



**Risk identification and management for the research use of
government administrative data**

Journal:	<i>Records Management Journal</i>
Manuscript ID	RMJ-03-2019-0016.R1
Manuscript Type:	Research Article
Keywords:	Government administrative data, Information governance, Risk Management, Sensitive data, Open data

SCHOLARONE™
Manuscripts

Risk identification and management for the research use of government administrative data

1. Introduction

Government administrative data have enormous potential for public and individual benefit, through improved educational and health services to citizens, medical research, environmental and climate interventions and exploitation of scarce energy resources. Administrative data is 'collected primarily for administrative (not research) purposes by government departments and other organisations for the purposes of registration, transaction and record keeping, usually during the delivery of a service', such as health care, vehicle licensing, tax and social security systems (<https://esrc.ukri.org/funding/guidance-for-applicants/research-ethics/useful-resources/key-terms-glossary/>). Administrative data are usually distinguished from data collected for statistical use, such as the census. Unlike administrative records, they do not provide evidence of activities and generally lack metadata and context relating to provenance. Administrative data, unlike open data, are not routinely made open or accessible, but only on request to named researchers for specified research projects, through research access protocols which often take months to negotiate and are subject to significant constraints around re-use, such as the use of safe havens. Researchers seldom make use of freedom of information or access to information protocols to access such data, as they need specific datasets and particular levels of granularity and an ability to re-process data, which are not made generally available. This article draws on research undertaken by the authors as part of the Administrative Data Research Centre in England (ADRC-E). The research examined perspectives on the sharing, linking and re-use (secondary use) of administrative data in England, viewed through three analytical themes: trust, consent and risk. The article presents the analysis of the identification and management of risk in the research use of government administrative data and presents a risk framework.

1
2
3 Risk management (that is, coordinated activities which allow organisations to control risks,
4 Lemieux, 2010) enables us to think about the balance between risk and benefit for the public
5 good and for other stakeholders. Mitigating activities or management mechanisms employed
6 to control the identified risks depend on the resources available to implement the options, on
7 the risk appetite or tolerance of the community and on the cost and likely effectiveness of the
8 mitigation. Mitigation and risk do not work in isolation and should be viewed holistically,
9 keeping the whole information infrastructure in balance, across the administrative data
10 system and between many different stakeholders.

11
12 This article seeks to establish a clearer picture of risk with regard to government
13 administrative data in England. It identifies and categorises the risks arising from the research
14 use of government administrative data. It identifies mitigating risk management activities,
15 linked to five key stakeholder communities, and discusses the locus of responsibility for risk
16 management actions. The identification of the risks and of mitigation strategies are derived
17 from the viewpoints of the interviewees and associated documentation and therefore reflect
18 their lived experience. The five stakeholder groups identified from the data are: individual
19 researchers; employers of researchers; wider research community; data creators and
20 providers; and data subjects and the broader public. The main sections of the article,
21 following the methodology and research context, set out the seven identified types of risk
22 events in the research use of administrative data, present a stakeholder mapping of the
23 communities in this research affected by the risks, and discuss the findings related to
24 managing and mitigating the risks identified. The conclusion presents the elements of a new
25 risk framework to inform future actions by the government data community and enable
26 researchers to exploit the power of administrative data for the public good.

26 **2. Methodology and research context**

27 **2.1 Methodology**

1
2
3 Between 2014 and 2017, the researchers conducted four case studies on government
4 administrative data for education, transport, energy and health. The purpose of the research
5 was to examine stakeholder perspectives about the sharing, linking and re-use (secondary
6 use) of administrative data. Following a scoping study, the qualitative research undertook 44
7 semi-structured interviews, plus one focus group, supported by documentary analysis and a
8 literature review. The secondary use of government administrative data by academic
9 researchers was the core of each case study. The research data was enriched by interviews
10 with government bodies as data providers, regulatory bodies, research funders and lobby
11 groups. In the education case study, we also interviewed data subjects. In spite of extensive
12 discussions with transport data providers, none were willing to be interviewed for the study.
13 The analysis is limited by the stakeholder views and may not present a complete picture: for
14 example, we did not interview many records and information managers in the research, and
15 we did focus on academic researchers. The choice of the four case studies allowed for some
16 exploration of different disciplinary perspectives (health, education, transport, energy) on the
17 same topic (government administrative data). The two years of the study did not allow for a
18 complete study of all classes of administrative data from all government providers, but by
19 considering four disciplinary perspectives, we were able to explore differences and avoid
20 mono-disciplinary approaches and assumptions. All interviewees were anonymised and
21 extracts are referenced by a code (eg. A10). Table 1 provides a summary of the research data
22 collection. The case studies, the key datasets accessed by academic researchers and the
23 interview protocol are explained in [anon] 2018, pp4-7. As explained there, interview
24 transcripts were thematically coded line-by-line, assigning a coding label to each component
25 and refining the codes into themes derived inductively from the data, in an iterative process
26 of analysis and assigning meaning. Most of the coding was undertaken by two of the authors,
27 reviewed by a third, and the analysis was discussed by the whole group. Three themes, trust,
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

risk and consent, framed the data analysis. The findings presented here arise from the coding analysis of the research data, concentrating on our analysis of the identification and management of risk in the research use of government administrative data: articles published elsewhere report on trust (anon, 2017) and consent (anon, 2018).

Case study	Number of interviewees	Types of interviewees	Coding range	Dates of data collection
Scoping study	5	Academic researchers	A1-A5	06/2014-07/2014
Education data	12, plus 4 in a Focus Group	Students (data subjects), HEI student data manager, academic researchers, data providers	FG1, A6-A17	12/2014-04/2015
Transport data	9	Academic researchers	A22-A25, A28, A29, A31-A33	09/2015-11/2015
Energy data	5	Academic researchers, data providers	A26, A27, A30, A34, A35	10/2015-01/2016
Health data	18	Academic researchers, data providers, policy advisors, information & data managers, lobby groups	A36-A53	05/2016-10/2016

Table 1: summary of interviewees

2.2 Research context

The identification and management of risk is discussed in the information compliance, governance and records management literature. ISO 31000 (2009b, revision 2017) defines risk as the ‘effect of uncertainty on objectives’, that is a deviation (positive or negative) from the expected outcome. ISO 31000 notes that risk is ‘a combination of the consequences of an event ... and the associated likelihood of occurrence’. In an information security context, AIC (risks which compromise the Availability, Integrity and Confidentiality of information) measures risk impact. ISO PD Guide 73 (2009a) describes a process of risk identification

1
2
3 comprising four elements (risk sources, events, causes and consequences, stakeholder needs).

4
5 Risk analysis includes calculations of the likelihood and potential consequences of a risk and
6
7 its severity or impact. Risk management and mitigation options can then be identified,
8
9 responsibility for action delegated, and costs and effects estimated. Routine hazard or
10
11 operational risks may be accepted and not do great damage. Control or uncertainty risks must
12
13 be constrained within acceptable limits. Opportunity or strategic risks may involve failure or
14
15 loss, but may lead to better outcomes.
16
17

18
19 Doty (2015) adopts Beck's characterisation of risk as a lens through which we seek to make
20
21 sense of an increasingly unpredictable society. Doty asserts the close relationship between
22
23 risk assessment and information, concluding that while risk assessment depends on
24
25 information, it is never possible to have a complete set of information and therefore, it is
26
27 impossible to eliminate risk altogether. Risk management is not always seen as a positive
28
29 activity: Sprehe (2005) warns against using management of risk as an argument in favour of
30
31 records management, describing it as 'an inadequate rationale for enterprise-wide records
32
33 management because it is essentially a defensive strategy'. We would argue for the positive
34
35 benefits of risk management as a proactive, not reactive, strategy, which enables an
36
37 organisation to assess risks, against an agreed risk appetite, and take action. Risk is also
38
39 related to benefit: a risk might be worth taking because of the benefit it will bring.
40
41
42

43
44 As Lemieux (2010) sets out, risks can be categorised in a variety of ways, and different
45
46 disciplines take different views on the best grouping. ISACA (2010) groups IT risks into
47
48 strategic risk; environmental risk; market risk; credit risk; operational risk; and compliance
49
50 risk; whereas ARMA International (2010) divides records and information risks into
51
52 administrative risks, records control risks, legal/regulatory risks, and technology risks. The
53
54 discussion of risk is often in the context of digital systems and cyber-security, which are
55
56 perceived to increase some information risks and therefore warrant further consideration,
57
58
59
60

such as that given by Quigley et al (2015) in their study of rhetorical analysis in cyber-security communications.

Our research is situated in the records management literature and adopts a risk management frame which broadly reflects information governance approaches. The *Records Management Journal* 'Special Issue: Information governance and ethics: information opportunities and challenges in a shifting world' (vol 29: 1-2, 2019) emphasised the interest in information governance in the records field. Several papers illustrate the relationships between information governance and the identification of risks to information security, which can then be addressed by a variety of rules and controls (eg Daneshmandnia, 2019, Xie, S L. 2019). Lemieux (2010) suggests that 'in contrast to varied definitions of risk management across distinct disciplines, the processes associated with its practice have evolved to become relatively standardized': thus a study of risk derived from administrative data management processes might have applicability across records and information management disciplines as well.

