

AUSTRALIAN DATA ARCHIVE
SUBMISSION TO THE PRODUCTIVITY COMMISSION
“DATA AVAILABILITY AND USE” PUBLIC INQUIRY
29 JULY 2016

INTRODUCTION

The Australian Data Archive (ADA) strongly supports the establishment of the public inquiry into Data Availability and Use, and welcomes the opportunity provided by the Productivity Commission to provide comment on these matters. ADA has provided access to data from government, academic and private sector sources for researchers since 1981, and currently supports 10 different Federal Government departments and agencies, among the 5000 datasets collected in over 1500 research projects going back to 1947. We believe that there are significant opportunities available through the improved provision of both public and private data to government, the academic community, industry and the public at large.

ABOUT ADA

The Australian Data Archive (ADA) provides a national service for the collection and preservation of digital research data and to make these data available for secondary analysis by academic researchers and other users. ADA was established at the ANU in 1981 with a brief to provide a national service for the collection and preservation of computer readable data relating to social, political and economic affairs and to make these data available for further analysis.

A team of professional data archivists provides both stewardship and outreach services to the Australian community. The archive:

- acquires, documents, preserves and disseminates data online to a broad range of social science researchers in the university, government, and other sectors
- provides the only comprehensive social science data collection in Australia, with a catalogue of over 5000 data sets
- holds data from Australian surveys, opinion polls and censuses and includes data from other countries within the Asia Pacific region
- supports access to datasets from 10 different Federal government departments and agencies
- provides specialist services within specific subject areas, including social sciences, health, indigenous studies, electoral behaviour, criminology and some humanities disciplines, and within specific data types, including quantitative, qualitative, time series and panel data, and historical statistics
- locates and manages access to overseas social science data sets required by Australian based researchers
- adopts, develops and applies standards in line with international best practice
- belongs to international organisations (such as the International Federation of Data Organizations and the International Association of Social Science Information Service and Technology) and plays a major role in cross-national collaborative projects

- provides support for the management and dissemination of grant-funded data collections such as ARC and NHMRC projects

GENERAL COMMENTS

As noted in ADA's mission, the importance of access to data for research purposes has been the primary reason for ADA's existence for over 35 years now. Many of the issues identified by the Commission in the issues paper for this inquiry have been long-running, and ADA supports the efforts to enable improved access to data that this inquiry should hopefully help to enable.

ADA notes in particular three issues highlighted in the paper:

1. *"Recurring data-related themes arising in Commission reports" (Box 2, page 8)*. Each of the six themes identified in this breakout box are consistent with issues commonly raised by ADA users, and are an accurate reflection of the potential improvements that could be made by improving the data access environment within Australia. Indeed, ADA provided advice to one of the PC inquiries mentioned in the report (on gambling in 2010) where similar concerns of a lack of data sharing and limited access to data for research and policy purposes also undermined the capacity to address problem gambling issues in the Australian context.
2. *Australia's commitment to public sector open data (p.12)*. The recent *Public Data Policy Statement* released by the Department of Prime Minister and Cabinet in December 2015 has provided a clear signal for Federal departments to give due consideration to access to data for the public good. While the Commission notes that they will need to make efforts "to put this commitment to the test" (p.13), ADA's recent experience in our work with government department would suggest that this policy statement has provided a degree of impetus within departments in terms of data access matters.
3. *Research data*. There has been a long term commitment since at least 2004 to support the management of and access to research data (that is, data generated and disseminated as an output of research conducted within Australia's university and broader research community), through the National Collaborative Research Infrastructure Strategy and related investments¹. The investments in research infrastructure, totalling over \$2.8 billion since 2004, particularly the Australian National Data Service established as part of this strategy, are not considered in the issues paper. The increasing volume of data outputs from federally funded research through the ARC and NHMRC (and other agencies) are also not considered. This data and infrastructure provide a significant complement to the public sector sources already under consideration in this review.

RESPONSES TO QUESTIONS

The remainder of our submission provides direct responses to the questions raised by the Commission in the published issues paper. We have not commented on all questions raised by the paper.

¹ <https://www.education.gov.au/national-collaborative-research-infrastructure-strategy-ncris>

A. QUESTIONS ON HIGH VALUE PUBLIC SECTOR DATA

What public sector datasets should be considered high-value data to the: business sector; research sector; academics; or the broader community?

