



University of the
West of England

Faculty of Business and Law

Can a change in attitudes improve effective access to administrative data for research?

Felix Ritchie

University of the West of England, Bristol, UK

Economics Working Paper Series
1607



University of the
West of England

bettertogether

Can a change in attitudes improve effective access to administrative data for research?

Felix Ritchie, University of the West of England, Bristol

Abstract

The re-use of administrative data for social research holds great potential. From a privacy perspective, administrative data present some additional challenges, including lack of consent, existence of matching databases, and the association with data breaches by administrative staff. Access to government data for research is currently undergoing a slow, small but significant transformation from the defensive strategies of the past. A key driver of this is attitudinal change; the new approach is characterised as evidence-based default-open, risk-managed, user-centred decision-making, and offers more security at lower cost with greater researcher value. Despite its apparent superiority, this approach is still a minority position. It fundamentally challenges the way decisions are made in the public sector, at an individual and institutional level, as well as making risks more explicit. The effective re-use of administrative data may then depend upon the degree to which attitudes to decision-making in the public sector can be changed.

Corresponding author: Felix Ritchie felix.ritchie@uwe.ac.uk
Bristol Business School, University of the West of England, Coldharbour Lane, Bristol BS16 1QY, UK

JEL codes: H11 M19

Key words: administrative data; confidentiality; data access; evidence-based decision-making; risk

Introduction

The use of administrative data for social policy research offers great possibilities. The data are often the population of interest – for example, all those who have engaged with a service in some way. The data are used for delivering the service, and so are likely to be checked for accuracy. Systems already exist for acquiring the data, so extracting the data can be done at a relatively low marginal cost. Finally, the data are likely to be fully identified at source, opening the possibility of datasets being linked through exact matching – for example, merging social security data with medical reports to evaluate the impact of interventions. The Wellcome Trust (2015) notes that this last aspect is particularly important for epidemiological (public health) studies, where the research value of linking contemporary medical records to longitudinal surveys is unparalleled.

Not everything is rosy in the garden of administrative delights. Data collected for administrative purposes may not be suitable for research purposes. Data cleaning for operational purposes may not meet the particular accuracy required for statistical analysis. Having access to the population of service users does not allow one to ask questions about non-service users. The dynamic nature of administrative data creates additional management and distribution problems (compared to, say, the defined sample for a periodic survey). Nevertheless, allowing the research use of administrative data opens up the possibility of large, low-cost, high-quality linkable data.

When the data under consideration are confidential (in other words, not suitable for public release), additional issues arise. Consent for re-use may not be given, and the use of statutory gateways to access information in such cases may not be clear. Sharing data may be seen as a breach of confidentiality by the service user irrespective of the legal basis of any access agreement, which may affect users' willingness to engage with the service. When considering the likelihood of breaches of confidentiality by a malicious party, the intruder is helped enormously by the fact that there always exists at least one perfect match for the data being targeted.

Overarching all of this is the attitude to risk in the public sector. Ritchie (2014b) argues that decisions which have a risky outcome are harder to accept in the public sector, for a variety of personal and institutional reasons. Such decisions are even harder when the benefit of such a decision is likely to be, at best, indirect, which is mostly the case for research data access (Moore, 2010; Ritchie and Welpton, 2012). Additionally, the resulting defensive decision-making is supported by almost all the academic literature on statistical protection of data.

Making administrative data available for research purposes has to balance the public benefit from research with the right to privacy of the individual to whom the data relate. There is however a growing concern that the current balance is tilted too far in favour of privacy. The counter-arguments in support of greater access include the lack of evidence for malicious ‘intruders’, a growing recognition of the role of mistakes, more understanding of the social elements of the researcher-data owner relationship, an acknowledgement of the indirect costs to society of over-protective activity, and the excellent record of the academic research community in looking after confidential data.

Two factors are driving this change in perspective. The first is operational improvements based on technological advances and psychological insights into researchers’ incentives. More important is attitudinal change: new ways of looking at the world provide the context within which the gains from operational advances can be maximised. Data access solutions in the public sector following these new approaches have been shown to deliver more security at lower cost with greater researcher value.

This new approach fundamentally challenges the way decisions are made in the public sector, at an individual and institutional level. Experience suggests that organisational attitudes can change, but only slowly – precedent is often the most powerful argument wielded in the public sector, and building precedents takes time. Making truly effective use of administrative data may therefore depend upon whether attitudes to decision-making in the public sector can be changed. This article explores that question.

The next section covers, briefly, the literature on public-sector decision-making and attitudes to risk; it also describes the current literature on statistical risk, showing how the conceptual homogeneity has made it extremely difficult for decision-makers in the public sector to challenge traditional perspectives. Section three considers the privacy challenges presented by administrative data, and their relationship to the statistical models of risk. Section four brings these elements together to discuss why ‘defensive’ decision-making dominates in the public sector.

Section 5 then describes how a different approach to evaluating risk is generating radically different perspectives, and the evidence for this new approach leading to better social and private outcomes. Given the apparent superiority of the new approach, Section 6 considers why its wider adoption has been limited. Section 7 concludes that to fully exploit the gains from administrative data, a much wider change in attitudes is needed.

Decision-making and risk

Decision-making in the public sector

Ritchie (2014b) reviews the literature on decision-making under uncertainty in the public sector. Much use is made of various forms of public choice theory, which adopts the rational, utility-maximising individual central to standard (neoclassical) economic models. Decision-making is carried out within a framework of known, fixed and consistent preferences, and the individual is only concerned with his or her own interests (which may include the interests of others).

This model has come under substantial attack from multiple directions. Authors such as Boyne (1998), Andrews et al (2011) and Carpenter and Krause (2015) argued that the statistical evidence for such models is profoundly flawed, Schillemans and Busuioc (2015) that the core theoretical assumptions underlying the models have no empirical support.