3. Results and Discussion

3.1 Identifying risks in the research use of administrative data

The thematic coding of the interview data identified seven different types of risk events in secondary use of administrative data, which can be further analysed for their sources, causes, impact, likelihood and potential consequences. Most of these are caused by unmanaged disclosure, whether accidental or deliberate, although research participants also distinguished risks caused by linkage and misinterpretation which may not entail disclosure. The seven risk events which were identified in the research data are summarised in Table 2 and in Figure 2.

Risk event #	Risk event description
1	<i>Identification of anonymised individuals</i> : leads to data misuse, harm (4, 5). Risk increases with data volume, making data more 'disclosive'. Some research requires identification. Risk appetite varies depending on sensitivity of data (eg health data), granularity of data (eg transport).

2	<i>Data linkage affordances</i> : increases other risks eg identification of subjects (1) or locations. Linkage errors, missing data, mismatches in data granularity and increased complexity may magnify risk and affect data reliability. Linked data can benefit individuals and groups, but low risk appetite may lead to opportunity risk (7).
3	<i>Misinterpretation of data</i> : lack of metadata, poor data documentation, linkage error, and researcher inexperience may lead to misinterpretation of data. Users misunderstand findings or misuse them for other purposes.
4	<i>Malicious misuse of data</i> : deliberate misuse leads to risk of harm (5).
5	<i>Harm to individuals and groups identified in the data</i> : results from deliberate misuse (4) or accidental re-identification, leads to potential psychological, physical, emotional, financial, reputational & other harm.
6	<i>Risk to commercial confidentiality</i> : privatisation of public functions, complexities of data production & ownership leads to data breaches, commercial risks, conflict between public policy benefit and commercial-in-confidence.
7	<i>Opportunity risk of not using data for research</i> : risk appetite varies between data providers, individual data subjects, and researchers, resulting in no agreed risk appetite. Risk aversion (eg legal) results in data not available for research in the public interest.

Table 2: Risk event identification

Risk event 1: identification of anonymised individuals

The most prominent risk is that of identification of individuals represented in the data, including re-identification from de-identified data. ‘De-anonymisation is a risk primarily to the person who is de-anonymised’ (A7) with consequences of other risks occurring, including malicious misuses (eg. identity theft) and harm to the individual. ‘Any use of data has disclosure risk, it’s impossible to reduce that to a zero chance.’ (A12).

Identification risk has long been a serious concern (and constraint) in the research use of educational administrative data, indicating a risk averse appetite by data providers.

‘Respondents were concerned that the level of anonymisation used in the National Pupil Database was not sufficient to prevent personal data about individuals held within the database from being discovered.’ (D1). One interviewee argued that the likelihood of identification varied with the data collected, ‘the risks in administrative data are often lower than in the risks in survey data, just because of the numbers’. (A12)

In several case studies, discussions of the risk of identification referenced the normative

1
2
3 legislative framework of UK data protection, with its content-based definitions of ‘personal’
4 and ‘sensitive’ data. Some participants sought to articulate how these might be understood,
5
6 for instance in an education data context, where date of birth information can link to other
7
8 personal information. However, A10 suggested that educational ‘attainment data, I would say
9
10 it’s not as disclosive as health data ... that’s very, highly sensitive.’ A16 compared accidental
11
12 release of ‘your GCSE [school exam] French result’ with ‘releasing all of your earnings
13
14 data’, describing ‘a different degree of privacy’.

15
16
17 Some suggested that the risk of identification may be re-assessed against shifting attitudes in
18
19 society at large, whereby privacy may no longer be considered an absolute right (and
20
21 consequently the researcher’s obligations may diminish). A15 remarked that ‘the younger
22
23 people are, the less worried they are about privacy, about data privacy... So, you know, lots
24
25 is changing here about public attitudes towards privacy’.

26
27
28 Others pointed out that some uses of administrative data in fact *require* identification of
29
30 individuals in the data. Not doing the research would risk loss of opportunity. A9 reported
31
32 ‘We’re doing a lot of randomised controlled trials. So for that, with the identifying, the whole
33
34 thing is we can identify children so we can track them over time.’

35
36
37 Interviewees perceived health data as highly sensitive and confidential, therefore the risks to
38
39 the data subjects (identification, misuse, misinterpretation) were significant. A51 articulated
40
41 this, alongside an acknowledgement that public perception of the uniquely sensitive nature of
42
43 health data is still a relative unknown.

44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000

Certainly where health data is concerned people are particularly sensitive. ... I think there are a lot of members of the public who have got things in their health record that they don’t want other people to know. Now you may say that is true of people’s financial information, and for some people that may be the case, but I would have thought there is a more general point around health and care data and people’s

1
2
3 medical histories and so on being more sensitive. (A51)
4

5 Perhaps as a result of the access which transport researchers have enjoyed to detailed
6
7 statistical datasets over many years, interviewees expressed the view that use of
8
9 administrative data in transport research carried fewer risks relating to disclosure than in
10
11 other domains. As A24 noted, 'with the transport data it's just oh, you know, oh you usually
12
13 cycle or you usually go by car'. Transport researchers recognised that sensitivity also relates
14
15 to increases in granularity or detail in the data. A25 considered that 'there's so much
16
17 sensitivity about the geo-location because when you get down to that unit of analysis it can
18
19 be easier to identify people', although A28 remarked on the lack of consistency, 'I couldn't
20
21 understand why something would be classified disclosive when something else isn't'.
22
23

24
25 The risk of disclosure, inadvertent or deliberate, increases in line with the volume of data
26
27 about each individual. As A32 said, 'the dataset is anonymised but at the same time there is a
28
29 lot of information in them, so you can narrow down what people are in the sample.' This was
30
31 reflected by data subjects in the focus group: 'I feel more comfortable with the idea of my
32
33 data being used for research at an aggregate level, whereas when you get down to a kind of
34
35 individual level, it kind of makes me feel that I'm losing control of how it's being used'(FG1)
36
37

38 *Risk event 2: data linkage affordances*

39

40
41 First, the likelihood of the risk of re-identification is increased by data linkage. Data subjects
42
43 recognise this: 'Yes, there's the idea of, not just one piece, but once you link all the pieces
44
45 together, you could be re-identified I think.' (FG1). As do researchers: 'As data become more
46
47 and more detailed through linkage, whether it's longitudinal linkage or linkage to other
48
49 datasets, then they become more disclosive' (A15). FG1 suggested that risks to individuals
50
51 increased by data linkage affordances might outweigh the public and private benefits.
52
53

54
55 Estimates of the increased risks through data linkage varied with different types of data, and
56
57 among different groups of interviewees. A12, using educational data, said that adding further
58
59
60

1
2
3 information to ‘datasets which are used by researchers in situ with people's educational
4 achievement, anonymised, ... wouldn't be any additional risk.’ Whereas linkage of transport
5 data increased the likelihood of identifiability of individuals and specific geographic
6 locations or areas, as A23 reported, through air pollution data.
7
8
9
10

11 Issues unique to linkage of health data emerged, including conflation in the minds of data
12 providers and policy makers between different types of secondary uses. A39 identified
13 significant differences in the risks of creating linked population data to support individual
14 patient management including ‘the use of population based linkages to identify individuals
15 who can be targeted [for preventive treatment]’, with linkages to support secondary uses for
16 ‘understanding patterns of service use, risk factors, outcomes in groups. ... public health
17 monitoring, service evaluation, performance management and research.’ (A39). This warrants
18 further exploration, particularly given the proposed models of consent and opt out that frame
19 the issues of privacy, data sharing and secondary use (Caldicott, 2013).
20
21
22
23
24
25
26
27
28
29
30
31
32

33 Secondly, interviewees identified the risk that the largely unknown degree of linkage error or
34 missing data affects the reliability of data and its analysis. Examples included identifying the
35 wrong individual and ‘missed links’ especially in ‘highly mobile populations’ (A39). A14 felt
36 that ‘if we link more datasets that the risk of something going wrong increases’. A12
37 suggested that academics might assume that data processed by a third party in a safe setting
38 was itself reliable, and not look out for missing data.
39
40
41
42
43
44
45
46

47 Opinions varied considerably however on the impact of linkage error as a risk to research
48 integrity. A14 remarked that ‘the way data is being linked and the methods and the reliability
49 is a black box’ and researchers need to maximise linkage rates and account for linkage error
50 better. A32 suggested that linking transport data, ‘because they’re in different forms and in
51 different scales, merging them into a dataset that we can use is a very complex task. ... it is
52 very easy to do the wrong type of merge’ (A32). A35 used energy data where datasets are
53
54
55
56
57
58
59
60

1
2
3 available at different levels of granularity, so the resolution of the linked data is low.
4

5 A third issue was increased complexity, leading to greater likelihood of disclosure and of
6 mismatches in data linkage. With energy data, where many different bodies provide, analyse
7 and process data for the NEED (National Energy Efficiency Data-Framework) dataset, which
8 holds ‘records for approximately 26 million households, and close to 30 million properties in
9 the data-framework. The quantity of data available gives rise to privacy risks.’ (D39). A30
10 noted that energy data was collected monthly, daily, even half hourly, and ‘the data privacy
11 concerns and risks do grow with the granularity of the data’.
12
13
14
15
16
17
18
19
20

21 *Risk event 3: misinterpretation of data*
22

23
24 Researchers identified several risks which might contribute to an increased risk of
25 misinterpretation of the data. Linkage based on population assumptions, described by A25,
26 leads to misinterpretation through unfounded assumptions or inappropriately generalised
27 data. Cavoli et al. (2015) give examples of assumptions made from a national dataset that
28 cannot be applied to a local population or *vice versa* and from datasets which are
29 insufficiently detailed to facilitate investigation of specific issues, such as unusual injuries or
30 effects upon particular minority communities.
31
32
33
34
35
36
37
38
39