There are numerous datasets that exist in the public sector that might be considered high-value to various communities, a number of which have been identified in other submissions to this inquiry. For example, the four databases detailed by the Department of Employment in their submission to the Inquiry² are all potentially of high value to end users, such as for policy and evaluation research, labour market analysis or social science research generally.

A broader consideration for the public and the research community is often the lack of information available about the datasets that do exist. It is difficult for those external to the public sector to identify data that may be of value when there is limited knowledge of what data exists. Even where there is information about the existence of data, it can often be difficult to discover that information, and to determine then how the data might be accessed. In this regard, the development of the data.gov.au portal is to be commended as a means for improving both discoverability and access of public sector data.

One potential framework for evaluating this element of the public sector data environment would be through the application of the FAIR principles recently established by the Force 11 group³ (also noted in the University of Sydney submission⁴). Application of the four FAIR principles – Findability, Accessibility, Interoperability and Re-usability – would provide a foundation for addressing some of the concerns noted above.

What characteristics define high-value datasets?

The characteristics of “high value” data are often in the eye of the beholder. For this reason, social science data archives like ADA have generally not sought to define what is “high-value”, but rather to provide the relevant support for our users to allow them to make an informed judgment about their potential value.

A key element of this approach has been the establishment of a common standard (known as a *metadata standard*) for the sharing of data for research use and between archives around the world. There is a long history of the development of these common standards in social science⁵, which has culminated in the establishment of the Data Documentation Initiative (DDI) standard in 1998, and of the DDI Alliance to govern the standard in 2003⁶. DDI-Lifecycle is designed to document and manage data across the entire life cycle, from conceptualization to data publication and analysis and beyond, while a lightweight version, DDI-Codebook, is intended primarily to document simple survey data.

DDI is now in use in over 90 countries around the world, and is used in the World Bank, World Health Organisation and Eurostat to support data and metadata dissemination internationally. In

² http://www.pc.gov.au/__data/assets/pdf_file/0012/202440/sub018-data-access.pdf

³ <https://www.force11.org/group/fairgroup/fairprinciples>

⁴ http://www.pc.gov.au/__data/assets/pdf_file/0011/202988/sub035-data-access.pdf

⁵ <http://www.ddialliance.org/system/files/DDI%20Timeline%20With%20Foundational%20Events-For-Website.pdf>

⁶ ADA is a current member of the Alliance and user of the DDI standard, and Steven McEachern (ADA Director and author of this submission) is current Chair of the DDI Alliance

Australia, the standard is used by ADA to support all of our data dissemination activities, and by the Australian Bureau of Statistics to facilitate the internal collection, management and dissemination of statistical data⁷.

The benefits of the use of a common metadata standard have been numerous, but are best reflected in the World Bank's International Household Survey Network⁸. The IHSN data catalog, through the use of a common software platform and common metadata standard for documenting data, have now enabled the findability and accessibility of 5629 datasets from 90 countries around the world through World Bank funding support.

The other core element of assessing potential value of a dataset lies in the assessment of data quality. Here a data quality framework such as the ABS's Data Quality Framework could be used as the basis for a common approach. The ABS framework (based on similar models in Europe and Canada) includes seven elements:

1. institutional environment,
2. relevance,
3. timeliness,
4. accuracy,
5. coherence,
6. interpretability, and
7. accessibility.

The importance of appropriate provision of information about the data is again critical here (through a standard such as DDI) to enable potential users to assess the quality of the data against such a framework.

What benefits would the community derive from increasing the availability and use of public sector data?

The benefits derived from improved availability and use of public sector data are numerous, and we believe that the summary included in the Issues Paper provide a good overview of the potential benefits. We would also recommend here two studies in this area that might provide additional insight:

1. The study by Houghton (2011)⁹ for the Australian National Data Service noted in the Issues Paper conducting cost-benefit analyses of the value of improved data access to public sector information
2. The study by Capgemini Consulting (2015)¹⁰ for the European Data Portal (EDP) estimating four benefits (direct market size, number of jobs created, cost savings, and efficiency gains) from the establishment of the EDP.