More importantly, others have argued that, to all intents and purposes, humans are irrational and are strongly influenced by aspects such as the starting position of any discussion, memory of recent events, perceptions of fairness, over- or under-confidence, and so on. For example, see Sutherland (1992) or Kahnemann (2012) for the behavioural psychology perspective; Feeney and DeHart-Davis (2009) for institutional influences; and Viscusi et al (2011), Linde and Sonnemans (2012), Kjeldsen and Jacobsen (2013) and Gazley (2014) for social effects.

Despite this, the rational decision-maker remains immensely important. A natural consequence of the rational model is the 'market-based governance model' (Somers, 2008; Adams, 2013), sometimes referred to as 'new public management' (Hartley, 2005), which encourages turning public services into private markets. This has been the dominant political perspective in recent decades, particularly in the Anglo-Saxon countries but also increasingly in other OECD economies. In other words, while the irrational, institutionalised, social individual is widely accepted in academic circles, in practical policymaking the rational model is still dominant.

Risky decision-making

A standard tenet in economics (and hence, public choice theory) is that people are risk averse; that is, given the choice between a certain outcome and an uncertain outcome with the same expected value, an individual would choose the certain outcome. Risk-aversion derives from 'expected utility theory', and provides the basis for much economic analysis (eg Gollier et al, 2013).

However, experimental work in the 1970s showed that this result was asymmetric and dependent upon the reference point and framing (Kahneman and Tversky, 1979): while risk-aversion is typical when considering gains, the same individual will display risk-loving behaviour when faced with a choice between losses. The result arises because people tend to regret losses more than they value gains (Kahneman, Knetsch and Thaler, 1991).

More generally, uncertainty is frequently associated with irrational decisions. The extensive surveys of Sutherland (1992) and Kahneman (2012) both note that humans are particularly poor at evaluating probabilities. Perceptions are influenced by the way options are presented, by the difficulty of any calculations, and by irrelevant information which is presented at the same time. People make use of heuristics (rules-of-thumb) to 'home in' on answers; the presence of uncertainty increases the importance of these heuristics.

Risky decision-making in the public sector

In economic theory, the collective rationality of the government makes it risk-neutral (Eckerd, 2014), as no individual bears the direct costs/benefits of any action. By design, the administrator's rewards are largely unaffected by performance; successful innovations may bring enhanced career prospects or some performance related pay, but in general the return to the administrator is expected to be independent of the success of the innovation. This symmetric approach to gains and losses ensures the rationality of the decision-making process.

This perspective has been challenged on the basis that decisions are still made by individuals who look to their own interests first. Public choice theory leads to a model of 'pork barrel' politics and administration, as decisions made by the bureaucrat leading to his or her specific benefit are diffused across a wide population. However, as Ritchie (2014b) notes, for most public servants the opposite is the case: when taking decisions they face specific costs (their continuing employment prospects if the decision turns out to be the wrong one) and diffuse benefits (many of the beneficiaries of the decision might not even recognise that they are beneficiaries). Various authors have noted the focus on 'avoiding failure' rather than a neutral assessment of options (eg Lofstedt, 2004; Moore, 2010; Mazzucato, 2014) as have practitioners and auditors (Office of the Auditor-General of Canada, 1998; House of Lords, 2006; Australian Government, 2015b).

Yang and Holzer (2013) suggest a potential reason for this in ‘grievance asymmetry’: users of government services tend to consume them out of necessity (for example, completing a tax return), and hence are only likely to reflect on the quality of service when that service fails. Bhatta (2003) notes that ‘government’ tends to be seen as indivisible: each failure is seen as symptomatic of a general failure of government. Finally, Cabantous et al (2012) show that different types of uncertainty are perceived differently: some bad calls are more acceptable than others; Ritchie (2014b) uses this argument to build a rationale for crises in public sector decision-making.

Overall, Ritchie (2014b) concludes that, while the theoretical rationale for risk-neutrality in the public sector seems sound, in practice the incentives faced by individuals making risky decisions encourage excessively cautious outcomes.

Attitudes to data risk in the statistical literature

There is a very large literature relating to ‘statistical disclosure control’ (SDC), the use of statistical techniques to analyse and address disclosure risks associated with the analysis of confidential data; see Hundepool et al (2010) for a review. SDC covers both the application of anonymisation techniques to source data before it is released to a researcher, and ensuring that outputs produced from confidential data do not lead to disclosure through, for example, a graph which highlights outliers.

The underlying techniques on which SDC theory is based have barely changed in recent decades. Hafner et al (2015) describe the ethos of any research paper or practice guide as:

- default-closed: data or statistics should not be released unless the protection method has been ‘proven’ to reduce confidentiality risk to a sufficiently low level
- intruder-driven: any release is assumed to be at risk of ‘intruders’ who are well-supplied with statistical knowledge and external data, and who have unlimited time, resources and malicious intent
- data-centred: discussions focus on protecting the data owner
- theory-driven: consideration of risks focuses on hypothetical possibilities rather than likely outcomes; little or no evidence is used to support attack scenarios
- worst-case planning: models typically define a ‘worst-case scenario’, which simplifies modelling; all other cases are included in this one

The orthodoxy of this position is such that it is rarely spelt out; for example, the manuals for the widespread SDC software tau-Argus and mu-Argus make no justification for this very specific world-view. Hundepool et al (2010) identify ‘Analysing data needs’ as one of the stages of anonymisation (albeit *after* identification of problematic variables and the data owner’s interests), but no further reference to this is made.

The SDC literature echoes the public-choice literature in its assumption of the objectivity and rationality of decisions being made. Risk and utility are measurable characteristics of the data; users and intruders evaluate information in the same way as data owners; user and intruder preferences and incentives are independent of the data protection measures taken.

