



Representativeness and optimal use of body mass index (BMI) in the UK Clinical Practice Research Datalink (CPRD)

Journal:	<i>BMJ Open</i>
Manuscript ID:	bmjopen-2013-003389
Article Type:	Research
Date Submitted by the Author:	11-Jun-2013
Complete List of Authors:	Bhaskaran, Krishnan; LSHTM, NCDE Forbes, Harriet; LSHTM, NCDE Douglas, Ian; LSHTM, NCDE Leon, David; LSHTM, NCDE Smeeth, Liam; LSHTM, NCDE
Primary Subject Heading:	Epidemiology
Secondary Subject Heading:	Epidemiology, Research methods
Keywords:	EPIDEMIOLOGY, PRIMARY CARE, STATISTICS & RESEARCH METHODS

SCHOLARONE™
Manuscripts

Peer Review Only

Representativeness and optimal use of body mass index (BMI) in the UK**Clinical Practice Research Datalink (CPRD)**

Authors: Krishnan Bhaskaran, Harriet J Forbes, Ian Douglas, David A Leon, Liam Smeeth

Faculty of Epidemiology & Population Health, London School of Hygiene and Tropical

Medicine, Keppel Street, London, WC1E 7HT, UK

Corresponding Author: Krishnan Bhaskaran, London School of Hygiene and Tropical

Medicine, Keppel Street, London, WC1E 7HT. Email: krishnan.bhaskaran@lshtm.ac.uk. Tel: +

44 (0) 20 7927 2268

Key words: CPRD, BMI, missing data, primary care databases, obesity.

Word Count: 3383

Abstract

Objectives: To assess the completeness and representativeness of body mass index (BMI) data in the Clinical Practice Research Datalink (CPRD), and determine an optimal strategy for their use.

Design: Descriptive study.

Setting: Electronic healthcare records from primary care.

Participants: A million patient random sample from the UK Clinical Practice Research Datalink (CPRD) primary care database, aged ≥ 16 years.

Primary and secondary outcome measures: BMI completeness in CPRD was evaluated by age, sex, and calendar period. CPRD-based summary BMI statistics for each calendar year (2003-10) were age- and sex-standardised and compared with equivalent statistics from the Health Survey for England (HSE).

Results: BMI completeness increased over calendar time from 37% in 1990-94 to 77% in 2005-11, was higher among females, and increased with age. When BMI at specific time points was assigned based on the most recent record, calendar year-specific mean BMI statistics underestimated equivalent HSE statistics by 0.75-1.1kg/m². Restricting to those with a recent (≤ 3 years) BMI resulted in mean BMI estimates closer to HSE (≤ 0.28 kg/m² underestimation), but excluded up to 47% of patients. An alternative strategy of imputing up-to-date BMI based on modelled changes in BMI over time since the last available record, also led to mean BMI estimates that were close to HSE (≤ 0.37 kg/m² underestimation).

Conclusions: Completeness of BMI in CPRD increased over time and varied by age and sex. At a given point in time, a large proportion of the most recent BMIs are unlikely to reflect current BMI; consequent BMI misclassification might be reduced by employing model-based imputation of current BMI.

Article summary

Article focus:

- Body mass index (BMI) data are frequently used in epidemiological analyses of primary care databases such as the UK Clinical Practice Research Datalink (CPRD), however their completeness and representativeness have not previously been assessed in detail.
- The aim of this article is to provide information on the completeness of BMI in CPRD primary care data, on their representativeness, and on the implications for their practical use in research.

Key messages:

- We found that completeness of BMI recordings in the Clinical Practice Research Datalink increased from 37% in 1990-4 to 77% in 2005-11 and differed by age and sex.
- At specific calendar time points, the most recent BMI recorded for a large proportion of patients was over 3 years old and was unlikely to reflect current BMI.
- The optimal strategy for assigning BMI status is likely to depend on the specific study population and research context. We suggest one possible approach that uses a model-based imputation of current BMI to reduce BMI misclassification.

Strengths and limitations of this study:

- Results presented here are based on a large random sample from the CPRD, therefore we can confidently generalise the findings to the whole CPRD database, and to similar databases based on UK primary care records.
- To assess the representativeness of CPRD BMI data, we compared with data from the Health Survey for England, which is based on a large nationally representative sample and includes BMI information measured by trained interviewers.
- Our study did not look at BMI recordings among children as this would require a different strategy.

Introduction

Overweight and obesity are major contributors to global disease burden[1] and are associated with substantial excess mortality[2]. The prevalence of obesity is increasing in both developed and developing countries[3, 4] and is a growing concern to policy makers. In England, the prevalence of obesity rose steadily from 1993 to 2010: from 13% to 26% in men, and from 16% to 26% in women[5]. Because of its association with various diseases and health outcomes, body mass index (BMI, the metric most widely used to classify overweight and obesity) is an important factor in many epidemiological studies, both as an exposure and as a potential confounder.

Databases of routinely collected electronic healthcare records are becoming an increasingly valuable resource in epidemiology, allowing population-level research on large, representative samples. The UK Clinical Practice Research Datalink (CPRD) (formerly the General Practice Research Database or GPRD) is widely used and contains medical records for approximately 8% of the UK population.[6] However, a shortcoming of these databases is that lifestyle data, such as BMI, tend to be opportunistically recorded and can be incomplete. Furthermore, those with non-missing lifestyle data may be unrepresentative of the general population. BMI has been an important covariate in many published studies based on CPRD[7-14] but the completeness and representativeness of the BMI data have not been previously documented.

Our aim was to undertake an in-depth investigation of BMI recordings in CPRD, including quantifying the completeness of BMI data, and assessing their representativeness by comparing summary statistics based on CPRD data with equivalent statistics from a representative general population survey.

Methods

Data sources

Clinical Practice Research Datalink (CPRD)

The Clinical Practice Research Datalink (CPRD) is a clinical database comprising anonymised computerised medical records from general practitioners (GPs) in the United Kingdom. Approximately 8% of the UK population are currently included and the database is broadly representative of the UK population.[15] CPRD contains demographic information, clinically relevant lifestyle data, prescription details, clinical events, preventive care provided, specialist referrals, and hospital admissions and their major outcomes. Data undergo quality checks and practices are designated as “up to standard” in CPRD from the date that they meet specified data entry quality criteria. For this study, we obtained a random sample of one million CPRD patients, because carrying out the analysis on the full CPRD database would be computationally difficult, and the reduction in precision of our estimates that would arise by restricting our analysis to a one million random sample is extremely small.

Body mass index data in CPRD

Height and weight measurements are recorded in CPRD whenever measured as part of routine care. We obtained all height and weight records and calculated BMI (BMI=weight/height²). Patient records without any measurements or with implausible measurements were excluded (Figure 1).

Health Survey for England

We obtained published Health Survey for England (HSE) data for BMI from the National Health Service (NHS) Information Centre.[16] The HSE is an annual survey designed to produce a representative sample of the adult population aged ≥16 years and living in private households. The methods are described in detail elsewhere.[17] Surveys were interviewer

1 administered with interviewers measuring the weight and height of all participants. Data
2
3 from 2003-10 were obtained, and these data have been weighted to reduce bias from non-
4
5 response, based on a logistic regression model incorporating age, sex, household type
6
7 (based on the number of adults and children living in a household), Strategic Health
8
9 Authority region, and social class (defined using the National Statistics Socio-economic
10
11 Classification system).
12
13

14 **Statistical methods**

15 **Completeness of BMI data in CPRD**

16
17 In the main analyses BMI completeness data in CPRD were estimated by calendar period
18
19 (1990-4, 1995-9, 2000-4, 2005-11). To calculate completeness for a particular calendar
20
21 period, all individuals from the one million sample who were registered, aged ≥ 16 years, and
22
23 under follow-up in “up to standard” practices on the mid-point of the period were identified
24
25 and included in the denominator. Among these individuals, the numerator comprised either
26
27 those with any previous BMI available in their electronic record regardless of how long ago
28
29 it was entered, or those with a BMI available up to 3 years prior to this date. Completeness
30
31 data were generated by age group, sex and among those whom, for clinical reasons, BMI
32
33 should be routinely monitored (those with type 2 diabetes, schizophrenia/other psychoses,
34
35 and ≥ 2 recent (last 6 months) statin prescriptions). We also investigated whether
36
37 completeness could be improved by searching for clinical codes (“Read codes”) indicating
38
39 BMI category. We have not presented confidence intervals for these descriptive statistics
40
41 because the sample size made sampling error negligible (for example, the standard errors
42
43 for the proportions with complete BMI data in age and calendar year subgroups were all
44
45 $< 0.5\%$).
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Comparison of CPRD BMI data with Health Survey for England data

We compared mean BMI over calendar time based on complete CPRD BMI data with equivalent HSE figures, for the period 2003-2010 (since, from 2003, HSE data were adjusted for non-response). CPRD mean BMI was based on patients registered and under up-to-standard follow-up at the mid-point of the calendar year. We produced two sets of CPRD mean BMI statistics: firstly we used last BMI observation carried forward (regardless of how long ago recorded); secondly we restricted to patients with a recent BMI available (up to 3 years before the mid-point of the calendar year). As above, confidence intervals are not presented because there was negligible sampling error (maximum standard error=0.02kg/m²). To make like-with-like comparisons with HSE, CPRD data were restricted to English practices, and mean BMI was age- and sex-standardised to the HSE population structure. Proportions classified as obese (BMI≥30kg/m²) over time based on CPRD and HSE data were also compared.

Model-based imputation of up-to-date BMI measures in CPRD

We explored whether outdated BMI measures in CPRD could be usefully updated by imputation based on a model predicting changes in individual-level BMI over time. We used data from individuals with multiple BMI records to model the expected change in BMI as a function of time since BMI recording (restricting to individuals with BMI records ≤ 10 years apart). We fitted a linear regression model with change in BMI as the outcome, and elapsed time included as a 3 knot cubic spline to allow for non-linearity; we also included interactions between the spline basis variables and indicator variables for age and sex.

Feasible weighted least squares estimation was used to allow for heteroskedasticity.[18]

Having specified a model for change in BMI over time, we first explored its performance among individuals with at least 2 BMIs entered in CPRD, by predicting the most recent BMI

1 based on the previous BMI record and the elapsed time; we compared the distribution of
2
3 the errors from this approach with the distribution of the errors from simply using the last
4
5 observation carried forward. We then repeated the comparison with the HSE mean BMI
6
7 data for each calendar year. This time we included all individuals with a BMI record in the
8
9 previous 10 years and used the model described above to impute current BMI at the mid-
10
11 point of the calendar year by predicting the change in BMI since the last available BMI
12
13 record. We did this within a multiple imputation framework (using 5 imputations) to
14
15 account for uncertainty in the modelled changes over time.[19]
16
17
18
19
20

21 The study was approved by the London School of Hygiene and Tropical Medicine Ethics
22
23 Committee.
24
25

26 **Results**

27 **Completeness of BMI data in CPRD**

28
29 In 1990-1994, 37% of individuals had at least one previously recorded BMI, and the
30
31 proportion increased to 77% by 2005-11 (Table 1). The proportion of individuals with a recent
32
33 BMI (recorded in the previous 3 years) was lower in each calendar period (35% in 1990-1994
34
35 rising to 51% in 2005-11). BMI completeness generally increased with age up to 75 years,
36
37 with a lower proportion in the oldest age group having data available. Data for single
38
39 calendar years are shown in Appendix Table A1 and illustrate similar patterns. BMI data
40
41 appeared to be consistently more widely available among women than men (Figure 2). As
42
43 expected, BMI completeness was higher in particular clinical subgroups: in total 97% of
44
45 patients with a record of type II diabetes had a recent BMI recorded, along with over 78% of
46
47 those with a diagnosis of schizophrenia/psychoses (Appendix Table A2). This is in line with
48
49 Quality and Outcomes Framework (QOF) which has encouraged BMI monitoring in these
50
51
52
53
54
55
56
57
58
59
60

1 conditions since 2004.[20] BMI completeness was also high among current statin users (82%
2
3
4 with a recent BMI available).

5
6
7 There was little extra information available in clinical (“Read”) codes relating to BMI. In the
8
9 most recent calendar period, out of 75518 individuals with no previous BMI record
10
11 available, only 1222 (1.6%) had ever had a clinical code that would enable classification into
12
13 BMI categories (underweight, normal, overweight/obese). Furthermore, for those with a
14
15 previous BMI, only a small proportion had more recent information related to BMI recorded
16
17 in a clinical code (7675/250430 = 3.0% in the most recent period).
18
19

20 21 22 **Summary statistics using complete CPRD BMI data and comparison with Health Survey for** 23 24 **England**

25
26 We found that age- and sex-standardised mean BMI based on CPRD data was consistently
27
28 and substantially lower (by up to 1.1kg/m²) than in the HSE data (mean BMI in CPRD =
29
30 25.7kg/m² in 2003 rising to 26.3 in 2010, compared with 26.8 kg/m² [95% CI 26.7 to 26.9]
31
32 and 27.3 [27.1 to 27.5] respectively in HSE; Figure 3).
33
34

35
36
37 When BMI entries more than 3 years old were discarded, between 33 to 47% of patients
38
39 were lost across calendar years. However, the estimated mean BMI in CPRD was
40
41 considerably closer to what would be expected based on the HSE data, with CPRD data
42
43 underestimating the HSE statistics by only between 0.04 to 0.28kg/m² in individual calendar
44
45 years, and the CPRD estimate falling within the HSE confidence interval for 2 of the most
46
47 recent 3 calendar years (mean BMI in CPRD = 26.9, 27.0 and 27.0 kg/m² compared with 27.0
48
49 [26.9 to 27.1], 27.0 [26.8 to 27.2] and 27.3 [27.1 to 27.5] in HSE, in 2008, 2009 and 2010
50
51 respectively). Age- and sex-stratified data demonstrated similar patterns, except that in the
52
53 eldest age group (75+ years), restricting to those with recent BMI measures did not bring
54
55 the estimated BMI substantially closer to HSE figures (Appendix Figure A1).
56
57
58
59
60

1
2
3
4 We also compared the proportions classified as obese between CPRD and HSE (Appendix
5
6 Figure A2). Consistent with the previous analysis, using any previous BMI reading to classify
7
8 individuals in CPRD resulted in lower obesity rates than expected based on HSE data, while
9
10 restricting to patients with a recent reading led to estimated obesity rates close to those in
11
12 HSE.
13
14

15 16 17 **Model-based imputation of up-to-date BMI measures in CPRD** 18

19 The contrast between BMI summary statistics based on recent measures and those based
20
21 on any previous measures suggested that older BMI records were tending to underestimate
22
23 current BMI. We therefore examined whether a model could be developed to impute
24
25 current BMI, taking into account elapsed time since the last measure. In a linear regression
26
27 model for change in BMI over time, we estimated that on average BMI increased over the
28
29 10-year period following a BMI record for those aged up to 69 years at the time of the
30
31 record and decreased over time in those aged 70 years or more (Appendix Figure A3). We
32
33 tested the predictive performance of our model by predicting the most recent BMI based on
34
35 the previous one, among CPRD patients with more than one recorded BMI available. When
36
37 the older BMI was less than 3 years old, there was little gain in applying the correction
38
39 compared with carrying the older observation forward (Figure 4). However, when there was
40
41 a longer gap, carrying the previous BMI forward tended to underestimate the later BMI,
42
43 while employing the model-based imputation removed the underestimation and led to
44
45 smaller errors on average (median error = -0.70kg/m^2 [IQR -2.18 to $+0.56$] using last
46
47 observation carried forward, compared with $+0.11\text{kg/m}^2$ [-1.29 to $+1.40$] using model-based
48
49 imputation).
50
51
52
53
54
55
56
57
58
59
60

1 We then repeated the comparison of mean BMI in CPRD versus HSE, this time using our
2
3 model for change in BMI over time as a basis for performing multiple imputations of current
4
5 BMI based on the latest available measure and the time since it was recorded. Estimated
6
7 mean BMIs were now in line with those based on only recent data in the earlier analysis,
8
9 and were only between 0.04 and 0.37kg/m² lower than HSE statistics in individual calendar
10
11 years (Figure 3, circles). Even with multiple imputation, confidence intervals remained
12
13 extremely narrow (<0.07kg/m²) due to the large sample size, so are not shown in the figure.
14
15 Of note, all patients with a BMI recorded up to 10 years before the midpoint of the calendar
16
17 year of interest were now included in the estimation of the “corrected” means; thus in
18
19 individual calendar years only 9 to 13% of patients were dropped, compared to 33-47% of
20
21 patients when dropping BMI records >3 years old.
22
23
24
25
26
27

28 Discussion

30 Main findings

31
32 BMI completeness has increased over calendar time (rising from 37% in 1990-94 to 77% in
33
34 2005-11). Completeness was higher among females, older age groups, and clinical
35
36 subgroups where recording BMI is encouraged. When BMI on the date of interest was
37
38 assigned to individual patients in CPRD using the last available record, regardless of how
39
40 long ago it was entered, we found that the resulting mean BMI statistics for the CPRD
41
42 population were consistently lower than equivalent HSE estimates (by up to 1.1kg/m²). This
43
44 appeared to be driven by older BMI records tending to systematically underestimate current
45
46 BMI: when only recent CPRD BMI records (≤3 years old) were used, mean BMI statistics
47
48 were closer to HSE estimates. However, a substantial number of patients were then
49
50 excluded altogether (33-47% across years). Finally, we suggested a process for modelling
51
52
53
54
55
56
57
58
59
60

1 changes in BMI after a BMI record, which could allow researchers to impute BMI on the date
2
3
4 of interest and avoid dropping large numbers without a recent measure from their analyses.
5
6

7 *Comparison with other studies*

8
9 There are very few comparable studies (Appendix Table A2). However, the proportion of
10
11 patients with a recent BMI recording in CPRD is in line with a summary of the QRESEARCH
12
13 database (a similar UK primary care database with data from over 530 general practices
14
15 using EMIS software rather than VISION software);[21] by March 2007, 58% of registered
16
17 patients aged 16+ years had their BMI recorded in the past 5 years; this compares with 51%
18
19 with a BMI recorded in the last 3 years in our analysis (for 2005-11). As in our study, the
20
21 QRESEARCH report showed an increase in completeness over time, rising from 42% in
22
23 2000/01 to 58% in 2007. In a third UK primary care database, THIN (The Health
24
25 Improvement Network), the proportion of newly registered patients between 2004 and
26
27 2006 with BMI data was in line with our findings; 62% of patients had a height recording and
28
29 66% had a weight recording within 12 months of registration.[22]
30
31
32
33
34
35

36 **Explanation of findings**

37 *Completeness*

38
39 Increasing completeness of BMI over time may reflect a general trend towards
40
41 encouragement to record BMI in primary care. Greater BMI completeness among females
42
43 and older age groups may have a number of explanations including higher consultation
44
45 rates in primary care [23, 24] and different prevalence's of diseases in which it is important
46
47 to monitor BMI.
48
49
50
51
52

53 *Comparison of CPRD BMI data with Health Survey for England data*

54
55 Mean BMI based on the CPRD population was lower in each calendar year than equivalent
56
57 HSE estimates when BMI in CPRD was assigned using the last available record; however,
58
59
60

