# Data Confidentiality Principles and Methods Report

October 2018

It is important to understand and apply confidentiality principles, rules, and methods to make sure that you:

- Do not release data that could identify people, households, or organisations unintentionally.

- Protect data provided by people and organisations, and ensure it isn't disclosed to anyone who is not authorised to access it.

- Use statistical methods to prevent data from being disclosed in a way that could identify a person, household, or organisation unintentionally.

Using statistical methods correctly protects the confidentiality of data. Methods such as perturbation, aggregation, suppression, limiting access, and building synthetic or confidential unit record files keep data confidential. When data is confidential, no individuals, households, or businesses can be identified, and no unauthorised people can access the data.

## Why do we have to protect data confidentiality?

Different organisations have different requirements relating to when they must or wish to protect the, privacy, security, and confidentiality of data so that people, households, and organisations can't be identified without their permission. This includes where we must or wish to protect the confidentiality of data throughout its life cycle — whenever we collect, use, store, and distribute it.

### What privacy, security, and confidentiality mean

The terms privacy, security, and confidentiality are often used interchangeably, but each term has a different meaning:

- **Privacy** refers to a person's ability to control the availability of data about themselves.

- **Security** refers to how an organisation stores and controls access to the data it holds.

- **Confidentiality** refers to the protection of data from, and about, individuals and organisations; and how we ensure that data is not made available or disclosed without authorisation.

## Degrees of identification in data

What do statisticians, data scientists and data analysts mean when they talk about confidentiality? How does identifiable data differ from de-identified or confidentialised information? Data identifiability is not binary. Data lies on a spectrum with multiple shades of identifiability. This is a primer on how to distinguish different categories of data in the NZ context.

### Identifiable
Data that directly or indirectly identifies an individual or business.

**Individual**

| Name | Hēni |
| Gender | Female |
| DOB | 31/01/1985 |
| Address | 28 My Road Postcode 6012 Wellington |

**Business**

| Name | Puzzles |
| Type | Paper Stationery Manufacturing |
| Employees | 34 |
| Expenditure | $398,000 |

Data that identifies a person without additional information or by linking to information in the public domain. Where an individual can be identified through connecting up information.

Personal, identifiable data like this are protected, and should only be released to the public providing we have explicit permission to do so.

*For example: Name, Date of birth, Gender.*

### De-identified
Data which has had information removed from it to reduce risk of spontaneous recognition.

**Individual**

| Name | *Unknown* |
| Gender | Female |
| DOB | 1985 |
| Address | Postcode 6012 Wellington |

**Business**

| Name | *Unknown* |
| Type | Manufacturing |
| Employees | 30 - 40 |
| Expenditure | $398,000 |

**De-identified:** Data which has had information removed from it to reduce risk of spontaneous recognition (likelihood of identifying a person, place or organisation without any effort).
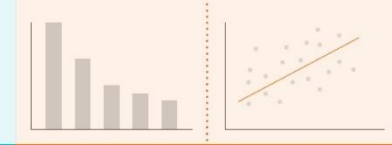
*For example: Data held within Stats NZ's Integrated Data Infrastructure and Longitudinal Business Database is de-identified before approved researchers can access in a secure data lab environment.*

**Partially confidentialised:** Data which has been modified to protect the confidentiality of respondents while also maintaining the integrity of data. Modification involves applying methods such as top-coding, data swapping, and collapsing categorical variables to the unit records.

### Confidentialised
Data which has had statistical methods applied to it to protect against disclosing unauthorised information.

**Individual**

| Name | *Unknown* |
| Gender | Female |
| Age | 30 - 40 years |
| Address | Wellington |

**Business**

| Name | *Unknown* |
| Type | Manufacturing |
| Employees | 10 - 100 |
| Expenditure | Under $500,000 |

Statistical methods include suppression, aggregation, perturbation, data swapping, top and bottom coding, etc. These prevent the unauthorised identification of individuals, households, or organisations. This data is publicly available.