40 Misunderstanding of characteristics in administrative data leads to misinterpretation. Cavoli
41 et al. (2015) give the example of cycling casualty rates growing in line with an increasing
42 popularity of cycling. Similarly, A29 noted differences between datasets: ‘if you have car
43 licence ... administrative sources always overstate the number ... if you die, you usually
44 don’t get round to telling the DVLA [Driver and Vehicle Licensing Agency] that you no
45 longer have a driving licence, so you’re counted in their set ... But you’ve dropped off the
46 population set.’
47
48
49
50
51
52
53
54

55 Misinterpretation can arise through variations in data capture, which may not be evident. A11
56 reported, ‘the UCAS [Universities and Colleges Admissions Service] data that I looked at on
57
58
59
60

1
2
3 schools, you find that certain schools code pupils in a certain way, other schools don't'. A34
4 reported the risk of the researcher not understanding the data creation process: 'The problem
5 is if you expect the data to be perfectly clean and suitable.' Misinterpretation (or
6
7
8 misunderstanding) also poses challenges to data providers. A13 commented, 'Sometimes
9
10
11 customers don't necessarily know what they're asking for. And then when they get it, they
12
13
14 don't quite understand what it is that they've been sent.'

15
16
17 A second aspect of the risk of misinterpretation is where the researcher's data analysis itself
18
19
20 is misunderstood or is misused, posing a risk to the researcher's professional reputation. A7,
21
22
23 A9 and A11 all gave examples of this type of risk. A9 worried that while his data analysis
24
25
26 could properly be used to evaluate educational outcomes over time, it could be misused
27
28
29 through ignorance, for instance in teacher performance reviews. While A11 reported not
30
31 publishing some findings for fear of misuse by parents in the school selection process.

32 33 *Risk event 4: malicious misuse of data*

34
35
36 The risk of misuse describes the possibility of researchers or third parties acting with
37
38
39 malicious intent towards an individual or group represented in the data. It results from re-
40
41
42 identification risk, but whereas the risk of disclosure is inadvertent or accidental, misuse is a
43
44
45 deliberate act. A16 gave a lurid example, 'if you are Hungarian and you're an abusive father,
46
47
48 you can then look up the name of every kid who took Hungarian GCSE [school exam], and
49
50
51 ... with a bit of work, you'd be able to find your ex-wife.'

52 53 *Risk event 5: harm to individuals and groups identified in the data*

54
55
56 Both deliberate misuse and accidental re-identification may result in harm to an individual.
57
58
59 Similar to 'damage and distress' in the UK Data Protection Act (s. 10), interviewees
60
61
62 recognised that the risk of causing harm to an individual extends beyond tangible damage
63
64
65 such as physical injury or financial loss to emotional or reputational harm. A15 commented
66
67
68 'by harm we don't just mean physical harm, because that's usually unlikely, it's reputational
69
70

1
2
3 harm or it's psychological harm.'

4
5 Groups or individuals within a particular group potentially risk harm. In education data, 'this
6
7 was a particular concern for respondents from bodies representing groups where the unique
8
9 characteristics of their representatives made them easily identifiable. Several respondents
10
11 were of the opinion that access to the National Pupil Database should remain restricted and/or
12
13 there should be an opt-out mechanism.' (D1). A17 suggested the risk was most likely to
14
15 occur in a population with similar characteristics where you could 'identify those people who
16
17 were different.'

18
19 The risk of prejudicing members of a particular group was noted by A15, 'let's say you come
20
21 to the conclusion that all people of the Muslim faith in Wolverhampton are likely to be
22
23 jihadist ... if you did something like that you would stigmatise a whole group
24
25 inappropriately.' This risk was also recognised by data subjects, such as A17 who said
26
27 'categorising people by their different nationalities or different religions. It ... give[s]
28
29 minorities a disadvantage, the way they're treated.'

30
31
32
33
34
35 *Risk event 6: risk to commercial confidentiality*

36
37 Interviews with energy researchers and data providers raised concerns for the largely
38
39 privatised energy industry about risks to commercial confidentiality. A government data
40
41 provider (A26) commented 'I just think their main worry is the data gets lost, it gets leaked
42
43 ... there are sort of two consequences, it would be embarrassing ... and also some of the
44
45 information might be considered to be commercially sensitive. So their competitors might use
46
47 it to their advantage.' Energy researchers, however, tended to be sceptical about the genuine
48
49 risks to commercial confidentiality. A35 commented that 'the reasons cited by the
50
51 suppliers... is commercial confidentiality, or ... the commercial worth of the data. ... in
52
53 reality all the suppliers know exactly what all the other suppliers are doing.'

54
55
56
57
58
59
60

1
2
3 Commercial risks were also raised by individual businesses being identifiable even in
4 aggregated data owing to geographic location or other distinguishing attributes, as noted by
5 A35, or data variables which might be commercially sensitive, such as property tenure,
6
7
8 A35, or data variables which might be commercially sensitive, such as property tenure,
9
10 ‘whether you own or merely occupy the building is a key piece of business information.’
11

12 *Risk event 7: opportunity risk of not using data for research*

13
14 Data providers varied in their risk appetite. A24 said that transport data was more easily
15 available with ‘less hoops to jump through’ than health data. A28 justified quite an
16 aggressive attitude towards the use of administrative data in transport research, the disclosure
17 risks notwithstanding, saying ‘I think we should always be pushing to use more. I think we
18 should always be testing, testing the boundary.’ A30 identified a conflict of interest between
19 the commercial interests of energy suppliers, who must also comply with government policy
20 aimed at better public outcomes (fuel poverty, energy efficiency, climate change).
21

22 Assessment of risk amongst data providers is, at best, inconsistent, evidently a source of
23 much frustration amongst researchers, including A28, ‘I don’t think they’re consistently
24 applied. ... I think a lot of it just comes down to the person on the day. ... certain people are
25 more pragmatic than others.’ A12 suggested that data reuse ‘is further up the list of priorities
26 in some Departments than others, and that reflects ... the extent to which the senior
27 politicians or policy makers feel that it’s an important issue.’
28

29 A balance between risks of using the data and the opportunity risk of not using it needs to be
30 struck. A31 reported meeting the data provider requirements, ‘the training and then I could
31 use the dataset only on a safe environment, a secure environment, to access the web, their
32 server. But... I abandoned this dataset.’ A24 warned that ‘being risk averse makes research
33 less useful.’
34

35 Risk aversion in broader public opinion is also an opportunity risk. A28 suggested that ‘the
36 biggest risk of course is this sort of public perception issue and, you know, the extent to
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 which the public understand the benefits of the data and also the safeguards that have been
4 put in place to ensure that there's nothing particularly disclosive.' Researcher A27 outlined
5 the opportunity risk if access to key administrative data sources on energy demand and use
6 was withheld, or if data are not provided in appropriate detail. Risks that future policy
7 development and industry-level evaluation may be hindered by restrictions on data access to
8 NEED is also recognised by DECC (Department of Energy and Climate Change) which
9 states 'the development of UK energy policy has required more detailed data to help deliver
10 and monitor reductions in energy use and emissions.' (D38)

11
12 One of the strongest emergent themes in the interviews with health researchers revolved
13 around access to data from NHS Digital following the Partridge Review (2014) and the
14 care.data programme (Digital Health Intelligence, 2014), which for a period of time resulted
15 in a lock down on the processing of data requests. The implementation of new procedures
16 with a stronger scrutiny and control led to prolonged delays to data access, and more access
17 request refusals. Researchers faced an increase in risk of lost opportunity. The interviews
18 provide evidence of how funded projects had to be significantly modified in quality and
19 scope to get around the lack of access to data, such as A40, forced to abandon a project
20 because of delays getting data, we 'simply haven't been able to answer the research questions
21 we set out with'. Coping with multiple data providers also risks access approvals and delays,
22 as A45 reported. Opportunity risk is increased by uncertainty and risk aversion.

23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47 **Figure 1: Risk events summary** [insert]

48
49 Research is put in jeopardy if the interdependencies between stakeholders, types of risk, and
50 risk management approaches are not balanced. Without a holistic understanding of this
51 interdependency, risk mitigation designed to protect the interests of one group may open up
52 more risk for others. Risk management must be viewed at a collective, societal level where
53 these interdependencies are balanced in order to best mitigate the risks for the benefit of all.
54
55
56
57
58
59
60

3.2 Stakeholder mapping

The research helped to identify stakeholder communities affected by the different risks who have a role in mitigation and risk management. As Figure 1 illustrates, we identified in the research data five entities with a stake in the beneficial use of government administrative data.