⁷ <http://www.amstat.org/meetings/ices/2012/papers/302207.pdf>

⁸ <http://ihsn.org/home/>

⁹ http://www.ands.org.au/__data/assets/pdf_file/0004/394285/houghton-cost-benefit-study.pdf

¹⁰ http://www.europeandataportal.eu/sites/default/files/edp_creating_value_through_open_data_0.pdf

QUESTIONS ON COLLECTION AND RELEASE OF PUBLIC SECTOR DATA

What are the main factors currently stopping government agencies from making their data available?

While it is difficult to know the reasons for specific agencies' limitations on enabling data access, some general reasons can be drawn from the work of the IHSN discussed earlier. In their 2010 working paper on "Dissemination of Microdata Files"¹¹, Dupriez and Boyko note six core areas of potential risk and/or cost:

1. Ethical issues and maintaining trust
2. Legal issues
3. Exposure to criticism and contradiction
4. Cost
5. Loss of exclusivity
6. Technical capacity

Notably however the paper also provides some suggestions for how these risks can be addressed in the situation of the dissemination of microdata / unit record data¹².

How could governments use their own data collections more efficiently and effectively?

There would appear to be significant potential for government departments to take advantage of their own data collections in a more efficient manner. For example, in updating their data and documentation in order to facilitate access to data for external use, agencies may well find that they are then able to make use of the improved data themselves – often known as "eating your own dog food" or "dogfooding" in software development¹³

Should the collection, sharing and release of public sector data be standardised? What would be the benefits and costs of standardising? What would standards that are 'fit for purpose' look like?

Given our experience with data and metadata standardisation¹⁴, we would strongly support the standardisation of public sector data for two reasons:

1. Consistency of experience for human users of the data: the use of common standards helps to enable users to compare between different data sources in order to assess their suitability for use and their quality.

¹¹ <http://ihsn.org/home/sites/default/files/resources/IHSN-WP005.pdf>

¹² Microdata is "the electronic data files containing the information about each unit of observation" (IHSN, 2010, p.3) collected in a data collection activity or administrative process – as opposed to aggregate statistics about the units.

¹³ <https://technet.microsoft.com/en-us/library/cc627315.aspx>

¹⁴ The DDI Alliance provides an overview of the benefits of structured and standardised metadata in systems such as public sector agencies: <http://www.ddialliance.org/training/why-use-ddi>

2. *“Machine actionability” of the data*: in order for information systems to be able to process data efficiently, they have to “know what to expect”. While there is on-going improvement in the development of methods for analysing “unstructured” data, the use of common standards enables developers of such systems to know what data their systems can expect, in order to perform relevant operations on those systems.

Several of the benefits of standardisation have been discussed elsewhere in this submission, and so we will not repeat them here. It should be noted however that the implementation of standards is not without cost, most notably in terms of data harmonisation and integration, administrative burden and technical infrastructure.

What criteria and decision-making tools do government agencies use to decide which public sector data to make publicly available and how much processing to undertake before it is released?

We cannot respond to this question directly, given that we are not a public sector agency. However, as an archive that supports access to public sector data through our facilities, we are often asked to provide information to public sector about options and issues in data access and dissemination, across a variety of research topics and areas. We would be happy to discuss further with the Commission if this is of interest.

What specific government initiatives (whether Australian Government, state, territory or local government, or overseas jurisdictions) have been particularly effective in improving data access and use?

Historically, the establishment of the agreement for microdata access between the then Australian Vice-Chancellors Committee (AVCC – now Universities Australia) and the Australian Bureau of Statistics in 1998 provided a model for improving access to unit record data from the Australian Bureau of Statistics across the Australian university sector, and provides a potential model for other public sector data providers for suggestions regarding improvements in the Australian situation.

A second major initiative by the Australian Bureau of Statistics was the establishment of the free access model for access to ABS publications and (most if not all) data, and the subsequent implementation of the Creative Commons licensing model for ABS public data and publications. Houghton (2011) estimates in his study cited in the Issues Paper, that the cost of moving to a free access model of approx. \$3.5 million per annum were outweighed by the benefits to users of over \$5 million per annum – and that these returns were relatively low compared to other projects which saw up to a fivefold difference between benefits and costs.

