



UNIVERSITY OF LEEDS

This is a repository copy of *Evidence from big data in obesity research: International case studies*.

White Rose Research Online URL for this paper:  
<http://eprints.whiterose.ac.uk/154851/>

Version: Accepted Version

---

**Article:**

Wilkins, E, Aravani, A [orcid.org/0000-0002-3976-7888](https://orcid.org/0000-0002-3976-7888), Downing, A [orcid.org/0000-0002-0335-7801](https://orcid.org/0000-0002-0335-7801) et al. (6 more authors) (2020) Evidence from big data in obesity research: International case studies. *International Journal of Obesity*, 44 (5). pp. 1028-1040. ISSN 0307-0565

<https://doi.org/10.1038/s41366-020-0532-8>

---

© The Author(s), under exclusive licence to Springer Nature Limited 2020. This is an author produced version of a paper published in *International Journal of Obesity*. Uploaded in accordance with the publisher's self-archiving policy.

**Reuse**

Items deposited in White Rose Research Online are protected by copyright, with all rights reserved unless indicated otherwise. They may be downloaded and/or printed for private study, or other acts as permitted by national copyright laws. The publisher or other rights holders may allow further reproduction and re-use of the full text version. This is indicated by the licence information on the White Rose Research Online record for the item.

**Takedown**

If you consider content in White Rose Research Online to be in breach of UK law, please notify us by emailing [eprints@whiterose.ac.uk](mailto:eprints@whiterose.ac.uk) including the URL of the record and the reason for the withdrawal request.



[eprints@whiterose.ac.uk](mailto:eprints@whiterose.ac.uk)  
<https://eprints.whiterose.ac.uk/>

1 **Title:** Evidence from big data in obesity research: International case studies

2 **Running title:** Evidence from big data in obesity research

3 **Authors:** Emma Wilkins<sup>1</sup>, Ariadni Aravani<sup>1</sup>, Amy Downing<sup>1</sup>, Adam Drewnowski<sup>2</sup>,  
4 Claire Griffiths<sup>3</sup>, Stephen Zwolinsky<sup>3</sup>, Mark Birkin<sup>4</sup>, Seraphim Alvanides<sup>5,6</sup>, Michelle A  
5 Morris<sup>1</sup>

6 <sup>1</sup> Leeds Institute for Data Analytics & School of Medicine, University of Leeds, United  
7 Kingdom

8 <sup>2</sup> Center for Public Health Nutrition, University of Washington, Seattle, USA

9 <sup>3</sup> School of Sport, Leeds Beckett University, United Kingdom

10 <sup>4</sup> Leeds Institute for Data Analytics & School of Geography, University of Leeds, United  
11 Kingdom

12 <sup>5</sup> Engineering & Environment, Northumbria University, United Kingdom

13 <sup>6</sup> GESIS – Leibniz Institute for the Social Sciences, Cologne, Germany  
14

15 **Corresponding author:** Michelle A Morris

16 **Email:** [m.morris@leeds.ac.uk](mailto:m.morris@leeds.ac.uk)

17 **Tel:** +44 113 343 0883

18 **Conflict of interest statement**

19 Authors report no conflict of interest.

20

21

22

1 **Abstract**

2 *Background/Objective:* Obesity is thought to be the product of over 100 different  
3 factors, interacting as a complex system over multiple levels. Understanding the  
4 drivers of obesity requires considerable data, which are challenging, costly and time-  
5 consuming to collect through traditional means. Use of 'big data' presents a potential  
6 solution to this challenge. Big data is defined by Delphi consensus as: "*always digital,*  
7 *has a large sample size, and a large volume or variety or velocity of variables that*  
8 *require additional computing power*<sup>1</sup>. 'Additional computing power' introduces the  
9 concept of Big Data Analytics. "The aim of this paper is to showcase international  
10 research case studies presented during a seminar series held by the Economic and  
11 Social Research Council (ESRC) Strategic Network for Obesity in the UK. These are  
12 intended to provide an in-depth view of how big data can be used in obesity research,  
13 and the specific benefits, limitations and challenges encountered.

14 *Methods and results:* Three case studies are presented. The first investigated the  
15 influence of the built environment on physical activity. It used spatial data on green  
16 spaces and exercise facilities alongside individual-level data on physical activity and  
17 swipe card entry to leisure centres, collected as part of a local authority exercise class  
18 initiative. The second used a variety of linked electronic health datasets to investigate  
19 associations between obesity surgery and the risk of developing cancer. The third  
20 used data on tax parcel values alongside data from the Seattle Obesity Study to  
21 investigate sociodemographic determinants of obesity in Seattle.

22 *Conclusions:* The case studies demonstrated how big data could be used to augment  
23 traditional data to capture a broader range of variables in the obesity system. They  
24 also showed that big data can present improvements over traditional data in relation

1 to size, coverage, temporality, and objectivity of measures. However, the case studies  
2 also encountered challenges or limitations; particularly in relation to  
3 hidden/unforeseen biases and lack of contextual information. Overall, despite  
4 challenges, big data presents a relatively untapped resource that shows promise in  
5 helping to understand drivers of obesity.

## 1 Introduction

2 Obesity is a complex health, social and economic challenge. It is widely recognised as  
3 a product of numerous multi-level factors, including individual, social, economic,  
4 environmental and political influences<sup>2-4</sup>. This complexity is represented in the  
5 Foresight Obesity System Map<sup>5</sup>, which lists 108 contributing factors, depicted as  
6 nodes in a system diagram. It is also reflected in the multi-disciplinary nature of obesity  
7 research, which covers disciplines as diverse as medicine, public health, geography  
8 and computer science. Whole systems approaches, which intervene across these  
9 multiple levels and domains, have been touted as a way to tackle the growing problem  
10 of obesity<sup>6</sup>. Understanding the drivers of obesity and responses to interventions within  
11 such a complex system requires considerable data. Use of ‘big data’ and associated  
12 analytics, presents a potential solution to this challenge.

