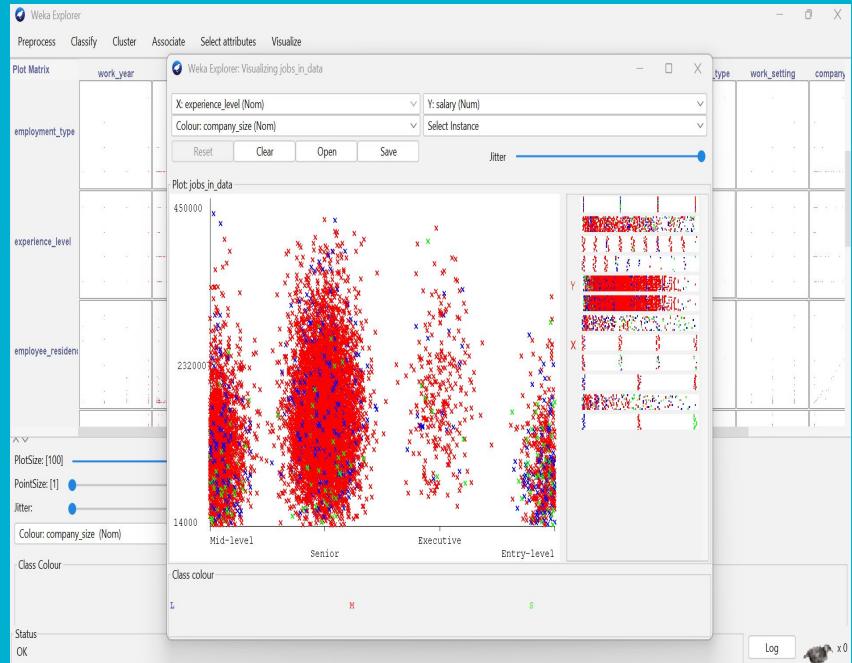


Pre-processing: Jobs and Salaries in Data Science

Team Members: Puja Shah & Sanjida Chowdhury

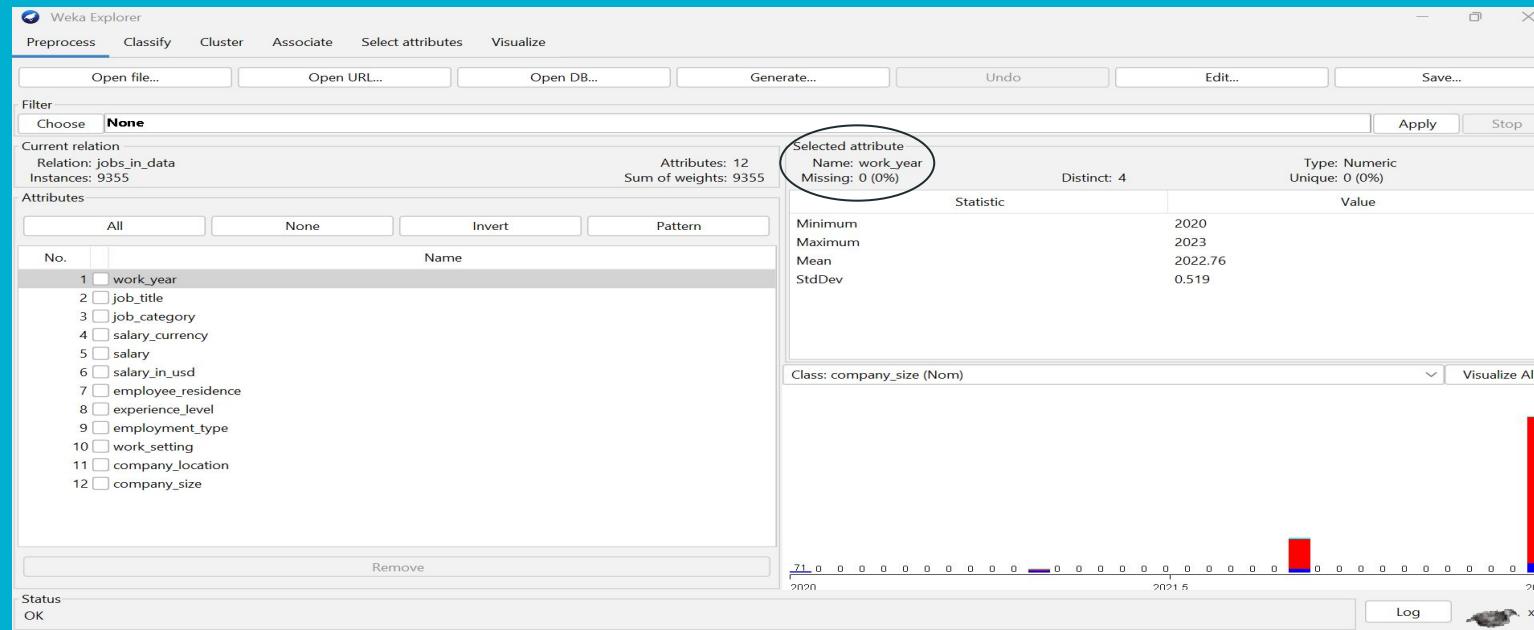
Dataset: “Jobs and Salaries in Data Science”

- The dataset captures a comprehensive snapshot of experienced level in jobs and salaries in the data science field for year 2023.
- The goal is to explore relationships, identify patterns, and understand factors influencing salary levels in data science.
- Attributes: work year, job title, job category, salary currency, salary, salary in USD, employee residence, experience level, employment type, work setting, company location, and company size.



Dealing with Missing Values

Original dataset without missing values:



Dealing with Missing Values

Intentionally created missing values (Replace with missing value):

Weka Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Open file... Open URL... Open DB... Generate... Undo Edit... Save...

Filter Choose ReplaceWithMissingValue -5-first-1st-5 1-P 0.1

Current relation Relation: jobs_in_data-weka.filters.unsupervised.attribute.ReplaceWithMissingValue-RB... Attributes: 12 Instances: 9355 Sum of weights: 9355

Attributes

No.	Name
1.	work_year
2.	job_title
3.	job_category
4.	salary_currency
5.	salary
6.	salary_in_usd
7.	employee_residence
8.	experience_level
9.	employment_type
10.	work_setting
11.	company_location
12.	company_size

Selected attribute Name: work_year Type: Numeric
Missing: 943 (10%) Distinct: 4 Unique: 0 (0%)

Statistic	Value
Minimum	2020
Maximum	2023
Mean	2022.758
StdDev	0.532

Class: company_size (Nom) Visualize All

Remove Log x0

Status OK

Viewer

Relation: jobs_in_data-weka.filters.unsupervised.attribute.ReplaceWithMissingValue-RB-first-last-51-P0.1