More recently, the investment of government departments in major data collection initiatives – such as the various longitudinal studies¹⁵ funded by the Department of Social Services, Department of Health and Department of Education – have provided a major resource for research and policy development in Australia. The establishment of these studies with a recognition of the policy and research potential of such data, and with an explicit emphasis on sharing the data for

¹⁵ For an overview of these studies, see: <https://www.dss.gov.au/about-the-department/publications-articles/research-publications/longitudinal-data-initiatives/guide-to-australian-longitudinal-studies>

these purposes, has resulted in major benefits for the Australian public and significant research output (for example as noted in the DSS Flosse research database of outputs from their studies¹⁶).

The other major recent initiative we would note is the establishment of the data.gov.au initiative and particularly the release of the Public Data Policy Statement in November 2015. Since its release, the Policy Statement has raised as a significant factor by a number of public sector agencies when considering disseminating their research data through ADA. It is apparent from our experience that departments are now giving serious consideration to data access considerations as a result of this policy, although as the Issues Paper notes this is a commitment that has not yet been put to the test. Nonetheless we believe the Policy Statement provides an important element of the overall framework for improving public sector data access in Australia.

QUESTIONS ON DATA LINKAGE

Which datasets, if linked or coordinated across public sector agencies, would be of high value to the community, and how would they be used?

Which rules, regulations or policies create unnecessary or excessive barriers to linking datasets?

ADA has no current direct involvement in data linkage at this time, and as such we have chosen not to respond to this question directly. Instead, we would acknowledge the major contributions that data linkage has provided for Australian academic and policy research, reflected in the submission by the Centre for Data Linkage at Curtin University¹⁷. We would also recommend further discussion with staff from ADA's parent centre, the ANU Centre for Social Research and Methods¹⁸. Several of the Centre staff have experience in the various datasets and procedures used in the Australian public sector environment.

How can Australia's government agencies improve their sharing and linking of public sector data? What lessons or examples from overseas should be considered?

Again, while ADA has no direct involvement in data linkage at this time, we would again defer to the Curtin University submission above regarding recommendations for improvements in Australian government agencies.

There are numerous examples of the value of data linkage, but we would point particularly to the experience of data linkage within the United Kingdom. Britain is unique in the world in having a portfolio of national birth cohort studies that follow individuals from birth through childhood and into adult life. These studies, the first of which was established in 1946, have already been instrumental in providing evidence relevant to a wide range of policy issues particularly in the areas of health, child development, education and employment.

A recent presentation¹⁹ by Prof. Jane Elliott from the Centre for Longitudinal Studies at University College, London (and now head of the Economic and Social Research Council in the UK), provided an overview of some of these, as well as an overview of recent initiatives for expanding data linkage in the UK environment. One example of these initiatives is the Administrative

¹⁶ <http://flosse.fahcsia.gov.au/>

¹⁷ http://www.pc.gov.au/__data/assets/pdf_file/0007/203101/sub041-data-access.pdf

¹⁸ <http://rssh.anu.edu.au/schools-centres/socialresearch/home>

¹⁹ <http://www.anu.edu.au/events/evidence-based-policy-the-value-of-longitudinal-studies-and-administrative-data-linkage>

Research Data Network funded by the ESRC, which facilitates access to administrative data and data linkage. The Network includes several recent examples of cases where linked data has been utilised to improve policy outcomes and public benefit in the United Kingdom²⁰. We would urge the Commission to explore some of the UK initiatives as exemplars that could be considered here in Australia.

QUESTIONS ON HIGH VALUE PRIVATE SECTOR DATA

What private sector datasets should be considered high-value data to: public policy; researchers and academics; other private sector entities; or the broader community? In each case cited, what characteristics define such datasets?

What would be the public policy rationale for any associated government intervention?

What benefits would the community derive from increasing the availability and use of private sector data?

ADA has no comment on these questions.

QUESTIONS ON ACCESS TO PRIVATE SECTOR DATA

As ADA has only a limited role in the provision of private sector data in Australia, we offer only general comments on these issues.

Are there any legislative or other impediments that maybe unnecessarily restricting the availability and use of private sector data? Should these impediments be reduced or removed?

What are the reasonable concerns that businesses have about increasing the availability of their data?

What principles, protocols or legislative requirements could manage the concerns of private sector data owners about increasing the availability of their data?

Should the collection, sharing and release of private sector data be standardised in some way? How could this be done and what would be the benefits and costs? What would standards that are 'fit for purpose' look like?