Figure 2: Stakeholder mapping [insert]

Individual researchers, usually working in an academic setting, seek as much unfettered access to relevant data for research, but want to avoid reputational risks for themselves, their research group, discipline and employer. Secondly, employers, mainly of academic researchers, have legal and ethical interests in information security and privacy, mitigating reputational risk, while enabling researchers to carry out high quality research. Thirdly, the wider research community, of individual researchers, research groups, their employers and representative bodies (eg subject associations and networks like Administrative Data Research Network), seeks to manage the risks of data re-use. The research community influences data provider and researcher behaviour and can act as a broker. The fourth stakeholder group are data creators and providers, mainly central government agencies but also trusted third party processors, commercial providers of data (especially in the energy data field) and regulatory oversight bodies, such as NHS Digital's IGARD (Independent Group Advising on the Release of Data). Data providers should give transparent and fair access through their protocols, to balance the risks associated with providing data to researchers and the opportunity risk if data was not provided. Data subjects and the broader public in whose interest research is carried out (including individuals, households, and businesses) provide the fifth stakeholder community. Once data subjects have provided their data, they cede control over its use. But they express views on governance mechanisms and on the risk appetite of public bodies. Promoting better public understanding of research use of

1
2
3 data and the risks and benefits attached to it, among individual citizens and collectively, will
4
5 strengthen the voice of this stakeholder group.
6

7
8 Risk appetite varies within the five stakeholder communities and between the different
9
10 stakeholders. It changes over time as part of a chain of interrelated consequences. One
11
12 interviewee (A30) compared ‘the risks from the perspective of an energy consumer, a
13
14 householder’ with ‘the perspectives of government departments and commercial
15
16 organisations, energy suppliers, National Grid’. Commenting on risk appetite, A30 said ‘I
17
18 think there is over-estimation of these risks. You know, government partners are extremely
19
20 risk-averse, and they fear, more than anything, the headline in the *Daily Mail*.’ A10 noted the
21
22 interdependent responsibilities for managing data risks dispersed across their research team,
23
24 the university employer’s ICT infrastructure and wider research community. Interviewees
25
26 A24 and A28 also, however, expressed frustration where risk management was excessively
27
28 burdensome given the data risks. ‘I’ve done safeguarding, ...ADRN [Administrative Data
29
30 Research Network] training, safe researcher training... lab induction, and ...background
31
32 checks.’(A28). While A24 noted that whatever was done, ‘the risk is not zero...’. Table 3
33
34 summarises the five stakeholder groups.
35
36
37
38
39

Stakeholder group	Risk focuses	Risk mitigation
Individual researchers	Access to relevant data for research; reputational risk; linked data	Researcher training; credentials; documenting data processing; data immersion; linkage protocols; statistical disclosure controls; compliance with ethical and provider frameworks.
Employers of researchers	Legal and ethical risks; organisational reputation risks; information security	Implement information security (eg safe harbours) & research data management practices; ethical research approval; legal compliance.
Wider research community	Reputational risks; legal and ethical risks.	Peer review and ethical frameworks; sanctions for researcher misbehaviour;

		influencing sector attitudes to data; developing linkage protocols and techniques; supporting researcher training and credentials; good data management practices.
Data creators & providers	High control over access and access conditions to data; balancing access risks v. opportunity risks.	Data supply agreements; data security; detailed data documentation; auditing and monitoring requirements (credentials, training, trusted third parties, information security); fair costs and charges for data access; publishing open data; commissioning research; project specific or time limited access to data.
Data subjects and the broader public	Public benefit v. private harm; anonymization risk;	Governance mechanisms; risk appetite; public engagement.

Table 3: Stakeholder groups, their risk focus and mitigations

3.3 Managing and mitigating the risks in the research use of administrative data

The analysis and coding of the research data revealed 22 risk mitigation actions considered or taken by the interviewees, which we aggregated into four clusters of risk mitigations. These are sectoral expectations of behaviour and professional reputation; data supply and data security protocols; controls of access to data and metadata; and changing attitudes to open government data. They are summarised in Table 4.

Risk mitigation cluster	Risk mitigation actions
<i>1: sectoral expectations of behaviour and professional reputation</i>	<ol style="list-style-type: none"> 1. Ethics and research integrity protocols 2. Sanctions for un-ethical behaviours 3. Researcher credentials 4. Researcher training 5. Peer review 6. Documenting the research process 7. Data immersion
<i>2: data supply and data security protocols</i>	<ol style="list-style-type: none"> 1. Data supply agreements 2. Data documentation 3. Data governance auditing and monitoring 4. Data security protocols

	<ol style="list-style-type: none"> 5. Good data practice 6. Linked data affordances
<i>3: controls of access to data and metadata</i>	<ol style="list-style-type: none"> 1. Limited time, limited data, project specific access 2. Statistical disclosure controls 3. Anonymization and aggregation of data 4. Granularity of accessible data 5. Data sharing 6. Commissioned research 7. Privatisation
<i>4: opening up data</i>	<ol style="list-style-type: none"> 1. Promoting a dialogue of public engagement and public understanding 2. Proactive publication of open data

Table 4: risk mitigation clusters and actions

Risk management 1: sectoral expectations of behaviour and professional reputation

The first group of mitigation actions focuses on ways of ensuring that researcher behaviour complies with the proper expectations of the wider community, as expressed in ethics and research integrity protocols, credentials and training, and peer review. A7 identified a conflict, saying ‘It’s a slightly tricky position of being my own policeman but you know I’ve got professional standards’. Education researcher A16 commented that individuals clearly understood and accepted their obligations. Risk management engendered less discussion in the transport case study, possibly because such data was less sensitive. A24 remarked, ‘It’s crazy to think that anyone who does it will be doing anything malicious with the data.’ While A28 said ‘I thoroughly agree that people will be responsible in the way they approach the analysis of the data and that kind of thing’.

Ethics and research integrity

Some risk mitigation responsibilities rest with the wider research community: local institutional infrastructure and across academia. For example, A15 stated, ‘primarily of course it’s the research community itself which has to satisfy any ethical considerations and requirements, which is done routinely through research ethics committees, or institutional review boards’. Some interviewees saw a separation between the practical management of data by, for example, the university IT infrastructure compliance with data security standards,

1
2
3 from the individual researcher's data use. For instance, while there might be an agreement to
4 destroy the data, the practical execution involved several parties, 'You can't leave it to the
5 researchers' (A14). Ethics processes do not appear to have evolved fully to accommodate the
6 use of unconsented personal data in research, to reflect the changing legal environment, or
7 distinguish between the risks of using aggregate and individual data. A7 noted 'my
8 institution... take[s] the view that aggregate data aren't particularly important, that it's all
9 about individual level data when ethics comes into play'.
10
11
12
13
14
15
16
17
18

19 *Sanctions*

20
21 The consequences of 'deliberately breaching confidentiality' are clear, including criminal
22 prosecution, or action by the data provider against both the researcher and their employer (eg
23 D10). A15 set this out in detail: 'most researchers work in a research institute or a university,
24 if they're getting access to administrative data ... their employer is bound in to a set of
25 conditions about how they must behave with respect to those data, and the penalties for
26 disobeying that or abusing that set of rules and guidelines is fairly extreme. ... a lifetime ban
27 on their receipt of [Research Council] funding,... a ban on funding to their organisation ...
28 the penalties which are now in place are quite draconian, so ... the institution has a very
29 strong incentive to ensure that their employees do not do anything inappropriate.' (A15)
30
31
32
33
34
35
36
37
38
39
40
41

42 *Researcher credentials*

43
44 However, sanctions do not adequately address risks in advance. A12 stated, 'it's better to
45 focus on credibility in the quality of the people who are accessing the data, ... checking who
46 is getting access to data, as much as what data they see.' A credentials-based approach is also
47 recognised by government data providers, for example, in handling of repeat requests for
48 educational data, which receive priority if 'The request is made by an accredited requester
49 who has already demonstrated their credentials' (D10 – National Pupil Database User Guide).
50
51
52
53
54
55
56
57

58 *Researcher training*

1
2
3 Researcher training helps to ensure users properly understand the data and avoid mistakes.

4
5 A12 explained, 'if you were training new researchers, it's in the descriptive statistics, the
6
7 table at the beginning of your paper, that you're most likely to accidentally make a mistake.'

8
9
10 Training is a requirement of data supply agreements, for example access to educational data,
11
12 'shall be restricted to only those persons who have received appropriate training regarding
13
14 data protection and security' (D2). Sometimes, the data provider delegated training to the
15
16 researcher or their employer. For example, 'The Requester [for NPD] shall ensure that each
17
18 proposed Permitted User receives appropriate training regarding data protection and security
19
20 to enable the Requester to comply with principle 7 of the Data Protection Act' (D6).

21
22
23 However, A8 described the increased requirement for researcher training and more stringent
24
25 data security protocols as 'a real balancing act'. 'Thinking, oh, we should put more processes
26
27 in place, ... I think that always seems sensible when you're foreseeing all these problems.

28
29 But then when it's practically implemented ... that will discourage people from using [data]'

30
31
32 Training requirements demanded of researchers seeking access to data seem out of line with
33
34 the requirements for data capture. For example, a data manager responsible for capturing
35
36 student data (A6), in response to a question about training or background checks, responded,
37
38 'No. I've been working on this kind of data for about fifteen years I suppose, but I don't
39
40 remember anything.'