13 Various definitions of ‘big data’ have been adopted<sup>7-9</sup>. In this paper, we adopt a  
14 definition of ‘big data’ established by a recent Delphi survey of international obesity  
15 and big data experts<sup>1</sup>, which agreed that, in contrast to traditional data, big data:

16 *“is always digital, has a large sample size, and a large volume or variety or velocity*  
17 *of variables that require additional computing power. It can include quantitative,*  
18 *qualitative, observational or interventional data from a wide range of sources (e.g.*  
19 *government, commercial, cohorts) that have been collected for research or other*  
20 *purposes, and may include one or several datasets. Specialist skills in computer*  
21 *programming, database management and data science analytics are usually*  
22 *required to analyse big data.”*

1 According to the Delphi survey, 'big data' can include not only 'novel' types of data  
2 such as social media, loyalty cards and sensors, but also routinely collected data, such  
3 as health records, government and census data.

4 The Economic and Social Research Council (ESRC) Strategic Network for Obesity  
5 ('the Network') was established to consider use of big data in obesity research<sup>10</sup>.  
6 Several outputs from the Network, which form part of this paper series, have  
7 demonstrated that research applications using big data, and associated analytics,  
8 within obesity research are rich and diverse. Timmins, Green *et al*<sup>11</sup> report a wide  
9 range of studies already using big data in obesity research. They reveal that big data  
10 could provide many benefits such as increased scope and objectivity, access to  
11 unreached populations, and the potential to evaluate real-world interventions. Big data  
12 and big data analytics have also been used to produce innovative data visualisation  
13 tools, with demonstrable policy impact<sup>12</sup>. Looking to the future, a mapping exercise<sup>13</sup>  
14 demonstrated that big data sources can provide information spanning almost 80% of  
15 the nodes in the Foresight Obesity System Map. The remainder of the nodes could be  
16 covered by more traditional data sources, demonstrating how synergy of big and  
17 traditional data can support whole systems approaches to obesity.

18 Big data also has limitations, such as concerns around data validity and  
19 representativeness<sup>11</sup>, which need to be balanced alongside benefits. Challenges exist  
20 around ethics, data governance, data handling and processing capabilities<sup>1, 7, 14</sup>.  
21 Consistent reporting of data sources, such as through use of the recently developed  
22 BEE-COAST framework<sup>13</sup>, better enables critique of these strengths and limitations.

23 Applications of big data in obesity research include use of retail sales data to evaluate  
24 the impact of obesity policy<sup>15</sup>, use of geotagged social media data to explore patterns  
25 in obesity-related behaviours<sup>16, 17</sup>, and use of smartphone data to assess physical

1 activity patterns over space and time<sup>18, 19</sup>. These examples draw on data from diverse  
2 sectors, highlighting again the multi-disciplinary nature of obesity research.

3 Uses of big data include both hypothesis generation ('exploratory analyses') and  
4 hypothesis testing. Recognising the distinction between these two forms of enquiry is  
5 important to avoid hypothesising after the results are known<sup>20</sup>. This may be  
6 particularly problematic in the case of big data research, as large sample sizes,  
7 coupled with repeated exploratory analyses, will lead to increased chance of detecting  
8 statistically significant associations that are of limited clinical and practical importance.

9 The aim of this paper is to showcase international research case studies presented  
10 during seminars held by the Network in the UK<sup>10</sup>. These are intended to complement  
11 existing high-level reviews of big data in obesity research<sup>11, 13</sup> by providing an in-depth  
12 view of how big data can be used in this field, and the specific benefits, limitations and  
13 challenges encountered.

## 14 **Methods and Results**

15 Three case studies presented at the Network Seminar Series<sup>10</sup> are reported. Each  
16 employed several sources of data, including 'big' and 'traditional' data to measure  
17 obesity-related exposures and/or outcomes. These data are reported using a  
18 standardised BEE-COAST framework<sup>13</sup> that cross references to the Foresight Obesity  
19 System Map nodes<sup>5</sup> highlighting the breadth of data coverage (Table 1).

20 Table 2 summarises all Network seminar presentations. Further information and  
21 seminar recordings can be found at <https://www.cdrc.ac.uk/research/obesity/>.

## 1 **Case Study 1: Uptake of physical activity in Leeds, UK**

2 *Griffiths and Zwolinsky, Seminar: May 2016, London School Hygiene Tropical*  
3 *Medicine*

4 *Background:* Physical activity can help prevent and manage a number of chronic  
5 health conditions, including obesity<sup>21, 22</sup>. The World Health Organisation<sup>23</sup>, and other  
6 bodies internationally<sup>24, 25</sup> have called upon authorities to increase opportunities for  
7 physical activity as a means to tackle obesity. Repurposing existing ‘big’ spatial data  
8 on the physical activity environment provides novel opportunities to support policy.

9 *Data:* Leeds Let’s Get Active Programme, Points of Interest (Table 1)

10 *Methods:* Links with Leeds City Council facilitated secondary analysis of data  
11 emerging from the Leeds Let’s Get Active (LLGA) programme; a council initiative to  
12 increase physical activity through exercise classes delivered at leisure centres.  
13 Exploratory, cross-sectional analyses investigated (i) the association between the  
14 number of neighbourhood physical activity opportunities and separate outcomes of  
15 sedentary behaviour and physical activity, controlling for neighbourhood deprivation,  
16 and (ii) whether residential proximity to participating leisure centres was related to  
17 attendance. Physical activity opportunities were derived from Points of Interest data;  
18 a large dataset detailing the locations of a wide range of features across the whole of  
19 Great Britain.