No.	1: work_year	2: job_title	3: job_category	4: salary_currency	5: salary	6: salary_in_usd	7: employee_residence	8: experience_level	9: employment_type	10: work_setting	11: company_location	12: company_size	
	Nominal	Nominal	Nominal	Nominal	Numeric	Numeric	Nominal	Nominal	Nominal	Nominal	Nominal	Nominal	
1	2023.0	Data Dev..	Data Engineer..	EUR	88000.0	Germany	Mid-level	Full-time	Hybrid	Germany	L		
2	2023.0	Data Arc..	Data Architect..	USD	186000.0	United States	Senior	Full-time	In-person	United States	M		
3		Data Arc..	Data Architect..	USD	81800.0	United States	Senior	Full-time	In-person	United States	M		
4	2023.0	Data Scie..	Data Science a..	USD	212000.0	United States	Senior	Full-time	In-person	United States	M		
5	2023.0	Data Scie..	Data Science a..	USD	93300.0	United States	Senior	Full-time	In-person	United States	M		
6	2023.0	Data Scie..	Data Science a..	USD	130000.0	United States	Senior	Full-time	Remote	United States	M		
7	2023.0	Data Scie..	Data Science a..	USD	100000.0	United States	Senior	Full-time	Remote	United States	M		
8	2023.0	Machine ...	Machine Learn..		224400.0	United States	Mid-level	Full-time	In-person	United States	M		
9	2023.0	Machine ...	Machine Learn..	USD	138700.0	United States	Mid-level	Full-time	In-person	United States	M		
10		Data Engi..	Data Engineer..	USD	210000.0	United States	Executive	Full-time	Remote	United States	M		
11	2023.0	Data Engi..	Data Engineer..	USD	168000.0	United States	Executive	Full-time	Remote	United States	M		
12	2023.0		Machine Learn..	USD	224400.0	United States	Senior	Full-time	In-person	United States	M		
13	2023.0		Machine Learn..	USD	138700.0	United States	Senior	Full-time	In-person	United States	M		
14	2023.0	Data Scie..	Data Science a..	GBP	350000.0	United Kingdom	Mid-level	Full-time	In-person	United States	M		
15	2023.0	Data Scie..	Data Science a..	GBP	30000.0	United Kingdom	Mid-level	Full-time	In-person	United Kingdom	M		
16	2023.0	Data Ana..	Data Analysis	USD	95000.0	95000.0	Entry-level	Full-time	In-person	United States	M		
17	2023.0	Data Ana..	Data Analysis	USD	75000.0	75000.0	United States	Entry-level	In-person	United States	M		
18	2023.0	Data Scie..	Data Science a..	USD	300000.0	300000.0	United States	Senior	Full-time	In-person	United States	M	
19		Data Scie..	Data Science a..	USD	234000.0	234000.0	United States	Senior	Full-time	In-person	United States	M	
20	2023.0		Leadership an..	USD	140000.0	United States	Mid-level	Full-time	In-person	United States	M		
21	2023.0	Analytics ..	Leadership an..	USD	120000.0	United States	Mid-level	Full-time	In-person	United States	M		
22	2023.0	Machine ...	Machine Learn..	USD	204500.0	United States	Mid-level	Full-time	In-person	United States	M		
23	2023.0	Machine ...	Machine Learn..	USD	142200.0	United States	Mid-level	Full-time	In-person	United States	M		
24	2023.0	Data Ana..	Data Analysis	USD	155000.0	United States	Mid-level	Full-time	In-person	United States	M		
25	2023.0	Data Ana..	Data Analysis	USD	110000.0	110000.0	United States	Mid-level	Full-time	In-person	United States	M	
26	2023.0	Machine ...	Machine Learn..	USD	266500.0	266500.0	United States	Senior	Full-time	United States	M		
27	2023.0	Machine ...	Machine Learn..		152000.0	United States	Senior	Full-time	In-person	United States	M		
28	2023.0	Applied S...	Data Science a..		222200.0	222200.0	United States	Mid-level	Full-time	United States	L		

Add instance Undo OK Cancel

Dealing with Missing Values

Replace missing values (Imputation with Mean):

Weka Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Open file... Open URL... Open DB... Generate... Undo Edit... Save... Filter Choose ReplaceMissingValues

Current relation: Relation: jobs_in_data-weka.filters.unsupervised.attribute.ReplaceWithMissingValue-Rfirst-last-1st-P0.1-weka.filters.unsupervised.attribute.ReplaceMissingValues Attributes: 12 Instances: 9355 Sum of weights: 9355

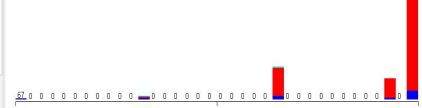
Attributes

No.	Name
1	work_year
2	job_title
3	job_category
4	salary_currency
5	salary
6	salary_in_usd
7	employee_residence
8	experience_level
9	employment_type
10	work_setting
11	company_location
12	company_size

Selected attribute: Name: work_year Type: Numeric
Missing: 0 (0%) Distinct: 5 Unique: 0 (0%)

Statistic	Value
Minimum	2020
Maximum	2023
Mean	2022.758
StdDev	0.496

Class: company_size (Nom) Visualize All



Status: OK

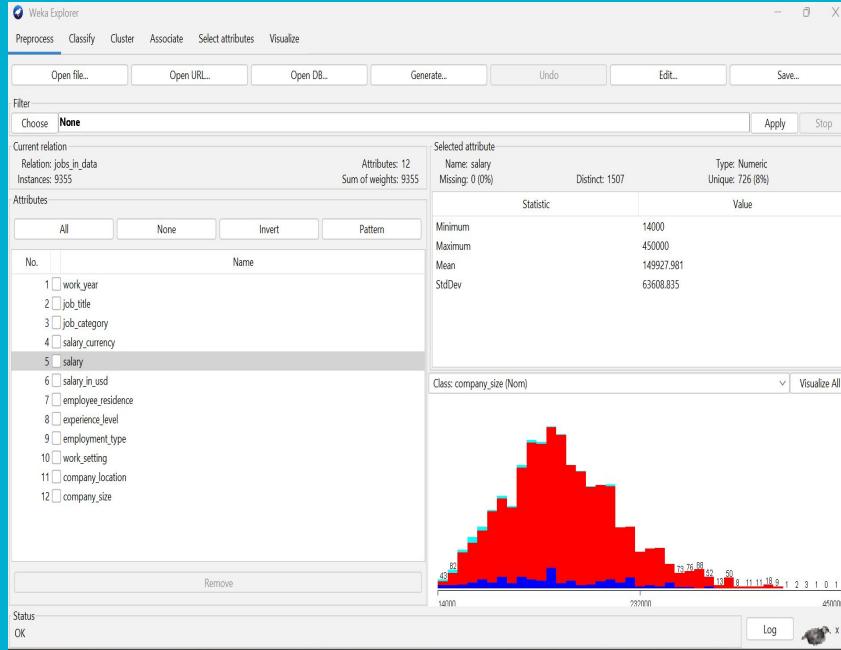
No.	1: work_year	2: job_title	3: job_category	4: salary_currency	5: salary	6: salary_in_usd	7: employee_residence	8: experience_level	9: employment_type	10: work_setting	11: company.location	12: company.size
	Numeric	Nominal	Nominal	Nominal	Numeric	Numeric	Nominal	Nominal	Nominal	Nominal	Nominal	Nominal
1	2023.0	Data Dev...	Data Engineer...	EUR	88000.0	150396.2287...	Germany	Mid-level	Full-time	Hybrid	Germany	L
2	2023.0	Data Arc...	Data Architect...	USD	186000.0	186000.0	United States	Senior	Full-time	In-person	United States	M
3	2022.7582...	Data Arc...	Data Architect...	USD	81800.0	81800.0	United States	Senior	Full-time	In-person	United States	M
4	2023.0	Data Sci...	Data Science a...	USD	212000.0	212000.0	United States	Senior	Full-time	In-person	United States	M
5	2023.0	Data Sci...	Data Science a...	USD	93300.0	93300.0	United States	Senior	Full-time	In-person	United States	M
6	2023.0	Data Sci...	Data Science a...	USD	130000.0	130000.0	United States	Senior	Full-time	Remote	United States	M
7	2023.0	Data Sci...	Data Science a...	USD	100000.0	100000.0	United States	Senior	Full-time	Remote	United States	M
8	2023.0	Machine ...	Machine Learn...	USD	224400.0	224400.0	United States	Mid-level	Full-time	In-person	United States	M
9	2023.0	Machine ...	Machine Learn...	USD	138700.0	138700.0	United States	Mid-level	Full-time	In-person	United States	M
10	2022.7582...	Data Engi...	Data Engineer...	USD	210000.0	210000.0	United States	Executive	Full-time	Remote	United States	M
11	2023.0	Data Engi...	Data Engineer...	USD	149462.8...	168000.0	United States	Executive	Full-time	Remote	United States	M
12	2023.0	Data Engi...	Machine Learn...	USD	224400.0	224400.0	United States	Senior	Full-time	In-person	United States	M
13	2023.0	Data Engi...	Machine Learn...	USD	138700.0	138700.0	United States	Senior	Full-time	In-person	United States	M
14	2023.0	Data Sci...	Data Science a...	GBP	35000.0	43064.0	United Kingdom	Mid-level	Full-time	In-person	United States	M
15	2023.0	Data Sci...	Data Science a...	GBP	30000.0	150396.2287...	United Kingdom	Mid-level	Full-time	In-person	United Kingdom	M
16	2023.0	Data Ana...	Data Analysis	USD	95000.0	95000.0	United States	Entry-level	Full-time	In-person	United States	M
17	2023.0	Data Ana...	Data Analysis	USD	75000.0	75000.0	United States	Entry-level	Full-time	In-person	United States	M
18	2023.0	Data Sci...	Data Science a...	USD	300000.0	300000.0	United States	Senior	Full-time	In-person	United States	M
19	2022.7582...	Data Sci...	Data Science a...	USD	234000.0	234000.0	United States	Senior	Full-time	In-person	United States	M
20	2023.0	Data Engi...	Leadership an...	USD	149462.8...	140000.0	United States	Mid-level	Full-time	In-person	United States	M
21	2023.0	Analytics ...	Leadership an...	USD	120000.0	150396.2287...	United States	Mid-level	Full-time	In-person	United States	M
22	2023.0	Machine ...	Machine Learn...	USD	149462.8...	204500.0	United States	Mid-level	Full-time	In-person	United States	M
23	2023.0	Machine ...	Machine Learn...	USD	149462.8...	142200.0	United States	Senior	Full-time	In-person	United States	M
24	2023.0	Data Ana...	Data Analysis	USD	149462.8...	155000.0	United States	Mid-level	Full-time	In-person	United States	M
25	2023.0	Data Ana...	Data Analysis	USD	110000.0	110000.0	United States	Senior	Full-time	In-person	United States	M
26	2023.0	Machine ...	Machine Learn...	USD	266500.0	266500.0	United States	Senior	Full-time	In-person	United States	M
27	2023.0	Machine ...	Machine Learn...	USD	149462.8...	152000.0	United States	Senior	Full-time	In-person	United States	M
28	2023.0	Anonymized S...	Data Science a...	USD	222200.0	222200.0	United States	Mid-level	Full-time	In-person	United States	I