1 when the analysis was restricted to those with a recent BMI record, estimates from CPRD
2
3 were close to HSE estimates. This suggests that the substantial proportion of BMI recordings
4
5 in CPRD that were outdated on the date of interest may have driven the apparent
6
7 underestimation of mean BMI in CPRD in the unrestricted analysis. This in turn would imply
8
9 that individual BMIs tend to increase over time, and indeed when we specifically modelled
10
11 changes in BMI over time, we found a pattern of increasing BMI with age for those <70
12
13 years old, consistent with prospective cohort studies with repeated BMI measurements [25-
14
15 27]. A simple adjustment of outdated BMIs based on these modelled changes over time
16
17 brought CPRD mean BMI statistics in line with HSE estimates, and when we validated the
18
19 adjustment in a subset of patients with repeated BMI measures, we found smaller errors on
20
21 average, compared with simply carrying outdated BMI records forwards.
22
23
24
25
26
27

28 Of note, we observed that CPRD consistently underestimated BMI compared to HSE among
29
30 those aged ≥ 75 years, even when only recent records were used; this may reflect the fact
31
32 that institutionalised patients are represented in CPRD but not in HSE: HSE may not be an
33
34 ideal comparison for this age group since elderly people in institutions (who are represented
35
36 in CPRD) may be more likely to be frail and have lower BMIs than those living in private
37
38 households.
39
40
41
42

43 **Implications**

44 First, our findings suggest that BMI completeness is likely to vary between studies
45
46 depending on the study population and study period. BMI data are not likely to be missing
47
48 completely at random (for example, missingness may vary by patient characteristics or
49
50 particular diseases). There may be information in the database, however, which predicts
51
52 missingness and which could satisfy the “missing at random” assumption required for
53
54 multiple imputation. A study exploring the potential of imputing missing data in THIN found
55
56
57
58
59
60

1 that after multiple imputation, summary statistics of height and weight were comparable
2
3 with data from nationally representative datasets.[22]
4

5
6 Second, our analyses suggest that the common practice of assigning BMI status based on
7
8 the nearest/most recent available record to the index date of interest might lead to
9
10 misclassification, given that a large number of patients have only substantially outdated BMI
11
12 records available at any particular time. Strategies to address this include restricting to
13
14 recent BMI, but this is likely to exclude a large numbers of patients. We have suggested an
15
16 alternative strategy based on updating the outdated BMIs by modelling changes in BMI over
17
18 time, though this is not without drawbacks: the approach requires an assumption that
19
20 individuals with ≥ 2 BMI records available (needed to estimate the model for changes over
21
22 time) are representative of the wider patient population, which may not be the case; it is
23
24 also a more complex strategy, particularly if done within a multiple imputation framework
25
26 to allow for uncertainty in the correction, which could be substantial in studies with smaller
27
28 sample sizes than considered here. Ultimately, the importance of these issues and the
29
30 optimal strategy to use is likely to depend on the particular study and the characteristics of
31
32 the study population.
33
34
35
36
37
38

39 **Strengths and Limitations**

40
41 Results presented here are based on a large random sample from the CPRD, therefore we
42
43 can confidently generalise the findings to the whole CPRD database. Although we cannot
44
45 assume these findings will relate to other routinely collected primary care databases in UK
46
47 based on other IT systems (CPRD is based on practices using VISION), they are likely to be
48
49 similar. This study did not look at BMI recordings among children as this would require a
50
51 different strategy. Completeness among 16-24 year age group may be artificially low
52
53 because weights recorded at age < 16 were excluded, so those at the lower end of the age
54
55
56
57
58
59
60

1 group will not have had as much time to accrue weight recordings. We believe HSE to be the
2
3 best available comparison for this study; it is a nationally representative, large sample
4
5 (sample size 14,836 in 2003 and 8,420 in 2010), utilising height and weight recordings
6
7 measured by a trained interviewer, and is weighted for non-response.[17, 28] However
8
9 there is a degree of missing data in HSE which is a limitation. In 2010 just over 85% of adults
10
11 interviewed provided valid height and weight recordings. [29] One of the most common
12
13 reasons for missing BMI was refusal (up to 8% were missing due to refusal),[17] which if
14
15 related to BMI status, may bias the estimates of mean BMI in HSE.
16
17
18
19

20 21 **Conclusions**

22
23 Completeness of BMI data in CPRD varies over time and by age and sex. BMI records may
24
25 become outdated over time and naive use could lead to misclassification of BMI status. The
26
27 optimal strategy for assigning BMI status to individuals in studies based on CPRD and similar
28
29 electronic healthcare databases is likely to depend on the specific study population and the
30
31 research context.
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

Conflicts of interest

The authors declare no conflicts of interest.

Funding

This report is independent research arising from a postdoctoral fellowship (for KB) supported by the National Institute for Health Research (PDF-2011-04-007). ID is supported by an MRC methodology research fellowship, LS is supported by a Wellcome Trust senior research fellowship in clinical science.

Data sharing statement

This analysis is based on a large random sample from the Clinical Practice Research Datalink, provided by the UK Medicines and Healthcare products Regulatory Agency. The authors' licence for using these data does not allow sharing of raw data with third parties.

peer review only

References

1. World Health Organisation. Global Health Risks: Mortality and burden of disease attributable to selected major risks. Geneva, Switzerland: 2009.
2. Flegal KM, Graubard BI, Williamson DF, Gail MH. Excess deaths associated with underweight, overweight, and obesity. *Jama-J Am Med Assoc* 2005;293(15):1861-7.
3. Swinburn BA, Sacks G, Hall KD, McPherson K, Finegood DT, Moodie ML, et al. Obesity 1 The global obesity pandemic: shaped by global drivers and local environments. *Lancet* 2011;378(9793):804-14.
4. Kelly T, Yang W, Chen CS, Reynolds K, He J. Global burden of obesity in 2005 and projections to 2030. *Int J Obesity* 2008;32(9):1431-7.
5. NHS Information Centre. Health survey for England - 2010: health and lifestyles 2011. Available from: <http://www.ic.nhs.uk/pubs/hse10report>.
6. CPRD. Clinical Practice Research Database (CPRD) website. Available from: <http://www.cprd.com/intro.asp>.
7. Delaney JA, Daskalopoulou SS, Brophy JM, Steele RJ, Opatrny L, Suissa S. Lifestyle variables and the risk of myocardial infarction in the general practice research database. *BMC Cardiovasc Disord* 2007;7:38.
8. Green J, Czanner G, Reeves G, Watson J, Wise L, Beral V. Oral bisphosphonates and risk of cancer of oesophagus, stomach, and colorectum: case-control analysis within a UK primary care cohort. *Brit Med J* 2010;341.
9. Tzoulaki I, Molokhia M, Curcin V, Little MP, Millett CJ, Ng A, et al. Risk of cardiovascular disease and all cause mortality among patients with type 2 diabetes prescribed oral antidiabetes drugs: retrospective cohort study using UK general practice research database. *Brit Med J* 2009;339:b4731.
10. Douglas I, Smeeth L, Irvine D. The use of antidepressants and the risk of haemorrhagic stroke: a nested case control study. *Br J Clin Pharmacol* 2011;71(1):116-20.
11. Andersohn F, Schade R, Suissa S, Garbe E. Long-Term Use of Antidepressants for Depressive Disorders and the Risk of Diabetes Mellitus. *Am J Psychiat* 2009;166(5):591-8.
12. Lawrenson R, Todd JC, Leydon GM, Williams TJ, Farmer RDT. Validation of the diagnosis of venous thromboembolism in general practice database studies. *Br J Clin Pharmacol* 2000;49(6):591-6.
13. Jick H, Zornberg GL, Jick SS, Seshadri S, Drachman DA. Statins and the risk of dementia. *Lancet* 2000;356(9242):1627-31.
14. van Staa TP, Wegman S, de Vries F, Leufkens B, Cooper C. Use of statins and risk of fractures. *Jama-J Am Med Assoc* 2001;285(14):1850-5.
15. Office for National Statistics. Key Health Statistics from General Practice 1998: Analyses of Morbidity and Treatment Data, Including Time Trends, England and Wales. Series MB6 No. 2, editor. London: Office for National Statistics; 2000.
16. NHS Information Centre. Health Survey for England - 2010: Trend tables: National Health Service; [cited on 19th June 2012]. Available from: <http://www.ic.nhs.uk/statistics-and-data-collections/health-and-lifestyles-related-surveys/health-survey-for-england/health-survey-for-england--2010-trend-tables>.
17. Aresu M, Boodhna G, Bryson A, Bridges S, Chaudhury M, Craig R, et al. Volume 2: Methods and documentation. In: Craig R, Mindell J, editors. Health Survey for England 2010. Leeds: NHS Information Centre for health and social care; 2011.
18. Greene WH. *Econometric Analysis*. Upper Saddle River, New Jersey: Prentice Hall.1997.

- 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
19. Rubin DB. Multiple Imputation for Nonresponse in Surveys: J. Wiley & Sons, New York. ; 1987.
 20. NHS. Quality and Outcomes Framework guidance for GMS contract 2011/12. Employers and British Medical Association; 2011.
 21. NHS Information Centre. A summary of public health indicators using electronic data from primary care 2008 [cited on 17th December 2012]. Available from: <http://www.ic.nhs.uk/article/2021/Website-Search?productid=4287&q=A+summary+of+public+health+indicators+using+electronic+data+from+primary+care.&sort=Relevance&size=10&page=1&area=both#top>.
 22. Marston L, Carpenter JR, Walters KR, Morris RW, Nazareth I, Petersen I. Issues in multiple imputation of missing data for large general practice clinical databases. *Pharmacoepidem Dr S* 2010;19(6):618-26.
 23. Rowlands S, Moser K. Consultation rates from the General Practice Research Database. *Br J Gen Pract* 2002;52(481):658-60.
 24. The Health and Social Care Information Centre. Trends in Consultation Rates in General Practice 1995 to 2009 2009 [cited on 17th December 2012]. Available from: <http://www.ic.nhs.uk/article/2021/Website-Search?productid=729&q=gresearch&sort=Relevance&size=10&page=1&area=both#top>.
 25. Li L, Law C, Power C. Body mass index throughout the life-course and blood pressure in mid-adult life: a birth cohort study. *J Hypertens* 2007;25(6):1215-23.
 26. Silverwood RJ, Pierce M, Thomas C, Hardy R, Ferro C, Sattar N, et al. Overweight across adult life and kidney function at age 60-4 years: the 1946 British birth cohort study. (awaiting publication) 2012.
 27. Tirosh A, Shai I, Afek A, Dubnov-Raz G, Ayalon N, Gordon B, et al. Adolescent BMI Trajectory and Risk of Diabetes versus Coronary Disease. *N Engl J Med* 2011;364(14):1315-25.
 28. Mindell J, Biddulph JP, Hirani V, Stamatakis E, Craig R, Nunn S, et al. Cohort Profile: The Health Survey for England. *Int J Epidemiol* 2012:1-9.

Table 1: Completeness of BMI data in the CPRD, by age and calendar period

Age group (yrs)	1990-4	1995-9	2000-4	2005-2011
16-24^a				
N registered	11423	17501	34452	42546
BMI in previous 3y (%)	26	28	25	32
BMI any previous (%)	26	37	30	37
25-34				
N registered	17477	29923	48659	50413
BMI in previous 3y (%)	37	39	36	49
BMI any previous (%)	38	66	67	72
35-44				
N registered	15953	28838	55991	61014
BMI in previous 3y (%)	36	36	31	46
BMI any previous (%)	39	67	71	80
45-54				
N registered	14507	27765	48093	55564
BMI in previous 3y (%)	39	37	32	50
BMI any previous (%)	42	70	73	84
55-64				
N registered	11680	20843	42258	49380
BMI in previous 3y (%)	42	40	37	57
BMI any previous (%)	44	74	77	87
65-74				
N registered	10678	17605	30997	34508
BMI in previous 3y (%)	36	37	40	67
BMI any previous (%)	38	71	79	91
75+				
N registered	8637	16005	29384	32523
BMI in previous 3y (%)	28	32	37	64
BMI any previous (%)	28	56	69	87
Total				
N registered	90355	158480	289834	325948
BMI in previous 3y (%)	35	36	34	51
BMI any previous (%)	37	64	67	77

N registered is all those under follow-up at mid-point of the period

^aNote, BMI measurements from age <16 years were not counted in this analysis, hence completeness in the 16-24 age group may be artificially low

Figure 1: Initial data processing to generate BMI for analysis

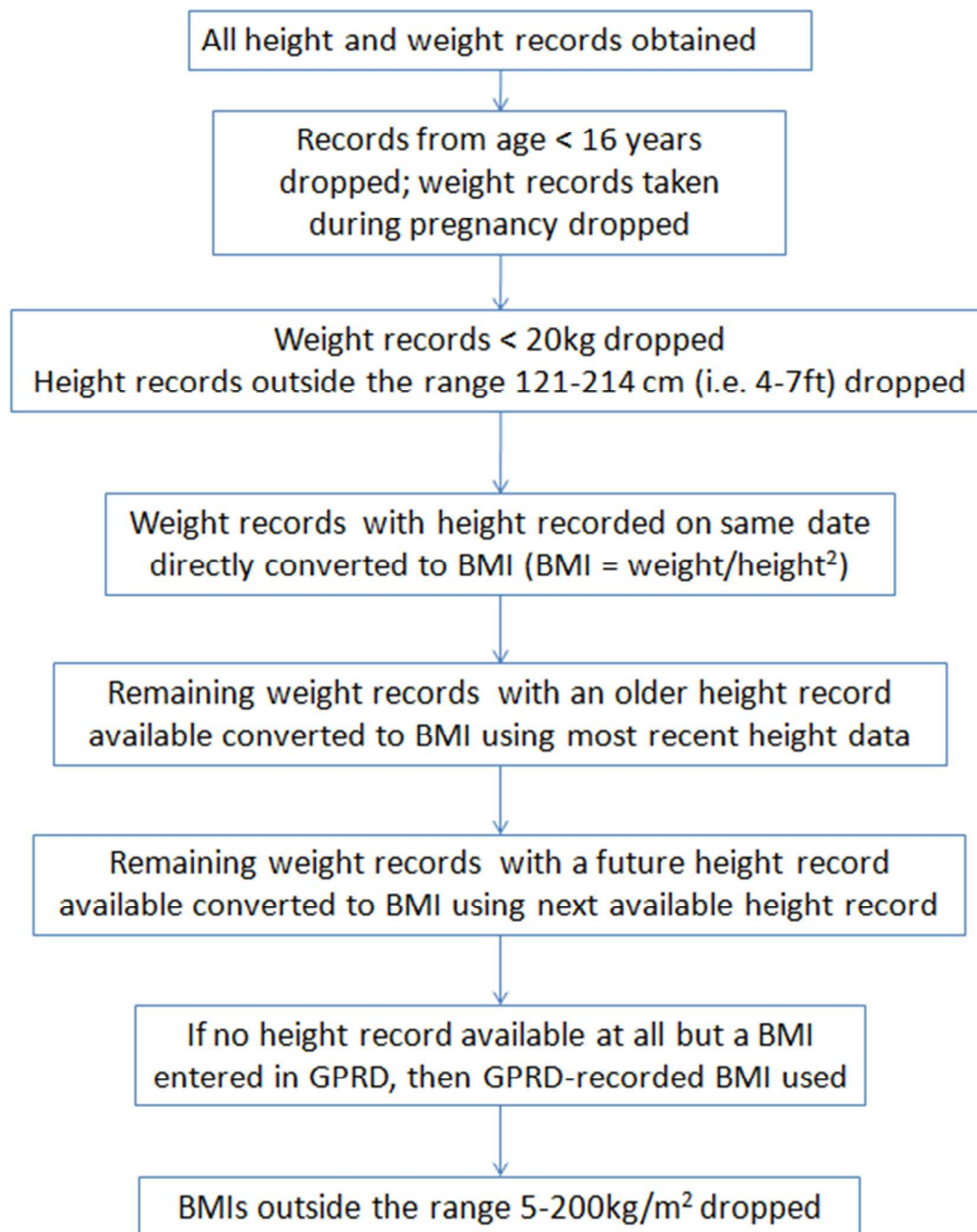
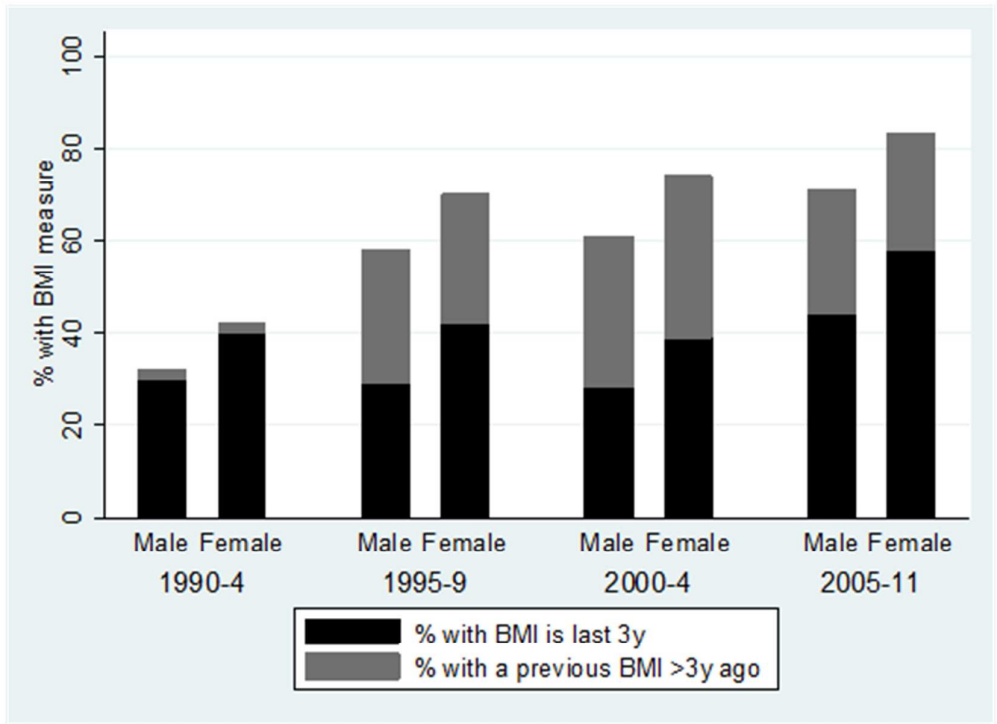


