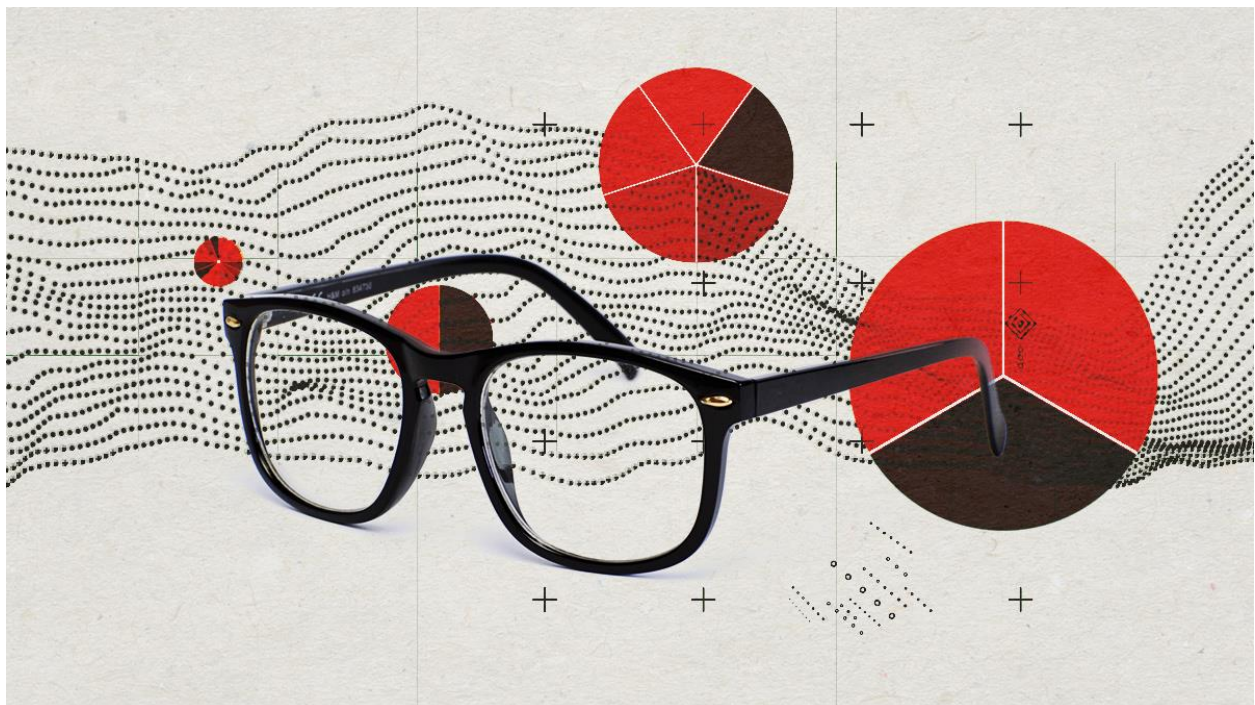


# Global Data Professional Salary

## Summary

This dataset contains survey responses from 882 data professionals from 46 countries who took part in the 2019 Global Data Professional Salary Survey.



*1HBR Staff/StudioM1/Moritz Otto/Getty Images*

## Data Source

The data is collected via surveys & provided by Brent Ozar in this [link](#).

## Data Contents

The dataset includes professional & personal details such as: location, salary, primary database, years of experience, education, hours of work per week, etc. from database administrators, data analysts, data architects, developers and data scientists in 2017-2019.

## Data Limitations

- **Missing information on benefits and perks:** The dataset does not include details about non-monetary benefits or perks such as health insurance, vacation days, retirement plans, stock options, and others, which are important elements of a data professional's total compensation.
- **Absence of company-specific data:** The dataset does not provide sufficient information about the companies the data professionals work for, such as their industry or location. This information could provide valuable context for understanding salary variations.
- **Missing answers:** The surveyors gradually asked more questions over time, so if a question wasn't asked in a year, the answers are populated with *Not Asked*.
- **Unreliable information:** The postal code field was optional and may be unreliable.