20 Participant postcodes were analysed in a Geographic Information System (GIS)  
21 together with data on the locations of physical activity opportunities from Points of  
22 Interest data and the locations of participating leisure centres. Physical activity  
23 opportunities separately included (i) green spaces and (ii) built facilities such as gyms,  
24 climbing facilities, and swimming pools. Neighbourhoods were defined using ‘Lower



1 Super Output Area' (LSOA) boundaries (a UK administrative geography containing  
2 ~1,500 people) and 2km circular buffers.

3 *Results:* LLGA data contained 29 796 self-reports of physical activity and sedentary  
4 behaviours, together with leisure centre attendance data from swipe cards. Analyses  
5 revealed no associations between any measure of physical activity opportunities and  
6 physical activity or sedentary behaviours, with the exception of counts of green spaces  
7 within LSOAs. Those with the highest counts of green spaces within LSOAs were  
8 more likely to meet physical activity guidelines.

9 Fewer than 50% of participants who registered for the programme attended a session.  
10 Of those that did, over one third did not attend the centre closest to them. Having a  
11 leisure centre within the residential Middle Super Output Area (an administrative  
12 geography containing ~8 000 people) or a 2km circular buffer only accounted for a  
13 small proportion of the variability in attendance rates. On further investigation, circular  
14 buffers of at least 4km around leisure centres were required to capture over 50% of  
15 attendees.

16 *Conclusion:* There is some indication that neighbourhood greenspace is related to  
17 physical activity. However, in agreement with other literature<sup>26, 27</sup>, this study shows  
18 different definitions of environment can produce different results. Future work must  
19 use measures that are relevant, consistent and transferable. Mere proximity to  
20 opportunities from home may not be a good indicator of actual exposure/opportunities.  
21 People frequently visit leisure centres that are not closest to home.

## 1 **Case Study 2: Obesity and Colorectal Cancer in England, UK**

2 *Aravani and Downing, Seminar: April 2017, Leeds Beckett University*

3 *Background:* Obesity is linked to an increased risk of several malignancies, including  
4 colorectal cancer<sup>28, 29</sup>. Counterintuitively, some research suggests surgery to reduce  
5 obesity ('obesity surgery') may increase the risk of colorectal cancer<sup>30-33</sup>. However,  
6 this association remains unclear, with the majority of studies having short follow-up  
7 time or lacking statistical power. This study tested the a-priori hypothesis that obesity  
8 surgery is associated with risk of colorectal cancer and also explored associations with  
9 other obesity-related cancers (breast, kidney or endometrial) across the English  
10 National Health Service (NHS).

11 *Data:* Hospital Episode Statistics (HES), National Cancer Registration and Analysis  
12 Service (NCRAS), Office for National Statistics (ONS) mortality data (Table 1).

13 *Methods:* This was a national population-based retrospective observational study.  
14 Individuals who underwent obesity surgery (the 'OS group') or had a hospital episode  
15 with a diagnosis of obesity but no obesity surgery (the 'no-OS group'), between April  
16 1997 and September 2013, were identified using HES data. HES data were obtained  
17 pre-linked with NCRAS and ONS mortality data. This allowed identification of  
18 individuals in the OS and no-OS groups who were subsequently diagnosed with  
19 colorectal cancer, or other obesity-related cancers. It also allowed identification of the  
20 time 'at risk' - the time from obesity diagnosis/surgery to development of a cancer,  
21 death or last follow-up (30<sup>th</sup> September 2013). Standardised incidence ratios (SIR)  
22 with 95% confidence intervals (CI) were calculated to define the risk of developing  
23 cancer in the OS and no-OS groups relative to the background English population,  
24 accounting for age and calendar year.

1 *Results:* A total of 1 002 607 obese patients were identified, of whom 4% (n=39 747)  
2 underwent obesity surgery. The OS group and no-OS groups had a median follow-up  
3 period of 3 years (range 1-16 years) and 2.5 years (range 1-16 years), respectively.  
4 In the no-OS cohort, 3 237 developed colorectal cancer with an SIR of 1.12 (95%CI  
5 1.08-1.16) relative to the background population. In the OS cohort, 43 developed  
6 colorectal cancer with an SIR of 1.26 (95%CI 0.92-1.71). There was a significantly  
7 increased risk of colorectal cancer among the oldest ( $\geq 50$  years) in the OS group (SIR:  
8 1.47, 95% CI: 1.02-2.06). High SIRs for renal and endometrial cancers were found in  
9 both the OS and non-OS groups<sup>34</sup>. By contrast, OS was associated with reduced  
10 breast cancer risk<sup>34</sup>.

11 *Conclusion:* Although the association between obesity surgery and subsequent  
12 colorectal cancer risk was limited by the small OS group size and short follow-up time,  
13 this study showed an elevated colorectal cancer risk continues after obesity surgery  
14 in individuals older than 50 years. The high SIRs for renal and endometrial cancers  
15 require further investigation.

16

# 1 **Case Study 3: Sociodemographic determinants of obesity in Seattle, USA**

2 *Drewnowski, Seminar March 2016, University of Cambridge*

3 *Background:* Socioeconomic status (SES), both at the individual and neighbourhood  
4 level is thought to contribute to obesity. However, studies of obesity and its  
5 determinants often do not contain important socioeconomic variables or include only  
6 self-reported measures, which are simplistic and subject to bias. Neighbourhood  
7 measures of SES are often only available for administrative geographies, which are  
8 subject to bias from the modifiable areal unit problem (MAUP)<sup>35</sup> and may not be suited  
9 to capturing neighbourhood effects on obesity<sup>36</sup>. A series of exploratory studies were  
10 conducted to examine whether residential property values – the second largest source  
11 of wealth in the US<sup>37</sup> - could be used as a proxy measure of individual and  
12 neighbourhood level SES, and to simulate obesity prevalence at a micro-scale.

13 *Data:* Seattle Obesity Study (SOS) I, II and III, King County Tax Parcel Values  
14 (Table1).