41 42 43 *Peer review*

44
45
46 The traditional apparatus of community-based quality control in academia, peer review, helps
47
48 to manage quality risks and, as A14, remarked, 'I think that's what the peer review process is
49
50 for. I think there'd be a danger if they [government data provider] were taking an interest that
51
52 it would be politically driven.' However, peer review is compromised if administrative
53
54 datasets are not available to reviewers or journal readers. A14 reported that 'a lot of journals
55
56 insist on making all the data you've used in your papers publicly available. So far I've got
57
58
59
60

1
2
3 away with saying that we could arrange for any independent third party to have permission to
4
5 access the data to check our results but we can't make it publicly available. ... But that, that
6
7 is a problem'. A1 identified conflicting priorities, 'the data provider's stipulations trump the
8
9 funder's stipulations. I think it's legitimate as a researcher to say that the work in this paper
10
11 was done on data with restricted access.' A related issue concerns the abilities of the
12
13 reviewers themselves, who, given the narrowness of access channels and novelty of research
14
15 using administrative data, may lack expertise to judge the data analysis. This may be a
16
17 systemic issue for quantitative social science research, as A8 reported 'I'm not sure the peer
18
19 review process would necessarily pick up the appropriate usage of any data, whether it be
20
21 administrative, large scale survey, census data or whatever, for sure.'

22 23 24 25 26 *Documenting the research process*

27
28 A34 argued that solutions lay in documenting the process of research analysis, as much as in
29
30 documenting the data. 'People need to be recording clearly how they created variables, for
31
32 instance, so what was the cleaning process, in our field, the energy field, it's just terrible in
33
34 terms of the lack of documentation, lack of clarity...'. But she noted the paucity of guidance
35
36 available to researchers about documenting data analysis, 'You can have a whole library just
37
38 on how to perform...various statistical computations etc., and then I've managed to find one
39
40 book, one book ... on how to actually document the steps as you go.' Transport and energy
41
42 researchers suggested that quantitative researchers could learn from qualitative research
43
44 methods of analysing secondary data sources, 'upping ...the amount of time that's just spent
45
46 on, you know, quality assurance of data, rather than just only run the stats' (A34).

47 48 49 50 51 *Data immersion*

52
53 A thorough understanding of administrative data is vital to mitigating the risk of
54
55 misinterpretation. As A22 asserted, 'it's your job as a researcher to understand how that [data
56
57 collection] may influence your results'. He gave an example of traffic accident data
58
59
60

1
2
3 (STATS19) which extended well beyond documentation, training, or desk-based processing.
4
5
6 ‘My boss actually rode along with one of them [the police] as sort of like an experience, you
7
8 know, this is how data is collected.’ (A22). He also emphasised the need for the researcher to
9
10 understand data variables, such as standard response categories, in detail. Data providers need
11
12 to make available sufficiently detailed documentation, including changes in data collection
13
14 mechanisms over time, to facilitate data immersion by researchers.
15

16 *Risk management 2: data supply and data security protocols*

17
18 The second group of mitigating activities focuses on protocols which govern and regulate the
19
20 supply and re-use of data.
21
22

23 *Data supply agreements*

24
25 Data providers have varied and detailed requirements embedded in their data supply
26
27 agreements. Some, such as the Department for Education, depend upon the researcher’s
28
29 individual declaration of compliance with the terms. A13 explained: ‘When they received the
30
31 data they signed up to the terms and conditions, and how they were going to be making the
32
33 data available, of outputs that were going to be arising out of that data, and ...they can let us
34
35 see what it is they’re going to be publishing, just so that we’re happy.’
36
37
38

39 *Data documentation*

40
41 Opinions varied as to the quality of the documentation supplied with energy data (A34, A35),
42
43 although A27 suggested certain advantages to raw state data. ‘So they don’t have a big book
44
45 of standards that are necessarily always coming with these datasets. So it does take quite an
46
47 effort, for the most part, to come to grips with all that could be in there, ... we want to know
48
49 how things are being recorded, and we want to be able to see the messiness of it, in order to
50
51 evaluate it’. It may help to understand how data anomalies arise, as A35 reported ‘where we
52
53 do have problems is they step outside their standard codings, and their standard descriptions.
54
55 ... they use that as a means of making the description more specific, but it does make our life
56
57
58
59
60

1
2
3 a bit more difficult because we have to clean it all up’.

4
5
6 *Data governance auditing and monitoring*

7
8 Data governance requirements have increased in recent years, as A15 suggested, ‘they’re
9
10 engaged in much more auditing and policing of what happens to data.’ Data providers have
11
12 established regulatory oversight bodies to advise, such as NHS Digital’s IGARD. Yet
13
14 potential loopholes remain over the auditing of conformity. A12 noted ‘they don’t check it
15
16 directly in the sense of coming down to check your computer, I guess they’re relying on
17
18 another external process to ensure that it’s safe.’ While A11 asked, ‘Do they ever audit? Er
19
20 no, not to my knowledge...they might ask for proof ... but they haven’t checked.’ Data
21
22 providers often rely on researcher good practice, thus reducing the compliance burden, as
23
24 A14 reported, ‘With the NPD you just use an ordinary computer, NPD-HESA [National Pupil
25
26 Database linked with Higher Education Statistics Agency] you just use ... on a secure
27
28 computer that only you can access, but the fact that there’s no controls over... what you take
29
30 away and nobody examines your outputs’. A24 suggested that researchers in turn rely on
31
32 institutional systems, ‘we sort of delegate it to our IT department. All the data are held on
33
34 computers, and ... I haven’t heard of any university datasets being hacked into.’ Delegation
35
36 of governance might result in an increase in risk.
37
38
39
40

41
42 *Data security protocols*

43
44 Data providers contribute to risk management by taking a balanced approach to data security.
45
46 For example, with education data, A12 suggested, ‘the kinds of information that the
47
48 government department is interested in, tends to be, can this data be inadvertently released
49
50 and reveal identity? That’s the overriding concern.’ Department for Education seeks to
51
52 ‘protect the privacy and confidentiality of individuals’ in education data (D10 4.1). They
53
54 warn that ‘NPD data users need to consider the risk of identifying individuals in their
55
56 analyses prior to publication / release’ and that researchers ‘should be on the alert for any rare
57
58
59
60

1
2
3 and justified and unintentional breach (D10 4.3). A12 put emphasis on ‘the security
4 questionnaire... We need to know that the organisation has got the relevant data security
5 base.’
6
7
8

9
10 Researchers accepted data security protocols, as A15 remarked, ‘I think data security is an
11 absolute given. You have to have data security if you’re going to persuade people that what
12 you are doing is not open to any kind of abuse of their personal information.’ ‘So gone are
13 the days of people just leaving data sitting on laptops or leaving laptops on trains, or leaving
14 diskettes, you know, disks and USB sticks and so on around’ (A15). Researchers felt that
15 they had not been the culprits in data breaches, ‘it’s always the government department that’s
16 lost the data, not the academics’ (A1), while A9 remarked ‘education researchers... we’re all
17 now suffering because of some mistakes made by other people.’
18
19
20
21
22
23
24
25
26
27

28 Some expressed concern that the balance between data security and access had tipped too far
29 in favour of security. A12 reported that, ‘When they decided that there would be individuals
30 in each government department responsible for the data, there was sea change in people’s
31 attitudes to releasing the data, and government departments became much, much more
32 careful... But the downside of that is ... never to release the data, that way you take no risk.’
33
34
35
36
37
38
39
40 A9 commented that data breaches ‘make the system harder for everybody.’
41

42 *Good data practice*

43
44 Even where cross-departmental best practice does exist, such as HM Treasury (2015) Aqua
45 Book, which resulted from a review of quality assurance of government analytical models,
46 departmental boundaries (and academic disciplines) limit adoption. Tools, guidance and
47 templates (GB. Department for Business, 2015) targeted at researchers internal to
48 government would also be useful to external researchers using government administrative data,
49 but seem not to be widely adopted: energy researcher A34 was the only interviewee to
50 mention the Aqua Book.
51
52
53
54
55
56
57
58
59
60

Linked data affordances

Linked data affordances can be a source of risk, but also a risk mitigation. For example, linking government datasets to existing cohort studies might expose flaws and biases in the data collection processes. A15 suggested that ‘one of the benefits that arises from giving researchers better access to administrative data is that they start to uncover all sorts of problems in the data, and also the documentation for data is often, you know, severely limited, lacking, or of poor quality. And it causes organisations to think about the way in which they document what they are doing. It also reveals what can happen when administrative processes change ... what data are generated, how they’re generated, how they are preserved, how they are linked through time’ (A15). A14 added, ‘having access to survey data, you can sort of unpick some of the biases in these administrative [data].’

Linking datasets from other providers allows beneficial sharing of research ethics best practice. A12 said, ‘they have research ethics committees, they have a very formal structure, they have a group collectively considering the ethics of the linkage... individuals who have expertise and competence, including lay members of these ethics committees to make judgements about whether linkage is a good idea.’

Too many different parties may, however, lead to unclear lines of accountability, lack of agreed procedures, confusion and delay, ‘if two datasets are linked, there has to be an explicit decision I guess as to who is responsible for the data’ (A12). A27 also identified confusion over responsibility, ‘I’ve been applying to join two different datasets together for a while now, and again the interest is there, everything’s there but actually there is not a clear recognition as to who says yes.’

Risk management 3: controls of access to data and metadata

A third group of risk mitigation actions focuses on controlling access to data and to metadata.

Data providers and third party data processors manage many of these mitigations, which

1
2
3 sometimes restrict or prevent beneficial research access to data.
4

5 *Limited time /limited data/ project specific access*
6

7
8 A common means of controlling researcher access to administrative data is to make it
9
10 available for a time-limited period only, for limited data types, and for a specific, pre-defined
11
12 research project. Data subjects would welcome an explanation of the limits around their data,
13
14 as a student asked, 'In what way have they accessed it and [who] is it given to...' (FG1).