Figure 2: Completeness of BMI data in CPRD, by gender and calendar period

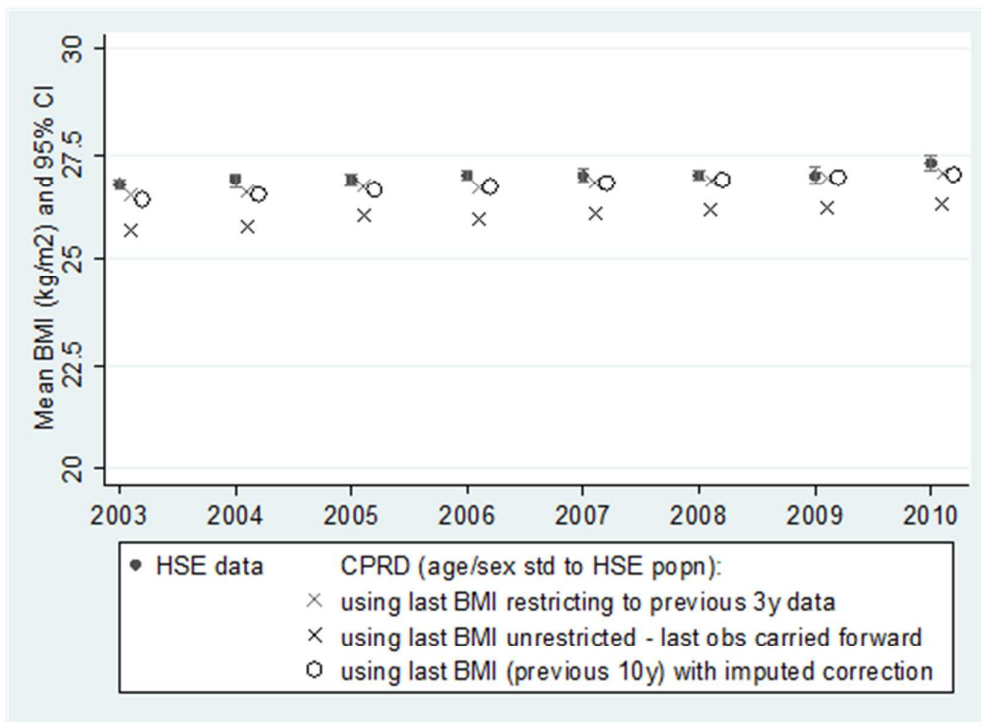


Note: Completeness data for each calendar period are based on all those under follow-up at mid-point of the period

For peer review only

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 3: Mean BMI over calendar time comparing those with BMI recorded in CPRD (English practices) with the Health Survey for England 2010 data

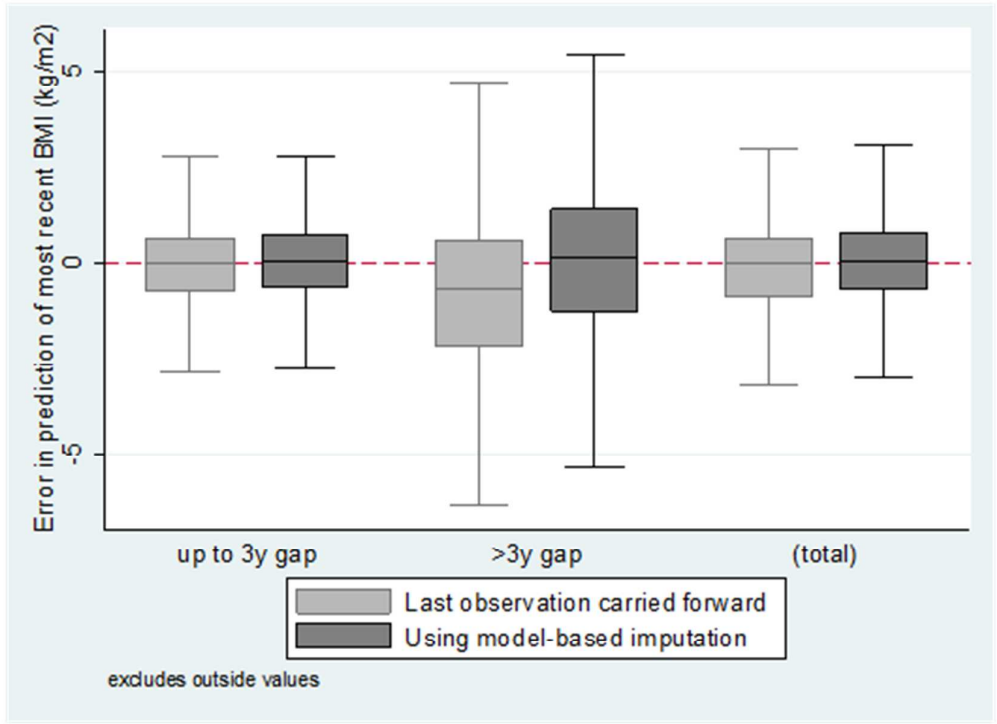


Note: CPRD figures are age- and sex- standardised to the Health Survey for England study population
 CPRD statistics are based on all patients registered at the mid-point of the calendar period and with a suitable previous BMI measure available (i.e. either any previous, or within the last 3 years)

Review only

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 4: Error in prediction of most recent BMI from older BMI, comparing simple last observation carried forward with model-based imputation of up to date BMI – stratified by time gap between readings



view only

Author contributions

I, Krishnan Bhaskaran, developed the analytical strategy for this paper, processed and analysed the data and wrote the paper.

I, Harriet Forbes, was involved in discussing the data processing and analysis of the data, as well as the writing of the paper.

I, Liam Smeeth, was involved in discussions of the analytical approach to this study and made comments on the analysis and the writing of the paper.

I, Ian Douglas, was involved in discussions of the analytical approach to this study and made comments on the analysis and the writing of the paper.

I, David Leon, was involved in discussions of the analytical approach to this study and made comments on the analysis and the writing of the paper.

STROBE Statement—Checklist of items that should be included in reports of *cohort studies*

	Item No	Recommendation
Title and abstract	1	(a) Indicate the study's design with a commonly used term in the title or the abstract n/a (we did not think there was an appropriate design keyword/term to describe this study as it is not a standard "exposure/outcome" study but is rather providing data quality information on a common exposure/covariate) (b) Provide in the abstract an informative and balanced summary of what was done and what was found P2
Introduction		
Background/rationale	2	Explain the scientific background and rationale for the investigation being reported P4
Objectives	3	State specific objectives, including any prespecified hypotheses P4
Methods		
Study design	4	Present key elements of study design early in the paper P6-7
Setting	5	Describe the setting, locations, and relevant dates, including periods of recruitment, exposure, follow-up, and data collection P5-6
Participants	6	(a) Give the eligibility criteria, and the sources and methods of selection of participants. Describe methods of follow-up P5-6 (b) For matched studies, give matching criteria and number of exposed and unexposed
Variables	7	Clearly define all outcomes, exposures, predictors, potential confounders, and effect modifiers. Give diagnostic criteria, if applicable P5-7
Data sources/ measurement	8*	For each variable of interest, give sources of data and details of methods of assessment (measurement). Describe comparability of assessment methods if there is more than one group P5-6
Bias	9	Describe any efforts to address potential sources of bias P6-7
Study size	10	Explain how the study size was arrived at P5
Quantitative variables	11	Explain how quantitative variables were handled in the analyses. If applicable, describe which groupings were chosen and why P6-7
Statistical methods	12	(a) Describe all statistical methods, including those used to control for confounding P6-7 (b) Describe any methods used to examine subgroups and interactions P6-7 (c) Explain how missing data were addressed P7 (d) If applicable, explain how loss to follow-up was addressed n/a (e) Describe any sensitivity analyses

		n/a
Results		
Participants	13*	(a) Report numbers of individuals at each stage of study—eg numbers potentially eligible, examined for eligibility, confirmed eligible, included in the study, completing follow-up, and analysed (b) Give reasons for non-participation at each stage (c) Consider use of a flow diagram
		FIG 1
Descriptive data	14*	(a) Give characteristics of study participants (eg demographic, clinical, social) and information on exposures and potential confounders (b) Indicate number of participants with missing data for each variable of interest (c) Summarise follow-up time (eg, average and total amount)
		P8-9 and FIG 2
Outcome data	15*	Report numbers of outcome events or summary measures over time n/a (no specific outcome)
Main results	16	(a) Give unadjusted estimates and, if applicable, confounder-adjusted estimates and their precision (eg, 95% confidence interval). Make clear which confounders were adjusted for and why they were included (b) Report category boundaries when continuous variables were categorized (c) If relevant, consider translating estimates of relative risk into absolute risk for a meaningful time period n/a (not an “exposure/outcome” study)
Other analyses	17	Report other analyses done—eg analyses of subgroups and interactions, and sensitivity analyses P9-11
Discussion		
Key results	18	Summarise key results with reference to study objectives P11-12
Limitations	19	Discuss limitations of the study, taking into account sources of potential bias or imprecision. Discuss both direction and magnitude of any potential bias P14-15
Interpretation	20	Give a cautious overall interpretation of results considering objectives, limitations, multiplicity of analyses, results from similar studies, and other relevant evidence P15
Generalisability	21	Discuss the generalisability (external validity) of the study results P14
Other information		
Funding	22	Give the source of funding and the role of the funders for the present study and, if applicable, for the original study on which the present article is based P16

*Give information separately for exposed and unexposed groups.

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

Note: An Explanation and Elaboration article discusses each checklist item and gives methodological background and published examples of transparent reporting. The STROBE checklist is best used in conjunction with this article (freely available on the Web sites of PLoS Medicine at <http://www.plosmedicine.org/>, Annals of Internal Medicine at <http://www.annals.org/>, and Epidemiology at <http://www.epidem.com/>). Information on the STROBE Initiative is available at <http://www.strobe-statement.org>.

For peer review only



Representativeness and optimal use of body mass index (BMI) in the UK Clinical Practice Research Datalink (CPRD)

Journal:	<i>BMJ Open</i>
Manuscript ID:	bmjopen-2013-003389.R1
Article Type:	Research
Date Submitted by the Author:	06-Aug-2013
Complete List of Authors:	Bhaskaran, Krishnan; LSHTM, NCDE Forbes, Harriet; LSHTM, NCDE Douglas, Ian; LSHTM, NCDE Leon, David; LSHTM, NCDE Smeeth, Liam; LSHTM, NCDE
Primary Subject Heading:	Epidemiology
Secondary Subject Heading:	Epidemiology, Research methods
Keywords:	EPIDEMIOLOGY, PRIMARY CARE, STATISTICS & RESEARCH METHODS

SCHOLARONE™
Manuscripts

Peer Review Only

Representativeness and optimal use of body mass index (BMI) in the UK**Clinical Practice Research Datalink (CPRD)**

Authors: Krishnan Bhaskaran, Harriet J Forbes, Ian Douglas, David A Leon, Liam Smeeth

Faculty of Epidemiology & Population Health, London School of Hygiene and Tropical

Medicine, Keppel Street, London, WC1E 7HT, UK

Corresponding Author: Krishnan Bhaskaran, London School of Hygiene and Tropical

Medicine, Keppel Street, London, WC1E 7HT. Email: krishnan.bhaskaran@lshtm.ac.uk. Tel: +

44 (0) 20 7927 2268

Key words: CPRD, BMI, missing data, primary care databases, obesity.

Word Count: 3680

Abstract

Objectives: To assess the completeness and representativeness of body mass index (BMI) data in the Clinical Practice Research Datalink (CPRD), and determine an optimal strategy for their use.

Design: Descriptive study.

Setting: Electronic healthcare records from primary care.

Participants: A million patient random sample from the UK Clinical Practice Research Datalink (CPRD) primary care database, aged ≥ 16 years.

Primary and secondary outcome measures: BMI completeness in CPRD was evaluated by age, sex, and calendar period. CPRD-based summary BMI statistics for each calendar year (2003-10) were age- and sex-standardised and compared with equivalent statistics from the Health Survey for England (HSE).

Results: BMI completeness increased over calendar time from 37% in 1990-94 to 77% in 2005-11, was higher among females, and increased with age. When BMI at specific time points was assigned based on the most recent record, calendar year-specific mean BMI statistics underestimated equivalent HSE statistics by 0.75-1.1kg/m². Restricting to those with a recent (≤ 3 years) BMI resulted in mean BMI estimates closer to HSE (≤ 0.28 kg/m² underestimation), but excluded up to 47% of patients. An alternative strategy of imputing up-to-date BMI based on modelled changes in BMI over time since the last available record, also led to mean BMI estimates that were close to HSE (≤ 0.37 kg/m² underestimation).

Conclusions: Completeness of BMI in CPRD increased over time and varied by age and sex. At a given point in time, a large proportion of the most recent BMIs are unlikely to reflect current BMI; consequent BMI misclassification might be reduced by employing model-based imputation of current BMI.

Article summary

Article focus:

- Body mass index (BMI) data are frequently used in epidemiological analyses of primary care databases such as the UK Clinical Practice Research Datalink (CPRD), however their completeness and representativeness have not previously been assessed in detail.
- The aim of this article is to provide information on the completeness of BMI in CPRD primary care data, on their representativeness, and on the implications for their practical use in research.

Key messages:

- We found that completeness of BMI recordings in the Clinical Practice Research Datalink increased from 37% in 1990-4 to 77% in 2005-11 and differed by age and sex.
- At specific calendar time points, the most recent BMI recorded for a large proportion of patients was over 3 years old and was unlikely to reflect current BMI.
- The optimal strategy for assigning BMI status is likely to depend on the specific study population and research context. We suggest one possible approach that uses a model-based imputation of current BMI to reduce BMI misclassification.

Strengths and limitations of this study:

- Results presented here are based on a large random sample from the CPRD, therefore we can confidently generalise the findings to the whole CPRD database, and to similar databases based on UK primary care records.
- To assess the representativeness of CPRD BMI data, we compared with data from the Health Survey for England, which is based on a large nationally representative sample and includes BMI information measured by trained interviewers.
- Our study did not look at BMI recordings among children as this would require a different strategy.

Introduction

Overweight and obesity are major contributors to global disease burden[1] and are associated with substantial excess mortality[2]. The prevalence of obesity is increasing in both developed and developing countries[3 4] and is a growing concern to policy makers. In England, the prevalence of obesity rose steadily from 1993 to 2010: from 13% to 26% in men, and from 16% to 26% in women[5]. Because of its association with various diseases and health outcomes, body mass index (BMI, the metric most widely used to classify overweight and obesity) is an important factor in many epidemiological studies, both as an exposure and as a potential confounder.

Databases of routinely collected electronic healthcare records are becoming an increasingly valuable resource in epidemiology, allowing population-level research on large, representative samples. The UK Clinical Practice Research Datalink (CPRD) (formerly the General Practice Research Database or GPRD) is widely used and contains comprehensive medical records for approximately 8% of the UK population,[6] allowing epidemiological studies to be carried out on a range of topics and with much greater statistical power than is typically available in traditional cohort studies. However, a shortcoming of these databases is that lifestyle data, such as BMI, tend to be opportunistically recorded (i.e. recorded when the patient is attending for other reasons, or when of direct clinical importance) and can be incomplete. Furthermore, those with non-missing lifestyle data may be unrepresentative of the general population. BMI has been an important covariate in many published studies based on CPRD[7-14] but the completeness and representativeness of the BMI data have not been previously documented.

Our aim was to undertake an in-depth investigation of BMI recordings in CPRD, including quantifying the completeness of BMI data, and assessing their representativeness by

1 comparing summary statistics based on CPRD data with equivalent statistics from a
2
3 representative general population survey. We also aimed to suggest and discuss how to deal
4
5 with the limitations of these routinely collected BMI data.
6
7

8 **Methods**

9 **Data sources**

10 **Clinical Practice Research Datalink (CPRD)**

11 The Clinical Practice Research Datalink (CPRD) is a clinical database comprising anonymised
12
13 computerised medical records from general practitioners (GPs) in the United Kingdom.
14
15 Approximately 8% of the UK population are currently included and the database is broadly
16
17 representative of the UK population.[15 16] Registration with a GP is near-universal in the
18
19 UK,[17] and GPs act as gatekeepers to the health system so that a CPRD data form a
20
21 comprehensive health record, comprising demographic information, clinically relevant
22
23 lifestyle data, prescription details, clinical events, preventive care provided, specialist
24
25 referrals, and hospital admissions and their major outcomes. Data undergo quality checks
26
27 and practices are designated as “up to standard” in CPRD from the date that they meet
28
29 specified data entry quality criteria. For this study, we obtained a random sample of one
30
31 million CPRD patients, because carrying out the analysis on the full CPRD database would be
32
33 computationally difficult, and the reduction in precision of our estimates that would arise by
34
35 restricting our analysis to a one million random sample is extremely small.
36
37
38
39
40
41
42
43
44
45
46

47 **Body mass index data in CPRD**

48 Height and weight measurements are recorded in CPRD whenever measured as part of
49
50 routine care. We obtained all height and weight records and calculated BMI
51
52 (BMI=weight/height²). Records without any measurements or with implausible
53
54 measurements were excluded (Figure 1).
55
56
57
58
59
60

Health Survey for England

We obtained published Health Survey for England (HSE) data for BMI from the National Health Service (NHS) Information Centre.[18] The HSE is an annual survey designed to produce a representative sample of the adult population aged ≥ 16 years and living in private households (sample size 14,836 in 2003 and 8,420 in 2010),. Surveys were interviewer administered with interviewers measuring the weight and height of all participants. Data from 2003-10 were obtained, and these data have been weighted to reduce bias from non-response, based on a logistic regression model incorporating age, sex, household type (based on the number of adults and children living in a household), Strategic Health Authority region, and social class (defined using the National Statistics Socio-economic Classification system). The methods are described in more detail elsewhere.[19]

Statistical methods

Completeness of BMI data in CPRD

In the main analyses BMI completeness data in CPRD were estimated by calendar period (1990-4, 1995-9, 2000-4, 2005-11). To calculate completeness for a particular calendar period, all individuals from the one million sample who were registered, aged ≥ 16 years, and under follow-up in “up to standard” practices on the mid-point of the period were identified and included in the denominator. Among these individuals, the numerator comprised either those with any previous BMI available in their electronic record regardless of how long ago it was entered, or those with a BMI available up to 3 years prior to this date. Completeness data were generated by age group ,sex and among those whom, for clinical reasons, BMI should be routinely monitored (those with type 2 diabetes, schizophrenia/other psychoses, and ≥ 2 recent (last 6 months) statin prescriptions). We also investigated whether completeness could be improved by searching for clinical codes (“Read codes”) indicating

1 BMI category. We have not presented confidence intervals for these descriptive statistics
2
3
4 because the sample size made sampling error negligible (for example, the standard errors
5
6 for the proportions with complete BMI data in age and calendar year subgroups were all
7
8 <0.5%).
9

10 **Comparison of CPRD BMI data with Health Survey for England data**

11
12 We compared mean BMI over calendar time based on complete CPRD BMI data with
13
14 equivalent HSE figures, for the period 2003-2010 (since, from 2003, HSE data were adjusted
15
16 for non-response). CPRD mean BMI was based on patients registered and under up-to-
17
18 standard follow-up at the mid-point of the calendar year. We produced two sets of CPRD
19
20 mean BMI statistics: firstly we used last BMI observation carried forward (regardless of how
21
22 long ago recorded); secondly we restricted to patients with a recent BMI available (up to 3
23
24 years before the mid-point of the calendar year). As above, confidence intervals are not
25
26 presented because there was negligible sampling error (maximum standard
27
28 error=0.02kg/m²). To make like-with-like comparisons with HSE, CPRD data were restricted
29
30 to English practices (for comparisons with HSE data only), and mean BMI was age- and sex-
31
32 standardised to the HSE population structure Proportions classified as obese (BMI≥30kg/m²)
33
34 over time based on CPRD and HSE data were also compared.
35
36
37
38
39
40
41
42

43 **Model-based imputation of up-to-date BMI measures in CPRD**

44
45 We explored whether outdated BMI measures in CPRD could be usefully updated by
46
47 imputation based on a model predicting changes in individual-level BMI over time. We used
48
49 data from individuals with multiple BMI records to model the expected change in BMI as a
50
51 function of time since BMI recording (restricting to individuals with BMI records ≤ 10 years
52
53 apart). We fitted a linear regression model with change in BMI as the outcome; the main
54
55 covariate predicting change in BMI was elapsed time, which was included as a 3 knot cubic
56
57
58
59
60