15 *Methods:* Data from the Seattle Obesity Study (SOS) I, II and III were used to assess  
16 associations between socioeconomic variables and health-related outcomes,  
17 including diet and obesity. Participants' residential addresses were geocoded to tax  
18 parcel centroids; plots of land owned by a single landowner and typically containing a  
19 single residential property or a block of properties e.g. flats. In a series of studies, SOS  
20 participants were ascribed individual and neighbourhood measures of SES based on  
21 King County Tax Parcel Values. Individual SES was operationalised as the average  
22 property value in the tax parcel of residence. Neighbourhood SES was operationalised  
23 as the average property value in the residential neighbourhood (various definitions  
24 including residential census tracts and home-centric buffers spanning multiple tax  
25 parcels). Multivariable linear regressions examined associations between these  
11

1 measures and obesity-related outcomes, including behaviours, diet quality (e.g.  
2 measures of soda and salad consumption) and obesity, controlling for age, gender,  
3 and race/ethnicity. This was contrasted with the performance of more traditional  
4 measures of SES, including education and income, at predicting obesity-related  
5 outcomes.

6 *Results:* Obesity-related outcomes were related both with property value measures of  
7 SES and more traditional SES measures. However, effect sizes for property value  
8 measures were typically equal to or greater than effect sizes for traditional measures.  
9 For example, among women, the prevalence ratio for obesity was 3.4 times greater  
10 among those having average residential tax parcel values in the lowest quartile  
11 compared to the highest (95% CI: 2.2-5.3)<sup>38</sup>. Contrastingly, education explained less  
12 variation in obesity rates (<high-school vs college, prevalence ratio: 1.7, 95% CI: 1.2-  
13 1.7). Average residential property values within residential census tracts were also  
14 associated with soda and salad consumption<sup>39</sup>, whereas income and education were  
15 not.

16 *Conclusion:* Residential property values present a convenient and readily-available  
17 measure of both individual and neighbourhood SES. They appear to better capture  
18 the multi-faceted nature of SES compared to single, self-reported measures such as  
19 education or income. They also have potential to be applied to spatial microsimulation  
20 models (a technique for estimating the characteristics of a population<sup>40</sup>) to model  
21 obesity rates at the micro-scale.

22

## 1 **Discussion**

2 These case studies demonstrate how big data and traditional data both have an  
3 important role in understanding the aetiology of obesity, alongside responses to  
4 obesity interventions. An earlier mapping exercise<sup>13</sup> demonstrated that combining big  
5 data with traditional data could provide information spanning over 82% of the 108  
6 nodes in the Foresight Obesity System Map. The data used in our three case studies  
7 spanned 34 nodes (31%), of which 59% were covered by big data sources. These  
8 case studies demonstrate that big data can successfully be used to augment  
9 traditional data to cover a wider scope of the obesity system, or to provide increased  
10 size, coverage, temporality, or objectivity of measures. The remaining discussion  
11 provides an in-depth review of the specific benefits, limitations and challenges  
12 encountered within these case studies.

## 13 **Benefits**

### 14 *Large size and coverage*

15 A key benefit, evident in all three case studies, was the potential size and coverage of  
16 the data. For example, by combining HES, NCRAS and ONS mortality data, Case  
17 Study 2 was able to assess cancer rates among over 1 million obese people.  
18 Moreover, the data were representative of the entire UK population with a recorded  
19 hospital admission since 1997, including populations that are often unreachable.  
20 Furthermore, as there was no option to opt in or out, recruitment and attrition biases,  
21 which hamper traditional cohort studies, were minimised.

22 While the data used in both Case Studies 1 and 3 were confined to relatively small  
23 geographic regions (Leeds, UK and King County, USA respectively), both had the  
24 potential to be extended nationally, or even internationally. For example, the Points of

1 Interest data used in Case Study 1 is available across the whole of Great Britain.  
2 Property values from county tax assessors are publicly available at the level of tax  
3 parcels for all US states, with alternate sources of property values (such as  
4 commercial property sales data) being available internationally<sup>41</sup>.

#### 5 *Better temporality*

6 Traditional epidemiologic obesity studies are largely cross-sectional or take repeated  
7 measures of exposures and/or outcomes at discrete time points<sup>42</sup>. The data used in  
8 these case studies provided improved temporality over traditional data in several  
9 respects. For example, the Points of Interest data used in Case Study 1 are updated  
10 quarterly, allowing fine-grained assessment of built environment dynamics, and close  
11 temporal linkage to obesity outcomes data. Financial and time constraints would make  
12 it unfeasible to collect environmental data at this frequency and scale through primary  
13 means. Historical Points of Interest data also allows older cohort studies to be linked  
14 with built environment variables. Data used in Case Study 2 currently span several  
15 decades and are updated continually, allowing tracking of health outcomes (hospital  
16 admissions, cancer incidences etc.) for an ever-growing cohort of people. The property  
17 values data used in Case Study 3, while only updated every 6 years, still has more  
18 frequent updates than decennial census data, which is typically used to measure  
19 SES<sup>43</sup>.

#### 20 *Objective measures*

21 The data used in all three case studies also provided the benefit of objective measures.  
22 Case Study 1 used spatial data from the UK's national mapping agency to objectively  
23 measure neighbourhood physical activity opportunities. This is in contrast with other  
24 studies, which have asked participants about perceptions of their local environment<sup>44</sup>.

1 Perception measures do not correlate well with objective measures, and both may be  
2 important to comprehensively capture built environment influences on obesity<sup>45</sup>. Case  
3 Study 2 used highly robust data from the NHS, Public Health England and ONS, which  
4 importantly included objective data on obesity diagnoses, surgery, cancer incidences  
5 and deaths. Finally, Case Study 3 demonstrated how property values could provide  
6 an objective proxy for individual socioeconomic status, which performs better than self-  
7 reported education or income at predicting obesity-related outcomes.

