



Scrape uncleaned data

Cleaning Pipeline

clean_general_pipeline

clean_company_pipeline

clean_rating_pipeline

clean_hq_pipeline

clean_salary_pipeline

clean_jobtype_pipeline

clean_size_pipeline

clean_type_pipeline

clean_sector_pipeline

clean_revenue_pipeline

EDA!



- split company and rating into separate columns
- aggregate/group company names that are spelled differently
- remove rows of non-company names

- remove rows with ≥ 5 NaNs
- remove specified columns: "job title", "founded", "job description", "industry"
- index remover (removes any index columns)

- fill NaN ratings with average company rating
- fill leftover NaN ratings by global average rating

- remove non-headquarter rows

- drop rows with NaN salary estimates
- parse out the float average salary estimate per instance/record

- remove NaN rows for job type
- parse out the job-type from string

- drop rows with NaNs