*For example: Stats NZ nz.stat datasets.*

## Why it is important to protect data confidentiality

New Zealand businesses, institutions, and organisations rely on high-quality, timely, and accurate data for planning, research, and information. Good data helps New Zealand grow and prosper.

The [New Zealand Data and Information Management Principles](#) mandate that government data and information should be open, readily available, well managed, reasonably priced and reusable unless there are necessary reasons for its protection. These principles include:

> "Open: Data and information held by government should be open for public access unless grounds for refusal or limitations exist under the Official Information Act or other government policy. In such cases they should be protected.
>
> Protected: Personal, confidential and classified data and information are protected."

### Data collection depends on goodwill and trust

Much of the data collected in New Zealand is about individual people, households, businesses, and organisations — including sensitive personal and commercial data. Data gatherers and users depend on the personal and commercial trust and goodwill of the people they collect data from. Maintaining confidentiality is crucial to the New Zealand data system.

### Data confidentiality is often a legal requirement

You're often required by law to keep data confidential. If you provide data to an unauthorised user, or provide identifiable information without consent, you may be breaking the law. If the information becomes public, the implications are more serious.

# What are the principles, laws, and ethics that govern data confidentiality?

Ways of keeping data confidentiality are governed by principles, laws, and ethics.

## Principles for managing data confidentiality

Principles and legislative requirements underpin the policies, standards and guidelines for data confidentiality. For example, Stats NZ's Microdata output guide describes the methods and rules researchers must use to confidentialise output produced from Stats NZ's microdata. The methods and rules are based on legislative requirements and four principles:

- **Utility** – Ensure data outputs are as rich, detailed, and unmodified as possible.

- **Safety** – Manage the risk of identifiable information being disclosed, down to the level required by law, ethical obligations, and the preservation of trust.

- **Simplicity** – Make rules as simple to apply and check as possible.

- **Consistency** – Aim for maximum consistency across outputs released over different channels, and across similar outputs from different sources of data.

Other sets of principles that are relevant to data confidentiality include:

- Open Data Charter 2015 – guides best practice for making data open (adopted by the NZ Government, March 2018)

- New Zealand Data and Information Management Principles 2011 – guides the management of data and information that the government holds on behalf of the public (agreed to in NZ Cabinet Minute (11) 29/12, August 2011).

## Data confidentiality required by law

Data users must comply with relevant legislation. Legislation with specific requirements about keeping data confidential include:

- Privacy Act 1993

- Official Information Act 1982

- Statistics Act 1975.

You may also need to comply with other legislative requirements when using specific types of data. For example, the Tax Administration Act 1994 sets out requirements for protecting tax data and the Health Information Privacy Code 1994 sets out rules for collecting, managing and using health information.

## Ethics influence the protection of data

An integral feature of any government data system is that it is underpinned by ethical principles, to ensure responsible data use and prevent harmful outcomes. Respect for people is about recognising the people behind the data and the interests of individuals and groups in how data is used.

Protecting confidentiality of data is an important way of showing respect for people. Whenever you release data you must take extra care with data that is personally or commercially sensitive.

Among the principles in the [International Statistical Institute's Declaration on Professional Ethics](#) are that, when statisticians produce statistics, they must guard privileged information, and protect the interests of individuals and organisations.

Government agencies and other producers of official statistics are also guided by the [United Nations Fundamental Principles for Official Statistics](#):

> *"… it is the utmost concern of official statistics, to secure the privacy of data providers (like households or enterprises) by assuring that no data is published that might be related to an identifiable person or business."*

Protecting personal identifying information and preserving security of any output is emphasised in the [Principles for safe and effective use of data and analytics](#) developed by the Chief Government Data Steward and the Privacy Commissioner.

Other ethical guidelines will be relevant for specific types of research. For example, the [National Ethics Advisory Committee's Ethical Guidelines for Observational Studies](#) covers research using health data.

# What are the methods used to keep data confidential?

It is essential to use confidentiality methods to protect individually identifiable information in microdata. You may also need to use them to protect larger datasets and data outputs.