The underlying assumptions are occasionally criticised. For example, Skinner (2012) noted that ‘worst-case’ is a subjective decision, driven by the concerns of the data owner. Hafner et al (2015) argue that ‘worst case’ in practice is simply the highest-risk mathematically tractable problem which can be fully parameterised. Overall, however, the SDC literature presents a unanimous front in the specification and analysis of confidentiality problems.

As Hafner et al (2015) point out, a common ethos and method is useful for the theoretical development of SDC: it provides consistency, familiarity, reproducibility of results, and a level playing field for comparing new results. The problem is that this perspective reinforces defensive attitudes to risk in the public sector.

Confidentiality challenges presented by administrative data

Ethical and legal challenges

Personal data collected for statistical purposes generally have well-defined legal and ethical gateways for access. In general, re-use of administrative data for research purposes does not. There may be an implicit or explicit duty of confidentiality; for example, in the UK an implicit duty of confidentiality always exists unless explicitly waived. Therefore, statutory gateways (laws) are used to provide a legal mechanism by which research use is enabled.

As it is not possible to define all possible states of the world in advance, laws necessarily leave some room for interpretation. For example, no data protection law provides definitions of 'confidential', 'personal' or 'anonymous' which are unambiguous. Thus when using statutory gateways, the data owners become responsible for the correct ethical and legal interpretation of the law ('data owner' is used as shorthand here for 'someone making decisions on the release of data').

Almost all statutory gateways have the same form:

- The data must only be used for statistical purposes, not for targeting individuals
- There must be a public benefit which outweighs the loss of privacy to the data subject
- Confidentiality measures must be put in place to minimise the loss of privacy

All three topics pose some conceptual problems. Consider the first: suppose research justifies a policy intervention targeted at a particular ethnic group in a particular location; at what point does the intervention become so precise that the data are being used to target individuals?

On the second topic, 'public benefit' is almost impossible to define, let alone measure, as is the potential privacy loss. Public benefit might be very specific (for example in the case of access to UK and US tax records where the tax department itself must benefit) or very general (gateways to data held by the Australian Department of Social Services only require that relevant ministers draw up 'guidelines' for what may be released).

In practice, both topics have precedents and experience, and generally are seen as ad hoc empirical problems rather than general barriers to research access.

Confidentiality also seems to be a 'solved' conceptual issue. There are three mechanisms for data access: fully anonymised data distributed without restriction (Public Use File, or PUF); partially anonymised data distributed to bona fide researchers with restrictions (Scientific Use File, or SUF); and highly detailed data limited to researchers using a secure facility under the data owner's control (Secure Use File, or SecUF). Although a minority of countries view PUFs as the only solution that maintains confidentiality, most OECD countries recognise that PUFs are only part of the answer; partial anonymisation combined with non-statistical measures also provides a sufficient guarantee of ongoing confidentiality. As with the first topics, this is seen often as an empirical problem: data anonymisation techniques have changed very little over the last twenty years, and there are now off-the-shelf tools such as mu-Argus and sdcMicro that implement all the standard techniques.

Ritchie (2014b), Hafner et al (2015) and Hafner et al (2016) strongly criticised the way that such techniques are used. It is clear that confidentiality can be maintained at any desired level through a combination of statistical and non-statistical measures, but the question of whether it should be is unanswered. As an analogy, traditional anonymisation can be likened to building a brick wall; there are different techniques, but the strengths and weaknesses of different methods are well-known. The critics of this approach argue that the builders fail to consider objectively (a) whether a wall is necessary and (b) whether it should be made of brick.

Public relations

The lack of consent can give rise to concerns over accountability: who knows what an anonymous Civil Servant is doing with your personal data? These concerns are likely to be heightened when the interaction with the government department is not voluntary (for example, when providing tax information or registering a change of address).

The Wellcome Trust (2015) reviewed barriers to data sharing for epidemiology, and argued that public attitudes to allowing research use of their administrative data are crucially dependent on three things:

- The way questions about access are asked (eg Haddow et al, 2011)
- The trust in the institution holding the data (GMC, 2007; Wellcome Trust, 2013)
- Anything data-related that might have recently occurred (eg canvassing public opinion in 2007, just after the UK tax department mislaid detailed financial data for ten million households)

The second of these, trust in the institution providing the access to data, seems to be the most important. Health organisations are generally seen as trustworthy; other government departments generally are seen as less trustworthy, although in the Nordic countries overall trust in government to manage data effectively is very high.

The Wellcome Trust (2015, p75) drew some general conclusions: *inter alia*, the general public are broadly

- concerned about the security of their data;
- unable to distinguish between operational and statistical use;
- unable to distinguish between levels of anonymisation;
- happy to change their minds when provided with more information.

This presents a problem for public organisations who have to justify their use of statutory gateways to override privacy concerns. Ethics would suggest keeping the public well-informed is essential to maintaining public confidence. However, practical experience suggests that raising the topic of research data access with the general public runs a high risk of creating active opposition. Accordingly, much of the public dialogue around research data sharing focuses on control of risks (“look at all the things we do to make sure your data isn’t misused”) rather than the benefits of use.

Disclosure risk

Almost all SDC theory is based around the concept of an ‘intruder’ who has access to some additional information (the ‘external dataset’) and uses it in combination with the confidential data to breach data privacy. Identification can be two-way: an intruder may use external information to help identify an individual in the data, and so uncover some information in the confidential data; or the intruder may use something in the confidential data to identify an individual in the external data.

In survey data, this concept of the external dataset is highly problematic. Data collected are de-identified as soon as possible. Data are collected for statistical analysis which may not have external relevance. Exact samples are generally not known, and even when collecting a (sub-) population there always exists some uncertainty over the specific form of inclusion. Finally, the entire data flow is usually managed within the data collecting organisation. Modern studies trying to match government survey data with external data sources have shown the difficulty in this, even in ideal circumstances (eg Evans and Ritchie, 2006; Hafner, 2008; both articles deal with business data which is generally assumed to be very identifiable).

