

Synthetic & De-identified Data in Healthcare Analytics



The rapid digitisation of healthcare and the integration of electronic health records (EHRs) have brought data to the forefront of healthcare analytics. Organisations striving for value-based care increasingly rely on data to make informed decisions, enhance patient outcomes and drive research. However, choosing the right type of data for analytics initiatives is crucial. Healthcare professionals primarily use three types of data—real-world, synthetic and de-identified—each offering unique benefits and challenges. Understanding when to use each data type can significantly impact the success of healthcare analytics projects.

The Importance of Real-world Data

Real-world data (RWD) refers to information collected from various sources reflecting patients' actual health status. EHRs, claims data, medical device registries, patient-reported outcomes and digital health devices all contribute to this dataset. RWD plays a crucial role in generating real-world evidence (RWE), which is instrumental in regulating and developing medical interventions. The evidence derived from RWD informs clinical trials and therapeutic advancements, particularly in areas such as cancer care and precision medicine.

However, despite its potential, RWD poses challenges. Data quality, availability and relevance to specific projects often undermine its practical use. Additionally, as healthcare organisations increasingly employ emerging technologies like artificial intelligence (Al) to process RWD, they must remain vigilant about data integrity and suitability for particular research goals. While RWD provides a wealth of insight, understanding when it is appropriate and how to manage its inherent limitations is crucial for stakeholders.

Advantages and Limitations of Synthetic Data

In contrast to RWD, synthetic data is artificially generated and designed to reflect the characteristics of real-world datasets without containing identifiable information. Synthetic data offers a compelling alternative where privacy and data harmonisation are critical. By simulating real-world scenarios, synthetic data allows researchers to train algorithms, develop applications and conduct clinical research while minimising privacy risks.

Despite its advantages, synthetic data has its drawbacks. For instance, the artificial nature of synthetic datasets can introduce biases or errors, compromising the quality of the analysis. Additionally, generating synthetic patient populations accurately can be challenging, limiting the dataset's applicability in large-scale studies. Issues such as data leakage, where information from a training set inadvertently influences the test set, can undermine AI model performance and reliability. Healthcare stakeholders must carefully assess these risks to determine whether synthetic data aligns with their analytics objectives.

De-identified Data and Privacy Concerns

As the name suggests, de-identified data involves masking or removing personal identifiers to ensure confidentiality while maintaining the dataset's utility. This data type is essential for adhering to the Health Insurance Portability and Accountability Act (HIPAA) regulations, allowing organisations to share information without compromising patient privacy. Researchers often use de-identified data to analyse demographic trends, evaluate healthcare disparities and improve patient care.

However, de-identification is not a foolproof solution. As Al and machine learning tools become more sophisticated, the risk of re-identification has increased. Even with direct identifiers removed, datasets can still be re-linked to individuals through other indirect variables, such as geographic data or treatment timelines. These challenges are prompting discussions about modernising HIPAA regulations to address the emerging privacy risks associated with advanced technologies. Healthcare organisations must adopt robust de-identification protocols that

extend beyond current regulations to safeguard patient data effectively.

Real-world, synthetic and de-identified data all serve distinct purposes in healthcare analytics. Real-world data offers unparalleled insight into patient outcomes and supports evidence-based decision-making, but it comes with concerns about quality and relevance. Synthetic data presents a privacy-friendly alternative but requires careful handling to prevent biases and errors. Meanwhile, de-identified data balances utility and confidentiality but remains vulnerable to re-identification risks. By recognising the advantages and limitations of each data type, healthcare organisations can strategically select the most suitable datasets for their initiatives. This understanding not only enhances research and patient care but also paves the way for safer and more effective use of emerging technologies in the healthcare sector.

Source: <u>TechTarget</u> Image Credit: <u>iStock</u>

Published on: Mon, 28 Oct 2024