Whenever we release data — to the public, a researcher, or any other kind of data user — we must make sure its confidentiality is appropriately protected.

We protect confidentiality by ensuring that details about individual people, households, businesses, or organisations are not identifiable, and cannot be deduced. Details must not be identifiable in the raw data, published statistics, or the data output users create.

Often you can release individually identifiable details, where you need to, provided you have received written authorisation from the individual to do so.

## Use statistical methods to protect microdata and larger sets of data

Unit record data and summary data — called microdata — is especially likely to be identifiable, as it is records of individual people, households, businesses, or organisations.

Statistical data that will be published needs to be organised in a way that prevents any individual details from being identified.

## Use these statistical methods for protecting the confidentiality of data

To protect the confidentiality of microdata — and where necessary, larger datasets — you can use one or more of these statistical methods:

- Perturbation – adding random noise to data outputs.

- Aggregation – combining and/or simplifying data outputs.

- Suppression – not reporting some data outputs.

- Limiting data access – putting conditions and/or limits on access to data.
- Synthesizing synthetic unit record files (SURFs) and confidential unit record files (CURFs) for general publication. These use a combination of perturbation, aggregation, suppression, and modelling until the data is confidential, but is also still a sufficiently accurate summary of the data to fulfil data users' needs.

## Review your confidentiality methods regularly

Review your confidentiality business rules, methods, and processes regularly – at least every three to five years. You need to ensure that new technology, or the public availability of additional data, has not increased the risk of disclosure. Introduce new measures for protecting confidentiality if you need to.

## What to do if there is a data breach

Even with protection in place, there is always a risk of disclosing identifiable data. A data breach or disclosure breach happens when data is released that identifies a person, household, business, or organisation.

You must acknowledge that there is always a risk of a data breach happening. The Office of the Privacy Commissioner's Data Safety Toolkit has guidelines on remedying, managing, and mitigating data breaches.

# How can we use perturbation to protect confidentiality?

Perturbation – adding random noise to data – is a widely used data confidentiality method. Perturbation works by adding a random value to the data, to mask the data. This is called adding 'random noise'.

Perturbation is a best-practice method. It is used by Stats NZ and by many international statistical agencies, including the US Census Bureau and the Australian Bureau of Statistics.

## Use a coordinated approach to count and magnitude tables

A count measures the number of individuals whose confidentiality is being protected.

A count magnitude (or value magnitude) measures a sum of counts (or sum of values) relating to the individual data you are protecting.

For instance:
- the human population in an area is a count
- how many television sets a population owns is a count magnitude
- how much they earn is a value magnitude.

Also:
- the number of businesses in an area is a count
- how many employees they employ is a count magnitude
- how much profit they make is a value magnitude.