### 8 *Augmentation of other data*

9 In Case Studies 1 and 3, big data were used to augment traditional data, illustrating  
10 the potential for big and traditional data to work in harmony. Both utilised location  
11 information (residential addresses) to link traditional data with built and socioeconomic  
12 environmental data. These represent important areas of the Foresight Obesity System  
13 Map frequently missing from traditional datasets. Case Study 3 also demonstrated that  
14 property values may provide improved measures of individual SES, even where  
15 alternate measures are included in traditional datasets. Moreover, measures of  
16 neighbourhood SES can be computed at a range of geographical scales, and  
17 unconstrained by administrative boundaries, minimising bias due to the MAUP<sup>35</sup>.  
18 These datasets also offer the potential for linkage with other big datasets such as  
19 electronic medical records. Indeed, an ongoing study ('Moving2Health') is seeking to  
20 link longitudinal electronic medical records with historical property values data<sup>46</sup> in an  
21 entirely new approach to studying built environment influences on health and disease.



## 1 **Limitations and challenges**

2 As well as the many benefits described above, limitations and challenges were also  
3 encountered. These can be divided into two categories: hidden/unforeseen biases and  
4 lack of contextual information.

### 5 *Hidden/unforeseen biases*

6 Bias within data is a concern for most research. Traditional studies seek to eliminate  
7 or reduce bias through design, with the well-established 'gold standard' being the  
8 randomised controlled trial. In epidemiological research, observational and case-  
9 control studies seek to minimise biases through methodological sampling techniques  
10 and rigorous data cleaning and handling procedures. "However, the process of  
11 collection, manipulation and extraction of value from big data - the big data analytics -  
12 is often opaque and may not follow expected research norms, making it challenging  
13 to identify and account for potential sources of bias."

14 As an example, while the data used in Case Study 2 was a national sample,  
15 differences in demographics between the general population and those (i) having a  
16 hospital episode and (ii) being eligible for obesity surgery, may lead to selection  
17 biases. In particular, people undergoing obesity surgery were required by the NHS to  
18 meet certain criteria ( $BMI \geq 40kg \cdot m^{-2}$  or  $35-40kg \cdot m^{-2}$  alongside at least one other  
19 obesity-related condition and inability to sustain weight loss through standard  
20 techniques). These factors may be associated with cancer risk independently of  
21 obesity treatment, confounding any observed associations. Indeed, in a negative  
22 control analysis, Case Study 2 found a higher incidence of lung cancer among those  
23 with obesity, and particularly those undergoing obesity treatment, compared to the  
24 background population<sup>34</sup>. This was unexpected given lung cancer is not an obesity-

1 related cancer and suggests residual confounding in the data; potentially due to  
2 increased smoking rates among those with obesity.

3 Another example of bias relates to systematic differences in the handling of data. Tax  
4 parcel values, as used in Case Study 3, are determined by independent counties  
5 according to state-level regulations. There may therefore be variability in valuation  
6 methods both at the county and state levels, leading to systematic biases in property  
7 valuations nationally. While not an issue in Case Study 3, as the study area was  
8 confined to one county, appropriate methods, such as multi-level modelling, would  
9 need to be considered in research spanning multiple counties or states. Comparability  
10 of house prices across large geographical areas also requires careful analysis in view  
11 of the known tendency towards spatial autocorrelation<sup>47</sup>.

12 Sources of bias can be hard to predict. A recent validation study showed that Points  
13 of Interest data, as used in Case Study 1, has variable completeness across different  
14 types of facilities (in this case, types of food outlets)<sup>48</sup>. This was thought to be due to  
15 differences in turnover/closure rates across outlet types, and the way Points of Interest  
16 data is sourced – with information on different outlet types being sourced from different  
17 data providers. Variability in data quality across outlet types led, in turn, to  
18 geographically varying errors due to differences in food outlet composition across  
19 environment types (e.g. deprived areas having more fast food outlets). It is unclear  
20 whether such bias would exist for listings of physical activity opportunities, as used in  
21 Case Study 1, but in any event, this example highlights how sources of bias may be  
22 difficult to anticipate.

1 *Lack of contextual information*

2 Lack of contextual information about the data was an additional challenge encountered  
3 across the case studies. This can lead to poorly performing predictive models and bias  
4 in causal models if confounders cannot be controlled for. Case Study 2 met a number  
5 of challenges in this respect. Firstly, the data did not include an earliest date of obesity  
6 diagnosis. This induces a time-related bias, with those undergoing surgery potentially  
7 having lived for longer with obesity than those not undergoing surgery.

8 Secondly, the HES data only classified procedures by type and not purpose, and it  
9 was not always clear whether procedure codes related to obesity surgery or to some  
10 other procedure (notably, some procedure codes could have encompassed both  
11 surgeries to treat obesity and surgeries to treat cancer). Procedural codes also  
12 changed over time. For example, prior to 2004 there were no codes for sleeve  
13 gastrectomy or gastric banding. It was unclear what coding was used to capture these  
14 surgeries prior to 2004 leading to further challenges in identifying obesity surgeries  
15 within the HES data.

16 A further 'missing information' challenge encountered in Case Study 2 was the  
17 absence of data on important covariates; notably BMI and other variables that may  
18 lead to increased cancer risk, and which may vary between the OS and no-OS groups.  
19 As mentioned above, using negative control analyses, the researchers detected  
20 potential residual confounding with the data. This highlights that even if sources of  
21 bias are identified, it may not be possible to control for them.