1 spline to allow for non-linearity; we also included interactions between the spline basis
2
3 variables and indicator variables for age and sex. Feasible weighted least squares estimation
4
5 was used to allow for heteroskedasticity.[20]
6
7

8
9 Having specified a model for change in BMI over time, we first explored its performance
10
11 among individuals with at least 2 BMIs entered in CPRD, by predicting the most recent BMI
12
13 based on the previous BMI record and the elapsed time; we compared the distribution of
14
15 the errors from this approach with the distribution of the errors from simply using the last
16
17 observation carried forward. We then repeated the comparison with the HSE mean BMI
18
19 data for each calendar year. This time we included all individuals with a BMI record in the
20
21 previous 10 years and used the model described above to impute current BMI at the mid-
22
23 point of the calendar year by predicting the change in BMI since the last available BMI
24
25 record. We did this within a multiple imputation framework (using 5 imputations) to
26
27 account for uncertainty in the modelled changes over time.[21]
28
29
30
31
32

33
34 The study was approved by the London School of Hygiene and Tropical Medicine Ethics
35
36 Committee.
37

38 39 **Results**

40 41 **Completeness of BMI data in CPRD**

42
43 In 1990-1994, 37% of individuals had at least one previously recorded BMI, and the
44
45 proportion increased to 77% by 2005-11 (Table 1). The proportion of individuals with a recent
46
47 BMI (recorded in the previous 3 years) was lower in each calendar period (35% in 1990-1994
48
49 rising to 51% in 2005-11). BMI completeness generally increased with age up to 75 years,
50
51 with a lower proportion in the oldest age group having data available. Data for single
52
53 calendar years are shown in Appendix Table A1 and illustrate similar patterns. BMI data
54
55
56
57
58
59
60

1 appeared to be consistently more widely available among women than men (Figure 2). As
2
3 expected, BMI completeness was higher in particular clinical subgroups: in total 97% of
4
5 patients with a record of type II diabetes had a recent BMI recorded, along with over 78% of
6
7 those with a diagnosis of schizophrenia/psychoses (Appendix Table A2). This is in line with
8
9 Quality and Outcomes Framework (QOF) which has encouraged BMI monitoring in these
10
11 conditions since 2004.[22] BMI completeness was also high among current statin users (82%
12
13 with a recent BMI available).
14
15

16
17
18 There was little extra information available in clinical (“Read”) codes relating to BMI. In the
19
20 most recent calendar period, out of 75518 individuals with no previous BMI record
21
22 available, only 1222 (1.6%) had ever had a clinical code that would enable classification into
23
24 BMI categories (underweight, normal, overweight/obese). Furthermore, for those with a
25
26 previous BMI, only a small proportion had more recent information related to BMI recorded
27
28 in a clinical code (7675/250430 = 3.0% in the most recent period).
29
30
31

32 33 **Summary statistics using complete CPRD BMI data and comparison with Health Survey for** 34 35 **England** 36

37
38 We found that age- and sex-standardised mean BMI based on CPRD data was consistently
39
40 and substantially lower (by up to 1.1kg/m²) than in the HSE data (mean BMI in CPRD =
41
42 25.7kg/m² in 2003 rising to 26.3 in 2010, compared with 26.8 kg/m² [95% CI 26.7 to 26.9]
43
44 and 27.3 [27.1 to 27.5] respectively in HSE; Figure 3).
45
46
47

48
49 When BMI entries more than 3 years old were discarded, between 33 to 47% of patients
50
51 were lost across calendar years. However, the estimated mean BMI in CPRD was
52
53 considerably closer to what would be expected based on the HSE data, with CPRD data
54
55 underestimating the HSE statistics by only between 0.04 to 0.28kg/m² in individual calendar
56
57 years, and the CPRD estimate falling within the HSE confidence interval for 2 of the most
58
59
60

1 recent 3 calendar years (mean BMI in CPRD = 26.9, 27.0 and 27.0 kg/m² compared with 27.0
2
3
4 [26.9 to 27.1], 27.0 [26.8 to 27.2] and 27.3 [27.1 to 27.5] in HSE, in 2008, 2009 and 2010
5
6 respectively). Age- and sex-stratified data demonstrated similar patterns, except that in the
7
8 eldest age group (75+ years), restricting to those with recent BMI measures did not bring
9
10 the estimated BMI substantially closer to HSE figures (Appendix Figure A1).
11
12

13
14
15
16
17 We also compared the proportions classified as obese between CPRD and HSE (Appendix
18
19 Figure A2). Consistent with the previous analysis, using any previous BMI reading to classify
20
21 individuals in CPRD resulted in lower obesity rates than expected based on HSE data, while
22
23 restricting to patients with a recent reading led to estimated obesity rates close to those in
24
25 HSE.
26
27

28 29 30 **Model-based imputation of up-to-date BMI measures in CPRD**

31
32 The contrast between BMI summary statistics based on recent measures and those based
33
34 on any previous measures suggested that older BMI records were tending to underestimate
35
36 current BMI. We therefore examined whether a model could be developed to impute
37
38 current BMI, taking into account elapsed time since the last measure. In a linear regression
39
40 model for change in BMI over time, we estimated that on average BMI increased over the
41
42 10-year period following a BMI record for those aged up to 69 years at the time of the
43
44 record and decreased over time in those aged 70 years or more (Appendix Figure A3). We
45
46 tested the predictive performance of our model by predicting the most recent BMI based on
47
48 the previous one, among CPRD patients with more than one recorded BMI available. When
49
50 the older BMI was less than 3 years old, there was little gain in applying the correction
51
52 compared with carrying the older observation forward (Figure 4). However, when there was
53
54 a longer gap, carrying the previous BMI forward tended to underestimate the later BMI,
55
56
57
58
59
60

1 while employing the model-based imputation removed the underestimation and led to
2
3 smaller errors on average (median error = $-0.70\text{kg}/\text{m}^2$ [IQR -2.18 to $+0.56$] using last
4
5 observation carried forward, compared with $+0.11\text{kg}/\text{m}^2$ [-1.29 to $+1.40$] using model-based
6
7 imputation).
8
9

10
11 We then repeated the comparison of mean BMI in CPRD versus HSE, this time using our
12
13 model for change in BMI over time as a basis for performing multiple imputations of current
14
15 BMI based on the latest available measure and the time since it was recorded. Estimated
16
17 mean BMIs were now in line with those based on only recent data in the earlier analysis,
18
19 and were only between 0.04 and $0.37\text{kg}/\text{m}^2$ lower than HSE statistics in individual calendar
20
21 years (Figure 3, circles). Even with multiple imputation, confidence intervals remained
22
23 extremely narrow ($<0.07\text{kg}/\text{m}^2$) due to the large sample size, so are not shown in the figure.
24
25 Of note, all patients with a BMI recorded up to 10 years before the midpoint of the calendar
26
27 year of interest were now included in the estimation of the “corrected” means; thus in
28
29 individual calendar years only 9 to 13% of patients were dropped, compared to 33-47% of
30
31 patients when dropping BMI records >3 years old.
32
33
34
35
36
37

38 Discussion

39 Main findings

40
41 BMI completeness has increased over calendar time (rising from 37% in 1990-94 to 77% in
42
43 2005-11). Completeness was higher among females, older age groups, and clinical
44
45 subgroups where recording BMI is encouraged. When BMI on the date of interest was
46
47 assigned to individual patients in CPRD using the last available record, regardless of how
48
49 long ago it was entered, we found that the resulting mean BMI statistics for the CPRD
50
51 population were consistently lower than equivalent HSE estimates (by up to $1.1\text{kg}/\text{m}^2$). This
52
53 appeared to be driven by older BMI records tending to systematically underestimate current
54
55
56
57
58
59
60

1 BMI: when only recent CPRD BMI records (≤ 3 years old) were used, mean BMI statistics
2
3 were closer to HSE estimates. However, a substantial number of patients were then
4
5 excluded altogether (33-47% across years). Finally, we suggested a process for modelling
6
7 changes in BMI after a BMI record, which could allow researchers to impute BMI on the date
8
9 of interest and avoid dropping large numbers without a recent measure from their analyses.
10
11

12 *Comparison with other studies*

13
14 There are very few comparable studies (Appendix Table A2). However, the proportion of
15
16 patients with a recent BMI recording in CPRD is in line with a summary of the QRESEARCH
17
18 database (a similar UK primary care database with data from over 530 general practices
19
20 using EMIS software rather than VISION software);[23] by March 2007, 58% of registered
21
22 patients aged 16+ years had their BMI recorded in the past 5 years; this compares with 51%
23
24 with a BMI recorded in the last 3 years in our analysis (for 2005-11). As in our study, the
25
26 QRESEARCH report showed an increase in completeness over time, rising from 42% in
27
28 2000/01 to 58% in 2007. In a third UK primary care database, THIN (The Health
29
30 Improvement Network), the proportion of newly registered patients between 2004 and
31
32 2006 with BMI data was in line with our findings; 62% of patients had a height recording and
33
34 66% had a weight recording within 12 months of registration.[24]
35
36
37
38
39
40
41
42

43 **Explanation of findings**

44 *Completeness*

45
46 Increasing completeness of BMI over time may reflect a general trend towards
47
48 encouragement to record BMI in primary care. Greater BMI completeness among females
49
50 and older age groups may have a number of explanations including higher consultation
51
52 rates in primary care [25 26] and different prevalences of diseases in which it is important to
53
54 monitor BMI.
55
56
57
58
59
60

Comparison of CPRD BMI data with Health Survey for England data

Mean BMI based on the CPRD population was lower in each calendar year than equivalent HSE estimates when BMI in CPRD was assigned using the last available record; however, when the analysis was restricted to those with a recent BMI record, estimates from CPRD were close to HSE estimates. This suggests that the substantial proportion of BMI recordings in CPRD that were outdated on the date of interest may have driven the apparent underestimation of mean BMI in CPRD in the unrestricted analysis. This in turn would imply that individual BMIs tend to increase over time, and indeed when we specifically modelled changes in BMI over time, we found a pattern of increasing BMI with age for those <70 years old, consistent with prospective cohort studies with repeated BMI measurements [27-29]; this pattern of increasing BMI over time is likely to be driven specifically by weight change, since adult height would not change substantially in this age range. A simple adjustment of outdated BMIs based on our modelled changes over time brought CPRD mean BMI statistics in line with HSE estimates, and when we validated the adjustment in a subset of patients with repeated BMI measures, we found smaller errors on average, compared with simply carrying outdated BMI records forwards.

Of note, we observed that CPRD consistently underestimated BMI compared to HSE among those aged ≥ 75 years, even when only recent records were used; this may reflect the fact that institutionalised patients are represented in CPRD but not in HSE: HSE may not be an ideal comparison for this age group since elderly people in institutions (who are represented in CPRD) may be more likely to be frail and have lower BMIs than those living in private households.

Implications

1 First, our findings suggest that BMI completeness is likely to vary between studies
2
3 depending on the study population and study period. BMI data are not likely to be missing
4
5 completely at random (for example, missingness may vary by patient characteristics or
6
7 particular diseases). There may be information in the database, however, which predicts
8
9 missingness and which could satisfy the “missing at random” assumption required for
10
11 multiple imputation. A study exploring the potential of imputing missing data in THIN found
12
13 that after multiple imputation, summary statistics of height and weight were comparable
14
15 with data from nationally representative datasets.[24]
16
17
18
19
20 Second, our analyses suggest that the common practice of assigning BMI status based on
21
22 the nearest/most recent available record to the index date of interest might lead to
23
24 misclassification, given that a large number of patients have only substantially outdated BMI
25
26 records available at any particular time. Strategies to address this include restricting to
27
28 recent BMI, but this is likely to exclude a large numbers of patients. We have suggested an
29
30 alternative strategy based on updating the outdated BMIs by modelling changes in BMI over
31
32 time, though this is not without drawbacks: the approach requires an assumption that
33
34 individuals with ≥ 2 BMI records available (needed to estimate the model for changes over
35
36 time) are representative of the wider patient population, which may not be the case; it is
37
38 also a more complex strategy, particularly if done within a multiple imputation framework
39
40 to allow for uncertainty in the correction, which could be substantial in studies with smaller
41
42 sample sizes than considered here. Other imputation strategies could also be considered in
43
44 certain contexts, such as the two-fold algorithm which imputes missing data from
45
46 longitudinal variables at particular time points by using adjacent data points.[30] Ultimately,
47
48 the pros and cons of various methods, and the optimal strategy to use is likely to depend on
49
50 the particular study and the characteristics of the study population.
51
52
53
54
55
56
57
58
59
60

Strengths and Limitations

Results presented here are based on a large random sample from the CPRD, therefore we can confidently generalise the findings to the whole CPRD database. Although we cannot assume these findings will relate to other routinely collected primary care databases in UK based on other IT systems (CPRD is based on practices using VISION), they are likely to be similar. This study did not look at BMI recordings among children as this would require a different strategy. Completeness among 16-24 year age group may be artificially low because weights recorded at age <16 were excluded, so those at the lower end of the age group will not have had as much time to accrue weight recordings. We believe HSE to be the best available comparison for this study; it is a nationally representative, large sample utilising height and weight recordings measured by a trained interviewer, and is weighted for non-response.[19 31] However there is a degree of missing data in HSE which is a limitation. In 2010 just over 85% of adults interviewed provided valid height and weight recordings. [29] One of the most common reasons for missing BMI was refusal (up to 8% were missing due to refusal),[19] which if related to BMI status, may bias the estimates of mean BMI in HSE. Our comparisons between CPRD-based and HSE-based BMI statistics focussed on the mean (and in the appendix, on the proportion classed as obese); these are the principal statistics published in HSE trend tables so we were not able to look at a broader range of measures of the BMI distribution that might be of interest to researchers using BMI data in the context of public health. Finally, we have not attempted to quantify or comment on the usefulness of BMI as a measure of adiposity, and researchers using BMI data should consider whether it is the best available measure for their purposes.

Conclusions

Completeness of BMI data in CPRD varies over time and by age and sex. BMI records may become outdated over time and naive use could lead to misclassification of BMI status. We used a 3-year cut-off to define a recent BMI; further research could include a systematic analysis of how long BMI records can be considered “up-to-date”, and whether this varies by patient characteristics. The optimal strategy for assigning BMI status to individuals in studies based on CPRD and similar electronic healthcare databases is likely to depend on the specific study population and the research context.

Conflicts of interest

The authors declare no conflicts of interest.

Funding

This report is independent research arising from a postdoctoral fellowship (for KB) supported by the National Institute for Health Research (PDF-2011-04-007). ID is supported by an MRC methodology research fellowship, LS is supported by a Wellcome Trust senior research fellowship in clinical science.

Data sharing statement

This analysis is based on a large random sample from the Clinical Practice Research Datalink, provided by the UK Medicines and Healthcare products Regulatory Agency. The authors' licence for using these data does not allow sharing of raw data with third parties.

peer review only

References

1. World Health Organisation. Global Health Risks: Mortality and burden of disease attributable to selected major risks. In: Organisation WH, ed. Geneva, Switzerland, 2009.
2. Flegal KM, Graubard BI, Williamson DF, et al. Excess deaths associated with underweight, overweight, and obesity. *Jama-J Am Med Assoc* 2005;**293**(15):1861-67
3. Swinburn BA, Sacks G, Hall KD, et al. Obesity 1 The global obesity pandemic: shaped by global drivers and local environments. *Lancet* 2011;**378**(9793):804-14
4. Kelly T, Yang W, Chen CS, et al. Global burden of obesity in 2005 and projections to 2030. *Int J Obesity* 2008;**32**(9):1431-37 doi: Doi 10.1038/ljo.2008.102[published Online First: Epub Date]].
5. NHS Information Centre. Health survey for England - 2010: health and lifestyles. Secondary Health survey for England - 2010: health and lifestyles 2011. <http://www.ic.nhs.uk/pubs/hse10report>.
6. CPRD. Clinical Practice Research Database (CPRD) website. Secondary Clinical Practice Research Database (CPRD) website. <http://www.cprd.com/intro.asp>.
7. Delaney JA, Daskalopoulou SS, Brophy JM, et al. Lifestyle variables and the risk of myocardial infarction in the general practice research database. *BMC Cardiovasc Disord* 2007;**7**:38 doi: 1471-2261-7-38 [pii]
- 10.1186/1471-2261-7-38[published Online First: Epub Date]].
8. Green J, Czanner G, Reeves G, et al. Oral bisphosphonates and risk of cancer of oesophagus, stomach, and colorectum: case-control analysis within a UK primary care cohort. *BMJ* 2010;**341**:c4444 doi: 10.1136/bmj.c4444[published Online First: Epub Date]].
9. Tzoulaki I, Molokhia M, Curcin V, et al. Risk of cardiovascular disease and all cause mortality among patients with type 2 diabetes prescribed oral antidiabetes drugs: retrospective cohort study using UK general practice research database. *British Medical Journal* 2009;**339**:b4731 doi: Artn B4731
- Doi 10.1136/Bmj.B4731[published Online First: Epub Date]].
10. Douglas I, Smeeth L, Irvine D. The use of antidepressants and the risk of haemorrhagic stroke: a nested case control study. *Brit J Clin Pharmacol* 2011;**71**(1):116-20 doi: DOI 10.1111/j.1365-2125.2010.03797.x[published Online First: Epub Date]].
11. Andersohn F, Schade R, Suissa S, et al. Long-Term Use of Antidepressants for Depressive Disorders and the Risk of Diabetes Mellitus. *Am J Psychiat* 2009;**166**(5):591-98 doi: DOI 10.1176/appi.ajp.2008.08071065[published Online First: Epub Date]].
12. Lawrenson R, Todd JC, Leydon GM, et al. Validation of the diagnosis of venous thromboembolism in general practice database studies. *Brit J Clin Pharmacol* 2000;**49**(6):591-96
13. Jick H, Zornberg GL, Jick SS, et al. Statins and the risk of dementia. *Lancet* 2000;**356**(9242):1627-31
14. van Staa TP, Wegman S, de Vries F, et al. Use of statins and risk of fractures. *Jama-J Am Med Assoc* 2001;**285**(14):1850-55
15. Office for National Statistics. *Key Health Statistics from General Practice 1998: Analyses of Morbidity and Treatment Data, Including Time Trends, England and Wales*. London: Office for National Statistics, 2000.
16. Parkinson JP, Davis S, Van Staa T. The General Practice Research Database: now and the future. In: Mann R, Andrews EB, eds. *Pharmacovigilance*. Chichester: John Wiley and Sons, 2007:341-48.
17. Schoonen WM, Thomas SL, Somers EC, et al. Do selected drugs increase the risk of lupus? A matched case-control study. *Br J Clin Pharmacol* 2010;**70**(4):588-96 doi: 10.1111/j.1365-2125.2010.03733.x[published Online First: Epub Date]].
18. NHS Information Centre. Health Survey for England - 2010: Trend tables. Secondary Health Survey for England - 2010: Trend tables. <http://www.ic.nhs.uk/statistics-and-data->