Stats NZ has developed a method which perturbs both count and magnitude tables: the Noised Counts and Magnitudes (NCM) method. NCM is part of Stats NZ's development of an Automated Confidentiality Service (ACS). The ACS includes software, applications, and expertise to help users automatically apply confidentiality methods and produce consistent results.
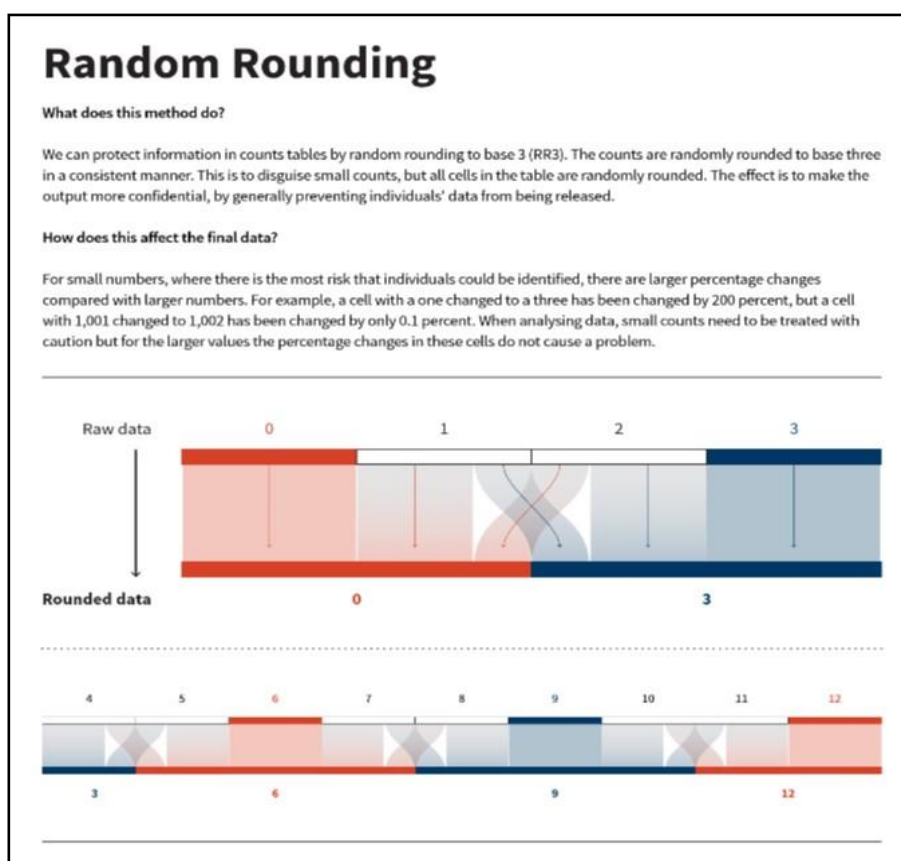
In the NCM method, each individual data record is assigned a uniformly distributed random number. These random numbers are fixed across time, to ensure the same degree of perturbation is applied to the individual over time.

## How to perturb counts

For count tables, random numbers generate a new random number for units grouped together in a cell. This is the basis for Fixed Random Rounding to base 3 (FRR3). It ensures the same group of individuals will always be rounded the same way in related tables.

In FRR3, you randomly round counts to base 3.

- Counts that are already multiples of three are left unchanged.

-  Those not a multiple of three, you round to the nearest multiple of three two-thirds of the time, or the next nearest multiple of three one-third of the time.



For example, a four will be rounded to either a three (2/3 likelihood) or a six (1/3 likelihood). This is to disguise small counts. But since all table data are rounded consistently, they are protected against both:

- differencing attacks, where closely related results might be subtracted from each other to discover underlying small counts

- Monte Carlo attacks, where attacks are run again and again, to discover the underlying raw numbers based on the distributions of results.

## Perturbation of magnitudes

Use an n% 'noise multiplier' to generate magnitude tables.

The noise protects sensitive data where there is a disclosure risk but cancels itself out in larger collections of data.

Individual values are protected by at least +/- n% for the most vulnerable data.

# How can we use aggregation to protect confidentiality?

Aggregation involves grouping categories together. You avoid disclosure by combining columns or rows into one new group. You combine or simplify data outputs. This reduces the amount of data available about individuals.

## Striking a balance between releasing data and saving labour

In the long run, aggregation is effective for striking a balance between releasing as much data as possible and limiting the work involved in producing tables.

Aggregation is useful when there are many cells with small numbers. By collapsing categories or combining data cells, you remove much of the sensitivity in the table.

You need subject matter knowledge to use this method. You need to know which values in the data are important for your data users, and how values have been aggregated in the past, so you can apply aggregation consistently.

Aggregation lowers the amount of detail in the final output data. You need to ensure that the resulting dataset is still useful for your users.

## Good data classifications and standards make aggregation easier

To maximise flexibility, code data at the lowest level of the classification possible.

Make sure that your data classifications and standards are relevant to your customer's needs.

Classifications and standards should:
- have an underlying conceptual basis
- fit within a statistical framework which is intuitive and easy to understand, navigate, and apply
- be internationally comparable when you need to compare data across countries
- be stable and comparable over time (balanced with the need to update classifications from time to time).