In contrast, with administrative data the external dataset definitely exists, and is probably fully-identified. Anyone who has access to the original data for administrative purpose is now a credible potential intruder. Unlike survey managers, who only need access to identifying information

for the data collection phase, and who have procedures in place to separate identifying and non-identifying information, administrators are likely to need permanent access to full-identified data. The source administrative data may not even be under the control of the department giving research access; it may have been delivered from a completely different organisation.

Finally, the research dataset is potentially more risky because every variable that can have a unique value is a potential identifier. Consider a researcher noting that “the highest weekly earnings in region X are £12,452”. With survey data, this is likely to be a ‘target’ variable: information an intruder would like to discover by finding out who the top earner is. With administrative data, anyone who has access to the source data can immediately identify that person, assuming that only one person in the region has that exact salary.

Confidentiality risks with administrative data: summary

The research use of administrative data presents a number of practical challenges, not least in the dynamic nature of the data: the data owners need to decide when, how and how often a research dataset is to be extracted from the administrative source. Whilst these are important problems and do have some privacy implications, in this section attention was focused on three specific issues which have a direct impact on the way data owners in the public sector make decisions:

- With statutory gateways rather than consent the basis for re-use, those making decisions about data access are required to interpret legal and ethical boundaries, and are responsible for their interpretations
- Without consent, public trust may be both more important and harder to gain
- Allowing researchers to re-use data which is accessible to others produces additional risks of disclosure

Defensive decision-making in the public sector

Default-open versus default-closed

Hafner et al (2015) report the effect of giving data professionals (that is, anyone engaged in the collection, management, dissemination or research use of data) in the public sector the following scenario.

Suppose you are responsible for making data release decisions, and have a mechanism for deciding whether such a release is ‘safe’. There are two ways that you could view your role in deciding on releases:

- a) Do not release data unless the release is shown to be safe
- b) Release data unless the release is shown to be unsafe

The audience is then asked two questions, with the second being posed only after the first has been answered:

Which of these should be your default perspective?

Which of these is your organisation’s default perspective in practice?

Hafner et al (2015) report that the experiment almost always reports the same results: 80%-90% say that (b) is their preferred perspective; the same proportion report that (a) is their organisation’s position.

Ritchie (2014a) argues that this difference in perspective is the key element limiting effective use of research data. Although the two statements are functionally identical (data releases either meet both rules or neither), psychologically they are worlds apart and will almost certainly lead to different outcomes.

Option (a), referred to as default-closed, has complete security but not access and hence no utility; if the security constraints can be relaxed without compromising confidentiality, utility is gained but security lost. Option (b), referred to as default-open, assumes that there is complete access to the data, and hence maximum usefulness; if the data has to be restricted to meet confidentiality criteria, then security is gained but utility is lost. Humans put more weight on losses than gains, and so a decision maker starting from 'no release' would come to a different conclusion to a decision-maker starting from 'no restrictions' (Ritchie, 2014a).

There is no conceptual reason for choosing default-open or default-closed. Both are epistemologically valid, and in the rational, objective world of public choice theory would produce the same recommendations.

In practice, almost all public sector data access is based on a default-closed model, for two reasons already discussed:

- Public servants are likely to be incentivised to avoid failure as a personal objective
- The SDC literature universally uses the language of the default-closed model

In these circumstances, the rational choice of a public servant who is not an expert in SDC is to choose the default-closed option. The purchasing-industry catchphrase "nobody ever got fired for buying IBM" has a direct, if windier, counterpart: "nobody ever got fired for saying that they thought the case for releasing data had not been made adequately to overcome the confidentiality risks."

Use of evidence

A second reason for the defensive stance taken by public-sector decision-makers is the lack of evidence in building release cases. As noted, the SDC literature focuses almost exclusively on worst-cases based on hypothetical scenarios, and rarely discusses the likelihood of different scenarios. From the perspective of a decision-maker without expert knowledge, this makes perfect sense. If data is protected against the worst case, then it is also protected against all other cases, and so there is no need to worry about the likelihood of other scenarios.

The problem with this is twofold. From a theoretical perspective, worst-case scenario planning has been criticised (Skinner, 2012; Hafner et al, 2015) as being convenient rather than objective or methodologically robust. More importantly, there is no evidence whatsoever to support these worst-case scenarios (Hafner et al, 2015; DSS, 2016); realised risks are almost always associated with non-malicious human behaviour (Desai and Ritchie, 2010; Ritchie and Welpton, 2014).

Lack of identification of benefits

Ritchie and Welpton (2012) argue that a key factor influencing public servants' willingness to focus on the downside of any data release decision is the frequent lack of any perceivable upside; The National Research Council (2014, p118), for example, notes that the US HIPAA guidelines explicitly exclude consideration of data utility when classifying datasets.

While advocates of greater research use of data claim that research use directly benefits the organisation through researcher feedback, in practice such benefits are small and rarely achieved (DSS, 2016), particularly relative to the perceived risk. This is the case even in organisations where the research community has had a significant input into data access arrangements, such the data laboratories run by the US Census Bureau, the UK Office for National Statistics, or the German Labour Department. As Ritchie (2004) notes, government data is typically collected to produce aggregate statistics; the insights into quality offered by external users of the data are minor, compared to the effort expended in getting key statistics right.

Crime and punishment in a rational world

Any confidentialisation of the data which does not lead to complete anonymisation must have some amount of residual risk. In a rational world, this is straightforward to manage: users are informed of their duty to maintain the confidentiality of the data and the necessary procedures; any resulting breach of confidentiality is the result of failure on the part of the user. Penalties ensure that the user acts appropriately; if the possibility of a breach is high, penalties should be raised until no rational user would breach confidentiality.