22 Challenges relating to missing contextual information were also evident, albeit to a  
23 lesser extent, in Case Studies 1 and 3. In Case Study 1, proprietary classifications  
24 were used to extract physical activity opportunities from Points of Interest data, but it

1 is unclear how these classifications were applied by the data provider, and how  
2 suitable they were for capturing physical activity opportunities relevant to obesity. For  
3 example, the classifications 'swimming pools' and 'tennis facilities' were likely to  
4 include both public and private (e.g. members-only) facilities. The data also did not  
5 include factors such as facility quality, cost or opening hours – all of which may  
6 influence facility utilisation. Similarly, while the property values data used in Case  
7 Study 3 appears to provide a good predictor of individual and neighbourhood  
8 socioeconomic context, it does not include information on other assets owned by  
9 people, and therefore may not perform well in areas where property represents only a  
10 small proportion of total assets.

## 11 **Future Directions and Conclusion**

12 The case studies presented in this paper highlight a variety of ways in which big data  
13 and associated analytics, have been used, alongside traditional data, in whole  
14 systems obesity research. They have provided detailed examples of how big data can  
15 present improvements over traditional data in relation to size, coverage, temporality,  
16 and objectivity of measures. Case study 3 also demonstrated that big data and big  
17 data analytics could be used to simulate data that is missing/unavailable from other  
18 datasets. For example, spatial microsimulation could be used to estimate  
19 neighbourhood obesity rates through combination of individual and area based  
20 characteristics<sup>40</sup>. However, these case studies also highlight that bigger data does not  
21 necessarily mean fewer challenges or limitations. Hidden/unforeseen biases and  
22 missing contextual information caused problems. Researchers should be mindful of  
23 these limitations, and look to mitigate them wherever possible, for example through  
24 using negative control analyses to test for biases, and linkage with additional datasets  
25 to provide additional contextual information.

1 The data used in the presented case studies, while meeting the definition of 'big data'  
2 as agreed by consensus of experts in a recent Delphi study <sup>49</sup>, may be regarded by  
3 some as being relatively simple, and perhaps not showcasing big data to its full  
4 potential. However, we feel the case studies presented here reflect the present state  
5 of big data and obesity research, which undoubtedly still has room for advancement  
6 in harnessing the full breadth and variety of big data. Other studies that are advancing  
7 the field of big data and obesity research in terms of the complexity of data and/or  
8 associated analyses have, for example, used loyalty card data to explore associations  
9 between objectively measured food purchases and individual characteristics <sup>50</sup>, or  
10 linked loyalty card food purchase data across the whole of London with medical  
11 prescription data to predict hypertension, high cholesterol, and diabetes at a fine  
12 spatial resolution <sup>51</sup>. Spatial microsimulation using census data has also been used to  
13 build a synthetic population for the UK, which has been linked via demographic  
14 characteristics to a nationally representative dietary survey (The National Diet and  
15 Nutrition Survey, allowing modelling of small-area variations in Body Mass Index,  
16 Calorie Intake and Physical Activity Level <sup>52</sup>. Nevertheless, there is still considerable  
17 scope for future innovation, such as through combining a greater number of diverse  
18 datasets to better capture the myriad of obesity drivers <sup>53</sup> and harnessing the temporal  
19 dimension of quickly-evolving datasets to track or predict changes over time.

20 Overall, in spite of challenges, big data and associated analytics, present a relatively  
21 untapped resource that shows promise in helping to understand obesity. We feel it is  
22 best utilised as a complement to traditional data, for example through data linkage or  
23 by providing a platform to test new methods to establish best practices in future  
24 research.

25

## 1 **Acknowledgements**

2 The ESRC Strategic Network for Obesity was funded via Economic and Social  
3 Research Council grant number ES/N00941X/1.

4 The authors would like to thank all of the network investigators<sup>a</sup> and members<sup>b</sup> for their  
5 participation in network meetings and discussion which contributed to the development  
6 of this paper.

7 <sup>a</sup> - [www.cdrc.ac.uk/research/obesity/investigators/](http://www.cdrc.ac.uk/research/obesity/investigators/)

8 <sup>b</sup> - [www.cdrc.ac.uk/research/obesity/network-members/](http://www.cdrc.ac.uk/research/obesity/network-members/)

9

## 10 **Conflict of interest statement**

11 Authors report no conflict of interest.

## 12 **References**

- 13 1. Vogel C, Zwolinsky S, Griffiths C, Hobbs M, Henderson E, Wilkins E. A Delphi  
14 study to build consensus on the definition and use of big data in obesity  
15 research. *International Journal of Obesity* 2019: e-pub ahead of print 17 Jan  
16 2019; doi:10.1038/s41366-018-0313-9.
- 17
- 18 2. Davison KK, Birch LL. Childhood overweight: A contextual model and  
19 recommendations for future research. *Obesity Reviews* 2001; **2**(3): 159-71.
- 20
- 21 3. Egger G, Swinburn B. An “ecological” approach to the obesity pandemic. *BMJ*  
22 1997; **315**(7106): 477-480.
- 23
- 24 4. Harrison K, Bost KK, McBride BA, Donovan SM, Grigsby-Toussaint DS, Kim J  
25 *et al.* Toward a developmental conceptualization of contributors to overweight  
26 and obesity in childhood: The six-Cs model. *Child Development Perspectives*  
27 2011; **5**(1): 50-58.
- 28
- 29 5. Butland B, Jebb S, Kopelman P, McPherson K, Thomas S, Mardell J *et al.*  
30 Foresight. Tackling obesities: future choices. Project report. 2007.