- 1 [collections/health-and-lifestyles-related-surveys/health-survey-for-england/health-survey-](#)
2 [for-england--2010-trend-tables.](#)
- 3
- 4 19. Aresu M, Boodhna G, Bryson A, et al. Volume 2: Methods and documentation. In: Craig R,
5 Mindell J, eds. Health Survey for England 2010. Leeds: NHS Information Centre for health
6 and social care, 2011.
- 7 20. Greene WH. *Econometric Analysis*. Upper Saddle River, New Jersey: Prentice Hall., 1997.
- 8 21. Rubin DB. *Multiple Imputation for Nonresponse in Surveys*: J. Wiley & Sons, New York. , 1987.
- 9 22. NHS. Quality and Outcomes Framework guidance for GMS contract 2011/12: Employers and
10 British Medical Association, 2011.
- 11 23. NHS Information Centre. A summary of public health indicators using electronic data from
12 primary care. Secondary A summary of public health indicators using electronic data from
13 primary care 2008. [http://www.ic.nhs.uk/article/2021/Website-](http://www.ic.nhs.uk/article/2021/Website-Search?productid=4287&q=A+summary+of+public+health+indicators+using+electronic+data+from+primary+care.&sort=Relevance&size=10&page=1&area=both#top)
14 [Search?productid=4287&q=A+summary+of+public+health+indicators+using+electronic+data](http://www.ic.nhs.uk/article/2021/Website-Search?productid=4287&q=A+summary+of+public+health+indicators+using+electronic+data+from+primary+care.&sort=Relevance&size=10&page=1&area=both#top)
15 [+from+primary+care.&sort=Relevance&size=10&page=1&area=both#top.](http://www.ic.nhs.uk/article/2021/Website-Search?productid=4287&q=A+summary+of+public+health+indicators+using+electronic+data+from+primary+care.&sort=Relevance&size=10&page=1&area=both#top)
- 16 24. Marston L, Carpenter JR, Walters KR, et al. Issues in multiple imputation of missing data for large
17 general practice clinical databases. *Pharmacoepidemiol Drug Saf* 2010;**19**(6):618-26 doi:
18 10.1002/pds.1934[published Online First: Epub Date]].
- 19 25. Rowlands S, Moser K. Consultation rates from the General Practice Research Database. *Brit J Gen*
20 *Pract* 2002;**52**(481):658-60
- 21 26. The Health and Social Care Information Centre. Trends in Consultation Rates in General Practice
22 1995 to 2009. Secondary Trends in Consultation Rates in General Practice 1995 to 2009
23 2009. [http://www.ic.nhs.uk/article/2021/Website-](http://www.ic.nhs.uk/article/2021/Website-Search?productid=729&q=qresearch&sort=Relevance&size=10&page=1&area=both#top)
24 [Search?productid=729&q=qresearch&sort=Relevance&size=10&page=1&area=both#top.](http://www.ic.nhs.uk/article/2021/Website-Search?productid=729&q=qresearch&sort=Relevance&size=10&page=1&area=both#top)
- 25 27. Li L, Law C, Power C. Body mass index throughout the life-course and blood pressure in mid-adult
26 life: a birth cohort study. *Journal of Hypertension* 2007;**25**(6):1215-23
- 27 28. Silverwood RJ, Pierce M, Thomas C, et al. Overweight across adult life and kidney function at age
28 60-4 years: the 1946 British birth cohort study. (awaiting publication) 2012
- 29 29. Tirosh A, Shai I, Afek A, et al. Adolescent BMI trajectory and risk of diabetes versus coronary
30 disease. *N Engl J Med* 2011;**364**(14):1315-25 doi: 10.1056/NEJMoa1006992[published Online
31 First: Epub Date]].
- 32 30. Welch C, Bartlett J, Petersen I. Application of multiple imputation using the two-fold fully
33 conditional specification algorithm in longitudinal clinical data. (In Press). *Stata Journal* 2013
- 34 31. Mindell J, Biddulph JP, Hirani V, et al. Cohort Profile: The Health Survey for England. *International*
35 *Journal of Epidemiology* 2012:1-9 doi: 10.1093/ije/dyr199[published Online First: Epub
36 Date]].
- 37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Table 1: Completeness of BMI data in the CPRD, by age and calendar period

Age group (yrs)	1990-4	1995-9	2000-4	2005-2011
16-24^a				
N registered	11423	17501	34452	42546
BMI in previous 3y (%)	26	28	25	32
BMI any previous (%)	26	37	30	37
25-34				
N registered	17477	29923	48659	50413
BMI in previous 3y (%)	37	39	36	49
BMI any previous (%)	38	66	67	72
35-44				
N registered	15953	28838	55991	61014
BMI in previous 3y (%)	36	36	31	46
BMI any previous (%)	39	67	71	80
45-54				
N registered	14507	27765	48093	55564
BMI in previous 3y (%)	39	37	32	50
BMI any previous (%)	42	70	73	84
55-64				
N registered	11680	20843	42258	49380
BMI in previous 3y (%)	42	40	37	57
BMI any previous (%)	44	74	77	87
65-74				
N registered	10678	17605	30997	34508
BMI in previous 3y (%)	36	37	40	67
BMI any previous (%)	38	71	79	91
75+				
N registered	8637	16005	29384	32523
BMI in previous 3y (%)	28	32	37	64
BMI any previous (%)	28	56	69	87
Total				
N registered	90355	158480	289834	325948
BMI in previous 3y (%)	35	36	34	51
BMI any previous (%)	37	64	67	77

N registered is all those under follow-up at mid-point of the period

^aNote, BMI measurements from age <16 years were not counted in this analysis, hence completeness in the 16-24 age group may be artificially low

Figure 1: Initial data processing to generate BMI for analysis

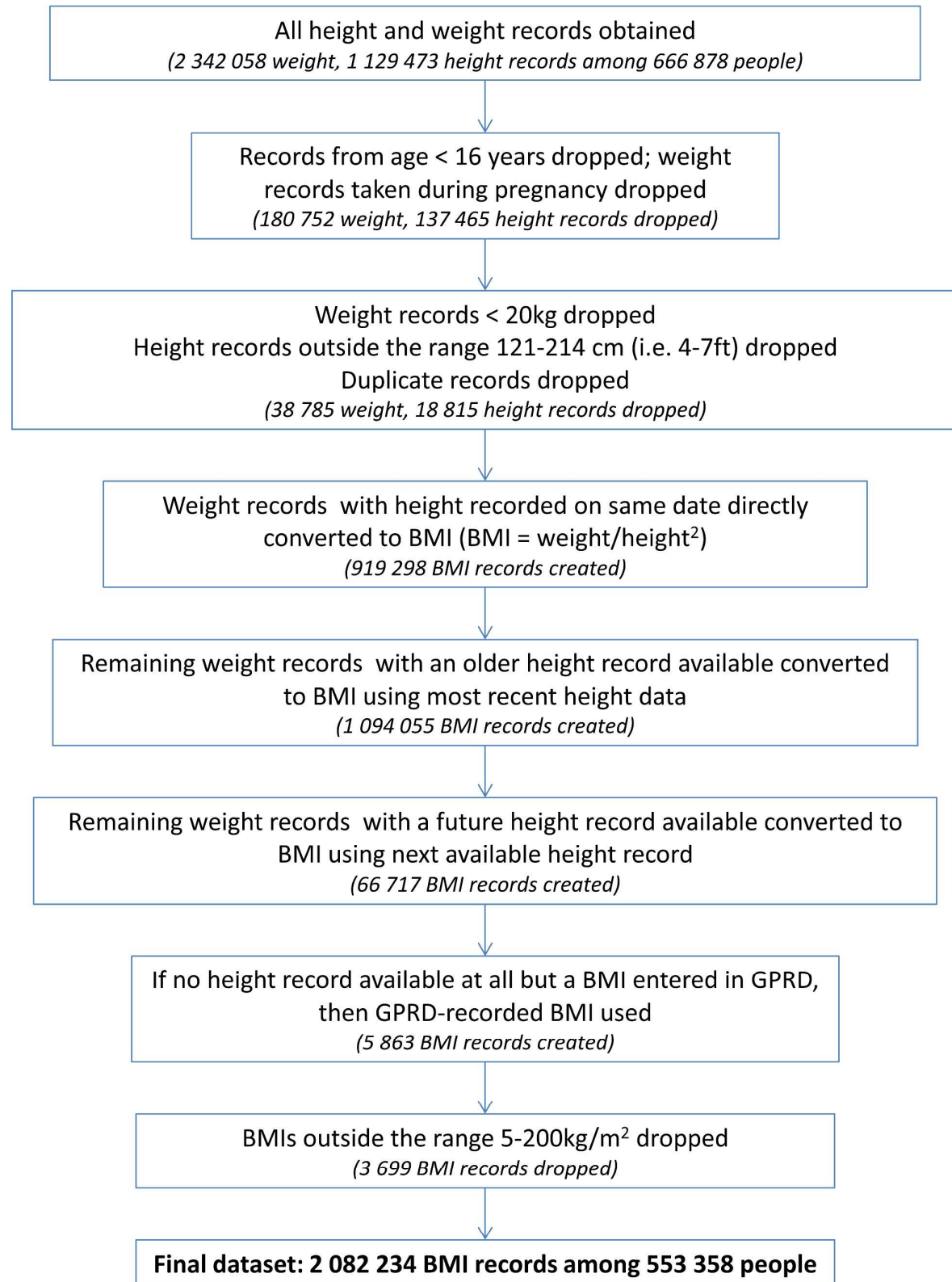
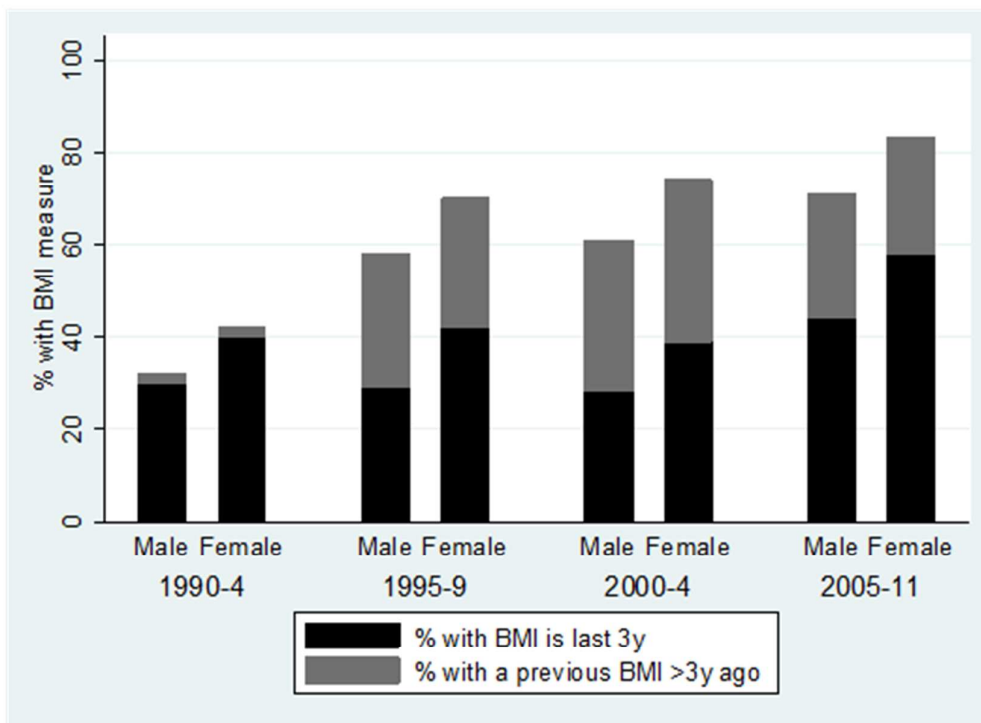


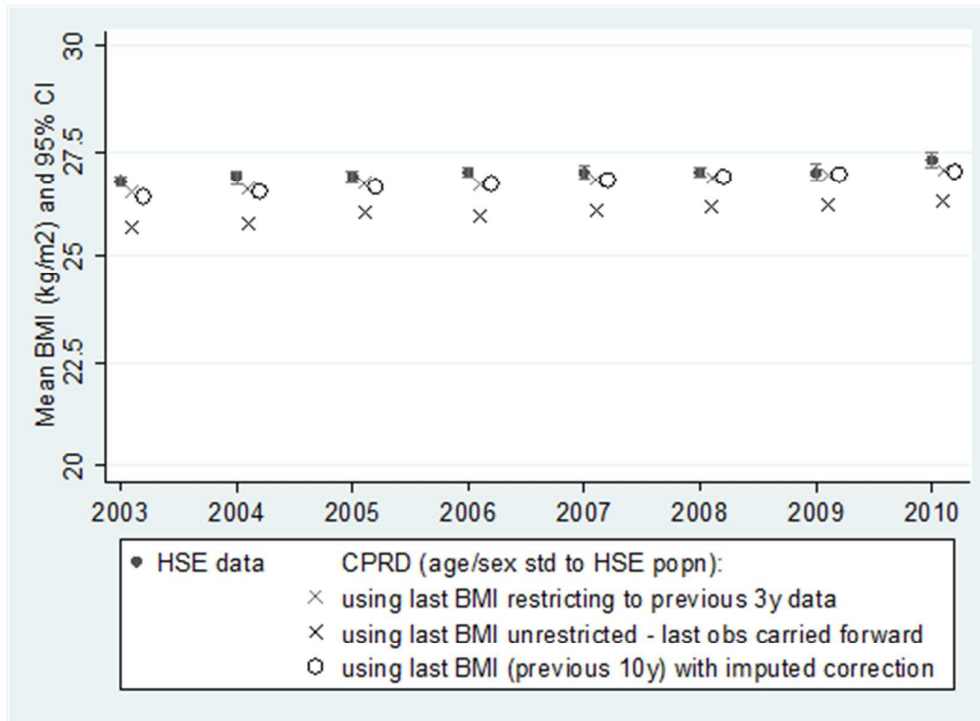
Figure 2: Completeness of BMI data in CPRD, by gender and calendar period



Note: Completeness data for each calendar period are based on all those under follow-up at mid-point of the period

For peer review only

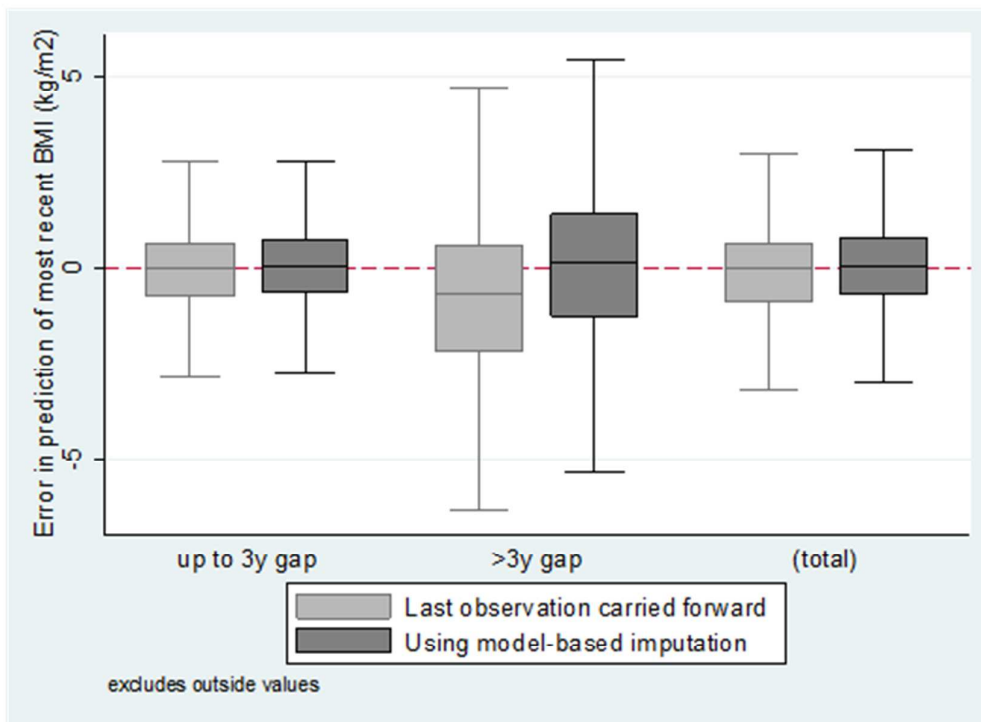
Figure 3: Mean BMI over calendar time comparing those with BMI recorded in CPRD (English practices) with the Health Survey for England 2010 data



Note: CPRD figures are age- and sex- standardised to the Health Survey for England study population

CPRD statistics are based on all patients registered at the mid-point of the calendar period and with a suitable previous BMI measure available (i.e. either any previous, or within the last 3 years)

Figure 4: Error in prediction of most recent BMI from older BMI, comparing simple last observation carried forward with model-based imputation of up to date BMI – stratified by time gap between readings



view only

Author contributions

I, Krishnan Bhaskaran, developed the analytical strategy for this paper, processed and analysed the data and wrote the paper.

I, Harriet Forbes, was involved in discussing the data processing and analysis of the data, as well as the writing of the paper.

I, Liam Smeeth, was involved in discussions of the analytical approach to this study and made comments on the analysis and the writing of the paper.

I, Ian Douglas, was involved in discussions of the analytical approach to this study and made comments on the analysis and the writing of the paper.

I, David Leon, was involved in discussions of the analytical approach to this study and made comments on the analysis and the writing of the paper.