This appeals to the defensive mindset: it absolves the data owner of responsibility for breaches, making clear that the user is at fault. However, it assumes that researchers come to the same 'rational' decision as the data owners would expect; but Desai and Ritchie (2010) point out that, there is very little evidence that users know or care about the data owners' interests. Moreover, training programmes which emphasise researchers' guilt are likely to be counter-productive (Desai and Ritchie, 2010; Jackson et al, 2012).

Impact on the release and re-use of administrative data

Those outside the SDC literature who have looked at data sharing (eg Research Councils UK, 2008; National Research Council, 2014, especially ch.5; Australian Government, 2015b; Academy of Social Sciences, 2016) have typically concluded that government re-use of data is overly defensive. It fails to take account of the full public benefit, is too concerned with hypothetical risks, places too much reliance on anonymisation without considering other risk control measures, and assumes that users are rational individuals who nevertheless care about data protection.

The re-use of administrative data exacerbates these problems. The lack of consent, and the consequent need for the data collector justify the re-use; the potential public relations disaster of re-using data supplied, possibly involuntarily; and the existence of matching data, all contribute to heighten the difficulty of allowing re-use of administrative data. In other words, given that government attitudes to data release are generally defensive, exploiting the research value of administrative data is likely to face even more objections than re-using datasets collected for statistical purposes.

The modern management approach

Since the mid 2000s, an alternative approach to microdata access has been developed. This model is sometimes referred to as the 'modern' approach (to distinguish it from the 'traditional' defensive approach), but here we adopt the acronym EDRU ('evidence-based, default-open, risk-managed, user-centred'; DSS, 2016). The model was originally developed to manage use of SecUFs through remote-access facilities, but since 2012 has been extended to all areas of data access.

The 'evidence-based, default-open, risk-managed, user-centred' ethos

The EDRU ethos directly challenges the problems identified above with data access; as the name implies, it has four components.

Central to the EDRU process is the use of evidence to assess risk, utility, and the desirability of options. If the data owner wishes to use worst-case planning, that worst-case should be based upon reasonable (even if low) probability, rather than hypothetical possibility. For example, Gregory (2014) uses expert panels to test anonymisation processes, rather than relying on statistical analysis.

The model is default-open: data are assumed to be available, and the only consideration is whether restrictions need to be placed on access to maintain confidentiality. Note that the

restrictions are placed on access, not on data; non-statistical solutions (such as licensing, or restricted access facilities) may be better solutions than reducing the value of data.

In the EDRU perspective risk is explicitly acknowledged to be subjective. This encourages joint acceptance of risk: if there is no 'right' answer, then the views of all interested parties have some validity. A unilateral decision by the data owner on the appropriate level of risk leaves the data owner exposed to criticism for being under- or over-cautious. Engaging with the users and other stakeholders avoids this criticism. Engaging the users in decisions about risk also helps to remind users that they are part of the solution, not just passive recipients of instructions.

The model is user-centred, recognising that the biggest risk in data access is the human factor (Ritchie and Welpton, 2014), but that a positive relationship can bring substantial security gains. This in turn leads to cost savings: designing an access system for users who can be trusted is much simpler than designing one for users who can't. A clear understanding of user needs helps to set objectives (Hafner et al, 2016). Finally, a positive relationship with users should increase the likelihood of the data owner reaping the direct benefits of data release, such as methodological and analytical feedback from engaged experts.

Implementation

The EDRU ethos leads to a number of important consequences for implementation.

First, personal and institutional relationships become much more important. The traditional approach to data management does not discuss attitudes as the process is seen as objective and the preferences of actors do not change. In contrast, in EDRU the default-open attitude is seen as the enabler for all the other elements, and therefore getting early agreement on the principles underlying the data solution is essential. For example, the Australian Department of Social Services drew up an attitudinal policy statement before carrying out a strategic review of options (DSS, 2016).

Second, relations with users come to the fore. A key factor is how users perceive that they are treated by the data owners (Desai and Ritchie, 2010). Users are an integral part of the design process: not just for identifying useful data, but for their contribution to the security model. Even if the solution is a PUF (and so the data owner intends no control over how the data is used, and by whom), the active appreciation of users encourages trust. For SUFs or SecUFs, the training of researchers becomes a key consideration. Unlike the static rational world of the traditional model, the EDRU model allows the preferences and incentives of all parties to change and develop.

Third, the consideration of non-statistical solutions to confidentiality challenges means that perturbation or restriction on the data becomes much less important. The popular 'Five Safes' framework for designing strategies (Desai et al, 2015) explicitly makes data limitation the residual, something to be done only when all the options relating to the users, the use and the environment have been considered and found wanting.

Finally, the emphasis on unavoidable human error and the certainty of mistakes encourages the development of a support framework involving all interested parties. Suppose Department X releases data to support research of interest to Department Y. Department Y needs to acknowledge both the value of the research and the risk taken by Department X so that, when something goes wrong, a fair assessment can be made of whether the risk was justified in the context of the overall benefit to society. A secondary consequence of this relationship management is that it can lower the perceived boundaries between producers and consumers of research data, increasing the sense of a research community.

Impact

The EDRU model was developed for access to SecUFs through remote access facilities, where the gains for non-statistical approaches to confidentiality are largest. Hafner et al (2015) note that the first EDRU-compliant facility, run by the UK Office for National Statistics, was able to support more

researchers with more data at costs significantly below those of similar facilities; yet it also achieved a class-leading level of security. The EDRU principles have been and are now generally acknowledged as the best practice for these sorts of facilities (for example Sullivan, 2011, or Webster, 2015).

Less progress has been made in areas where less control is possible. Hafner et al (2016) describe the first SUFs built on EDRU principles. The data were previously confidentialised using the defensive approach: the SUFs was widely perceived to have teaching value only, and Hafner et al (2016) show that identified security risks were not addressed. Applying the EDRU approach led to minimal perturbation of the data while still addressing the security risks omitted in the previous solution.