To what extent can voluntary data sharing arrangements—between businesses /between businesses and consumers/ involving third party intermediaries—improve outcomes for the availability and use of private data? How could participation levels be increased?

Would such voluntary arrangements raise competition issues? How might this change if private sector information sharing were mandated? Is authorisation (under the Competition and Consumer Act 2010 (Cth)) relevant?

²⁰ <https://adrn.ac.uk/research-projects/case-studies>

ADA has no comment on these questions, other than to again recommend the use of common metadata standards if such an approach to private sector data sharing is adopted.

What role can governments usefully play in promoting the wider availability of private datasets that have the potential to deliver substantial spillover benefits?

How can the sharing and linking of private sector data be improved in Australia? What lessons or examples from overseas should be considered?

There may be some potential benefits to government assisting in enabling access to private sector data in Australia. Again here we would point to the United Kingdom's Administrative Data Research Data Network as an example. The ADRN is in fact the first phase of a broader "Big Data Network" program funded by the Economic and Social Research Council. The second phase of this program is the "Business and Local Government Data" network²¹, established in 2014. Along with researcher support, this network also includes support for small and medium sized enterprises to access and use new data sources for their own business.

Who should have the ownership rights to data that is generated by individuals but collected by businesses? For which data does unclear ownership inhibit its availability and use?

ADA has no comment on this question.

QUESTIONS ON CONSUMER ACCESS TO, AND CONTROL OVER, DATA

What impediments currently restrict consumers' access to and use of public and private sector data about themselves? Is there scope to streamline individuals' access to such data and, if there is, how should this be achieved?

Are regulatory solutions of value in giving consumers more access to and control over their own data?

Are there other ways to encourage greater cultural acceptance amongst businesses of consumer access to data about them?

What role do third party intermediaries currently play in assisting consumers to access and use data about themselves? What barriers impede the availability (and take-up) of services offered by third party intermediaries?

What datasets, including datasets of aggregated data on consumer outcomes at the product or provider level, would provide high value to consumers in helping them make informed decisions? What criteria should be used to identify such datasets? What, if any, barriers are impeding consumers' access to, and use of, such data?

ADA has no comment on these questions.

²¹ <http://www.blgdataresearch.org/>

QUESTIONS ON RESOURCE COSTS OF ACCESS

How should the costs associated with making more public sector data widely available be funded?

ADA strongly supports the provision of access to public sector at minimal or no costs to the public at large where possible. The evidence from the study by Houghton (2010) also notes that the costs of supporting data access are often outweighed by the gains in benefits that are associated with the provision of free access.

As such, while we do not recommend a specific cost model, we would recommend to the Commission that it prioritise consideration of models that minimise the need for external charging of users as the means of funding the costs of access. One recent project studying cost models for the provision of digital research data that the Commission may wish to consider was the Collaboration to Consider the Costs of Curation (4C) Project²², which included 13 European repositories headed by the Joint Information Systems Committee in the UK.

To what extent are data-related resources in agencies being directed towards dealing with data management and access issues versus data analysis and use?

ADA is unable to comment on this question as we have no knowledge of agencies current resourcing in these areas.

What pricing principles should be applied to different datasets? What role should price signals play in the provision of public sector data?

As noted above, ADA strongly supports the provision of access to public sector at minimal or no costs to the public at large where possible. As such, we believe that price signals should not play a role in the provision of public sector data.

Is availability of skilled labour an issue in areas such as data science or other data-specific occupations? Is there a role for government in improving the skills base in this area?

ADA has no comment on this question.

QUESTIONS ON PRIVACY PROTECTION

What types of data and data applications (public sector and private sector) pose the greatest concerns for privacy protection?

ADA is unable to comment on this question in detail, given the broad types and forms of data available. However we would be happy to discuss further with the commission regarding the various types of data in use in the social sciences and related disciplines.

How can individuals' and businesses' confidence and trust in the way data is used be maintained and enhanced?

²² <http://4cproject.eu/>

A recent development in the provision of research data in the UK and Europe more broadly has been the adoption of appropriate trust models to support access to data for research and broader use. A potential model for this is the Five Safes model developed at the UK Data Archive in conjunction with the Office of National Statistics in the United Kingdom. An implementation of this model, known as the “Trusted Access” model, has recently been adopted by the ABS as the reference model for assessing new and existing data access methods for ABS products and services²³, and also at Statistics NZ²⁴.