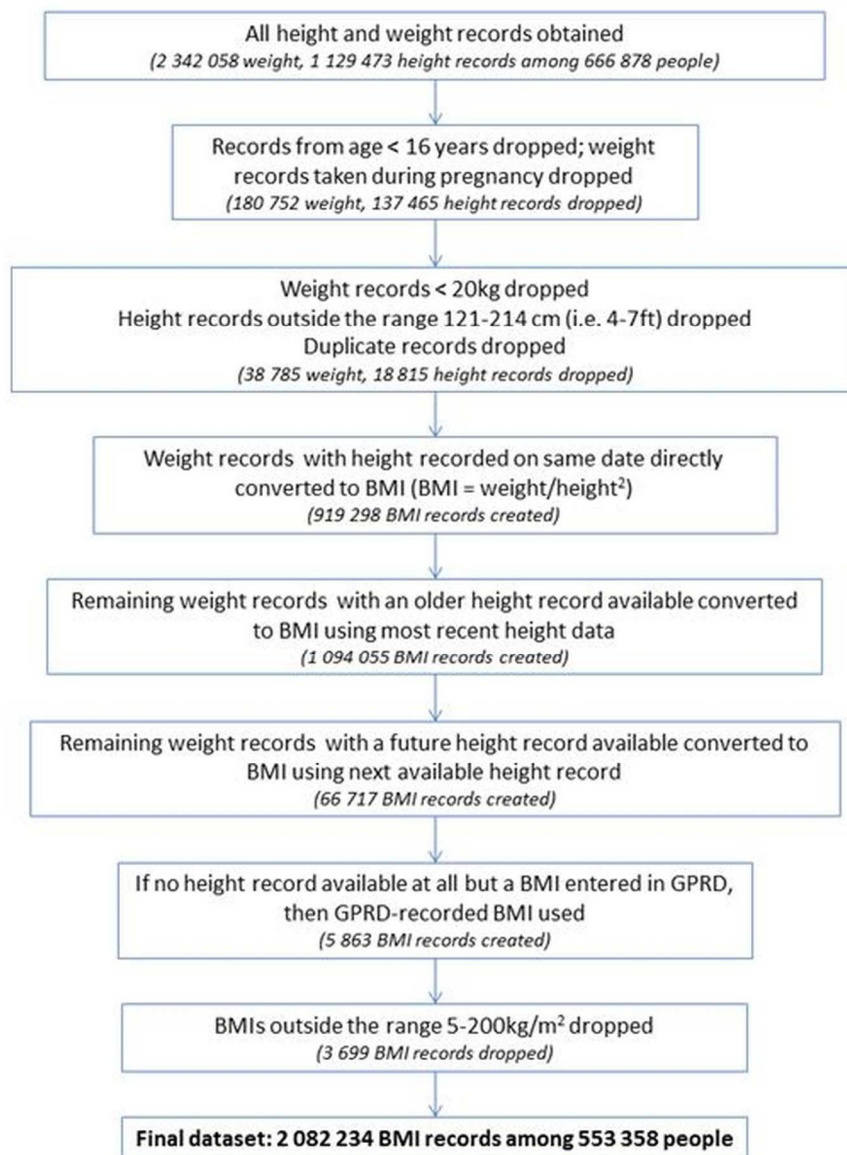


Figure 1: Initial data processing to generate BMI for analysis
141x187mm (300 x 300 DPI)

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

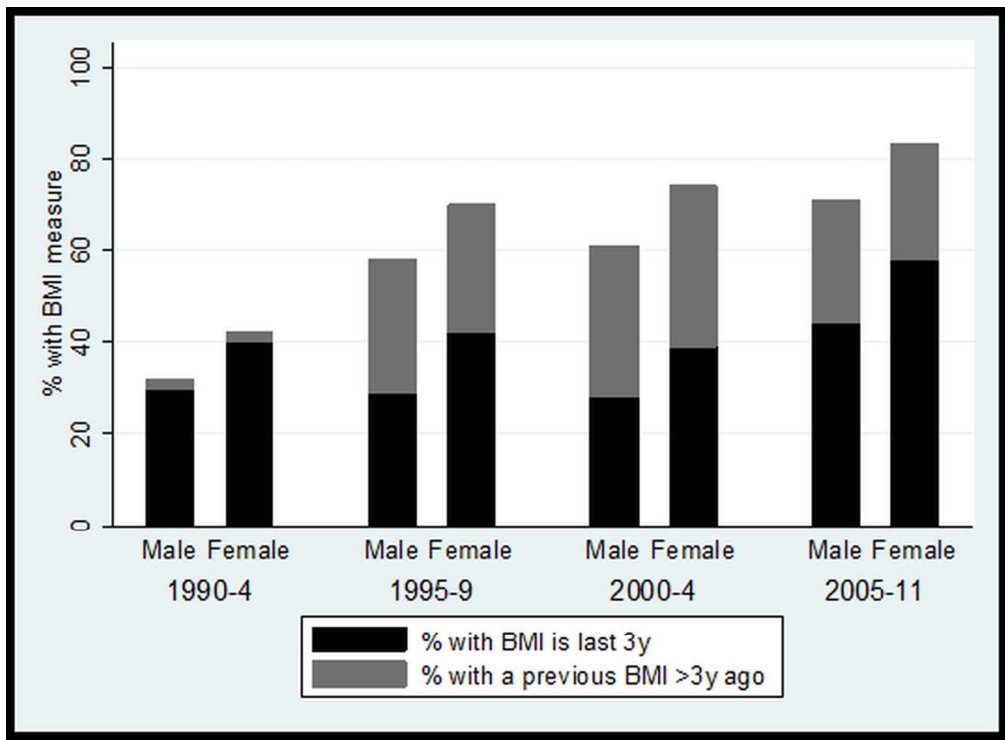


Figure 2: Completeness of BMI data in CPRD, by gender and calendar period
141x103mm (300 x 300 DPI)

Review only

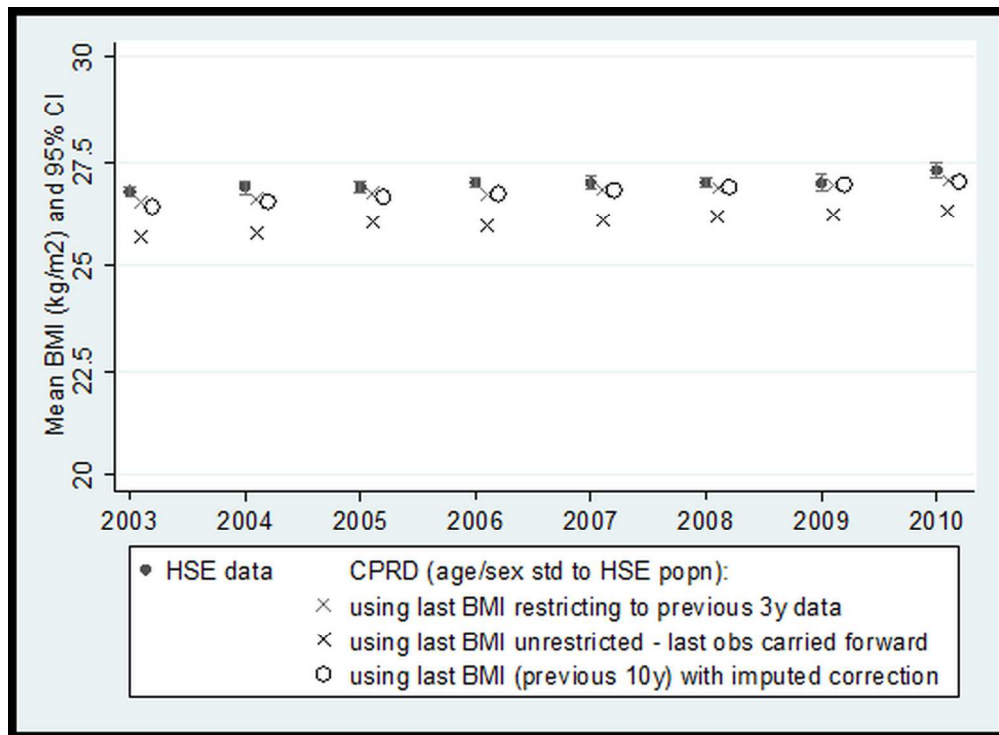


Figure 3: Mean BMI over calendar time comparing those with BMI recorded in CPRD (English practices) with the Health Survey for England 2010 data
 141x103mm (300 x 300 DPI)

View only

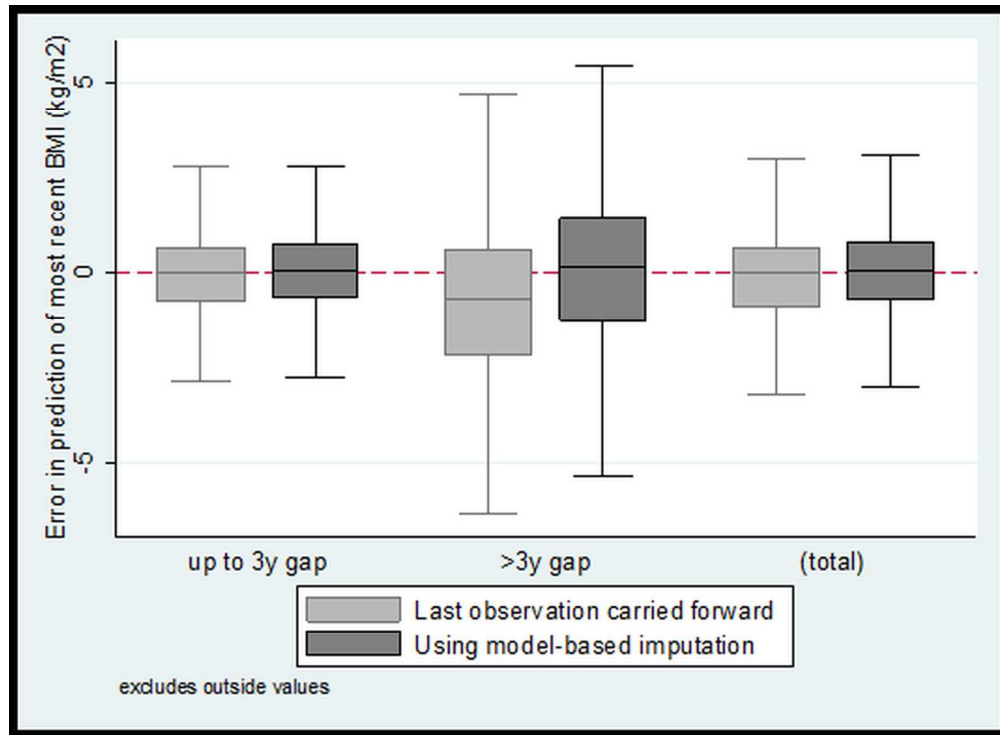


Figure 4: Error in prediction of most recent BMI from older BMI, comparing simple last observation carried forward with model-based imputation of up to date BMI – stratified by time gap between readings
141x103mm (300 x 300 DPI)

Review only

1
2
3
4
5
6
7
8 **Representativeness and optimal use of body mass index (BMI) in the UK**

9
10 **Clinical Practice Research Datalink (CPRD)**

11
12
13
14 **Authors:** Krishnan Bhaskaran, Harriet J Forbes, Ian Douglas, David A Leon, Liam Smeeth

15
16 Faculty of Epidemiology & Population Health, London School of Hygiene and Tropical

17
18 Medicine, Keppel Street, London, WC1E 7HT, UK

19
20
21
22 **Corresponding Author:** Krishnan Bhaskaran, London School of Hygiene and Tropical

23
24 Medicine, Keppel Street, London, WC1E 7HT. Email: krishnan.bhaskaran@lshtm.ac.uk. Tel: +

25
26 44 (0) 20 7927 2268

27
28
29
30
31 **Key words:** CPRD, BMI, missing data, primary care databases, obesity.

32
33 **Word Count:** ~~3383~~3680

Abstract

Objectives: To assess the completeness and representativeness of body mass index (BMI) data in the Clinical Practice Research Datalink (CPRD), and determine an optimal strategy for their use.

Design: Descriptive study.

Setting: Electronic healthcare records from primary care.

Participants: A million patient random sample from the UK Clinical Practice Research Datalink (CPRD) primary care database, aged ≥ 16 years.

Primary and secondary outcome measures: BMI completeness in CPRD was evaluated by age, sex, and calendar period. CPRD-based summary BMI statistics for each calendar year (2003-10) were age- and sex-standardised and compared with equivalent statistics from the Health Survey for England (HSE).

Results: BMI completeness increased over calendar time from 37% in 1990-94 to 77% in 2005-11, was higher among females, and increased with age. When BMI at specific time points was assigned based on the most recent record, calendar year-specific mean BMI statistics underestimated equivalent HSE statistics by $0.75\text{-}1.1\text{kg/m}^2$. Restricting to those with a recent (≤ 3 years) BMI resulted in mean BMI estimates closer to HSE ($\leq 0.28\text{kg/m}^2$ underestimation), but excluded up to 47% of patients. An alternative strategy of imputing up-to-date BMI based on modelled changes in BMI over time since the last available record, also led to mean BMI estimates that were close to HSE ($\leq 0.37\text{kg/m}^2$ underestimation).

Conclusions: Completeness of BMI in CPRD increased over time and varied by age and sex. At a given point in time, a large proportion of the most recent BMIs are unlikely to reflect current BMI; consequent BMI misclassification might be reduced by employing model-based imputation of current BMI.

Article summary

Article focus:

- Body mass index (BMI) data are frequently used in epidemiological analyses of primary care databases such as the UK Clinical Practice Research Datalink (CPRD), however their completeness and representativeness have not previously been assessed in detail.
- The aim of this article is to provide information on the completeness of BMI in CPRD primary care data, on their representativeness, and on the implications for their practical use in research.

Key messages:

- We found that completeness of BMI recordings in the Clinical Practice Research Datalink increased from 37% in 1990-4 to 77% in 2005-11 and differed by age and sex.
- At specific calendar time points, the most recent BMI recorded for a large proportion of patients was over 3 years old and was unlikely to reflect current BMI.
- The optimal strategy for assigning BMI status is likely to depend on the specific study population and research context. We suggest one possible approach that uses a model-based imputation of current BMI to reduce BMI misclassification.

Strengths and limitations of this study:

- Results presented here are based on a large random sample from the CPRD, therefore we can confidently generalise the findings to the whole CPRD database, and to similar databases based on UK primary care records.
- To assess the representativeness of CPRD BMI data, we compared with data from the Health Survey for England, which is based on a large nationally representative sample and includes BMI information measured by trained interviewers.
- Our study did not look at BMI recordings among children as this would require a different strategy.

Introduction

Overweight and obesity are major contributors to global disease burden[1] and are associated with substantial excess mortality[2]. The prevalence of obesity is increasing in both developed and developing countries[3 4] and is a growing concern to policy makers. In England, the prevalence of obesity rose steadily from 1993 to 2010: from 13% to 26% in men, and from 16% to 26% in women[5]. Because of its association with various diseases and health outcomes, body mass index (BMI, the metric most widely used to classify overweight and obesity) is an important factor in many epidemiological studies, both as an exposure and as a potential confounder.

Databases of routinely collected electronic healthcare records are becoming an increasingly valuable resource in epidemiology, allowing population-level research on large, representative samples. The UK Clinical Practice Research Datalink (CPRD) (formerly the General Practice Research Database or GPRD) is widely used and contains [comprehensive](#) medical records for approximately 8% of the UK population, [6] [allowing epidemiological studies to be carried out on a range of topics and with much greater statistical power than is typically available in traditional cohort studies](#). However, a shortcoming of these databases is that lifestyle data, such as BMI, tend to be opportunistically recorded [\(i.e. recorded when the patient is attending for other reasons, or when of direct clinical importance\)](#) and can be incomplete. Furthermore, those with non-missing lifestyle data may be unrepresentative of the general population. BMI has been an important covariate in many published studies based on CPRD[7-14] but the completeness and representativeness of the BMI data have not been previously documented.

Our aim was to undertake an in-depth investigation of BMI recordings in CPRD, including quantifying the completeness of BMI data, and assessing their representativeness by

1
2
3
4
5 comparing summary statistics based on CPRD data with equivalent statistics from a
6
7 representative general population survey. [We also aimed to suggest and discuss how to deal](#)
8
9 [with the limitations of these routinely collected BMI data.](#)
10

11 **Methods**

12 **Data sources**

13 **Clinical Practice Research Datalink (CPRD)**

14
15
16
17 The Clinical Practice Research Datalink (CPRD) is a clinical database comprising anonymised
18
19 computerised medical records from general practitioners (GPs) in the United Kingdom.

20
21 Approximately 8% of the UK population are currently included and the database is broadly
22
23 representative of the UK population.[15 16] [Registration with a GP is near-universal in the](#)
24
25 [UK,\[17\] and GPs act as gatekeepers to the health system so that a CPRD data contains form a](#)
26
27 [comprehensive health record, comprising](#) demographic information, clinically relevant
28
29 lifestyle data, prescription details, clinical events, preventive care provided, specialist
30
31 referrals, and hospital admissions and their major outcomes. Data undergo quality checks
32
33 and practices are designated as “up to standard” in CPRD from the date that they meet
34
35 specified data entry quality criteria. For this study, we obtained a random sample of one
36
37 million CPRD patients, because carrying out the analysis on the full CPRD database would be
38
39 computationally difficult, and the reduction in precision of our estimates that would arise by
40
41 restricting our analysis to a one million random sample is extremely small.
42
43
44

45 **Body mass index data in CPRD**

46
47 Height and weight measurements are recorded in CPRD whenever measured as part of
48
49 routine care. We obtained all height and weight records and calculated BMI
50
51 (BMI=weight/height²). [Patient rR](#)records without any measurements or with implausible
52
53 measurements were excluded (Figure 1).
54

Health Survey for England

We obtained published Health Survey for England (HSE) data for BMI from the National Health Service (NHS) Information Centre.[18] The HSE is an annual survey designed to produce a representative sample of the adult population aged ≥ 16 years and living in private households ([sample size 14,836 in 2003 and 8,420 in 2010](#)). ~~The methods are described in detail elsewhere.~~[19] Surveys were interviewer administered with interviewers measuring the weight and height of all participants. Data from 2003-10 were obtained, and these data have been weighted to reduce bias from non-response, based on a logistic regression model incorporating age, sex, household type (based on the number of adults and children living in a household), Strategic Health Authority region, and social class (defined using the National Statistics Socio-economic Classification system). [The methods are described in more detail elsewhere.](#)[19]

Statistical methods

Completeness of BMI data in CPRD

In the main analyses BMI completeness data in CPRD were estimated by calendar period (1990-4, 1995-9, 2000-4, 2005-11). To calculate completeness for a particular calendar period, all individuals from the one million sample who were registered, aged ≥ 16 years, and under follow-up in "up to standard" practices on the mid-point of the period were identified and included in the denominator. Among these individuals, the numerator comprised either those with any previous BMI available in their electronic record regardless of how long ago it was entered, or those with a BMI available up to 3 years prior to this date. Completeness data were generated by age group, sex and among those whom, for clinical reasons, BMI should be routinely monitored (those with type 2 diabetes, schizophrenia/other psychoses, and ≥ 2 recent (last 6 months) statin prescriptions). We also investigated whether

1
2
3
4
5 completeness could be improved by searching for clinical codes (“Read codes”) indicating
6
7 BMI category. We have not presented confidence intervals for these descriptive statistics
8
9 because the sample size made sampling error negligible (for example, the standard errors
10
11 for the proportions with complete BMI data in age and calendar year subgroups were all
12
13 <0.5%).
14

15 16 **Comparison of CPRD BMI data with Health Survey for England data**

17
18 We compared mean BMI over calendar time based on complete CPRD BMI data with
19
20 equivalent HSE figures, for the period 2003-2010 (since, from 2003, HSE data were adjusted
21
22 for non-response). CPRD mean BMI was based on patients registered and under up-to-
23
24 standard follow-up at the mid-point of the calendar year. We produced two sets of CPRD
25
26 mean BMI statistics: firstly we used last BMI observation carried forward (regardless of how
27
28 long ago recorded); secondly we restricted to patients with a recent BMI available (up to 3
29
30 years before the mid-point of the calendar year). As above, confidence intervals are not
31
32 presented because there was negligible sampling error (maximum standard
33
34 error=0.02kg/m²). To make like-with-like comparisons with HSE, CPRD data were restricted
35
36 to English practices (for comparisons with HSE data only), and mean BMI was age- and sex-
37
38 standardised to the HSE population structure Proportions classified as obese (BMI≥30kg/m²)
39
40 over time based on CPRD and HSE data were also compared.
41
42

43 44 **Model-based imputation of up-to-date BMI measures in CPRD**

45
46 We explored whether outdated BMI measures in CPRD could be usefully updated by
47
48 imputation based on a model predicting changes in individual-level BMI over time. We used
49
50 data from individuals with multiple BMI records to model the expected change in BMI as a
51
52 function of time since BMI recording (restricting to individuals with BMI records ≤ 10 years
53
54 apart). We fitted a linear regression model with change in BMI as the outcome, ~~and; the~~
55