Dealing with Missing Values

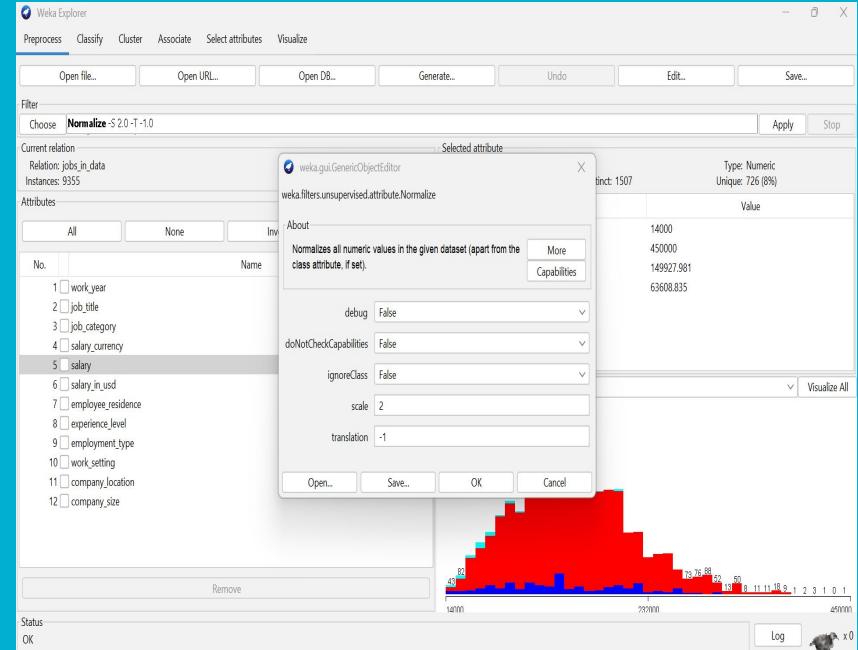
- The intentional creation of missing values allowed us to simulate scenarios where data may be incomplete.
- By deliberately introducing missing values in critical attributes such as salary and experience level, we test the dataset's resilience and its ability to handle gaps in information.
- Missing values were replaced using imputation with the mean of the respective attributes.
- This process ensures robustness testing and prepares the dataset for subsequent preprocessing steps.

Normalization (min-max)

Original data

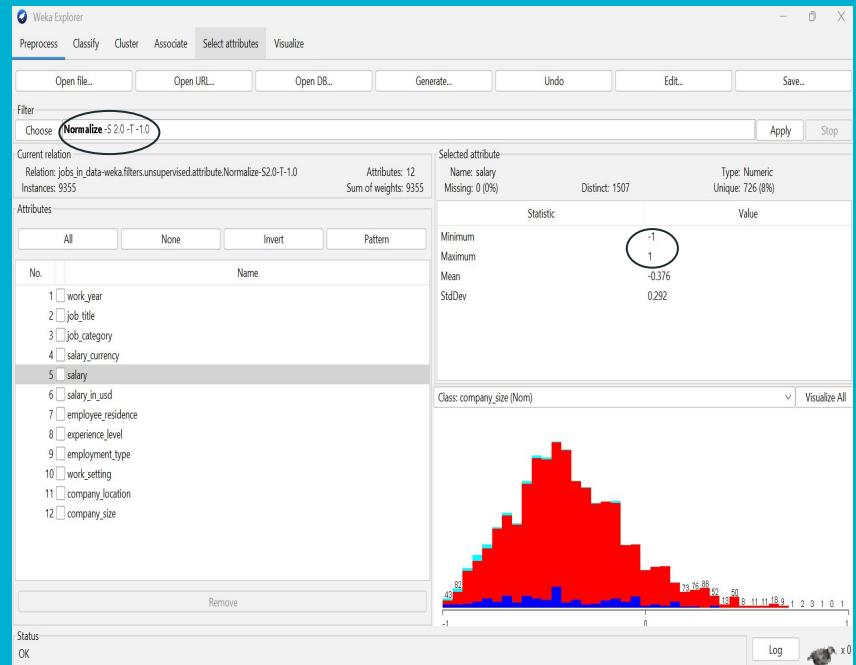


Applying 'Normalize' to data



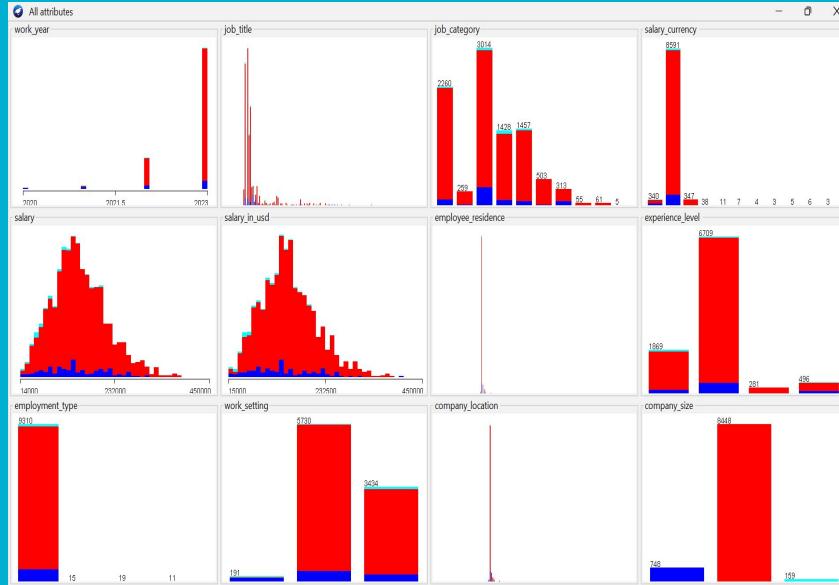
Normalization (min-max)

- Application of min-max scaling to standardize numerical attributes between -1 to 2.

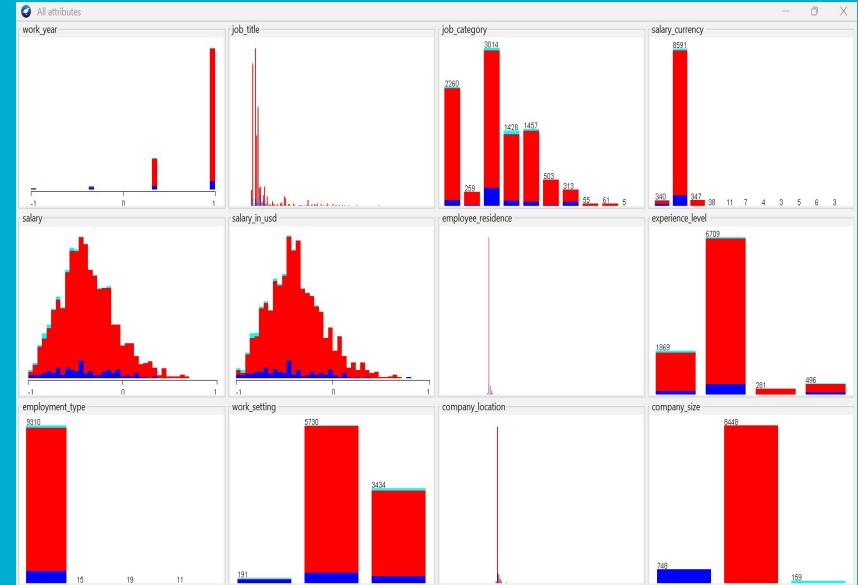


Normalization (min-max)

Orignal



After normalization



Normalization (min-max)