The Five Safes model has five related elements:

1. Safe People: Can the researcher(s) be trusted to use the data in an appropriate manner?
2. Safe Projects: Is the data to be used for an appropriate purpose?
3. Safe Settings: Does the access environment prevent unauthorised use?
4. Safe Data: Is there a disclosure risk in the data itself?
5. Safe Output: Are the statistical results non-disclosive²⁵

ADA supports the adoption of these principles in the ABS environment and would recommend that the Commission give consideration to such models as a means for supporting public and business trust and confidence.

What weight should be given to privacy protection relative to the benefits of greater data availability and use, particularly given the rate of change in the capabilities of technology?

Are further changes to the privacy-related policy framework needed? What are these specific changes and how would they improve outcomes? Have such approaches been tried in other jurisdictions?

How could coordination across the different jurisdictions in regard to privacy protection and legislation be improved?

ADA has no comment on these issues.

How effective are existing approaches to confidentialisation and data security in facilitating data sharing while protecting privacy?

²³ <http://www.abs.gov.au/AUSSTATS/abs@.nsf/be4aa82cd8cf7f07ca2570d60018da27/e4d483bab4e1ad93ca257f4c00170bb6!OpenDocument>

²⁴ http://www.stats.govt.nz/browse_for_stats/snapshots-of-nz/integrated-data-infrastructure/keep-data-safe.aspx

²⁵ Per the ABS “Trusted Access” model cited above

In our experience, the existing approaches to confidentialisation and security have provided a reasonably effective mechanism for supporting access to data for research purposes. However, we would again highlight the Five Safes/Trusted Access model discussed above as a potential basis for evaluating the data security and confidentiality characteristics of data relative to other elements of a data access system.

What lessons from overseas jurisdictions can Australia learn from regarding the use of individuals' and businesses' data, particularly in regard to protecting privacy and commercially sensitive or commercially valuable information?

What are the benefits and costs of allowing an individual to request deletion of personal information about themselves? In what circumstances and for what types of information should this apply?

What competing interests (such as the public interest) or practical requirements would indicate that the ability to request deletion should not apply?

ADA does not wish to comment on these issues.

QUESTIONS ON OTHER RESTRICTIONS

Having regard to current legislation and practice, are further protocols or other measures required to facilitate the disclosure and use of data about individuals while protecting privacy interests? What form should any such protocols or other measures take?

ADA would recommend the use of a model such as the Five Safes/Trusted Access model for facilitating an assessment of the processes involved in data access.

ADA also has significant experience in the development of protocols to support dissemination of data for research and other purposes. We would be happy to discuss this further with the Commission.

Is there need for a more uniform treatment of commercial in confidence data held by the Australian Government and state and territory governments?

Are there merits in codifying the treatment and classification of business data for privacy or security purposes? What would this mean in practice?

ADA has no comment on these issues.

QUESTIONS ON DATA SECURITY

Are security measures for public sector data too prescriptive? Do they need to be more flexible to adapt to changing circumstances and technologies?

In ADA's experience, the emergence of both new technologies for access (such as the ABS TableBuilder environment²⁶) and the development of confidentialisation and data security methods, suggest that a more flexible approach to data security is needed going forward.

The Five Safes model is informative here. In past experience, there has been a heavy dependence on the use of confidentialisation procedures (a "safe data" approach) to manage concerns with the dissemination of data in unexpected and insecure ways ("unsafe settings"). However the increasing availability and viability of remote access ("safe settings") may make it possible to reduce the need for excessive confidentialisation in environments where both the data is unable to be removed ("safe setting") and the output of any analysis can be vetted ("safe outputs").

How do data security measures interact with the Privacy Act?

How should the risks and consequences of public sector and private sector data breaches be assessed and managed? Is data breach notification an appropriate and sufficient response?

ADA has no comment on these issues.

CONTACT FOR FURTHER INFORMATION

For further information on any of the comments in this submission, please contact and author of this submission, Dr. Steven McEachern, Director of the Australian Data Archive at the Australian National University.

²⁶ <http://www.abs.gov.au/websitedbs/censushome.nsf/home/tablebuilder>