

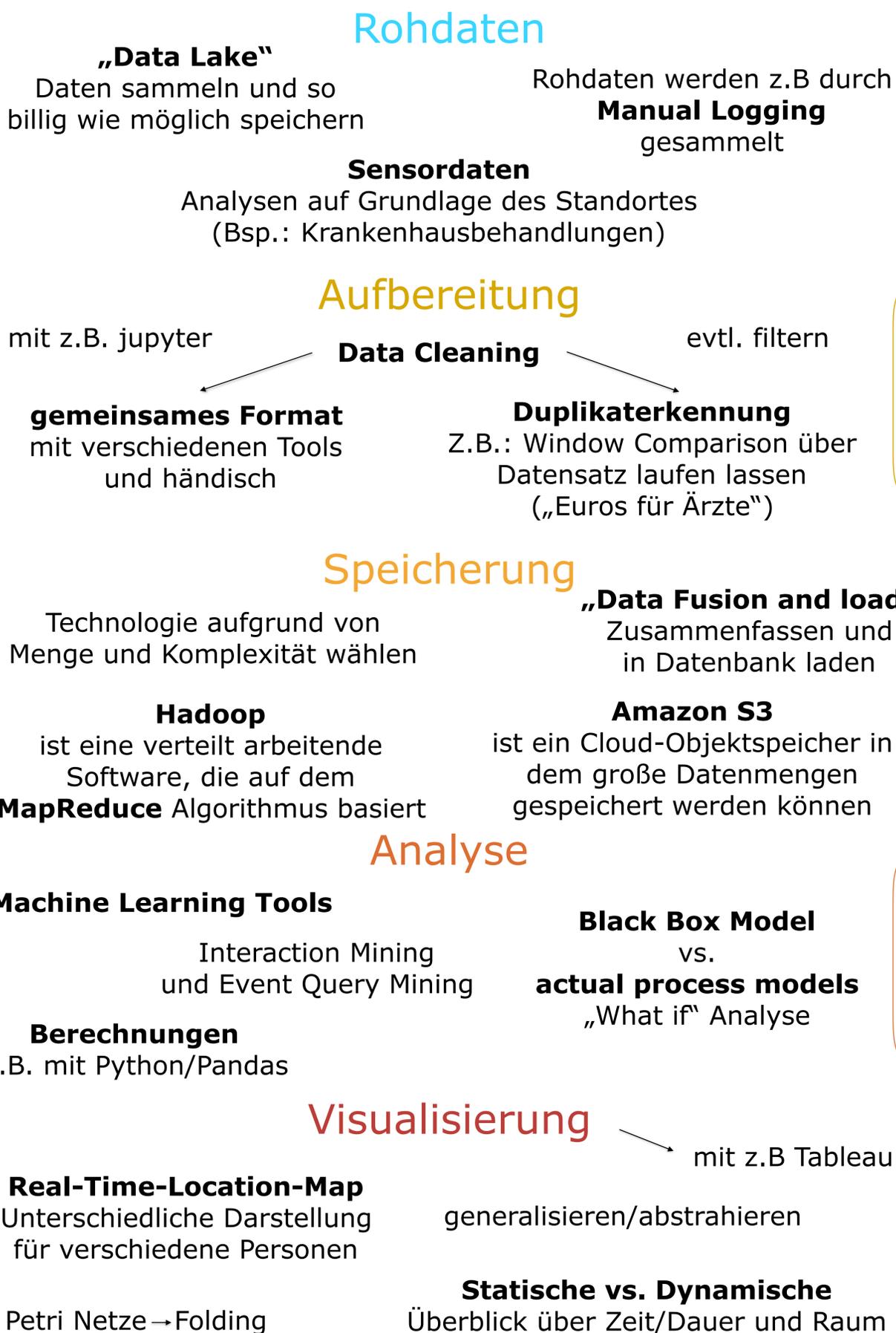
Data Engineering Prozesse

Cold Warm Hot ist eine Strategie um Data Engineering Prozesse anschaulich darzustellen zu können. Es gibt kalte Daten mit denen noch nichts gemacht wurde, durch die Aufbereitung werden sie zu warmen Daten und schließlich nach der Analyse liegen heiße Daten vor. Zuerst müssen die **Rohdaten** gesammelt werden, im nächsten Schritt werden sie vereinheitlicht, gefiltert und Duplikate werden eliminiert, das nennt man **Aufbereitung**. Im dritten Schritt werden die Daten mittels einer entsprechenden Technologie **gespeichert**, damit im nächsten Schritt die **Analyse** erleichtert wird. Hier sind kreative Lösungen gefragt. Zuletzt werden die Analyseergebnisse **visualisiert** wobei auf eine zielgruppengerechte Darstellungsart geachtet werden muss.

keine Fragen keine Antworten

Fragen stellen

Fragen beantworten



verschiedene Datenquellen

z.B.: PDF, XLS, Web
→ nicht einheitlich
„Euros für Ärzte“: Daten aus 40 verschiedene Quellen

Duplikat-erkennung

false negative: Duplikat-Paar wird fälschlicherweise nicht erkannt → ist z.B. im Datenjournalismus die besser Alternative

Was ist MapReduce?

Programmiermodell für nebenläufige Berechnungen über große Datenmengen auf Computerclustern

Vereinigung von Wissen

Es muss Informatiker und z.B Journalistenwissen vereinigt werden um sinnvolle Antworten zu finden/erkennen

Warum Daten-visualisierung?

Entscheidungsträger besitzen weniger Wissen als Entwickler
→ Mehrwert in möglichst kurzer Zeit darstellen