Orignal

Viewer
Relation: jobs_in_data

No.	1:work_year	2:job_title	3:job_category	4:salary_currency	5:salary_Nominal	6:salary_in_usd_Nominal	7:employee_residence_Nominal	8:experience_level_Nominal	9:employment_type_Nominal	10:work_setting_Nominal	11:company_location_Nominal	12:company_size_Nominal
1	2023.0	Data Dev..	Data Engineer..	EUR	88000.0	95012.0	Germany	Mid-level	Full-time	Hybrid	Germany	L
2	2023.0	Data Arc..	Data Architect..	USD	186000.0	186000.0	United States	Senior	Full-time	In-person	United States	M
3	2023.0	Data Arc..	Data Architect..	USD	81800.0	81800.0	United States	Senior	Full-time	In-person	United States	M
4	2023.0	Data Scie..	Data Science a..	USD	212000.0	212000.0	United States	Senior	Full-time	In-person	United States	M
5	2023.0	Data Scie..	Data Science a..	USD	93300.0	93300.0	United States	Senior	Full-time	In-person	United States	M
6	2023.0	Data Scie..	Data Science a..	USD	130000.0	130000.0	United States	Senior	Full-time	Remote	United States	M
7	2023.0	Data Scie..	Data Science a..	USD	100000.0	100000.0	United States	Senior	Full-time	Remote	United States	M
8	2023.0	Machine ..	Machine Learn..	USD	224400.0	224400.0	United States	Mid-level	Full-time	In-person	United States	M
9	2023.0	Machine ..	Machine Learn..	USD	138700.0	138700.0	United States	Mid-level	Full-time	In-person	United States	M
10	2023.0	Data Engi..	Data Engineer..	USD	210000.0	210000.0	United States	Executive	Full-time	Remote	United States	M
11	2023.0	Data Engi..	Data Engineer..	USD	168000.0	168000.0	United States	Executive	Full-time	Remote	United States	M
12	2023.0	Machine ..	Machine Learn..	USD	224400.0	224400.0	United States	Senior	Full-time	In-person	United States	M
13	2023.0	Machine ..	Machine Learn..	USD	138700.0	138700.0	United States	Senior	Full-time	In-person	United States	M
14	2023.0	Data Scie..	Data Science a..	GBP	35000.0	43064.0	United Kingdom	Mid-level	Full-time	In-person	United Kingdom	M
15	2023.0	Data Scie..	Data Science a..	GBP	30000.0	36912.0	United Kingdom	Mid-level	Full-time	In-person	United Kingdom	M
16	2023.0	Data Ana..	Data Analysis	USD	95000.0	95000.0	United States	Entry-level	Full-time	In-person	United States	M
17	2023.0	Data Ana..	Data Analysis	USD	75000.0	75000.0	United States	Entry-level	Full-time	In-person	United States	M
18	2023.0	Data Scie..	Data Science a..	USD	300000.0	300000.0	United States	Senior	Full-time	In-person	United States	M
19	2023.0	Data Scie..	Data Science a..	USD	234000.0	234000.0	United States	Senior	Full-time	In-person	United States	M
20	2023.0	Analytics ..	Leadership an..	USD	140000.0	140000.0	United States	Mid-level	Full-time	In-person	United States	M
21	2023.0	Analytics ..	Leadership an..	USD	120000.0	120000.0	United States	Mid-level	Full-time	In-person	United States	M
22	2023.0	Machine ..	Machine Learn..	USD	204500.0	204500.0	United States	Mid-level	Full-time	In-person	United States	M
23	2023.0	Machine ..	Machine Learn..	USD	142200.0	142200.0	United States	Mid-level	Full-time	In-person	United States	M
24	2023.0	Data Ana..	Data Analysis	USD	155000.0	155000.0	United States	Mid-level	Full-time	In-person	United States	M
25	2023.0	Data Ana..	Data Analysis	USD	110000.0	110000.0	United States	Mid-level	Full-time	In-person	United States	M
26	2023.0	Machine ..	Machine Learn..	USD	266500.0	266500.0	United States	Senior	Full-time	In-person	United States	M
27	2023.0	Machine ..	Machine Learn..	USD	152000.0	152000.0	United States	Senior	Full-time	In-person	United States	M
28	2023.0	Applied S...	Data Science a..	USD	222700.0	222700.0	United States	Mid-level	Full-time	In-person	United States	I

Add instance Undo OK Cancel

After normalization

Viewer
Relation: jobs_in_data-weka.filters.unsupervised.attribute.Normalizer-S2.0-T.1.0

No.	1:work_year	2:job_title	3:job_category	4:salary_currency	5:salary_Nominal	6:salary_in_usd_Nominal	7:employee_residence_Nominal	8:experience_level_Nominal	9:employment_type_Nominal	10:work_setting_Nominal	11:company_location_Nominal	12:company_size_Nominal
1	1.0	Data Dev..	Data Engineer..	EUR	-0.66050458175..	-0.632128735..	Germany	Mid-level	Full-time	Hybrid	Germany	L
2	1.0	Data Arc..	Data Architect..	USD	-0.21109174311..	-0.213793103..	United States	Senior	Full-time	In-person	United States	M
3	1.0	Data Arc..	Data Architect..	USD	-0.688990255880..	-0.6923753563..	United States	Senior	Full-time	In-person	United States	M
4	1.0	Data Scie..	Data Science a..	USD	-0.091743119260..	-0.094252873..	United States	Senior	Full-time	In-person	United States	M
5	1.0	Data Scie..	Data Science a..	USD	-0.6362365321100..	-0.64..	United States	Senior	Full-time	In-person	United States	M
6	1.0	Data Scie..	Data Science a..	USD	-0.467889802568..	-0.471254567..	United States	Senior	Full-time	Remote	United States	M
7	1.0	Data Scie..	Data Science a..	USD	-0.605504581559..	-0.60919502..	United States	Senior	Full-time	Remote	United States	M
8	1.0	Machine ..	Machine Learn..	USD	-0.0348623853211..	-0.037241379..	United States	Mid-level	Full-time	In-person	United States	M
9	1.0	Machine ..	Machine Learn..	USD	-0.42791651376..	-0.431694567..	United States	Mid-level	Full-time	In-person	United States	M
10	1.0	Data Engi..	Data Engineer..	USD	-0.10097431926..	-0.10348267..	United States	Executive	Full-time	Remote	United States	M
11	1.0	Data Engi..	Data Engineer..	USD	-0.2935779816513..	-0.296551724..	United States	Executive	Full-time	Remote	United States	M
12	1.0	Machine ..	Machine Learn..	USD	-0.0348623853211..	-0.037241379..	United States	Senior	Full-time	In-person	United States	M
13	1.0	Machine ..	Machine Learn..	USD	-0.42791651376..	-0.431694567..	United States	Senior	Full-time	In-person	United States	M
14	1.0	Data Scie..	Data Science a..	GBP	-0.903669724770..	-0.87097014..	United Kingdom	Mid-level	Full-time	In-person	United Kingdom	M
15	1.0	Data Scie..	Data Science a..	GBP	-0.92605504587156..	-0.899255172..	United Kingdom	Mid-level	Full-time	In-person	United Kingdom	M
16	1.0	Data Ana..	Data Analysis	USD	-0.6284403669724..	-0.632183908..	United States	Entry-level	Full-time	In-person	United States	M
17	1.0	Data Ana..	Data Analysis	USD	-0.720183468238..	-0.72413791..	United States	Entry-level	Full-time	In-person	United States	M
18	1.0	Data Scie..	Data Science a..	USD	-0.311926055045..	-0.310944827..	United States	Senior	Full-time	In-person	United States	M
19	1.0	Data Scie..	Data Science a..	USD	-0.0091743119266..	-0.00696551..	United States	Senior	Full-time	In-person	United States	M
20	1.0	Analytics ..	Leadership an..	USD	-0.4220183468238..	-0.425287356..	United States	Mid-level	Full-time	In-person	United States	M
21	1.0	Analytics ..	Leadership an..	USD	-0.517167467889..	-0.517241379..	United States	Mid-level	Full-time	In-person	United States	M
22	1.0	Machine ..	Machine Learn..	USD	-0.126146788990..	-0.128735632..	United States	Mid-level	Full-time	In-person	United States	M
23	1.0	Machine ..	Machine Learn..	USD	-0.411926055045..	-0.415172413..	United States	Mid-level	Full-time	In-person	United States	M
24	1.0	Data Ana..	Data Analysis	USD	-0.352210091743..	-0.356321839..	United States	Mid-level	Full-time	In-person	United States	M
25	1.0	Data Ana..	Data Analysis	USD	-0.5596330275229..	-0.563218390..	United States	Mid-level	Full-time	In-person	United States	M
26	1.0	Machine ..	Machine Learn..	USD	-0.1502568807339..	-0.156321839..	United States	Senior	Full-time	In-person	United States	M
27	1.0	Machine ..	Machine Learn..	USD	-0.366974770642..	-0.37014942..	United States	Senior	Full-time	In-person	United States	M
28	1.0	Applied S...	Data Science a..	USD	-0.044951284403..	-0.047356321..	United States	Mid-level	Full-time	In-person	United States	I