Classifications and standards must be unambiguous, exhaustive, and mutually exclusive:

- **unambiguous** – observations can be clearly classified into a certain group based on defined classification principles and criteria

- **exhaustive** – all cases of the observation data can be classified

- **exclusive** – groups are clearly defined so data can't be classified into more than one group.

Classifications and standards must be systematic and operationally feasible. To achieve this:

- classify observations consistently using agreed criteria

- define concepts and variables related to the classification

- make sure unspecified or residual groups like 'not elsewhere classified' contain few cases. If the size of the residual group grows considerably, you need to revise the classification system

- to minimise bias in the data, use automated processes and methods, such as coding tools (where practical)

- ensure classifications are hierarchical, with a main group level which you break down further into lower classification levels.

Use a common collapsing strategy for aggregations. Give classifications names that reflect both the most detailed and the collapsed levels.

# How can we use suppression to protect data confidentiality?

When you suppress data, you do not report selected data. Suppression is removing data from an output that reveals individualised information.

## Suppress data by not reporting some data outputs

If a data value reveals too much data about a person, household, or business, you can remove the data value from the output by *suppressing*. You replace its number value with another value, such as an empty space, a zero, or a character like 'S' or 'C'. This is *primary suppression*.

But if you decide a data value is at risk, suppressing only that value is not enough. If you give subtotals or marginal totals, it is still possible to determine the suppressed cell's actual value. You need to suppress other data values too, to protect the primary data value. Suppressing these other data values, in the same way, is *secondary suppression*.

You need to suppress other cells, so the value of the cell you first suppressed can't be determined. To suppress the fewest cells possible, complete a square of suppressions:

$2^N$ total suppressions for an N dimensional table (for example, $2^2$ = 4 total suppressions for a 2-dimensional table).

## Using secondary suppression

Secondary suppression is often not an easy task. To do it, you need:

- your criteria for performing secondary suppression (for example, minimising data loss, or sticking with previous cell suppression – refer to the criteria listed below)

- methods for identifying the best suppression pattern.

Use these criteria to decide how to apply secondary suppression.

### Historical criteria

It is important to keep track of your publication history of primary and secondary suppressions, and to take care not to disclose data where you change which cells are suppressed, over time. Changing previous cell suppression trends might cause either:

- a disclosure for the normally suppressed cell

- a problem for another cell that is now suppressed.

### Cost function criteria

You might want to:

- suppress cells that have small values (but do not suppress cells containing zeros)

- suppress a minimum number of cells

- minimise the number of values you suppress.

### Availability criteria

You might need to publish certain cells for statistical information reasons, so you cannot use them for secondary cell suppression; this might give you problems finding enough, or appropriate, cells to suppress.

To test if a suppression pattern is effective enough, make sure that in each row or column you suppress, there are at least two suppressions. For a 2-way table, each suppression should be the corner of a square or rectangle of suppressions.

## Use an automated tool for suppression

Primary and secondary suppression can be a time-consuming manual process. Some automated tools to help include *Tau-ARGUS*, *G-Confid*, and *sdcTable*.

# How can we limit access to data to protect confidentiality?

Unit record datasets that contain information about specific people, households, and organisations (microdata) are most likely to reveal identifiable information. Protect confidentiality by imposing strict limitations on access to it.

## Put conditions around the access to sensitive data

Only grant access to microdata to researchers who state the statistical purposes for wanting access.

Where you approve access, consider drawing up a legally binding contract to control access to the data.

## Follow Stats NZ's best practice principles for data access

Stats NZ assesses research proposals to access microdata using the following principles:

- access to microdata must be for statistical purposes and/or bona fide research purposes

- access to microdata must be consistent with relevant legislation

- access to microdata is at the discretion of the data custodian

- access to microdata must protect respondents' confidentiality

- access to microdata must not adversely affect data collection
- decisions on requests for access to microdata will be provided through transparent processes.

Sometimes, negotiations for researcher access involve multiple data custodians. Each custodian should consider and grant access individually.