## Data Wrangling

Original column	Change	New column	Reason
<i>None</i>	<b>Created</b>	<b>id</b>	<b>Column contains unique identifiers</b>
Survey Year	Renamed	survey_year	Conformity to Python naming convention
Timestamp	Renamed	timestamp	Conformity to naming convention
SalaryUSD	Renamed	salary_in_usd	Conformity to naming convention
Country	Renamed	country	Conformity to naming convention
<b>PostalCode</b>	<b>Removed</b>	<i>None</i>	<b>Deemed unreliable by the author</b>
PrimaryDatabase	Renamed	primary_db	Conformity to naming convention
YearsWithThisDatabase	Renamed	years_of_exp_with_primary_db	Conformity to naming convention; Clarity
OtherDatabases	Renamed	other_dbs	Conformity to naming convention
EmploymentStatus	Renamed	employment_status	Conformity to naming convention
JobTitle	Renamed	job_title	Conformity to naming convention

ManageStaff	Renamed	manage_staff	Conformity to naming convention
YearsWithThisTypeOfJob	Renamed	years_of_exp_with_data_jobs	Conformity to naming convention; Clarity
HowManyCompanies	Renamed	number_of_companies_worked_for	Conformity to naming convention; Clarity
OtherPeopleOnYourTeam	Renamed	number_of_team_members	Conformity to naming convention; Clarity
CompanyEmployeesOverall	Renamed	number_of_company_employees	Conformity to naming convention; Clarity
DatabaseServers	Renamed	number_of_db_servers	Conformity to naming convention; Clarity
Education	Renamed	education	Conformity to naming convention
EducationIsComputerRelated	Renamed	education_is_computer_related	Conformity to naming convention
Certifications	Renamed	certifications	Conformity to naming convention
HoursWorkedPerWeek	Renamed	hours_worked_per_week	Conformity to naming convention
TelecommuteDaysPerWeek	Renamed	wfh_days_per_week	Conformity to naming convention; Clarity

PopulationOfLargestCityWithin20Miles	Renamed	pop_of_largest_city_within_20_miles	Conformity to naming convention
EmploymentSector	Renamed	employment_sector	Conformity to naming convention
LookingForAnotherJob	Renamed	looking_for_another_job	Conformity to naming convention
CareerPlansThisYear	Renamed	career_plans_this_year	Conformity to naming convention
Gender	Renamed	gender	Conformity to naming convention
OtherJobDuties	Renamed	other_job_duties	Conformity to naming convention
KindsOfTasksPerformed	Renamed	tasks_performed	Conformity to naming convention
<b>Counter</b>	<b>Removed</b>	<b>None</b>	<b>Column is redundant</b>

## Data Consistency Checks

Column	Consistency issue	Possible reasons	Solution
id	<i>None</i>	<i>None</i>	<i>None</i>
survey_year	<i>None</i>	<i>None</i>	<i>None</i>
Timestamp	<i>None</i>	<i>None</i>	<i>None</i>

salary_in_usd	<b>Non-numeric values</b> <b>Special characters</b> <b>Abnormally high &amp; low values</b>	<b>Incorrect input</b> <b>Hourly rates</b> <b>Part-time salaries</b> <b>Low salaries</b>	<b>Values below 1000 or above 1000000 are removed</b>
country	<i>None</i>	<i>None</i>	<i>None</i>
primary_db	<i>None</i>	<i>None</i>	<i>None</i>
years_of_exp_with_primary_db	<b>Abnormally high values</b>	<b>Incorrect input</b>	<b>Outliers are removed</b>
other_dbs	<b>Missing data</b>	<b>Respondents not sharing info</b>	<b>Missing values are replaced with “Not Provided”</b>
employment_status	<i>None</i>	<i>None</i>	<i>None</i>
job_title	<i>None</i>	<i>None</i>	<i>None</i>
manage_staff	<i>None</i>	<i>None</i>	<i>None</i>
years_of_exp_with_data_jobs	<i>None</i>	<i>None</i>	<i>None</i>
number_of_companies_worked_for	<i>None</i>	<i>None</i>	<i>None</i>
number_of_team_members	<i>None</i>	<i>None</i>	<i>None</i>