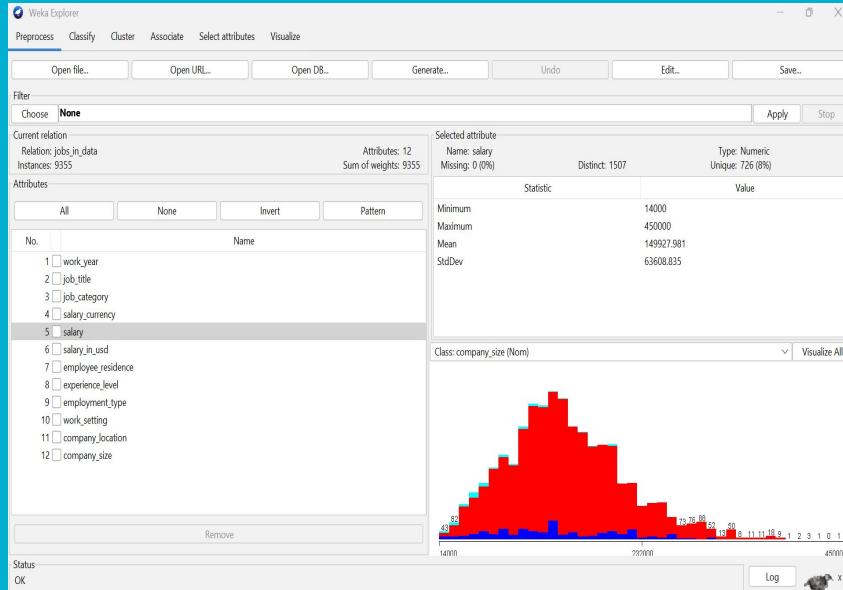
Add instance Undo OK Cancel

Normalization (min-max)

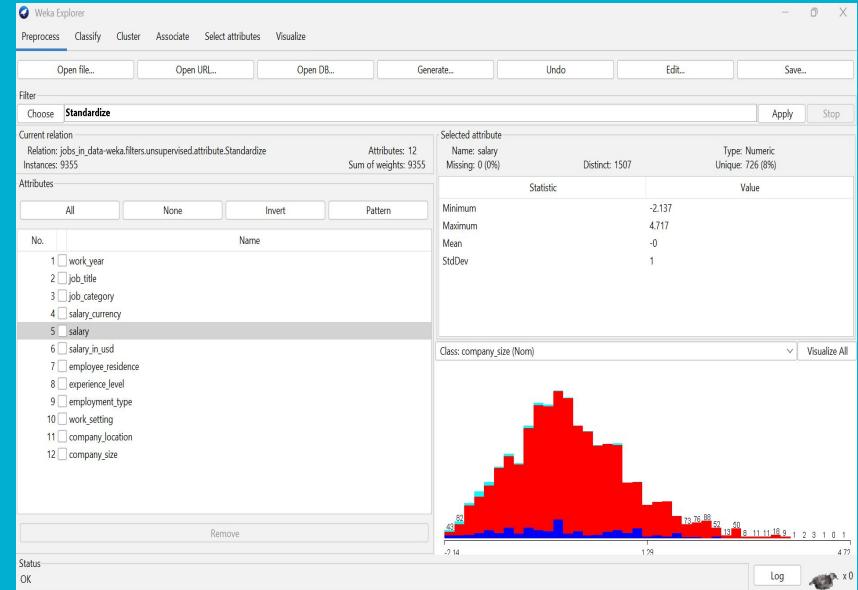
- The original dataset contains numerical features with varying scales, which may impact the performance of certain algorithms.
- Min-Max normalization is applied to ensure that numerical features are standardized between -1 and 1. This helps prevent attributes with larger scales from dominating the analysis, making the dataset more suitable for a wide range of algorithms.

Normalization (z-score)

Original

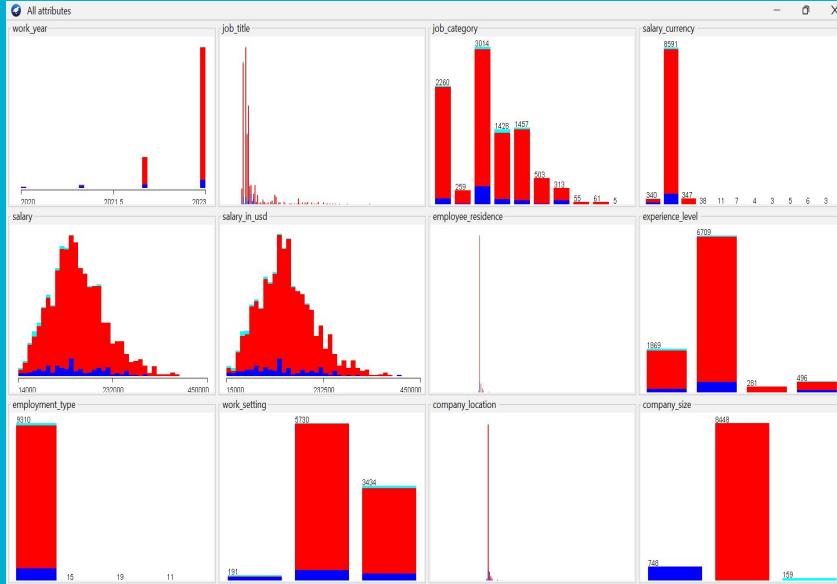


After applying 'Standardize'

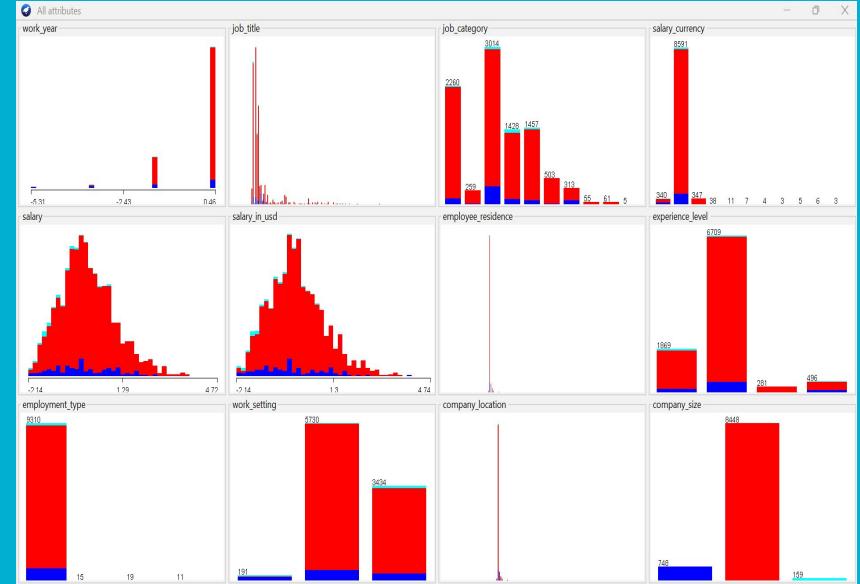


Normalization (z-score)

Original



After applying 'Standardize'



Normalization (z-score)