15
16 Researchers were largely accepting of the consequent need to destroy data at the end of a
17
18 project, although they expressed frustration at the impact this has upon their ability to carry
19
20 out follow-up research. A9 said, 'In theory, I've got the data to do it, for another project...
21
22 I'm constrained from doing it. Now I could re-apply, but I don't know how long it would
23
24 take.' An alternative approach is to replace project-specific permissions with a licence for
25
26 accredited persons to use the data. A9 explained, 'you have a licence to have this data, right,
27
28 and within a set of guidelines about anonymity and protection and security, and not asking
29
30 for data which is more sensitive than you need, but then you have a licence to basically play
31
32 with the data then and come up with research questions and research ideas... to find out why
33
34 this happened or what you could do about it'.
35
36
37
38

39
40 Some data providers manage risks by providing only 'a sample of data would be used for a
41
42 number of reasons, including the reduction in the risk of data disclosure or loss' (D39).
43

44
45 Similarly, educational research does not usually need access to a full set of attributes held in a
46
47 particular dataset. Sensitive data could be withheld, as reported by A16 and A9: 'it's got
48
49 nothing to do with individuals. I am not interested. ... In fact, when I ran the analysis the first
50
51 time, I deleted every column except the local authority, the school, the free school meal and
52
53 the ever-free school meal...I chop out everything that I'm not going to need'.
54

55
56 *Statistical disclosure controls*
57

58
59 Administrative and other data is sometimes subject to statistical disclosure controls to ensure
60

1
2
3 that anonymised individuals are not identifiable. Linking data which includes unusual
4 characteristics or clusters can allow anonymised people or places to be re-identified. Data
5 supply agreements often require the mechanism of statistical disclosure control or rounding
6 mechanisms 'to ensure that all statistics published are at a level of anonymisation and
7 aggregation which will ensure that no Personal Data or Sensitive Personal Data are
8 published, and thereby ensure the confidentiality of individuals' (D2, D6). A15 explained,
9
10 'As you create a longitudinal record, in other words you link the information that people have
11 given you through time, whether per person or the family, then you're creating a ... richer
12 source of information but you're also creating something which is more and more disclosive
13 through time, even if the identity of the individual is removed before those data are made
14 available to the research community. And therefore you've got to engage in disclosure
15 control mechanisms, so that ... you can say that, you know, these data have been used for
16 research but you were not identified and your identity cannot possibly have been revealed in
17 the process of conducting that research' (A15). A11 followed a similar protocol, 'my models
18 are all aggregate, results are displayed in the aggregate, and typically sort of cloaked in a sort
19 of geo-demographic classifier which, I mean you could argue that's another level of sort of
20 abstraction from the individual.'

21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 *Anonymisation and aggregation of data*

43
44 Anonymisation can go beyond individuals to entities, such as businesses or schools. A14
45 suggested a way to 'anonymise the school data, so, so you could just, you know, make sure
46 the schools all had the same identifier but anonymise it, effectively that would be sufficient to
47 completely anonymise the data and you know, make it publicly available. ... I've always said
48 to journals we can't make this data available ... but I can see that changing' (A14).

49
50
51
52
53
54
55
56 Researchers were keen to stress that no identifying data about pupils would be published:

57
58 'We're not even identifying schools, let alone individuals' (A9).
59
60

1
2
3 Energy data from the linked data sources in NEED is anonymised. A26 reported, ‘The whole
4 idea of anonymising the dataset was to prevent individuals from being identified and
5 individual households from being identified, so we had to take out all of the data that would
6 allow it to be linked.’ Multiple layers of aggregation also make re-identification difficult.
7
8
9

10 11 12 *Granularity of accessible data*

13
14 Decreasing granularity of data makes it less disclosive, but as energy researcher A30
15 suggested, less useful for some research. ‘So there is a fairly reasonable quality and quantity
16 of aggregated energy statistics available... at the national level or kind of sub-national level,
17 the sort of going down to local authority level and maybe a little bit lower, you can get
18 reasonable energy data. Going beyond that is extremely difficult’. In contrast, the Valuation
19 Office Agency data on business rates was cited as a world-leading example of the research
20 potential of raw administrative data. ‘It’s completely disaggregate, which is one of the huge
21 values of it. ... As a resource, it’s essentially unparalleled in the world’ (A35).
22
23
24
25
26
27
28
29
30
31
32

33 *Data sharing*

34
35 The research identified a need for government agencies to collaborate on data access
36 requirements and criteria, and apply lessons learned from data sharing initiatives to requests
37 for access and linkage of data for re-use. Interdepartmental collaboration within government
38 (or the absence thereof) was raised in the transport case study and came through strongly in
39 interviews with energy researchers, such as A35, ‘There’s no joined up thinking between
40 government departments in terms of data. It just isn’t there.’ Different approaches to data
41 sharing between government bodies leads to inconsistent risk assessment about similar types
42 of data, and consequently their release to researchers. A27 reported ‘the argument of not
43 allowing detailed individual access ... I’ve requested it before and have been turned down,
44 ‘cause they treat everything like an FoI, and they say, well, it’s to do with privacy.’
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 In spite of the work of ADRN, negotiations for data access generally proceed on a case-by-
4 case basis and in isolation from each other, such that decisions reached do not provide
5 precedents when similar requests arise in other parts of government. Access is re-negotiated
6 from first principles, resulting in protracted negotiations, such as researcher, A30's 3 year
7 negotiations with energy suppliers over access to data on energy efficiency interventions
8 under government-funded schemes. Whilst DVLA has been at the forefront of transport data
9 sharing initiatives for administrative purposes within government and with industry, there is
10 no evidence that this has smoothed the path for researchers wishing to use the same data for
11 secondary analysis.
12
13
14
15
16
17
18
19
20
21
22

23 *Commissioned research*

24
25
26 Delays and inconsistencies in the approval process for obtaining administrative data led some
27 researchers to seek other routes to data access. One is via government-sponsored research,
28 directly or through third party organisations. A10 hoped that saying, 'we can access NPD
29 data from other funders, we'll say we can get access' would facilitate her education research.
30 However, the need to specify the purposes for which the data will be used makes this a
31 frustrating experience, particularly where linked data is concerned. A9 reflected on 'a tactical
32 error' running three projects together in an application for pupil data resulting in confusion,
33 questions and delays, 'three different things here and we need different things for different
34 reasons, and then there were three parties involved'.
35
36
37
38
39
40
41
42
43
44
45
46

47 *Privatisation*

48
49 Many researchers were concerned about government data going behind paywalls, where
50 publicly funded researchers had to pay government departments for access to data to
51 undertake research in the public interest, and the effect of privatisation of government
52 functions curbing access to data. A35 gave an example of privatisation of government
53 Executive Agencies, such as Ordnance Survey Ltd, fearing another 'Building Research
54
55
56
57
58
59
60

1
2
3 Establishment ... a government research department on buildings which got privatised and
4 then they said, oh well, we're keeping all of our data then, despite the fact that it was
5 collected with public money, you can't get it.' Privatisation is a rather extreme example of
6 control of access to data which at least has the benefit of reducing the risks associated with
7 data sharing and re-use.
8
9

10 11 12 13 14 15 *Risk management 4: opening up data*

16
17 The final group of mitigating activities relate to changing the risk attitudes about data access,
18 partly through encouraging more proactive publication of open data, and by promoting a
19 dialogue of public engagement and public understanding of research access to government
20 data.
21
22
23
24

25 26 27 *Public engagement*

28
29 The researcher's individual responsibility extends well beyond managing the risks in data and
30 academic outputs into proactive public engagement activities to demonstrate the balance of
31 risks and the public benefits of secondary research using administrative data. A14 felt that
32 'researchers have a responsibility to actually show the public worth, this involves, you know,
33 getting support from the groups, and making a case.' Public engagement would also help to
34 mitigate risks such as misinterpretation, allowing the researcher to explain 'why these data
35 may or may not show what they appear to show and so on' (A7) and to demonstrate the
36 public value of the research.
37
38
39
40
41
42
43
44
45

46 47 48 *Open Data*

49
50 Making data as widely available as possible is one radical approach to managing risk. A12
51 exclaimed, 'You know, the best way to counter that is to have data accessible to as many
52 people as possible in education that's exactly what's happened, so making this data
53 available, particularly to other disciplines has meant that some of the debate around certain
54 controversial issues that in the past have been very dominated by individuals or small
55
56
57
58
59
60

1
2
3 datasets that are not replicable, now it's much harder because the data is publicly available.'
4
5 A16 said, 'some people think we should just put the National Pupil Database on the web so
6
7 anyone who wants to use it... and obviously in a world of open policy-making and open data,
8
9 there's much to be said for that.' A compromise between confidentiality and open access to a
10
11 dataset might be to look for new ways to make parts of the dataset publicly available. For
12
13 instance, A16 suggested issuing a 1% or 2% open sample of NPD data.
14
15

16
17 Open data helps researchers to understand what kinds of administrative data might be useful
18
19 for their research work, before applying for data access. A1 gave an example from HESA, the
20
21 Higher Education Statistics Agency, 'they publish some open tables on their website and they
22
23 also have non-open data, but it's fairly clear from the open data what sort of things they have,
24
25 and you can imagine what research work you might be able to do using that.'
26
27