number_of_company_employees	<i>None</i>	<i>None</i>	<i>None</i>
number_of_db_servers	<i>None</i>	<i>None</i>	<i>None</i>
education	<i>None</i>	<i>None</i>	<i>None</i>
<b>education_is_computer_related</b>	<b>Missing data</b>	<b>Respondents not sharing info</b>	<b>Replaced with “Not Provided”</b>
certifications	<i>None</i>	<i>None</i>	<i>None</i>
hours_worked_per_week	Abnormally high values	Incorrect input	Outliers are removed
wfh_days_per_week	<i>None</i>	<i>None</i>	<i>None</i>
pop_of_largest_city_within_20_miles	<i>None</i>	<i>None</i>	<i>None</i>
employment_sector	<i>None</i>	<i>None</i>	<i>None</i>
looking_for_another_job	<i>None</i>	<i>None</i>	<i>None</i>
career_plans_this_year	<i>None</i>	<i>None</i>	<i>None</i>
gender	<i>None</i>	<i>None</i>	<i>None</i>
<b>other_job_duties</b>	<b>Missing data</b>	<b>Respondents not sharing info</b>	<b>Replaced with “Not Provided”</b>
<b>tasks_performed</b>	<b>Missing data</b>	<b>Respondents not sharing info</b>	<b>Replaced with “Not Provided”</b>

## Column Details

Column	Column Description
id	Unique identifier
survey_year	The year the survey was conducted
Timestamp	Timestamp when the survey answer was submitted
salary_in_usd	Salary in USD
country	The country where the respondent lives
primary_db	The primary database that the respondent works with
years_of_exp_with_primary_db	The number of years of experience with the primary database
other_dbs	Other databases which the respondent also works with
employment_status	Employment status
job_title	Job title
manage_staff	Whether the respondent is working as a manager
years_of_exp_with_data_jobs	The number of years of experience with data-related jobs
number_of_companies_worked_for	The number of companies the respondent has worked for



number_of_team_members	The number of people working in the same team
number_of_company_employees	The number of people working in the same company
number_of_db_servers	The number of database servers that the respondents works with
education	Education level
education_is_computer_related	Whether the education is related to computer
certifications	Data-related certifications
hours_worked_per_week	Number of work hours per week
wfh_days_per_week	Number of work-from-home days per week
pop_of_largest_city_within_20_miles	The population of the largest city within 20 miles of the respondent's living location
employment_sector	The sector the respondent works in
looking_for_another_job	Whether the respondent is looking for as job
career_plans_this_year	Any career plans for the current year
other_job_duties	Any other job duties
tasks_perfomed	Any tasks performed at work

## Summary Statistics

Variables	time -variant/ invariant	structured/unstru ctured	qualitative/ quanti tative	qualitative	quantitative
				nominal/ordinal/b inary	discrete/continuo us
id	time-invariant	structured	qualitative	nominal	
survey_year	time-invariant	structured	qualitative	ordinal	
salary_in_usd	time-variant	structured	quantitative		continuous
country	time-invariant	structured	qualitative	nominal	
primary_db	time-invariant	structured	qualitative	nominal	
years_of_exp_with_primary_db	time-invariant	structured	quantitative		continuous
other_dbs	time-invariant	unstructured	qualitative	nominal	
employment_status	time-invariant	structured	qualitative	nominal	
job_title	time-invariant	structured	qualitative	nominal	
manage_staff	time-invariant	structured	qualitative	binary	
years_of_exp_with_data_jobs	time-invariant	structured	quantitative		continuous

number_of_companier_worked_for	time-invariant	structured	qualitative	ordinal	
number_of_team_members	time-invariant	structured	qualitative	ordinal	
number_of_company_employees	time-invariant	structured	qualitative	ordinal	
number_of_db_servers	time-invariant	structured	quantitative		continuous
education	time-invariant	structured	qualitative	ordinal	
education_is_computer_related	time-invariant	structured	qualitative	nominal	
certifications	time-invariant	structured	qualitative	nominal	
hours_worked_per_week	time-invariant	structured	quantitative		continuous
wfh_days_per_week	time-invariant	structured	qualitative	ordinal	
pop_of_largest_city_within_20_miles	time-invariant	structured	qualitative	ordinal	
employment_sector	time-invariant	structured	qualitative	nominal	
looking_for_another_job	time-invariant	structured	qualitative	nominal	
career_plans_this_year	time-invariant	structured	qualitative	nominal	
other_job_duties	time-invariant	structured	qualitative	nominal	
tasks_perfomed	time-invariant	structured	qualitative	nominal	

## Research Questions

- What is the median salary level by profession?
- What is the average income of data professionals in each country?
- How does experience level affect income?
- Are there salary gaps amongst data professionals between countries?