Viewer
Relation: jobs_in_data

No.	1: work_year	2: job_title	3: job_category	4: salary_currency	5: salary	6: salary_in_usd	7: employee_residence	8: experience_level	9: employment_type	10: work_setting	11: company_location	12: company_size
	Numeric	Nominal	Nominal	Nominal	Numeric	Numeric	Nominal	Nominal	Nominal	Nominal	Nominal	Nominal
1	2023.0	Job Dev..	Data Engineer..	EUR	88000.0	95012.0	Germany	Mid-level	Full-time	Hybrid	Germany	L
2	2023.0	Data Archit..	Data Architect..	USD	186000.0	186000.0	United States	Senior	Full-time	In-person	United States	M
3	2023.0	Data Archit..	Data Architect..	USD	81800.0	81800.0	United States	Senior	Full-time	In-person	United States	M
4	2023.0	Data Scie..	Data Science a..	USD	212000.0	212000.0	United States	Senior	Full-time	In-person	United States	M
5	2023.0	Data Scie..	Data Science a..	USD	93300.0	93300.0	United States	Senior	Full-time	In-person	United States	M
6	2023.0	Data Scie..	Data Science a..	USD	130000.0	130000.0	United States	Senior	Full-time	Remote	United States	M
7	2023.0	Data Scie..	Data Science a..	USD	100000.0	100000.0	United States	Senior	Full-time	Remote	United States	M
8	2023.0	Machine ..	Machine Learn..	USD	224400.0	224400.0	United States	Mid-level	Full-time	In-person	United States	M
9	2023.0	Machine ..	Machine Learn..	USD	138700.0	138700.0	United States	Mid-level	Full-time	In-person	United States	M
10	2023.0	Data Engin..	Data Engineer..	USD	210000.0	210000.0	United States	Executive	Full-time	Remote	United States	M
11	2023.0	Data Engin..	Data Engineer..	USD	168000.0	168000.0	United States	Executive	Full-time	Remote	United States	M
12	2023.0	Machine ..	Machine Learn..	USD	224400.0	224400.0	United States	Senior	Full-time	In-person	United States	M
13	2023.0	Machine ..	Machine Learn..	USD	138700.0	138700.0	United States	Senior	Full-time	In-person	United States	M
14	2023.0	Data Scie..	Data Science a..	GBP	35000.0	43064.0	United Kingdom	Mid-level	Full-time	In-person	United Kingdom	M
15	2023.0	Data Scie..	Data Science a..	GBP	30000.0	36912.0	United Kingdom	Mid-level	Full-time	In-person	United Kingdom	M
16	2023.0	Data Anal..	Data Analysis	USD	95000.0	95000.0	United States	Entry-level	Full-time	In-person	United States	M
17	2023.0	Data Anal..	Data Analysis	USD	75000.0	75000.0	United States	Entry-level	Full-time	In-person	United States	M
18	2023.0	Data Scie..	Data Science a..	USD	300000.0	300000.0	United States	Senior	Full-time	In-person	United States	M
19	2023.0	Data Scie..	Data Science a..	USD	234000.0	234000.0	United States	Senior	Full-time	In-person	United States	M
20	2023.0	Analytics ..	Leadership an..	USD	140000.0	140000.0	United States	Mid-level	Full-time	In-person	United States	M
21	2023.0	Analytics ..	Leadership an..	USD	120000.0	120000.0	United States	Mid-level	Full-time	In-person	United States	M
22	2023.0	Machine ..	Machine Learn..	USD	204500.0	204500.0	United States	Mid-level	Full-time	In-person	United States	M
23	2023.0	Machine ..	Machine Learn..	USD	142200.0	142200.0	United States	Mid-level	Full-time	In-person	United States	M
24	2023.0	Data Anal..	Data Analysis	USD	155000.0	155000.0	United States	Mid-level	Full-time	In-person	United States	M
25	2023.0	Data Anal..	Data Analysis	USD	110000.0	110000.0	United States	Mid-level	Full-time	In-person	United States	M
26	2023.0	Machine ..	Machine Learn..	USD	266500.0	266500.0	United States	Senior	Full-time	In-person	United States	M
27	2023.0	Machine ..	Machine Learn..	USD	152000.0	152000.0	United States	Senior	Full-time	In-person	United States	M
28	2023.0	Applied S...	Data Science a..	USD	222700.0	222700.0	United States	Mid-level	Full-time	In-person	United States	I

Add instance Undo OK Cancel

Viewer
Relation: jobs_in_data-weka.filters.unsupervised.attribute.Standardize

No.	1: work_year	2: job_title	3: job_category	4: salary_currency	5: salary	6: salary_in_usd	7: employee_residence	8: experience_level	9: employment_type	10: work_setting	11: company_location	12: company_size
	Numeric	Nominal	Nominal	Nominal	Numeric	Numeric	Nominal	Nominal	Nominal	Nominal	Nominal	Nominal
1	0.4611452..	Data Dev..	Data Engineer..	EUR	-0.9735715147482..	-0.875115469..	Germany	Mid-level	Full-time	Hybrid	Germany	L
2	0.4611452..	Data Archit..	Data Architect..	USD	0.567913245057..	0.5655803720..	United States	Senior	Full-time	In-person	United States	M
3	0.4611452..	Data Archit..	Data Architect..	USD	-1.071045883462..	-1.08420976..	United States	Senior	Full-time	In-person	United States	M
4	0.4611452..	Data Scie..	Data Science a..	USD	0.97858971065..	0.976623472..	United States	Senior	Full-time	In-person	United States	M
5	0.4611452..	Data Scie..	Data Science a..	USD	-0.89203389306..	-0.90221379..	United States	Senior	Full-time	In-person	United States	M
6	0.4611452..	Data Scie..	Data Science a..	USD	-0.313289516654..	-0.321309590..	United States	Senior	Full-time	Remote	United States	M
7	0.4611452..	Data Scie..	Data Science a..	USD	-0.78492210103..	-0.796163150..	United States	Senior	Full-time	Remote	United States	M
8	0.4611452..	Machine ..	Machine Learn..	USD	1.1707810440655..	1.172896276..	United States	Mid-level	Full-time	In-person	United States	M
9	0.4611452..	Machine ..	Machine Learn..	USD	-0.176516064553..	-0.183602058..	United States	Mid-level	Full-time	In-person	United States	M
10	0.4611452..	Data Engin..	Data Engineer..	USD	0.94493739921007..	0.9449665568..	United States	Executive	Full-time	Remote	United States	M
11	0.4611452..	Data Engin..	Data Engineer..	USD	0.2841117683183..	0.280171584..	United States	Executive	Full-time	Remote	United States	M
12	0.4611452..	Machine ..	Machine Learn..	USD	1.1707810440655..	1.172896276..	United States	Senior	Full-time	In-person	United States	M
13	0.4611452..	Machine ..	Machine Learn..	USD	-0.176516064553..	-0.183602058..	United States	Senior	Full-time	In-person	United States	M
14	0.4611452..	Data Scie..	Data Science a..	GBP	-1.80679273942..	-1.697317993..	United Kingdom	Mid-level	Full-time	In-person	United Kingdom	M
15	0.4611452..	Data Scie..	Data Science a..	GBP	-1.885380161484..	-1.794748529..	United Kingdom	Mid-level	Full-time	In-person	United Kingdom	M
16	0.4611452..	Data Anal..	Data Analysis	USD	-0.865352754234..	-0.87530510..	United States	Entry-level	Full-time	In-person	United States	M
17	0.4611452..	Data Anal..	Data Analysis	USD	-1.177949271310..	-1.191674450..	United States	Entry-level	Full-time	In-person	United States	M
18	0.4611452..	Data Scie..	Data Science a..	USD	2.35925179557125..	2.369527247..	United States	Senior	Full-time	In-person	United States	M
19	0.4611452..	Data Scie..	Data Science a..	USD	1.32170347389924..	1.324894916..	United States	Senior	Full-time	In-person	United States	M
20	0.4611452..	Analytics ..	Leadership an..	USD	-0.156078652171..	-0.163025071..	United States	Mid-level	Full-time	In-person	United States	M
21	0.4611452..	Analytics ..	Leadership an..	USD	-0.47050309137..	-0.479594110..	United States	Mid-level	Full-time	In-person	United States	M
22	0.4611452..	Machine ..	Machine Learn..	USD	0.857913427443..	0.857910082..	United States	Mid-level	Full-time	In-person	United States	M
23	0.4611452..	Machine ..	Machine Learn..	USD	-0.12492261984..	-0.128202476..	United States	Mid-level	Full-time	In-person	United States	M
24	0.4611452..	Data Anal..	Data Analysis	USD	0.079376445534..	0.074401708..	United States	Mid-level	Full-time	In-person	United States	M
25	0.4611452..	Data Anal..	Data Analysis	USD	-0.627711245620..	-0.637878630..	United States	Mid-level	Full-time	In-person	United States	M
26	0.4611452..	Machine ..	Machine Learn..	USD	1.8326387835390..	1.839747105..	United States	Senior	Full-time	In-person	United States	M
27	0.4611452..	Machine ..	Machine Learn..	USD	0.0325743852085..	0.026916352..	United States	Senior	Full-time	In-person	United States	M
28	0.4611452..	Applied S...	Data Science a..	USD	1.136196538792..	1.138073682..	United States	Mid-level	Full-time	In-person	United States	L

