl. dataset pipeline)								Datensaufbereitung (engl. dataset preparation)			
															
Datensatzrecherche (engl. dataset research)	Datensammlung (engl. dataset collection)	Datensatzprüfung (engl. dataset check)	Datensatzauswahl (engl. dataset selection)	Explorative Datenanalyse (engl. exploratory data analysis)								Datensatzbereinigung (engl. dataset cleaning)	Datensatzvalidierung (engl. dataset validation)		
Suchstrategie (engl. search strategy)	Download (engl. data download)	Speicherung (engl. data storage)			Datenstrukturanalyse (engl. data structure analysis)	Analyse der strukturierten Daten (engl. analysis of structured data)			Analyse der unstrukturierten Daten (engl. analysis of unstructured data)			Fehlwertbehandlung (engl. missing value handling)	Duplikatentfernung (engl. duplicate removal)		
					Ortsdaten (engl. location data)	Zeitdaten (engl. time data)	Bewertungen (engl. ratings)		Fehlwerterkennung (engl. missing value detection)	Duplikaterkennung (engl. duplicate detection)	Textlängenanalyse (engl. text length analysis)				
															
Rohzeichenfolgen (engl. raw strings)															
Suchbegriffe <ul style="list-style-type: none"><li>„Beschwerden“</li><li>„Kommentar“</li><li>„NLP“</li><li>„Kundenbeschwerden“</li><li>„Online-Kommentare“</li><li>„Produktbewertungen“</li></ul>		+ lokale Sprachierung				+ Bundesstaaten	+ Zeitraumbestimmung		+ NaNs (Not a Number)			+ Imputation			
Anforderungen <ul style="list-style-type: none"><li>organischer Datensatzes</li></ul>							Datumsanalyse		+ NaTs (Not a Text)				„Strategies for handling missing data in text include imputation where missing values are estimated based on existing data and removal where incomplete observations are excluded from the analysis. These strategies help in ensuring the completeness and integrity of the text data [12].“ (Upadhye, 2020, p. 207)		
													# Zellen ohne Text entfernen (0,4% ist vermaschlicht) df_clean = df.dropna(subset=[...])		
													# Oder Mindestänge erzwingen df_clean = df[df['text'].str.len() > 20]		