

Brief Article

The Author

March 5, 2024

Differential privacy (DP) is a privacy-preserving concept in data analysis and statistics. The goal of differential privacy is to allow the inclusion of an individual's data in a dataset for analysis while protecting the privacy of that individual.

So, how does DP achieve this? Let's take a simple example:

Suppose you are part of a company that wants to calculate the average salary of its employees without revealing the salary of any specific employee. The company collects the following data:

Employee A: \$60,000

Employee B: \$70,000

Employee C: \$80,000

Employee D: \$90,000

Non-Differentially Private Calculation:

In a traditional average calculation, you might sum all salaries and divide by the number of employees: $\text{Average} = \frac{60,000 + 70,000 + 80,000 + 90,000}{4} = 75,000$

Differentially Private Calculation:

In a differentially private approach, you add random noise to the calculation to protect individual privacy. Let's say you add random noise between -1,000 and +1,000 to each employee's salary:

$$\text{Average} = \frac{60,000 + 70,000 + 80,000 + 90,000 + \text{noise}}{4}$$

Here, the noise is a random value chosen from the range -1,000 to +1,000. This ensures that even if an individual's salary changes slightly, it is challenging to determine which specific employee's salary contributed to the final result.

In essence, differential privacy aims to introduce uncertainty in the output distribution. Considering two datasets that differ only in one data record, a differentially private algorithm will generate output distributions that are slightly different for these two datasets. In other words, observers cannot discern which dataset the algorithm used from the output, and thus, they cannot determine whether a specific data record is in the dataset or not.