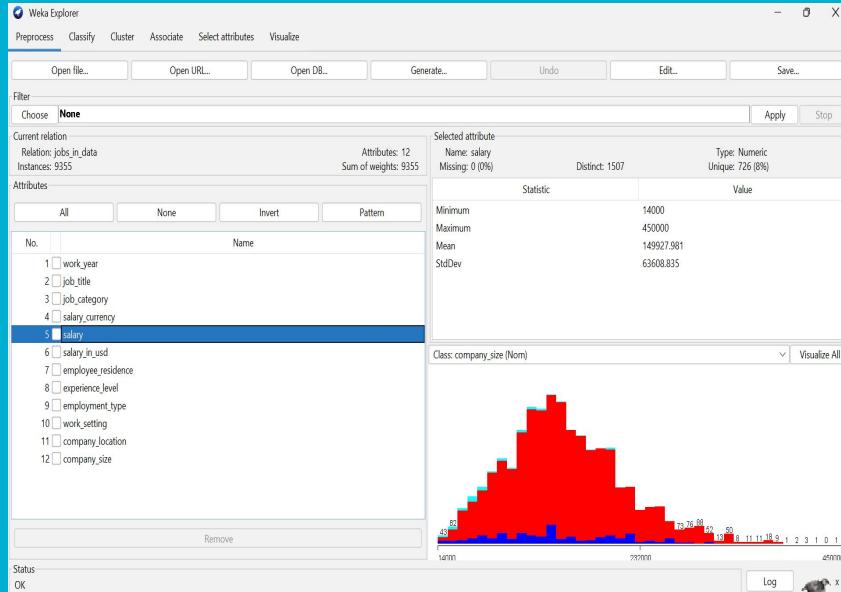
Add instance Undo OK Cancel

Normalization (z-score)

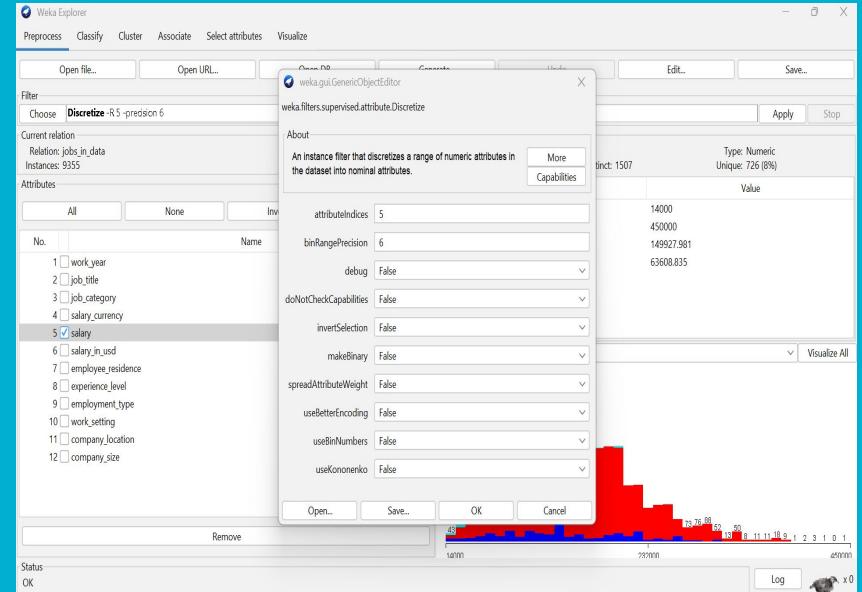
- Z-Score normalization (Standardize) is applied to ensure that numerical features have a mean of 0 and a standard deviation of 1. This helps in centering and scaling the data, making it suitable for algorithms that assume normally distributed features.

Discretization

Original

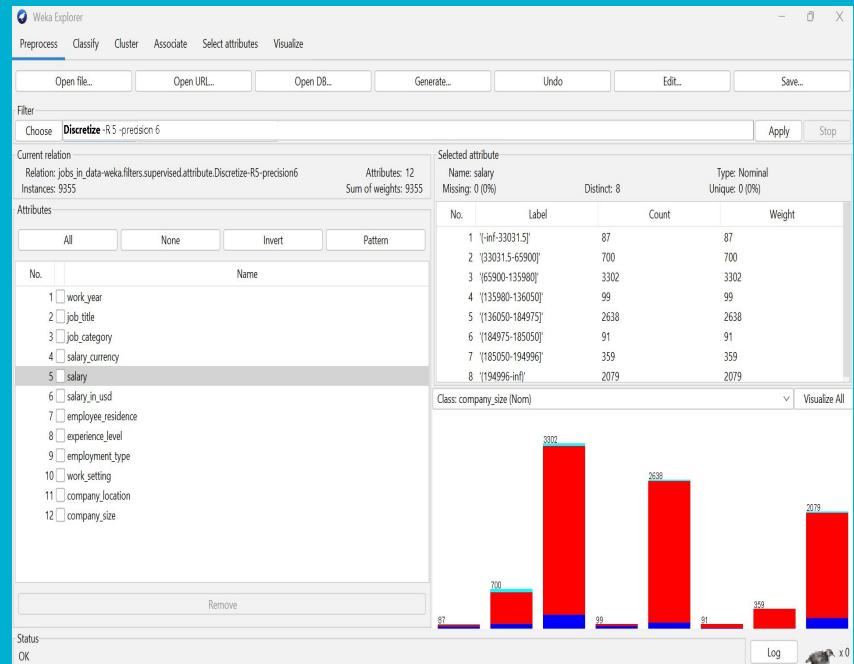


Applying 'Discretize'



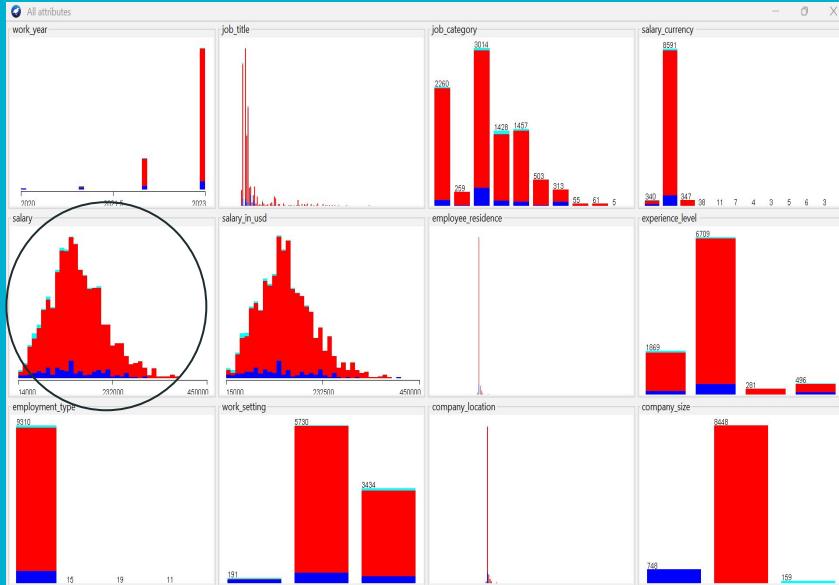
Discretization

- Implementation of supervised discretization techniques to transform continuous variables into discrete bins.
- Conversion of numerical attributes into nominal categories to simplify analysis and enhance interpretability.

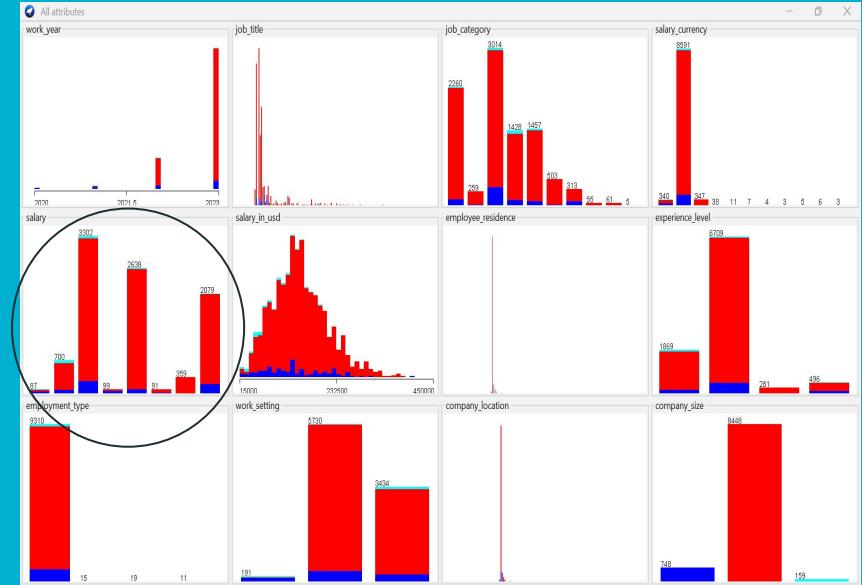


Discretization

Original



After 'Discretize'