- 1  
2 6. Rutter HR, Bes-Rastrollo M, de Henauw S, Lahti-Koski M, Lehtinen-Jacks S,  
3 Mullerova D *et al.* Balancing upstream and downstream measures to tackle the  
4 obesity epidemic: A position statement from the European association for the  
5 study of obesity. *Obesity Facts* 2017; **10**(1): 61-63.
- 6  
7 7. Mittelstadt BD, Floridi L. The ethics of big data: Current and foreseeable issues  
8 in biomedical contexts. *Science and Engineering Ethics* 2016; **22**(2): 303-41.
- 9  
10 8. Kaisler S, Armour F, Espinosa JA, Money W. Big data: Issues and challenges  
11 moving forward. *2013 46th Hawaii International Conference on System*  
12 *Sciences* 2013: 995-1004.
- 13  
14 9. Herland M, Khoshgoftaar TM, Wald R. A review of data mining using big data  
15 in health informatics. *Journal of Big Data* 2014; **1**(2): e-pub 24 June 2014;  
16 <https://doi.org/10.1186/2196-1115-1-2>.
- 17  
18 10. Morris M, Birkin M. The ESRC Strategic Network for Obesity: Tackling obesity  
19 with big data. *International Journal of Obesity* 2018; **42**(12): 1948-1950.
- 20  
21 11. Timmins K, Green M, Radley D, Morris M, Pearce J. How has big data  
22 contributed to obesity research? A review of the literature. *International Journal*  
23 *of Obesity* 2018; **42**: 1951–1962.
- 24  
25 12. Monsivais P, Francis O, Lovelace R, Chang M, Strachan E, Burgoine T. Data  
26 visualisation to support obesity policy: case studies of data tools for planning  
27 and transport policy in the UK. *International Journal of Obesity* 2018; **42**(12):  
28 1977-1986.
- 29  
30 13. Morris M, Wilkins E, Timmins K, Bryant M, Birkin M, Griffiths C. Can big data  
31 solve a big problem? Reporting the obesity data landscape in line with the  
32 Foresight obesity system map. *International Journal of Obesity* 2018; **42**(12):  
33 1963-1976.
- 34  
35 14. Vayena E, Salathé M, Madoff LC, Brownstein JS. Ethical challenges of big data  
36 in public health. *PLOS Computational Biology* 2015; **11**(2): e1003904.
- 37  
38 15. Silver LD, Ng SW, Ryan-Ibarra S, Taillie LS, Induni M, Miles DR *et al.* Changes  
39 in prices, sales, consumer spending, and beverage consumption one year after  
40 a tax on sugar-sweetened beverages in Berkeley, California, US: A before-and-  
41 after study. *PLoS Medicine* 2017; **14**(4): e1002283.
- 42

- 1 16. Gore RJ, Diallo S, Padilla J. You are what you tweet: connecting the geographic  
2 variation in america's obesity rate to Twitter content. *PloS One* 2015; **10**(9):  
3 e0133505.
- 4
- 5 17. Nguyen QC, Li D, Meng H-W, Kath S, Nsoesie E, Li F *et al.* Building a national  
6 neighborhood dataset from geotagged Twitter data for indicators of happiness,  
7 diet, and physical activity. *JMIR Public Health and Surveillance* 2016; **2**(2):  
8 e158.
- 9
- 10 18. Hirsch JA, James P, Robinson JR, Eastman KM, Conley KD, Evenson KR *et*  
11 *al.* Using MapMyFitness to place physical activity into neighborhood context.  
12 *Frontiers in Public Health* 2014; **2**(19): 1-9.
- 13
- 14 19. Althoff T, Hicks JL, King AC, Delp SL, Leskovec J. Large-scale physical activity  
15 data reveal worldwide activity inequality. *Nature* 2017; **547**(7663): 336–339.
- 16
- 17 20. Kerr NL. HARKing: hypothesizing after the results are known. *Pers Soc Psychol*  
18 *Rev* 1998; **2**(3): 196-217.
- 19
- 20 21. Lee IM, Shiroma EJ, Lobelo F, Puska P, Blair SN, Katzmarzyk PT *et al.* Effect  
21 of physical inactivity on major non-communicable diseases worldwide: an  
22 analysis of burden of disease and life expectancy. *Lancet* 2012; **380**(9838):  
23 219-29.
- 24
- 25 22. Bennett JE, Li G, Foreman K, Best N, Kontis V, Pearson C *et al.* The future of  
26 life expectancy and life expectancy inequalities in England and Wales:  
27 Bayesian spatiotemporal forecasting. *Lancet* 2015; **386**(9989): 163-70.
- 28
- 29 23. World Health Organisation. Report of the Commission on ending childhood  
30 obesity. 2016.
- 31
- 32 24. Centers for Disease Control and Prevention. Recommended community  
33 strategies and measurements to prevent obesity in the United States. 2009.
- 34
- 35 25. Local Government Association. Building the foundations: Tackling obesity  
36 through planning and development. 2016.
- 37
- 38 26. Burgoine T, Alvanides S, Lake AA. Creating 'obesogenic realities'; Do our  
39 methodological choices make a difference when measuring the food  
40 environment? *International journal of health geographics* 2013; **12**(33): e-pub  
41 3 July 2013; doi: 10.1186/1476-072X-12-33.
- 42