Gregory (2014) independently conceived the EDRU principles, and applied them to the creation of PUF for administrative records of electricity use. The assumption was made that the data should be available; user need and acceptable restrictions were identified through active consultation; and the disclosure risks were tested not by theoretical models, but by setting up panels of users who were offered prizes for successful re-identification (none were successful). This approach received substantial positive support from the user community, and is seen as potential model for wider adoption across the UK government.

The most radical developments are currently underway in Australia. A formal statement on data strategy was published (Australian Government, 2015b) to the effect that data access would be default open, consultative, and seeking to maximise the benefits of re-use and linking. Webster (2015) and DSS (2016) describe complete data strategies for the Australian Bureau of Statistics and Australian Department of Social Services, respectively, based on the EDRU ethos.

Challenges

If EDRU is a superior and cheaper approach to data access why is it not universally adopted? There are three possible reasons for this.

First, EDRU directly challenges a data access paradigm which has developed over decades. Changing such modes of thought is a slow process and typically requires extensive negotiation. The preponderance of the traditional method also means that EDRU advocates are in a very small minority, and struggle to make an impact. It is noticeable that EDRU has had the most impact on the remote-access community which is relatively small and which shares ideas on a regular basis.

Second, EDRU is relatively new and therefore has few precedents. For decision-makers wanting reassurance that a defensible decision has been made, precedents are invaluable. In the case of administrative data this is even more so: with the exception of Gregory (2014) and DSS(2016), EDRU has largely been applied to survey data and so the risk posed by external matching datasets has not been demonstrable assessed.

Third, EDRU is time-consuming, require investment in relations between data owners, users, and other interested parties. Part of this relationship building also involves getting agreement to share risk, which may not be appealing to those who are just used to seeing the policy benefits of research, for example.

One challenge which has not been widely addressed is public perceptions. EDRU assumes that mistakes will be made; therefore, plans should be put in place to minimise impact. Such plans should include dealing with public criticism of risky behaviour. In theory a sensible way to avoid this is to pro-actively manage public expectations; for example, the UK Administrative Data Research Network has produced publicity material designed for the general public, and has privacy campaigners on its advisory board. However, there are also examples of badly designed public information campaigns which have led to a public backlash against the proposed development; O'Dowd (2014) discusses the ill-fated care.data health data linkage programme in the UK. At present, there is no consensus over whether greater public engagement increases or lowers the acceptability of projects.

Conclusion

Governments collect very large amounts of data. There is enormous value in making these data available for research use, and the evidence of the last half-century shows that the re-use of these data by academics is valuable for public policy and extremely low risk. Nevertheless, decision-makers typically adopt a defensive attitude to the release of data, as a result of personal and institutional concerns. Administrative data presents more confidentiality challenges compared to data collected solely for statistical purposes, and so is even less likely to be released for research.

A small but growing number of data access solutions take a different path, leading to more access to better quality data at a lower cost and with lower risk. Central to the model is the change in attitudes to risk, relationships, and reference points. This approach has been usefully implemented in restricted access facilities over the last data, and is increasingly being adopted in a wider range of data access solutions.

If this is so successful, why is adoption slow? The answer lies in the change in mindset needed to make use of these new approaches: risk is negotiated rather than estimated, and cannot be transferred to the user; user objectives determine outcomes instead of data owner objectives; and the assessment of hard evidence replaces hypothetical speculation on worst cases. All of these are difficult for data owners to accept, particularly when this is a very minority view, there are few specialists, and the available precedents to copy are few and far between.

Nevertheless, experience shows that it is possible to change the defensive mindset, leading to gains for all parties. In answer to the question “can we make effective use of administrative data for social research?”, the answer is “yes; if we really want to.”

References

- Academy of Social Sciences (2016) *Response to Cabinet Office Consultation on Better Use of Data in Government* <https://www.acss.org.uk/wp-content/uploads/2016/04/Academy-of-Social-Sciences-response-to-Cabinet-Office-consultation-on-Better-Use-of-Data-in-Government.pdf>
- Adams V. (2013) *Markets of sorrow, labours of faith*. Durham: Duke University Press.
- Andrews R., Boyne G.A. and Walker R.M. (2011) “Dimensions of Publicness and Organizational Performance: A Review of the Evidence” *J Public Adm Res Theory* v21:3 ppi301-i319. doi: 10.1093/jopart/mur026
- Australian Government (2015a) *Public Sector Data Management*. Department of the Prime Minister and Cabinet. https://www.dpmc.gov.au/sites/default/files/publications/public_sector_data_mgt_project.pdf
- Australian Government (2015b) *Public Data Policy Statement*. Department of the Prime Minister and Cabinet. https://www.dpmc.gov.au/sites/default/files/publications/aust_govt_public_data_policy_statement_1.pdf
- Bhatta G. (2003) “Don't just do something, stand there! Revisiting the issue of risks in innovation in the public sector”. *The Innovation Journal*.
- Boyne, G.A. (1998) “Bureaucratic Theory Meets Reality: Public Choice and Service Contracting in U.S. Local Government”. *Public Administration Review*. v58:6 pp474-484
- Cabantous L., Hilton D., Kunreuther H. and Michel-Kerjan E. (2012) “Is imprecise knowledge better than conflicting expertise? Evidence from insurers’ decisions in the United States” *J. Risk and Uncertainty* v42 pp211–232
- Carpenter D. and Krause G. (2015) “Transactional Authority and Bureaucratic Politics”. *Journal of Public Administration Research and Theory* v25 pp5-25