1
2
3
4
5 | main covariate predicting change in BMI was elapsed time, which was included as a 3 knot
6
7 cubic spline to allow for non-linearity; we also included interactions between the spline
8
9 basis variables and indicator variables for age and sex. Feasible weighted least squares
10
11 estimation was used to allow for heteroskedasticity.[20]
12
13

14 Having specified a model for change in BMI over time, we first explored its performance
15
16 among individuals with at least 2 BMIs entered in CPRD, by predicting the most recent BMI
17
18 based on the previous BMI record and the elapsed time; we compared the distribution of
19
20 the errors from this approach with the distribution of the errors from simply using the last
21
22 observation carried forward. We then repeated the comparison with the HSE mean BMI
23
24 data for each calendar year. This time we included all individuals with a BMI record in the
25
26 previous 10 years and used the model described above to impute current BMI at the mid-
27
28 point of the calendar year by predicting the change in BMI since the last available BMI
29
30 record. We did this within a multiple imputation framework (using 5 imputations) to
31
32 account for uncertainty in the modelled changes over time.[21]
33
34

35 The study was approved by the London School of Hygiene and Tropical Medicine Ethics
36
37 Committee.
38
39

40 **Results**

41 **Completeness of BMI data in CPRD**

42
43 In 1990-1994, 37% of individuals had at least one previously recorded BMI, and the
44
45 proportion increased to 77% by 2005-11(Table 1).The proportion of individuals with a recent
46
47 BMI (recorded in the previous 3 years) was lower in each calendar period (35% in 1990-1994
48
49 rising to 51% in 2005-11). BMI completeness generally increased with age up to 75 years,
50
51 with a lower proportion in the oldest age group having data available. Data for single
52
53
54
55
56
57
58
59
60

1
2
3
4
5 calendar years are shown in Appendix Table A1 and illustrate similar patterns. BMI data
6
7 appeared to be consistently more widely available among women than men (Figure 2). As
8
9 expected, BMI completeness was higher in particular clinical subgroups: in total 97% of
10
11 patients with a record of type II diabetes had a recent BMI recorded, along with over 78% of
12
13 those with a diagnosis of schizophrenia/psychoses (Appendix Table A2). This is in line with
14
15 Quality and Outcomes Framework (QOF) which has encouraged BMI monitoring in these
16
17 conditions since 2004.[22] BMI completeness was also high among current statin users (82%
18
19 with a recent BMI available).
20

21
22 There was little extra information available in clinical (“Read”) codes relating to BMI. In the
23
24 most recent calendar period, out of 75518 individuals with no previous BMI record
25
26 available, only 1222 (1.6%) had ever had a clinical code that would enable classification into
27
28 BMI categories (underweight, normal, overweight/obese). Furthermore, for those with a
29
30 previous BMI, only a small proportion had more recent information related to BMI recorded
31
32 in a clinical code (7675/250430 = 3.0% in the most recent period).
33
34

35 **Summary statistics using complete CPRD BMI data and comparison with Health Survey for** 36 37 **England**

38
39 We found that age- and sex-standardised mean BMI based on CPRD data was consistently
40
41 and substantially lower (by up to 1.1kg/m²) than in the HSE data (mean BMI in CPRD =
42
43 25.7kg/m² in 2003 rising to 26.3 in 2010, compared with 26.8 kg/m² [95% CI 26.7 to 26.9]
44
45 and 27.3 [27.1 to 27.5] respectively in HSE; Figure 3).
46
47

48
49 When BMI entries more than 3 years old were discarded, between 33 to 47% of patients
50
51 were lost across calendar years. However, the estimated mean BMI in CPRD was
52
53 considerably closer to what would be expected based on the HSE data, with CPRD data
54
55 underestimating the HSE statistics by only between 0.04 to 0.28kg/m² in individual calendar
56
57
58
59
60

1
2
3
4
5 years, and the CPRD estimate falling within the HSE confidence interval for 2 of the most
6
7 recent 3 calendar years (mean BMI in CPRD = 26.9, 27.0 and 27.0 kg/m² compared with 27.0
8
9 [26.9 to 27.1], 27.0 [26.8 to 27.2] and 27.3 [27.1 to 27.5] in HSE, in 2008, 2009 and 2010
10
11 respectively). Age- and sex-stratified data demonstrated similar patterns, except that in the
12
13 eldest age group (75+ years), restricting to those with recent BMI measures did not bring
14
15 the estimated BMI substantially closer to HSE figures (Appendix Figure A1).
16
17

18
19
20
21 We also compared the proportions classified as obese between CPRD and HSE (Appendix
22
23 Figure A2). Consistent with the previous analysis, using any previous BMI reading to classify
24
25 individuals in CPRD resulted in lower obesity rates than expected based on HSE data, while
26
27 restricting to patients with a recent reading led to estimated obesity rates close to those in
28
29 HSE.
30
31

32 **Model-based imputation of up-to-date BMI measures in CPRD**

33
34 The contrast between BMI summary statistics based on recent measures and those based
35
36 on any previous measures suggested that older BMI records were tending to underestimate
37
38 current BMI. We therefore examined whether a model could be developed to impute
39
40 current BMI, taking into account elapsed time since the last measure. In a linear regression
41
42 model for change in BMI over time, we estimated that on average BMI increased over the
43
44 10-year period following a BMI record for those aged up to 69 years at the time of the
45
46 record and decreased over time in those aged 70 years or more (Appendix Figure A3). We
47
48 tested the predictive performance of our model by predicting the most recent BMI based on
49
50 the previous one, among CPRD patients with more than one recorded BMI available. When
51
52 the older BMI was less than 3 years old, there was little gain in applying the correction
53
54 compared with carrying the older observation forward (Figure 4). However, when there was
55
56

1
2
3
4
5 a longer gap, carrying the previous BMI forward tended to underestimate the later BMI,
6
7 while employing the model-based imputation removed the underestimation and led to
8
9 smaller errors on average (median error = -0.70kg/m^2 [IQR -2.18 to $+0.56$] using last
10
11 observation carried forward, compared with $+0.11\text{kg/m}^2$ [-1.29 to $+1.40$] using model-based
12
13 imputation).
14

15
16 We then repeated the comparison of mean BMI in CPRD versus HSE, this time using our
17
18 model for change in BMI over time as a basis for performing multiple imputations of current
19
20 BMI based on the latest available measure and the time since it was recorded. Estimated
21
22 mean BMIs were now in line with those based on only recent data in the earlier analysis,
23
24 and were only between 0.04 and 0.37kg/m^2 lower than HSE statistics in individual calendar
25
26 years (Figure 3, circles). Even with multiple imputation, confidence intervals remained
27
28 extremely narrow ($<0.07\text{kg/m}^2$) due to the large sample size, so are not shown in the figure.
29
30 Of note, all patients with a BMI recorded up to 10 years before the midpoint of the calendar
31
32 year of interest were now included in the estimation of the “corrected” means; thus in
33
34 individual calendar years only 9 to 13% of patients were dropped, compared to 33-47% of
35
36 patients when dropping BMI records >3 years old.
37
38

39 Discussion

40 Main findings

41
42 BMI completeness has increased over calendar time (rising from 37% in 1990-94 to 77% in
43
44 2005-11). Completeness was higher among females, older age groups, and clinical
45
46 subgroups where recording BMI is encouraged. When BMI on the date of interest was
47
48 assigned to individual patients in CPRD using the last available record, regardless of how
49
50 long ago it was entered, we found that the resulting mean BMI statistics for the CPRD
51
52 population were consistently lower than equivalent HSE estimates (by up to 1.1kg/m^2). This
53
54
55

1
2
3
4
5 appeared to be driven by older BMI records tending to systematically underestimate current
6
7 BMI: when only recent CPRD BMI records (≤ 3 years old) were used, mean BMI statistics
8
9 were closer to HSE estimates. However, a substantial number of patients were then
10
11 excluded altogether (33-47% across years). Finally, we suggested a process for modelling
12
13 changes in BMI after a BMI record, which could allow researchers to impute BMI on the date
14
15 of interest and avoid dropping large numbers without a recent measure from their analyses.
16
17

18 *Comparison with other studies*

19
20 There are very few comparable studies (Appendix Table A2). However, the proportion of
21
22 patients with a recent BMI recording in CPRD is in line with a summary of the QRESEARCH
23
24 database (a similar UK primary care database with data from over 530 general practices
25
26 using EMIS software rather than VISION software);[23] by March 2007, 58% of registered
27
28 patients aged 16+ years had their BMI recorded in the past 5 years; this compares with 51%
29
30 with a BMI recorded in the last 3 years in our analysis (for 2005-11). As in our study, the
31
32 QRESEARCH report showed an increase in completeness over time, rising from 42% in
33
34 2000/01 to 58% in 2007. In a third UK primary care database, THIN (The Health
35
36 Improvement Network), the proportion of newly registered patients between 2004 and
37
38 2006 with BMI data was in line with our findings; 62% of patients had a height recording and
39
40 66% had a weight recording within 12 months of registration.[24]
41
42

43 **Explanation of findings**

44 *Completeness*

45
46
47 Increasing completeness of BMI over time may reflect a general trend towards
48
49 encouragement to record BMI in primary care. Greater BMI completeness among females
50
51 and older age groups may have a number of explanations including higher consultation
52
53
54
55
56
57
58
59
60

1
2
3
4
5 rates in primary care [25 26] and different prevalence^s of diseases in which it is important
6
7 to monitor BMI.
8
9

10 *Comparison of CPRD BMI data with Health Survey for England data*

11
12 Mean BMI based on the CPRD population was lower in each calendar year than equivalent
13
14 HSE estimates when BMI in CPRD was assigned using the last available record; however,
15
16 when the analysis was restricted to those with a recent BMI record, estimates from CPRD
17
18 were close to HSE estimates. This suggests that the substantial proportion of BMI recordings
19
20 in CPRD that were outdated on the date of interest may have driven the apparent
21
22 underestimation of mean BMI in CPRD in the unrestricted analysis. This in turn would imply
23
24 that individual BMIs tend to increase over time, and indeed when we specifically modelled
25
26 changes in BMI over time, we found a pattern of increasing BMI with age for those <70
27
28 years old, consistent with prospective cohort studies with repeated BMI measurements [27-
29
30 29]; this pattern of increasing BMI over time is likely to be driven specifically by weight
31
32 change, since adult height would not change substantially in this age range. A simple
33
34 adjustment of outdated BMIs based on ~~these~~our modelled changes over time brought
35
36 CPRD mean BMI statistics in line with HSE estimates, and when we validated the adjustment
37
38 in a subset of patients with repeated BMI measures, we found smaller errors on average,
39
40 compared with simply carrying outdated BMI records forwards.
41
42

43
44 Of note, we observed that CPRD consistently underestimated BMI compared to HSE among
45
46 those aged ≥75 years, even when only recent records were used; this may reflect the fact
47
48 that institutionalised patients are represented in CPRD but not in HSE: HSE may not be an
49
50 ideal comparison for this age group since elderly people in institutions (who are represented
51
52 in CPRD) may be more likely to be frail and have lower BMIs than those living in private
53
54 households.
55

Implications

1
2
3
4
5
6
7 First, our findings suggest that BMI completeness is likely to vary between studies
8
9 depending on the study population and study period. BMI data are not likely to be missing
10
11 completely at random (for example, missingness may vary by patient characteristics or
12
13 particular diseases). There may be information in the database, however, which predicts
14
15 missingness and which could satisfy the “missing at random” assumption required for
16
17 multiple imputation. A study exploring the potential of imputing missing data in THIN found
18
19 that after multiple imputation, summary statistics of height and weight were comparable
20
21 with data from nationally representative datasets.[24]
22

23
24 Second, our analyses suggest that the common practice of assigning BMI status based on
25
26 the nearest/most recent available record to the index date of interest might lead to
27
28 misclassification, given that a large number of patients have only substantially outdated BMI
29
30 records available at any particular time. Strategies to address this include restricting to
31
32 recent BMI, but this is likely to exclude a large numbers of patients. We have suggested an
33
34 alternative strategy based on updating the outdated BMIs by modelling changes in BMI over
35
36 time, though this is not without drawbacks: the approach requires an assumption that
37
38 individuals with ≥ 2 BMI records available (needed to estimate the model for changes over
39
40 time) are representative of the wider patient population, which may not be the case; it is
41
42 also a more complex strategy, particularly if done within a multiple imputation framework
43
44 to allow for uncertainty in the correction, which could be substantial in studies with smaller
45
46 sample sizes than considered here. Other imputation strategies could also be considered in
47
48 certain contexts, such as the two-fold algorithm which imputes missing data from
49
50 longitudinal variables at particular time points by using adjacent data points.[30] Ultimately,
51
52 the ~~importance of these issues~~ pros and cons of various methods, and the optimal strategy
53
54
55
56
57
58
59
60

1
2
3
4
5 to use is likely to depend on the particular study and the characteristics of the study
6
7 population.
8
9

10 **Strengths and Limitations**

11
12 Results presented here are based on a large random sample from the CPRD, therefore we
13
14 can confidently generalise the findings to the whole CPRD database. Although we cannot
15
16 assume these findings will relate to other routinely collected primary care databases in UK
17
18 based on other IT systems (CPRD is based on practices using VISION), they are likely to be
19
20 similar. This study did not look at BMI recordings among children as this would require a
21
22 different strategy. Completeness among 16-24 year age group may be artificially low
23
24 because weights recorded at age <16 were excluded, so those at the lower end of the age
25
26 group will not have had as much time to accrue weight recordings. We believe HSE to be the
27
28 best available comparison for this study; it is a nationally representative, large sample
29
30 ~~(sample size 14,836 in 2003 and 8,420 in 2010),~~ utilising height and weight recordings
31
32 measured by a trained interviewer, and is weighted for non-response.[19 31] However there
33
34 is a degree of missing data in HSE which is a limitation. In 2010 just over 85% of adults
35
36 interviewed provided valid height and weight recordings. [29] One of the most common
37
38 reasons for missing BMI was refusal (up to 8% were missing due to refusal),[19] which if
39
40 related to BMI status, may bias the estimates of mean BMI in HSE. [Our comparisons](#)
41
42 [between CPRD-based and HSE-based BMI statistics focussed on the mean \(and in the](#)
43
44 [appendix, on the proportion classed as obese\); these are the principal statistics published in](#)
45
46 [HSE trend tables so we were not able to look at a broader range of measures of the BMI](#)
47
48 [distribution that might be of interest to researchers using BMI data in the context of public](#)
49
50 [health. Finally, we have not attempted to quantify or comment on the usefulness of BMI as](#)
51
52
53
54
55
56
57
58
59
60

1
2
3
4
5 a measure of adiposity, and researchers using BMI data should consider whether it is the
6 best available measure for their purposes.
7
8
9

10 **Conclusions**

11
12 Completeness of BMI data in CPRD varies over time and by age and sex. BMI records may
13
14 become outdated over time and naive use could lead to misclassification of BMI status. We
15 used a 3-year cut-off to define a recent BMI; further research could include a systematic
16 analysis of how long BMI records can be considered “up-to-date”, and whether this varies
17 by patient characteristics. The optimal strategy for assigning BMI status to individuals in
18
19
20
21
22 studies based on CPRD and similar electronic healthcare databases is likely to depend on the
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

Conflicts of interest

The authors declare no conflicts of interest.

Funding

This report is independent research arising from a postdoctoral fellowship (for KB) supported by the National Institute for Health Research (PDF-2011-04-007). ID is supported by an MRC methodology research fellowship, LS is supported by a Wellcome Trust senior research fellowship in clinical science.

Data sharing statement

This analysis is based on a large random sample from the Clinical Practice Research Datalink, provided by the UK Medicines and Healthcare products Regulatory Agency. The authors' licence for using these data does not allow sharing of raw data with third parties.

References

1. World Health Organisation. Global Health Risks: Mortality and burden of disease attributable to selected major risks. In: Organisation WH, ed. Geneva, Switzerland, 2009.
2. Flegal KM, Graubard BI, Williamson DF, et al. Excess deaths associated with underweight, overweight, and obesity. *Jama-J Am Med Assoc* 2005;**293**(15):1861-67
3. Swinburn BA, Sacks G, Hall KD, et al. Obesity 1 The global obesity pandemic: shaped by global drivers and local environments. *Lancet* 2011;**378**(9793):804-14
4. Kelly T, Yang W, Chen CS, et al. Global burden of obesity in 2005 and projections to 2030. *Int J Obesity* 2008;**32**(9):1431-37 doi: Doi 10.1038/ljo.2008.102[published Online First: Epub Date]].
5. NHS Information Centre. Health survey for England - 2010: health and lifestyles. Secondary Health survey for England - 2010: health and lifestyles 2011. <http://www.ic.nhs.uk/pubs/hse10report>.
6. CPRD. Clinical Practice Research Database (CPRD) website. Secondary Clinical Practice Research Database (CPRD) website. <http://www.cprd.com/intro.asp>.
7. Delaney JA, Daskalopoulou SS, Brophy JM, et al. Lifestyle variables and the risk of myocardial infarction in the general practice research database. *BMC Cardiovasc Disord* 2007;**7**:38 doi: 1471-2261-7-38 [pii] 10.1186/1471-2261-7-38[published Online First: Epub Date]].
8. Green J, Czanner G, Reeves G, et al. Oral bisphosphonates and risk of cancer of oesophagus, stomach, and colorectum: case-control analysis within a UK primary care cohort. *BMJ* 2010;**341**:c4444 doi: 10.1136/bmj.c4444[published Online First: Epub Date]].
9. Tzoulaki I, Molokhia M, Curcin V, et al. Risk of cardiovascular disease and all cause mortality among patients with type 2 diabetes prescribed oral antidiabetes drugs: retrospective cohort study using UK general practice research database. *British Medical Journal* 2009;**339**:b4731 doi: Artn B4731 Doi 10.1136/Bmj.B4731[published Online First: Epub Date]].
10. Douglas I, Smeeth L, Irvine D. The use of antidepressants and the risk of haemorrhagic stroke: a nested case control study. *Brit J Clin Pharmacol* 2011;**71**(1):116-20 doi: DOI 10.1111/j.1365-2125.2010.03797.x[published Online First: Epub Date]].
11. Andersohn F, Schade R, Suissa S, et al. Long-Term Use of Antidepressants for Depressive Disorders and the Risk of Diabetes Mellitus. *Am J Psychiat* 2009;**166**(5):591-98 doi: DOI 10.1176/appi.ajp.2008.08071065[published Online First: Epub Date]].
12. Lawrenson R, Todd JC, Leydon GM, et al. Validation of the diagnosis of venous thromboembolism in general practice database studies. *Brit J Clin Pharmacol* 2000;**49**(6):591-96
13. Jick H, Zornberg GL, Jick SS, et al. Statins and the risk of dementia. *Lancet* 2000;**356**(9242):1627-31
14. van Staa TP, Wegman S, de Vries F, et al. Use of statins and risk of fractures. *Jama-J Am Med Assoc* 2001;**285**(14):1850-55
15. Office for National Statistics. *Key Health Statistics from General Practice 1998: Analyses of Morbidity and Treatment Data, Including Time Trends, England and Wales*. London: Office for National Statistics, 2000.
16. Parkinson JP, Davis S, Van Staa T. The General Practice Research Database: now and the future. In: Mann R, Andrews EB, eds. *Pharmacovigilance*. Chichester: John Wiley and Sons, 2007:341-48.
17. Schoonen WM, Thomas SL, Somers EC, et al. Do selected drugs increase the risk of lupus? A matched case-control study. *Br J Clin Pharmacol* 2010;**70**(4):588-96 doi: 10.1111/j.1365-2125.2010.03733.x[published Online First: Epub Date]].
18. NHS Information Centre. Health Survey for England - 2010: Trend tables. Secondary Health Survey for England - 2010: Trend tables. <http://www.ic.nhs.uk/statistics-and-data->