28
29 Transport researchers appear to be enthusiastic advocates of open data as a means of
30
31 validating data and reducing the risks of misinterpretation. A24, suggested that 'If data are
32
33 publicly available ... somebody who's ... queuing to cross check and validate every dataset,
34
35 and they might say, I've found this dataset, I can prove it is all inaccurate. ... So if
36
37 everything's publicly available, it means they can be scrutinised more easily.'
38
39

40 **4. Conclusions**

41
42 The discussion has analysed elements of the risk framework associated with research use of
43
44 government administrative data, brought together in Figure 3. Table 2 and Figure 1 identifies
45
46 and describes seven risk events, including re-identification, data linkage, misinterpretation,
47
48 misuse, harm, commercial and opportunity risks. Table 4 summarises 22 risk mitigations in
49
50 four clusters which can be adopted by the five groups of stakeholders (researchers,
51
52 employers, research community, data providers, data subjects) identified in Figure 2 and
53
54 Table 3.
55
56

57
58 **Figure 3 Risk framework elements [insert]**
59
60

1
2
3 The analysis derives from the experiences of those interviewed for the study, captured in the
4 interview transcripts and association documentation. It is therefore reflective of that
5
6 experience and does not attempt to present all possible risks or all possible mitigations. By
7
8 interviewing stakeholders in four case studies, we sought to avoid mono-disciplinary views
9
10 and assumptions on risk. We believe that the analysis presented here will be of value to those
11
12 working in other disciplinary areas and to those such as records and information managers
13
14 who may find it helpful to better understand the perspectives presented here. It was surprising
15
16 that records managers were not more visible either in data provider or research user
17
18 organisations as they have relevant skills in data management, however they were seldom
19
20 directly involved in administrative data in our case studies. We do not claim that the risk
21
22 framework will be applicable in all data settings, rather that it adds richness to our
23
24 understanding of data risks and how best to manage them, especially by highlighting the
25
26 views of academic researchers who use administrative data.
27
28
29
30
31

32
33 Further discussion is needed between stakeholders to develop an agreed framework for
34
35 managing risks collaboratively, which can then be implemented, monitored, reviewed and
36
37 continually improved. At present the data community lacks mechanisms to agree and
38
39 implement risk mitigation and risk management frameworks. As a result, risks are
40
41 experienced locally, lessons are not learned, and there is no opportunity for the research
42
43 community to better manage risks. No forum or community of practice exists in which
44
45 discussion of best practice in administrative data management, and specifically risks and risk
46
47 mitigation, can take place. A balance is required to create a risk management environment
48
49 that benefits all stakeholders, since the exercise of risk mitigation by one group while
50
51 decreasing the risk for that group, may in fact shift the burden and increase the risks for
52
53 another. The approach to risk must reflect the nature of the data (sensitivity), the conditions
54
55 of its creation (confidentiality, legal frameworks), its depth (granularity) and its breadth
56
57
58
59
60

1
2
3 (degree of linkage) which combine to influence the weighting of risk factors. An
4
5 understanding of risk weighting, not discussed in this article, should inform the approach to
6
7 risk. Further research is needed into cultures of risk management or risk appetite: perceptions
8
9 of what is too risky vary across the different stakeholder groups. Risk mitigation cannot
10
11 simply focus inwards on effective systems and processes but must also look outwards and
12
13 include public awareness of data sharing, and the public perception of risk. It is in the
14
15 interests of all those responsible for the creation, management and re-use of government
16
17 administrative data to make progress on a common mechanism, in order to ensure that the
18
19 power of administrative data is fully exploited for the public good.
20
21
22
23
24
25

26 **References**

- 27
28 ARMA International (2010). *Evaluating and Mitigating Records and Information*
29
30 *Risk*, ARMA, Overland Park, KS.
31
32
33 Caldicott, F. (2013). *Information: To Share Or Not To Share? The Information Governance*
34
35 *Review*. London, Department of Health.
36
37
38 Cavoli, C., Christie, N., Mindell, J., Titheridge, H. (2015). Linking transport, health and
39
40 sustainability: Better datasets for better policy-making. *Journal of Transport & Health*, 2:2,
41
42 111-119.
43
44
45 Daneshmandnia, A. (2019). The influence of organizational culture on information
46
47 governance effectiveness. *Records Management Journal*, 29: 1-2, 18-41.
48
49
50 Digital Health Intelligence. (2014). *Care.data: a row waiting to happen*.
51
52 <https://www.digitalhealth.net/2014/01/care-data-a-row-waiting-to-happen/> Accessed
53
54 31.01.18.
55
56
57 Doty, P. (2015). U.S. homeland security and risk assessment. *Government Information*
58
59 *Quarterly*. 32: 3, 342-352.
60

1
2
3 GB. Department for Business, Energy & Industrial Strategy (2015). Aqua Book resources.

4
5 <https://www.gov.uk/government/collections/aqua-book-resources>.

6
7
8 Accessed 31.01.18.

9
10 GB. HM Treasury. (2015). *Aqua Book: guidance on producing quality analysis for*

11
12 *government*. [https://www.gov.uk/government/publications/the-aqua-book-guidance-on-](https://www.gov.uk/government/publications/the-aqua-book-guidance-on-producing-quality-analysis-for-government)

13
14 [producing-quality-analysis-for-government](https://www.gov.uk/government/publications/the-aqua-book-guidance-on-producing-quality-analysis-for-government) Accessed 31.01.18.

15
16
17 International Standards Organisation. (2009a). *PD ISO Guide 73:2009 Risk management —*
18
19 *Vocabulary*.

20
21 International Standards Organisation. (2009b). *BS ISO 31000:2009 Risk management —*
22
23 *Principles and guidelines (revision consultation ISO/DIS 31000 :2017(E))*.

24
25 ISACA (2010). *The Risk IT Framework*, ISACA, Rolling Meadows, IL.

26
27 Lemieux, V L. (2010). The records-risk nexus: exploring the relationship between records
28
29 and risk. *Records Management Journal*, 20: 2, 199-216.

30
31
32 Partridge, N. (2014). *Data Release Review*. Leeds: Health and Social Care Information
33
34
35 Centre.

36
37 Quigley, K., Burns, C., Stallard, K. (2015) 'Cyber Gurus': A rhetorical analysis of the
38
39 language of cybersecurity specialists and the implications for security policy and critical
40
41 infrastructure protection. *Government Information Quarterly*. 32: 2, p108-117.

42
43
44 Sexton, A., Shepherd, E., Duke-Williams, O., Eveleigh, A. (2017). A balance of trust in the
45
46 use of government administrative data. *Archival Science*. 17:4, 305-330.

47
48
49 Sexton, A., Shepherd, E., Duke-Williams, O., & Eveleigh, A. (2018). The role and nature of
50
51 consent in government administrative data. *Big Data & Society*. 5:2, 1-17.

52
53
54 Sprehe, J.T. (2005). The positive benefits of electronic records management in the context of
55
56 enterprise content management. *Government Information Quarterly*. 22: 2, 297-303.

1
2
3 Xie, S L. (2019). A must for agencies or a candidate for deletion: A grounded theory
4 investigation of the relationships between records management and information security.
5
6
7
8 *Records Management Journal*. 29: 1-2, 57-85.
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Records Management Journal