- Desai T. and Ritchie F. (2010) "Effective researcher management", in *Work session on statistical data confidentiality 2009*; Eurostat.
- Desai T., Ritchie F., and Welpton R. (2016) *The Five Safes: designing data access for research*. Working papers in Economics no. 1601, University of the West of England, Bristol. January
- DSS(2016) *Data access strategy: final report*. Australian Department of Social Services, June.
- Eckerd A. (2014) "Risk Management and Risk Avoidance in Agency Decision Making". *Public Administration Review*, v74:5 pp616-629
- Evans P. and Ritchie F. (2009) *UK Company Statistics Reconciliation Project*, Department of Business Enterprise and Regulatory Reform; URN 09/599
- Feeney M.K. and DeHart-Davis L. (2009) "Bureaucracy and public employee behaviour: a case of local government". *Review of Public Personnel Administration* v29:4 pp 311-326
- Gazley B. (2014) "Good Governance Practices in Professional Associations for Public Employees: Evidence of a Public Service Ethos?". *Public Administration Review*, v74:6 pp736–747
- GMC (2007) *Public and Professional attitudes to privacy of healthcare data: A Survey of the Literature*. General Medical Council.
- Gollier C., Hammit J.K. and Treich N. (2013) "Risk and choice: a research saga". *J. Risk and Uncertainty* v47 pp129-148
- Gregory M. (2014) "DECC's National Energy Efficiency Data-Framework – Anonymised dataset". <http://www.turing-gateway.cam.ac.uk/documents/Gregory.pptx>
- Haddow G., Bruce A., Sathanandam S. & Wyatt J. (2011) "'Nothing is really safe': a focus group study on the processes of anonymizing and sharing of health data for research purposes" *Journal of Evaluation in Clinical Practice*, v17:6, pages 1140–1146
- Hafner, H.-P. (2008). *Die Qualität der Angriffsdatenbank für die Matchingexperimente mit den Daten des KSE-Panels 1999 – 2002*. Mimeo. IAB
- Hafner H.-P., Ritchie F. and Lenz R. (2016) "User-centred threat identification for anonymized microdata". Working papers in Economics no. 1503, University of the West of England, Bristol. March 2014. Forthcoming in JoS
- Hafner H-P., Lenz R., Ritchie F., and Welpton R. (2015) "Evidence-based, context-sensitive, user-centred, risk-managed SDC planning: designing data access solutions for scientific use", in *UNECE/Eurostat Worksession on Statistical Data Confidentiality 2015*, Helsinki.
- Hartley J. (2005) "Innovation in Governance and Public Services: Past and Present", *Public Money & Management*. v25:1 pp27-43
- House of Lords (2006) *Government Policy on the Management of Risk*. Select Committee on Economics Affairs, 5th Report of Session 2005-6 Volume 1:Report
- Hundepool, A., Domingo-Ferrer, J., Franconi, L., Giessing, S., Lenz, R., Longhurst, J., Schulte Nordholt, E., Seri, G. and De Wolf, P. (2010). *Handbook on Statistical Disclosure Control*, ESSNet SDC http://neon.vb.cbs.nl/casc/.SDC_Handbook.pdf
- Jackson, J., Hough, M., Bradford, B., Hohl, K. and Kuha, J. (2012) Policing by consent: UK evidence on legitimate power and influence. *ESS Country Specific Topline Results Series Issue 1*. London: City University
- Kahneman D. (2012) *Thinking, fast and slow*. London: Penguin Books.
- Kahneman D., Knetsch J. and Thaler R. (1991) "Anomalies: the endowment effect, loss aversion and status quo bias", *Journal of Economic Perspectives* v5:1 pp193-206.
- Kahneman, D.; Tversky, A. (1979). "Prospect theory: An analysis of decisions under risk". *Econometrica* **47** (2): 263–291.
- Kjeldsen A.M. and Jacobsen C.B. (2013) "Public Service Motivation and Employment Sector: Attraction or Socialization?". *Journal of Public Administration Research and Theory* v23 pp899-926
- Linde J. and Sonnemans J. (2011) "Social comparison and risk choices". *J. Risk and Uncertainty*. V44 pp45-72

- Lofstedt R.E. (2004) "The swing of the regulatory pendulum i Europe: from precautionary principle to (regulatory) impact analysis". *J. Risk and Uncertainty* v28:3 pp237-260
- Mazzucato M. (2014) *The Entrepreneurial State*. London: Anthem Press
- Moore M.H. (2010) "Break-Through Innovations and Continuous Improvement: Two Different Models of Innovative Processes in the Public Sector". *Public Money & Management* January 2005 v44
- National Research Council (2014). *Proposed Revisions to the Common Rule for the Protection of Human Subjects in the Behavioral and Social Sciences*. National Research Council, Washington, DC: The National Academies Press.
- O' Dowd A., (2014) "Patients could withhold information from GPs because of confusion over care.data scheme, doctors warn", *British Medical Journal*, February
- Office of the Auditor-General of Canada (1998) *Innovation in the Federal Government: the Risk Not Taken*. Public Policy Forum Discussion Paper
- Research Councils UK (2008) Response to the ICO consultation on anonymisation. <http://www.rcuk.ac.uk/RCUK-prod/assets/documents/submissions/200812RCUKResponseICOconsultationanonymisation.pdf>
- Ritchie F. (2004) "Business Data Linking – Recent UK experience", *Austrian Journal of Statistics* v33:1-2 pp89-97
- Ritchie F. (2014a) "Access to sensitive data: satisfying objectives, not constraints". *Journal of Official Statistics* v30:3 pp533-54
- Ritchie F. (2014b) "Resistance to change in government: risk, inertia and incentives". Working papers in Economics no. 1412, University of the West of England, Bristol. December
- Ritchie F. and Welpton R. (2012) "Data access as a public good" in *Work session on statistical data confidentiality 2011*, UNECE/Eurostat.
- Ritchie F. and Welpton R. (2014) "Addressing the human factor". Working papers in Economics no. 1413, University of the West of England, Bristol. December
- Schillemans, T. and Busuioc M. (2015) "Predicting Public Sector Accountability: From Agency Drift to Forum Drift". *Journal of Public Administration Research and Theory* v25 pp191-215
- Somers M. (2008) *Genealogies of Citizenship: markets, statelessness and the right to have rights*. Cambridge: CUP.
- Skinner C. (2012) Statistical Disclosure Risk: Separating Potential and Harm, *Int. Stat. Rev.* v80:3 pp349–368
- Sullivan F. (2011) The Scottish Health Informatics Programme, presentation to Health Statistics User Group. <http://www.rss.org.uk/uploadedfiles/userfiles/files/Frank-Sullivan-linkage.ppt>
- Sutherland S. (1992) *Irrationality*. Pinter and Martin; London
- Viscusi W.K., Phillips O.R. and Kroll S. (2011) "Risky investment decisions: how are individuals influenced by their groups?" *J. Risk and Uncertainty* v43 pp81-106
- Webster A. (2015), "Expanding access to public data: background paper", ASSA Policy Roundtable
- Wellcome Trust (2013) *Summary Report of Qualitative Research into Public Attitudes to Personal Data and Linking Personal Data*. July.
- Wellcome Trust (2015) *Enabling data linkage to maximise the value of public health research data*. London, March
- Yang K. and Holzer M. (2006) "The Performance–Trust Link: Implications for Performance Measurement" *Public Administration Review*. v66:1 pp114–126