- 1 27. Wilkins E, Morris M, Radley D, Griffiths C. Methods of measuring associations  
2 between the Retail Food Environment and weight status: Importance of  
3 classifications and metrics. *SSM - population health* 2019; e-pub ahead of print  
4 4 May 2019; doi: <https://doi.org/10.1016/j.ssmph.2019.100404>.
- 5
- 6 28. Bardou M, Barkun AN, Martel M. Obesity and colorectal cancer. *Gut* 2013;  
7 **62**(6): 933-947.
- 8
- 9 29. Siegel R, Desantis C, Jemal A. Colorectal cancer statistics, 2014. *CA Cancer*  
10 *Journal for Clinicians* 2014; **64**(2): 104-117.
- 11
- 12 30. Derogar M, Hull MA, Kant P, Östlund M, Lu Y, Lagergren J. Increased risk of  
13 colorectal cancer after obesity surgery. *Annals of Surgery* 2013; **258**(6): 983-  
14 988.
- 15
- 16 31. Kant P, Hull MA. Excess body weight and obesity—the link with gastrointestinal  
17 and hepatobiliary cancer. *Nature Reviews Gastroenterology and Hepatology*  
18 2011; **8**(4): 224-238.
- 19
- 20 32. Östlund MP, Lu Y, Lagergren J. Risk of obesity-related cancer after obesity  
21 surgery in a population-based cohort study. *Annals of Surgery* 2010; **252**(6):  
22 972-976.
- 23
- 24 33. Sainsbury A, Goodlad RA, Perry SL, Pollard SG, Robins GG, Hull MA.  
25 Increased colorectal epithelial cell proliferation and crypt fission associated with  
26 obesity and roux-en-Y gastric bypass. *Cancer Epidemiology Biomarkers &*  
27 *Prevention* 2008; **17**(6): 1401-1410.
- 28
- 29 34. Aravani A, Downing A, Thomas JD, Lagergren J, Morris EJA, Hull MA. Obesity  
30 surgery and risk of colorectal and other obesity-related cancers: An English  
31 population-based cohort study. *Cancer Epidemiology* 2018; **53**: 99-104.
- 32
- 33 35. Openshaw S. The modifiable areal unit problem. In: *Concepts and Techniques*  
34 *in Modern Geography*. Geo Books: Norwich, 1984, pp 1-41.
- 35
- 36 36. Kwan M-P. The uncertain geographic context problem. *Annals of the*  
37 *Association of American Geographers* 2012; **102**(5): 958-968.
- 38
- 39 37. Di Zhu X, Yang Y, Liu X. The importance of housing to the accumulation of  
40 household net wealth. In: Joint Center for Housing Studies, Harvard University,  
41 2003.
- 42

- 1 38. Rehm CD, Moudon AV, Hurvitz PM, Drewnowski A. Residential property values  
2 are associated with obesity among women in King County, WA, USA. *Social*  
3 *Science & Medicine* 2012; **75**(3): 491-495.
- 4
- 5 39. Drewnowski A, Buszkiewicz J, Aggarwal A. Soda, salad, and socioeconomic  
6 status: Findings from the Seattle Obesity Study (SOS). *SSM-Population Health*  
7 2019; **7**: e100339.
- 8
- 9 40. Birkin M, Morris MA, Birkin TM, Lovelace R. Using census data in  
10 microsimulation modelling. In: Stillwell J, Duke-Williams O (eds). *The*  
11 *Routledge Handbook of Census Resources, Methods and Applications*, 1 edn.  
12 Routledge, 2018.
- 13
- 14 41. Jiao J, Drewnowski A, Moudon AV, Aggarwal A, Oppert J-M, Charreire H *et al.*  
15 The impact of area residential property values on self-rated health: A cross-  
16 sectional comparative study of Seattle and Paris. *Preventive Medicine Reports*  
17 2016; **4**: 68-74.
- 18
- 19 42. Nguyen DM, El-Serag HB. The epidemiology of obesity. *Gastroenterology*  
20 *Clinics* 2010; **39**(1): 1-7.
- 21
- 22 43. Pickett KE, Pearl M. Multilevel analyses of neighbourhood socioeconomic  
23 context and health outcomes: a critical review. *Journal of Epidemiology &*  
24 *Community Health* 2001; **55**(2): 111-122.
- 25
- 26 44. Timperio A, Salmon J, Telford A, Crawford D. Perceptions of local  
27 neighbourhood environments and their relationship to childhood overweight  
28 and obesity. *International Journal of Obesity* 2005; **29**(2): 170-175.
- 29
- 30 45. Roda C, Charreire H, Feuillet T, Mackenbach JD, Compernelle S, Glonti K *et*  
31 *al.* Mismatch between perceived and objectively measured environmental  
32 obesogenic features in European neighbourhoods. *Obesity Reviews* 2016;  
33 **17**(S1): 31-41.
- 34
- 35 46. Drewnowski A, Arterburn D, Zane J, Aggarwal A, Gupta S, Hurvitz PM *et al.*  
36 The Moving to Health (M2H) approach to natural experiment research: A  
37 paradigm shift for studies on built environment and health. *SSM-Population*  
38 *Health* 2019; **7**: 100345.
- 39
- 40 47. Bourassa SC, Cantoni E, Hoesli M. Predicting house prices with spatial  
41 dependence a comparison of alternative methods. *The Journal of Real Estate*  
42 *Research* 2010; **32**(2): 139-160.

43

- 1 48. Wilkins EL, Radley D, Morris MA, Griffiths C. Examining the validity and utility  
2 of two secondary sources of food environment data against street audits in  
3 England. *Nutr J* 2017; **16**(82): 1-13.
- 4
- 5 49. Vogel C, Zwolinsky S, Griffiths C, Hobbs M, Henderson E, Wilkins E. A Delphi  
6 study to build consensus on the definition and use of big data in obesity  
7 research. *Int J Obes (Lond)* 2019.
- 8
- 9 50. Nevalainen J, Erkkola M, Saarijarvi H, Nappila T, Fogelholm M. Large-scale  
10 loyalty card data in health research. *Digit Health* 2018; **4**: 2055207618816898.
- 11
- 12 51. Aiello L, Schifanello R, Quercia D, Del Prete L. Large-scale and high-resolution  
13 analysis of food purchases and health outcomes. *EPJ Data Science* 2019;  
14 **8**(14).
- 15
- 16 52. Birkin M, Morris M, Birkin T, Lovelce R. Using census data in microsimulation  
17 modelling. In: Stillwell J, Duke-Williams O (eds). *The Routledge Handbook of*  
18 *Census Resources, Methods and Applications*, 2018.
- 19
- 20 53. Morris MA, Wilkins E, Timmins KA, Bryant M, Birkin M, Griffiths C. Can big data  
21 solve a big problem? Reporting the obesity data landscape in line with the  
22 Foresight obesity system map. *Int J Obes (Lond)* 2018; **42**(12): 1963-1976.
- 23
- 24