When you consider granting access to data, also consider the Privacy Act. The Act governs the use of data beyond the purpose for which it was originally collected.
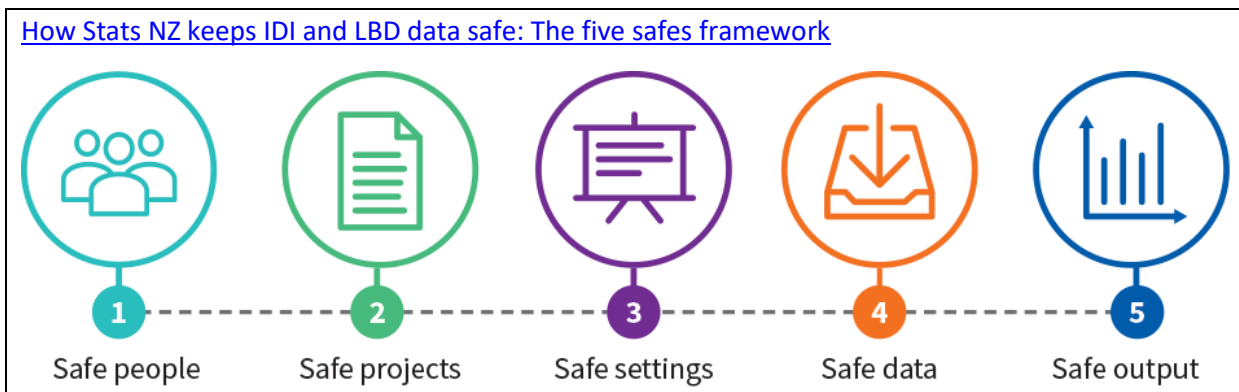
In some situations, you may need to consult the Privacy Commissioner. For example, you may have a case where a legal provision parallels or constrains the relevant legislation. Or the privacy implications of the research may not be clear.

## The 'five safes' framework

At Stats NZ, microdata researchers operate within the 'five safes' framework. We only grant access to microdata if all the following conditions are met:

- **safe people** – researchers can be trusted to use the data appropriately and follow procedures
- **safe projects** – the project has a statistical purpose and is in the public interest
- **safe settings** – security arrangements prevent unauthorised access to the data
- **safe data** – identifiers are removed before data is made available
- **safe output** – the statistical results produced do not contain any results that disclose individually identifiable information.

How Stats NZ keeps IDI and LBD data safe: The five safes framework



| 1 Safe people | 2 Safe projects | 3 Safe settings | 4 Safe data | 5 Safe output |

## The microdata output guide goes into more detail

The Microdata output guide is Stats NZ's best-practice guide for ensuring confidentiality in outputs from microdata. It covers how to use the statistical methods in greater detail, with examples.

# How can we build synthetic and confidential unit record files to support the general publication of microdata?

You can publish open microdata once its confidentiality is protected. You use statistical methods to prepare synthetic unit record files (SURFs) and confidential unit record files (CURFs) that are suitable for general publication.

You use the methods of perturbation, aggregation, and suppression to process microdata so individual people, households, businesses, and organisations cannot be identified.

## Overseas precedents for publishing open microdata

Publishing CURFs is done overseas, for example, the Integrated Public Use Microdata Series (IPUMS) published by the US Census Data for Social, Economic and Health Research. Open government initiatives have pioneered the release of CURFS, rather than national statistics organisations.

## Building confidential unit record files (CURFs) for general publication as open microdata

You build CURFs by perturbing, aggregating, and supressing microdata, until the data no longer discloses identifiable information about individuals, but is also still an accurate enough summary estimate of the data to meet the customers' needs.

## The role of synthetic data in CURFs

When you create CURFs, you may confidentialise data by replacing the real data with data you have processed or modelled. Lightly confidentialised CURFs are called partly synthetic data. Heavily confidentialised CURFs are called fully synthetic data, or SURFs.