Recent UWE Economics Papers

See <http://www1.uwe.ac.uk/bl/research/bristoleconomicanalysis> for a full list.

2016

- 1607 **Can a change in attitudes improve effective access to administrative data for research?**
Felix Ritchie
- 1606 **Application of ethical concerns for the natural environment into business design: a novel business model framework**
Peter Bradley, Glenn Parry and Nicholas O'Regan
- 1605 **Refining the application of the FLQ Formula for estimating regional input coefficients: an empirical study for South Korean regions**
Anthony T. Flegg and Timo Tohmo
- 1604 **Higher education in Uzbekistan: reforms and the changing landscape since independence**
Kobil Ruziev and Davron Rustamov
- 1603 **Circular economy**
Peter Bradley
- 1602 **Do shadow banks create money? 'Financialisation' and the monetary circuit**
Jo Michell
- 1601 **Five Safes: designing data access for research**
Tanvi Desai, Felix Ritchie and Richard Welpton

2015

- 1509 **Debt cycles, instability and fiscal rules: a Godley-Minsky model**
Yannis Dafermos
- 1508 **Evaluating the FLQ and AFLQ formulae for estimating regional input coefficients: empirical evidence for the province of Córdoba, Argentina**
Anthony T. Flegg, Leonardo J. Mastronardi and Carlos A. Romero
- 1507 **Effects of preferential trade agreements in the presence of zero trade flows: the cases of China and India**
Rahul Sen, Sadhana Srivastava and Don J Webber
- 1506 **Using CHARM to adjust for cross-hauling: the case of the Province of Hubei, China**
Anthony T. Flegg, Yongming Huang and Timo Tohmo
- 1505 **University entrepreneurship education experiences: enhancing the entrepreneurial ecosystems in a UK city-region**
Fumi Kitagawa, Don J. Webber, Anthony Plumridge and Susan Robertson
- 1504 **Can indeterminacy and self-fulfilling expectations help explain international business cycles?**
Stephen McKnight and Laura Povoledo
- 1503 **User-focused threat identification for anonymised microdata**
Hans-Peter Hafner, Felix Ritchie and Rainer Lenz
- 1502 **Reflections on the one-minute paper**
Damian Whittard
- 1501 **Principles- versus rules-based output statistical disclosure control in remote access environments**
Felix Ritchie and Mark Elliot

2014

- 1413 **Addressing the human factor in data access: incentive compatibility, legitimacy and cost-effectiveness in public data resources**
Felix Ritchie and Richard Welpton
- 1412 **Resistance to change in government: risk, inertia and incentives**
Felix Ritchie
- 1411 **Emigration, remittances and corruption experience of those staying behind**
Artjoms Ivlevs and Roswitha M. King
- 1410 **Operationalising ‘safe statistics’: the case of linear regression**
Felix Ritchie
- 1409 **Is temporary employment a cause or consequence of poor mental health?**
Chris Dawson, Michail Veliziotis, Gail Pacheco and Don J Webber
- 1408 **Regional productivity in a multi-speed Europe**
Don J. Webber, Min Hua Jen and Eoin O’Leary
- 1407 **Assimilation of the migrant work ethic**
Chris Dawson, Michail Veliziotis, Benjamin Hopkins
- 1406 **Empirical evidence on the use of the FLQ formula for regionalizing national input-output tables: the case of the Province of Córdoba, Argentina**
Anthony T. Flegg, Leonardo J. Mastronardi and Carlos A. Romero
- 1405 **Can the one minute paper breathe life back into the economics lecture?**
Damian Whittard
- 1404 **The role of social norms in incentivising energy reduction in organisations**
Peter Bradley, Matthew Leach and Shane Fudge
- 1403 **How do knowledge brokers work? The case of WERS**
Hilary Drew, Felix Ritchie and Anna King
- 1402 **Happy moves? Assessing the impact of subjective well-being on the emigration decision**
Artjoms Ivlevs
- 1401 **Communist party membership and bribe paying in transitional economies**
Timothy Hinks and Artjoms Ivlevs

2013

- 1315 **Global economic crisis and corruption experience: Evidence from transition economies**
Artjoms Ivlevs and Timothy Hinks
- 1314 **A two-state Markov-switching distinctive conditional variance application for tanker freight returns**
Wessam Abouarghoub, Iris Biefang-Frisancho Mariscal and Peter Howells
- 1313 **Measuring the level of risk exposure in tanker shipping freight markets**
Wessam Abouarghoub and Iris Biefang-Frisancho Mariscal
- 1312 **Modelling the sectoral allocation of labour in open economy models**
Laura Povoledo