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

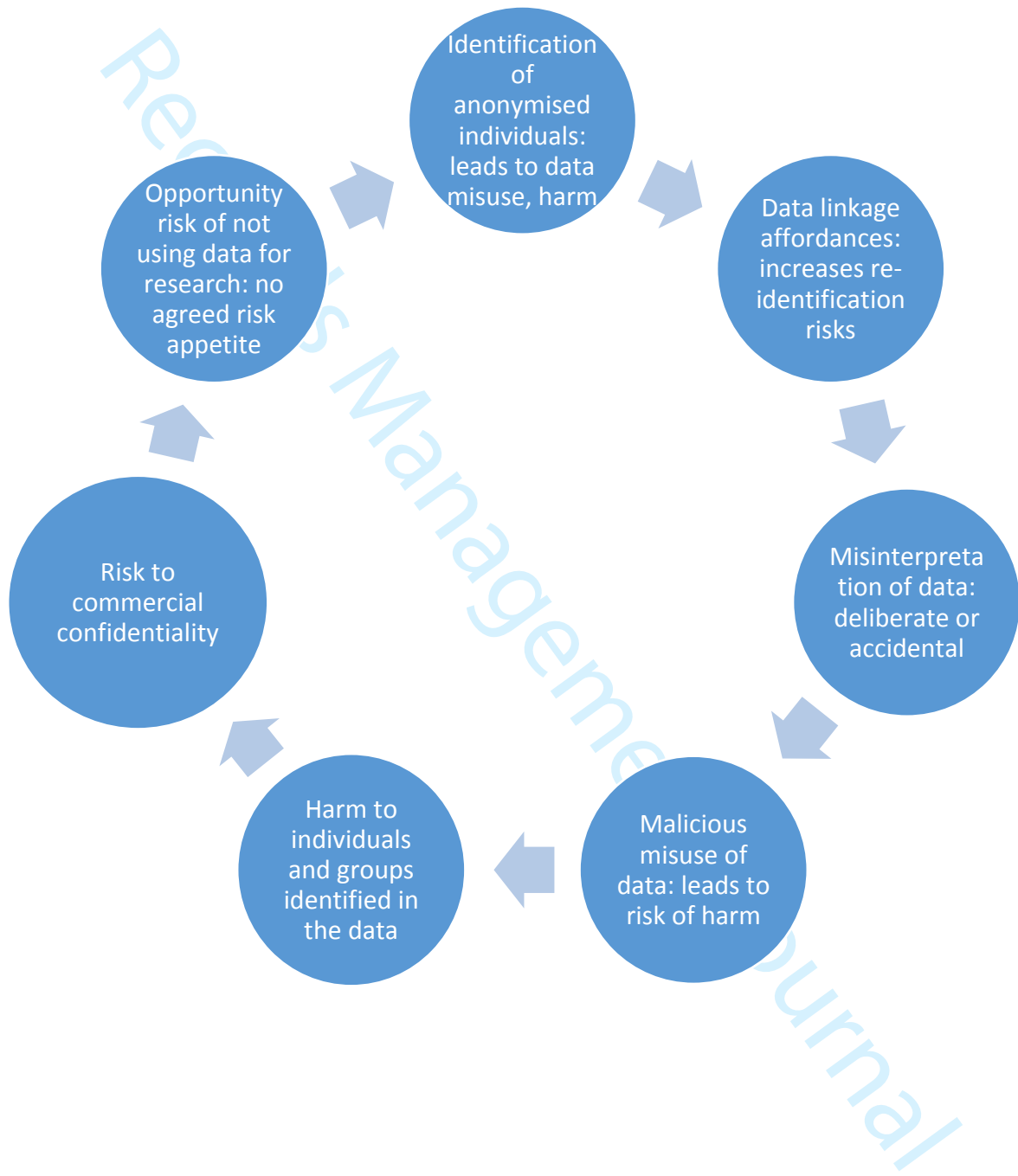


Figure 1: Risk events summary

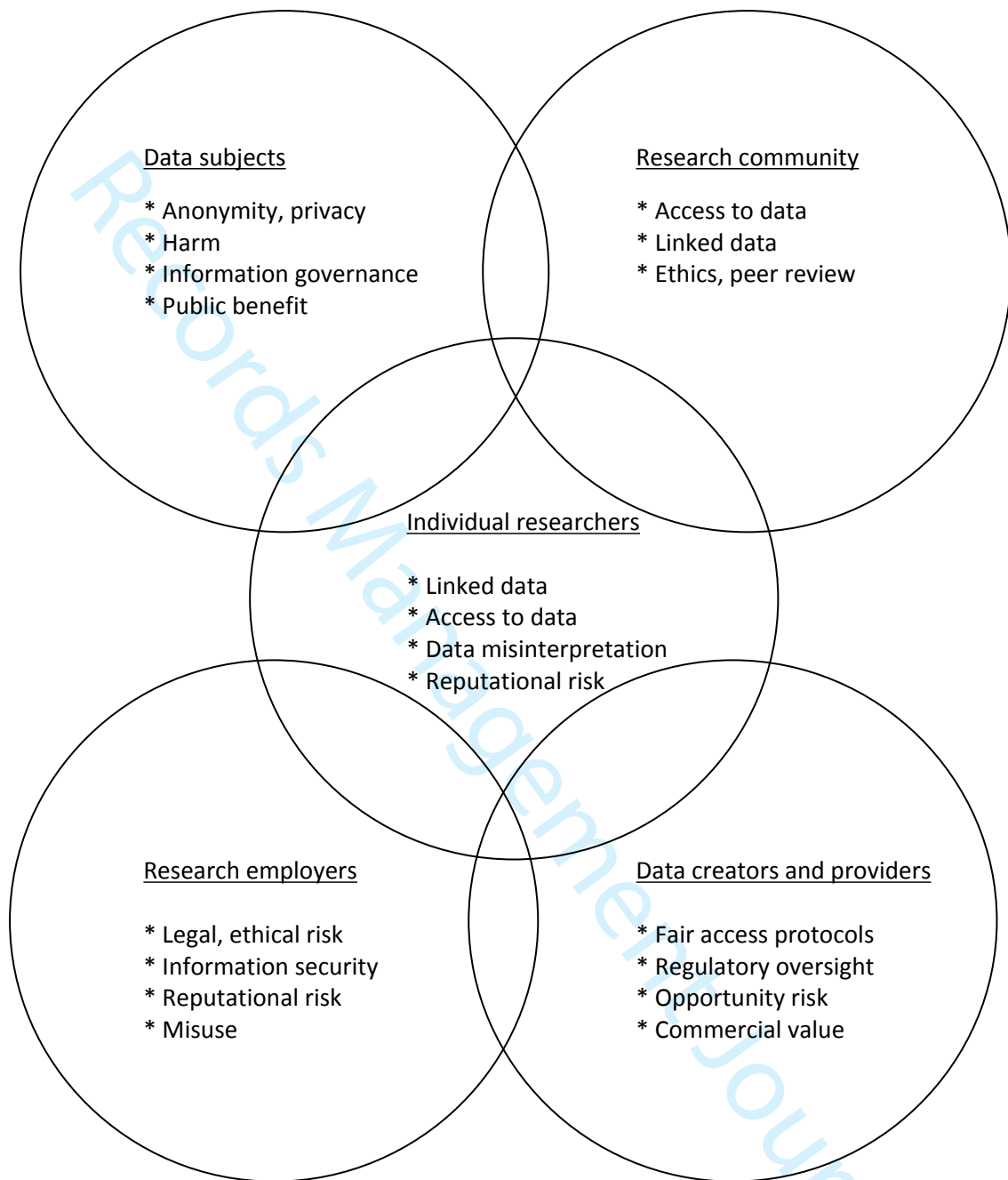


Figure 2: Stakeholder mapping

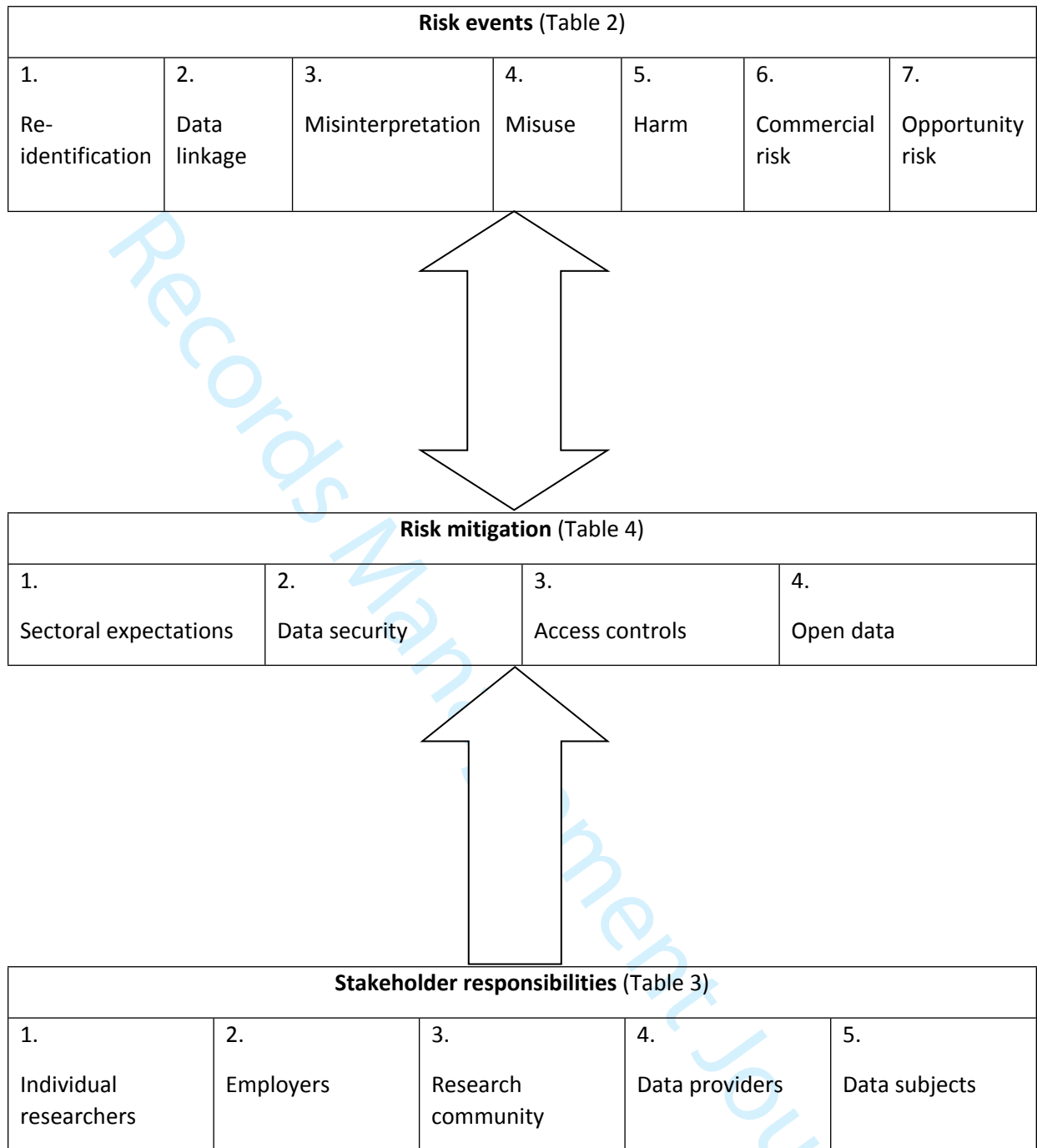


Figure 3: Risk framework elements