## Creating CURFs and SURFs requires expertise and resources

Creating CURFs and SURFS is challenging and requires technical expertise. Research continues into how to automate the work. Current techniques can quantify the confidentiality importance of each variable and mitigate the risk for each variable. You can use k-anonymity testing, and Special Unique Detection Algorithms (SUDA), within automated tools like *sdcMicro*.

## The trade-off between confidentiality and utility

Often, the more heavily you confidentialise a record, the less useful it is to your customers or end-users. You need to strike a balance between confidentiality and usefulness.

If you cannot ensure data is confidential, you may need to withhold it.

# References

Future of Privacy Forum (2016). *A visual guide to practical data de-identification.*
Retrieved from
https://fpf.org/2016/04/25/a-visual-guide-to-practical-data-de-identification/.

ICT.govt.nz (2011). *New Zealand Data and Information Management Principles*.
Retrieved from
https://www.ict.govt.nz/guidance-and-resources/open-government/new-zealand-data-and-
information-management-principles/.

International Statistical Institute (2010). *Declaration on professional ethics*.
Retrieved from
https://www.isi-web.org/index.php/news-from-isi/34-professional-ethics/296-
declarationprofessionalethics-2010uk.

National Ethics Advisory Committee (2012) *National Ethics Advisory Committee's Ethical Guidelines
for Observational Studies*.
Retrieved from
https://neac.health.govt.nz/publications-and-resources/neac-publications/streamlined-ethical-
guidelines-health-and-disability.

OECD (2007) *OECD Glossary of Statistical Terms*.
Retrieved from
https://stats.oecd.org/glossary/index.htm

Open Data Charter (2015). *International Open Data Charter Principles*.
Retrieved from
https://opendatacharter.net/principles/.

Simson Garfinkel, National Institute of Standards and Technology (NIST) (2015). *De-identification of
personal information (NISTIR 8053)*.
Retrieved from
https://csrc.nist.gov/publications/detail/nistir/8053/final.

Privacy Commissioner, Stats NZ (2018). *Principles for the safe and effective use of data and analytics*.
Retrieved from
https://www.stats.govt.nz/assets/Uploads/Data-leadership-fact-sheets/Principles-safe-and-
effective-data-and-analytics-May-2018.pdf.

Stats NZ (2007). *Principles and protocols for producers of Tier 1 statistics*.
Retrieved from
http://archive.stats.govt.nz/about_us/who-we-are/home-statisphere/tier-1/principles-
protocols.aspx

Stats NZ (2015). *Privacy, security, and confidentiality of information supplied to Statistics NZ*.
Retrieved from
http://archive.stats.govt.nz/about_us/legisln-policies-protocols/confidentiality-of-info-supplied-to-
snz/safeguarding-confidentiality.aspx.

Stats NZ (2016a). *Microdata output guide (Fourth edition)*.
Retrieved from
http://archive.stats.govt.nz/tools_and_services/microdata-access/data-lab/microdata-output-guide.aspx.

Stats NZ (2016b). *Introducing new method for confidentialising business demography tables*.
Retrieved from
http://archive.stats.govt.nz/browse_for_stats/businesses/business_characteristics/new-method-for-confidentialising-tables.aspx.

Stats NZ (2017b). *Information, Privacy, Security and Confidentiality Policy*.
Retrieved from
http://archive.stats.govt.nz/about_us/legisln-policies-protocols/info-priv-sec-confid-policy.aspx.

Stats NZ (2017c). *How we keep IDI and LBD data safe: The five safes*.
Retrieved from
http://archive.stats.govt.nz/browse_for_stats/snapshots-of-nz/integrated-data-infrastructure/keep-data-safe.aspx#safes

Stats NZ (2018c). *Microdata access protocols*.
Retrieved from
http://archive.stats.govt.nz/about_us/legisln-policies-protocols/microdata-access-protocols.aspx.

United Nations Statistics Division (2015) *UN Fundamental Principles of Official Statistics – Implementation guidelines.*
Retrieved from
https://unstats.un.org/unsd/dnss/gp/Implementation_Guidelines_FINAL_without_edit.pdf