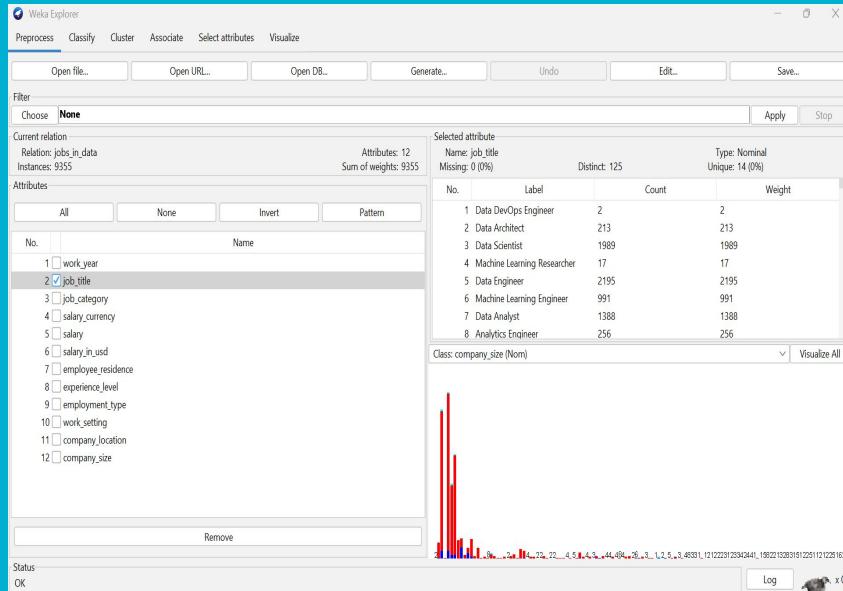


Discretization

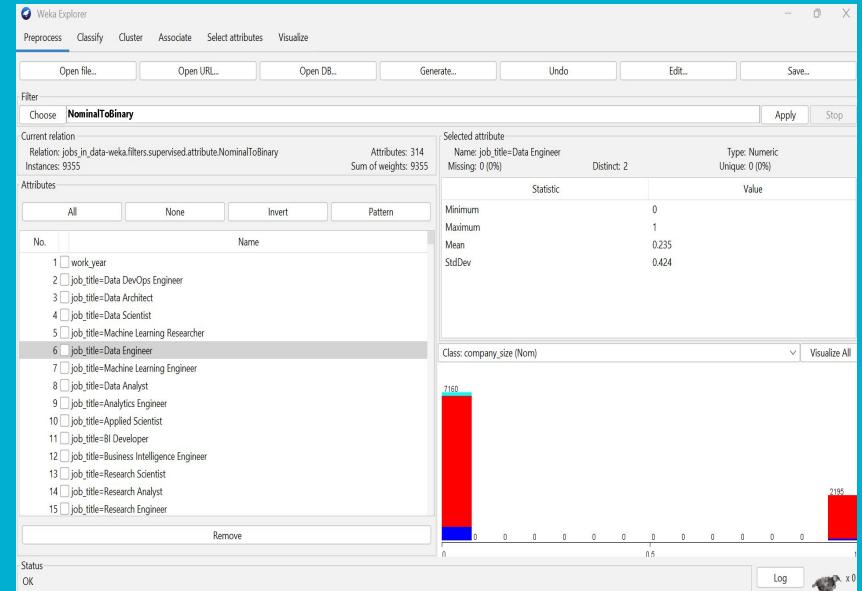
- Continuous variables, such as salary, may benefit from discretization to simplify the dataset and identify patterns.
- Supervised discretization methods allow us to use class information to guide the binning process, leading to more meaningful and interpretable results.
- Transformation of numerical attributes into nominal categories enhances the dataset's interpretability and facilitates analysis, particularly when exploring relationships between different attributes.

Change Nominal to Binary

Original



Applied Nominal to Binary

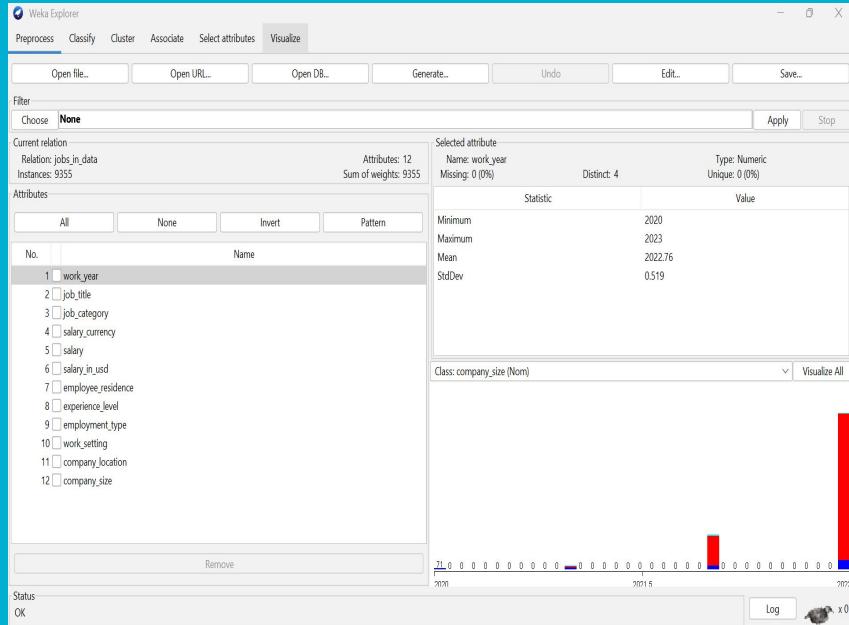


Change Nominal to Binary

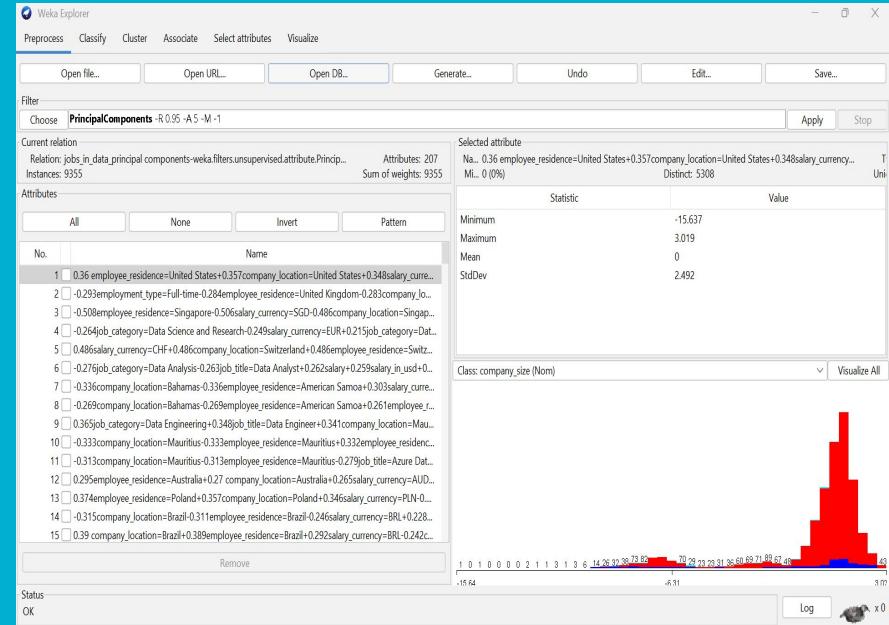
- Nominal attributes, such as job title, job category, often require transformation for better algorithm compatibility.
- The supervised approach ensures that the transformation is guided by class information, preserving the underlying relationships within the data.

Data Reduction-PCA

Original



Application of PCA



Data Reduction-PCA

- High-dimensional datasets can pose challenges for computational efficiency and interpretability.
- Principal Component Analysis (PCA) is applied to identify the principal components, effectively reducing the dataset's dimensionality. This facilitates faster computation and aids in visualizing essential patterns and relationships within the data.

Summary of Operations

- The preprocessing operations collectively enhance dataset robustness by addressing missing values, standardize numerical features to prevent scale dominance, and simplify the dataset through supervised discretization. Transformation of nominal attributes into binary format improves algorithm compatibility, while PCA reduced dimensionality for improved efficiency and interpretability.
- Together, these preprocessing steps prepare the dataset for analysis, ensuring it is more robust, standardized, and suitable for a wide range of analytical techniques.

Thank You