# **BIG DATA BASICS**



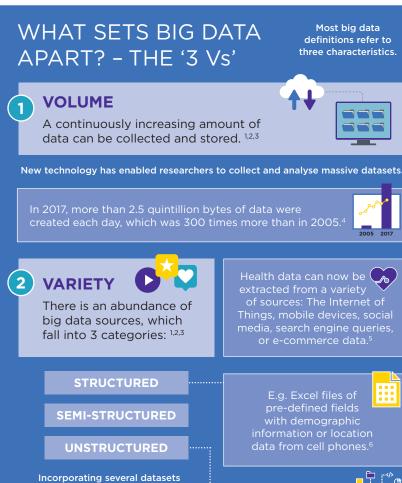
National Collaborating Centre for Infectious Diseases Centre de collaboration nationale

In public health, the use of big data for planning is still emerging and questions are common. Here, we begin with basic questions: what is 'big data', what sets these data apart from other data, and how are big data being used?

## WHAT IS 'BIG DATA'?

Big data are not new but continue to change with advancements in technology. This has led to many different definitions and a lack of conceptual clarity. One common definition for big data is this:

Very large and diverse datasets that are analysed computationally at high velocity to reveal patterns, trends and associations.<sup>1,2</sup>



within an analysis allows researchers E.g. Photos, videos, audio files, social to enrich their findings and identify new trends and associations.<sup>2</sup> media content, websites, atellite imagery, open-ended survey responses.<sup>6</sup>

Cell phones and wearable devices provide high velocity, continuously produced data collected in real time.<sup>7</sup>

Less than 20% of all data are structured and unstructured data make up an increasing proportion.6

VELOCITY The speed at which data are generated is high.<sup>1,2,3</sup>

Data are generated at vastly different speeds, depending on the source.

## HOW ARE BIG DATA USED?

### **HIGH VOLUME DATA**

A recent study that analysed the content of 2.8 million tweets related

to the COVID-19 pandemic helps public health to understand the public's information needs and respond faster and

and mobility of people, and



## more appropriately.8

**MULTIPLE DATA SETS** Studies that bring together data on commercial air travel, location

travel restrictions help predict the spread of COVID-19 and assess the effects of travel restrictions compared to other public health measures.9,10

A study that analysed cell phone data helps public health unravel the relationship between peoples' movements and socioeconomic factors, and to understand how public health restrictions may aff<u>ect different</u> populations differently and at different times.<sup>11</sup>



REFERENCES

3

- Mooney SJ, Pejaver V. Big data in public health: terminology, machine learning, and privacy. Annu. Rev. Public Health. 2018 Apr 1;39:95-112.
- Fuller D, Buote R, Stanley K. A glossary for big data in population and public health: discussion and commentary on terminology and research methods. J Epidemiol Community Health. 2017 Nov 1;71(11):1113-7.
- Ylijoki O, Porras J. Perspectives to definition of big data: a mapping study and discussion. J. Innov. Manag. 2016 May 4;4(1):69-91.
- 4. Herschel R, Miori VM. Ethics & big data. Technology in Society. 2017;49:31-36.
- Jia Q, Guo Y, Wang G & Barnes SJ. Big data analytics in the fight against major public health incidents (including COVID-19): a conceptual framework. Int. J. Environ. Res. Public Health. 2020;17(17): 6161.
- Marr B. What's the difference between structured, semi-structured and unstructured data? Forbes. 2019 Oct 18.
- 7. Badr HS, Du H, Marshall M, Dong E, Squire MM, & Gardner LM. Association between mobility patterns and COVID-19 transmission in the USA: a mathematical modelling study. The Lancet Infect. Dis., 2020; 20(11):1247-1254.
- 8. Abd-Alrazaq A, Alhuwail D, Househ M, Hamdi M, Shah Z. Top concerns of tweeters during the COVID-19 pandemic: infoveillance study. J.Med. Internet Res. 2020;22(4):e19016.
- 9. Watts A, Au NH, Thomas-Bachli A, Forsyth J, Mayah O, Popescu S. & Bogoch II. Potential for inter-state spread of Covid-19 from Arizona, USA: analysis of mobile device location and commercial flight data. J. Travel Med. 2020, 1–3.
- 10. Chinazzi M, Davis JT, Ajelli M, et al. The effect of travel restrictions on the spread of the 2019 novel coronavirus (COVID-19) outbreak, Science, 2020;368(6489);395-400,
- Long JA, Ren C. Associations between mobility and socio-economic indicators vary across the timeline of the Covid-19 pandemic. Computers, environment and urban systems. 2022 (91)1.