- 1
2
3
4
5 [collections/health-and-lifestyles-related-surveys/health-survey-for-england/health-survey-](#)
6 [for-england--2010-trend-tables.](#)
- 7 19. Aresu M, Boodhna G, Bryson A, et al. Volume 2: Methods and documentation. In: Craig R,
8 Mindell J, eds. Health Survey for England 2010. Leeds: NHS Information Centre for health
9 and social care, 2011.
- 10 20. Greene WH. *Econometric Analysis*. Upper Saddle River, New Jersey: Prentice Hall., 1997.
- 11 21. Rubin DB. *Multiple Imputation for Nonresponse in Surveys*: J. Wiley & Sons, New York. , 1987.
- 12 22. NHS. Quality and Outcomes Framework guidance for GMS contract 2011/12: Employers and
13 British Medical Association, 2011.
- 14 23. NHS Information Centre. A summary of public health indicators using electronic data from
15 primary care. Secondary A summary of public health indicators using electronic data from
16 primary care 2008. [http://www.ic.nhs.uk/article/2021/Website-](http://www.ic.nhs.uk/article/2021/Website-Search?productid=4287&q=A+summary+of+public+health+indicators+using+electronic+data+from+primary+care.&sort=Relevance&size=10&page=1&area=both#top)
17 [Search?productid=4287&q=A+summary+of+public+health+indicators+using+electronic+data](http://www.ic.nhs.uk/article/2021/Website-Search?productid=4287&q=A+summary+of+public+health+indicators+using+electronic+data+from+primary+care.&sort=Relevance&size=10&page=1&area=both#top)
18 [+from+primary+care.&sort=Relevance&size=10&page=1&area=both#top.](http://www.ic.nhs.uk/article/2021/Website-Search?productid=4287&q=A+summary+of+public+health+indicators+using+electronic+data+from+primary+care.&sort=Relevance&size=10&page=1&area=both#top)
- 19 24. Marston L, Carpenter JR, Walters KR, et al. Issues in multiple imputation of missing data for large
20 general practice clinical databases. *Pharmacoepidemiol Drug Saf* 2010;**19**(6):618-26 doi:
21 10.1002/pds.1934[published Online First: Epub Date]].
- 22 25. Rowlands S, Moser K. Consultation rates from the General Practice Research Database. *Brit J Gen*
23 *Pract* 2002;**52**(481):658-60
- 24 26. The Health and Social Care Information Centre. Trends in Consultation Rates in General Practice
25 1995 to 2009. Secondary Trends in Consultation Rates in General Practice 1995 to 2009
26 2009. [http://www.ic.nhs.uk/article/2021/Website-](http://www.ic.nhs.uk/article/2021/Website-Search?productid=729&q=gresearch&sort=Relevance&size=10&page=1&area=both#top)
27 [Search?productid=729&q=gresearch&sort=Relevance&size=10&page=1&area=both#top.](http://www.ic.nhs.uk/article/2021/Website-Search?productid=729&q=gresearch&sort=Relevance&size=10&page=1&area=both#top)
- 28 27. Li L, Law C, Power C. Body mass index throughout the life-course and blood pressure in mid-adult
29 life: a birth cohort study. *Journal of Hypertension* 2007;**25**(6):1215-23
- 30 28. Silverwood RJ, Pierce M, Thomas C, et al. Overweight across adult life and kidney function at age
31 60-4 years: the 1946 British birth cohort study. (awaiting publication) 2012
- 32 29. Tirosh A, Shai I, Afek A, et al. Adolescent BMI trajectory and risk of diabetes versus coronary
33 disease. *N Engl J Med* 2011;**364**(14):1315-25 doi: 10.1056/NEJMoa1006992[published Online
34 First: Epub Date]].
- 35 30. Welch C, Bartlett J, Petersen I. Application of multiple imputation using the two-fold fully
36 conditional specification algorithm in longitudinal clinical data. (In Press). *Stata Journal* 2013
- 37 31. Mindell J, Biddulph JP, Hirani V, et al. Cohort Profile: The Health Survey for England. *International*
38 *Journal of Epidemiology* 2012;1-9 doi: 10.1093/ije/dyr199[published Online First: Epub
39 Date]].
- 40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

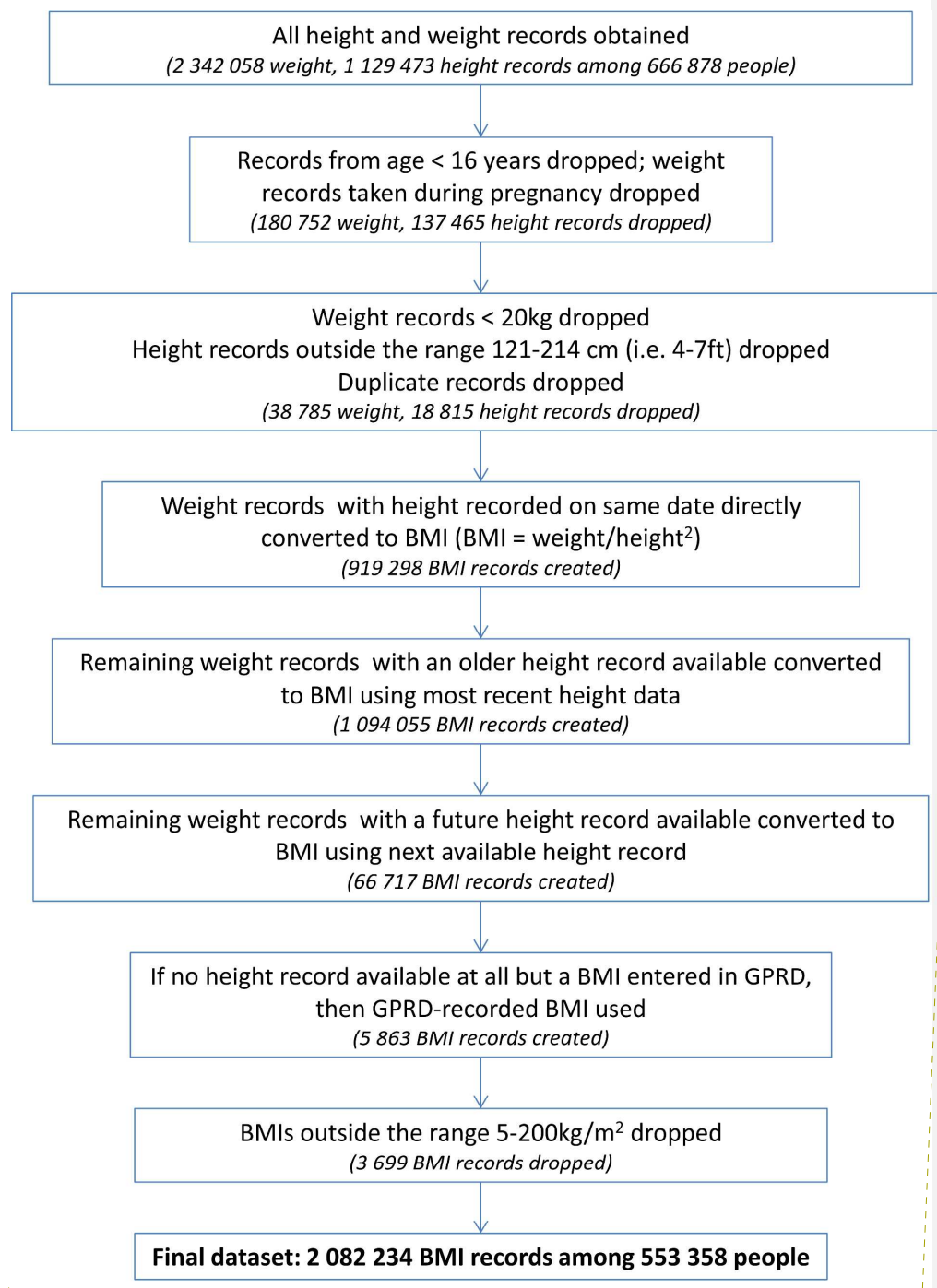
Table 1: Completeness of BMI data in the CPRD, by age and calendar period

Age group (yrs)	1990-4	1995-9	2000-4	2005-2011
16-24^a				
N registered	11423	17501	34452	42546
BMI in previous 3y (%)	26	28	25	32
BMI any previous (%)	26	37	30	37
25-34				
N registered	17477	29923	48659	50413
BMI in previous 3y (%)	37	39	36	49
BMI any previous (%)	38	66	67	72
35-44				
N registered	15953	28838	55991	61014
BMI in previous 3y (%)	36	36	31	46
BMI any previous (%)	39	67	71	80
45-54				
N registered	14507	27765	48093	55564
BMI in previous 3y (%)	39	37	32	50
BMI any previous (%)	42	70	73	84
55-64				
N registered	11680	20843	42258	49380
BMI in previous 3y (%)	42	40	37	57
BMI any previous (%)	44	74	77	87
65-74				
N registered	10678	17605	30997	34508
BMI in previous 3y (%)	36	37	40	67
BMI any previous (%)	38	71	79	91
75+				
N registered	8637	16005	29384	32523
BMI in previous 3y (%)	28	32	37	64
BMI any previous (%)	28	56	69	87
Total				
N registered	90355	158480	289834	325948
BMI in previous 3y (%)	35	36	34	51
BMI any previous (%)	37	64	67	77

N registered is all those under follow-up at mid-point of the period

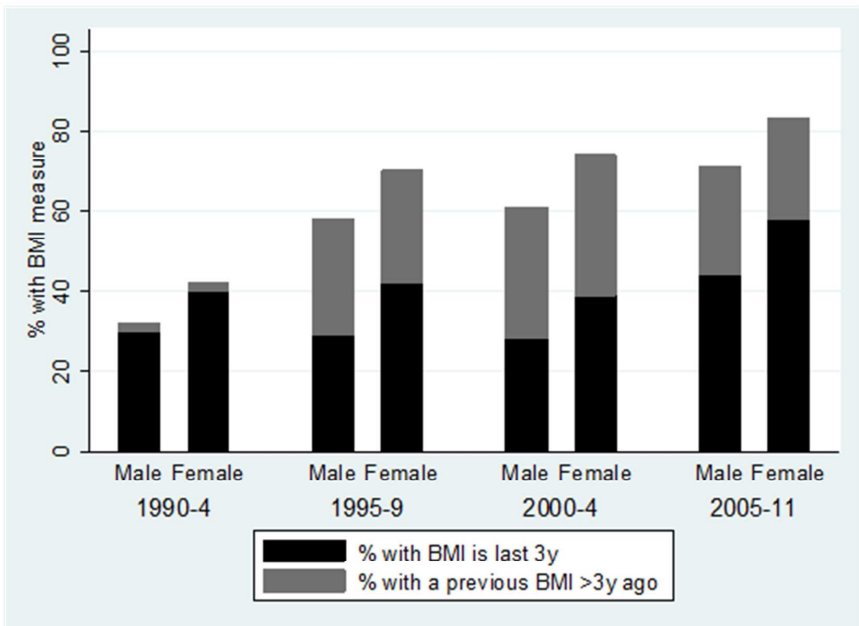
^aNote, BMI measurements from age <16 years were not counted in this analysis, hence completeness in the 16-24 age group may be artificially low

Figure 1: Initial data processing to generate BMI for analysis



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

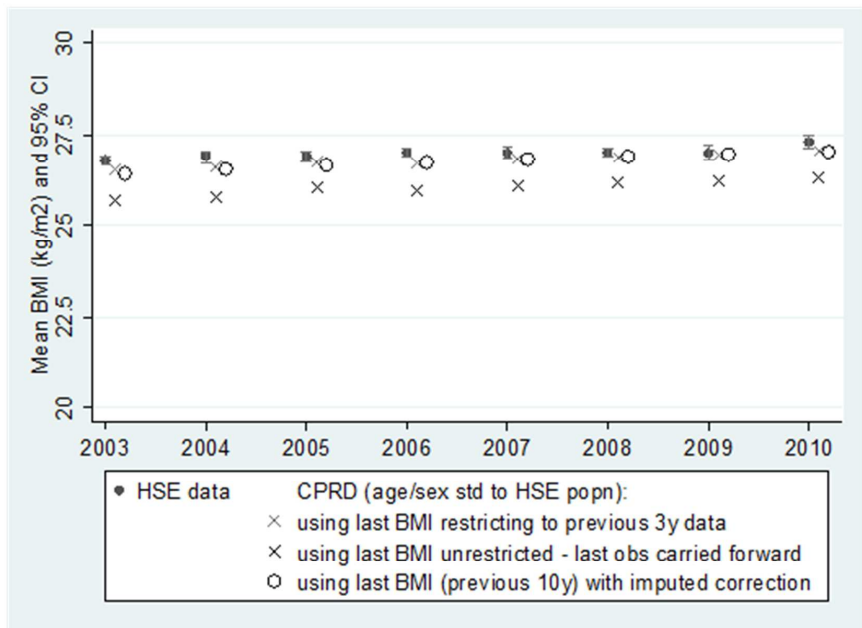
Figure 2: Completeness of BMI data in CPRD, by gender and calendar period



Note: Completeness data for each calendar period are based on all those under follow-up at mid-point of the period

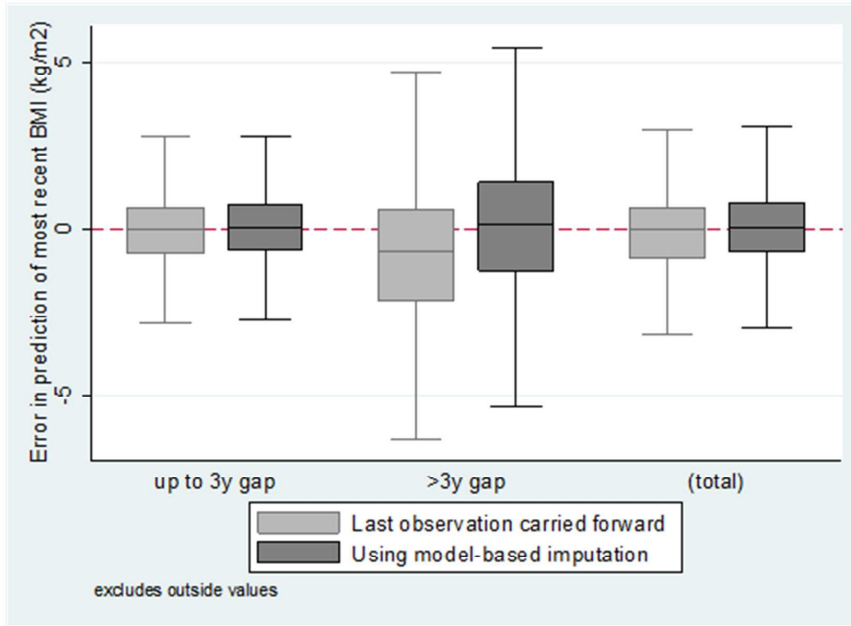
Review only

Figure 3: Mean BMI over calendar time comparing those with BMI recorded in CPRD (English practices) with the Health Survey for England 2010 data



Note: CPRD figures are age- and sex- standardised to the Health Survey for England study population
 CPRD statistics are based on all patients registered at the mid-point of the calendar period and with a suitable previous BMI measure available (i.e. either any previous, or within the last 3 years)

Figure 4: Error in prediction of most recent BMI from older BMI, comparing simple last observation carried forward with model-based imputation of up to date BMI – stratified by time gap between readings



Author contributions

I, Krishnan Bhaskaran, developed the analytical strategy for this paper, processed and analysed the data and wrote the paper.

I, Harriet Forbes, was involved in discussing the data processing and analysis of the data, as well as the writing of the paper.

I, Liam Smeeth, was involved in discussions of the analytical approach to this study and made comments on the analysis and the writing of the paper.

I, Ian Douglas, was involved in discussions of the analytical approach to this study and made comments on the analysis and the writing of the paper.

I, David Leon, was involved in discussions of the analytical approach to this study and made comments on the analysis and the writing of the paper.

STROBE Statement—Checklist of items that should be included in reports of *cohort studies*

NOTE Page numbers refer to revised manuscript, tracked changes version

	Item No	Recommendation
Title and abstract	1	(a) Indicate the study's design with a commonly used term in the title or the abstract n/a (we did not think there was an appropriate design keyword/term to describe this study as it is not a standard "exposure/outcome" study but is rather providing data quality information on a common exposure/covariate) (b) Provide in the abstract an informative and balanced summary of what was done and what was found P2
Introduction		
Background/rationale	2	Explain the scientific background and rationale for the investigation being reported P4-5
Objectives	3	State specific objectives, including any prespecified hypotheses P4-5
Methods		
Study design	4	Present key elements of study design early in the paper P6-7
Setting	5	Describe the setting, locations, and relevant dates, including periods of recruitment, exposure, follow-up, and data collection P5-6
Participants	6	(a) Give the eligibility criteria, and the sources and methods of selection of participants. Describe methods of follow-up P5-6 (b) For matched studies, give matching criteria and number of exposed and unexposed
Variables	7	Clearly define all outcomes, exposures, predictors, potential confounders, and effect modifiers. Give diagnostic criteria, if applicable P5-7
Data sources/ measurement	8*	For each variable of interest, give sources of data and details of methods of assessment (measurement). Describe comparability of assessment methods if there is more than one group P5-6
Bias	9	Describe any efforts to address potential sources of bias P6-8
Study size	10	Explain how the study size was arrived at P5
Quantitative variables	11	Explain how quantitative variables were handled in the analyses. If applicable, describe which groupings were chosen and why P6-7
Statistical methods	12	(a) Describe all statistical methods, including those used to control for confounding P6-7 (b) Describe any methods used to examine subgroups and interactions P6-7 (c) Explain how missing data were addressed

		P7
		(d) If applicable, explain how loss to follow-up was addressed
		n/a
		(e) Describe any sensitivity analyses
		n/a
Results		
Participants	13*	(a) Report numbers of individuals at each stage of study—eg numbers potentially eligible, examined for eligibility, confirmed eligible, included in the study, completing follow-up, and analysed
		(b) Give reasons for non-participation at each stage
		(c) Consider use of a flow diagram
		FIG 1
Descriptive data	14*	(a) Give characteristics of study participants (eg demographic, clinical, social) and information on exposures and potential confounders
		(b) Indicate number of participants with missing data for each variable of interest
		(c) Summarise follow-up time (eg, average and total amount)
		P8-9 and FIG 2
Outcome data	15*	Report numbers of outcome events or summary measures over time
		n/a (no specific outcome)
Main results	16	(a) Give unadjusted estimates and, if applicable, confounder-adjusted estimates and their precision (eg, 95% confidence interval). Make clear which confounders were adjusted for and why they were included
		(b) Report category boundaries when continuous variables were categorized
		(c) If relevant, consider translating estimates of relative risk into absolute risk for a meaningful time period
		n/a (not an “exposure/outcome” study)
Other analyses	17	Report other analyses done—eg analyses of subgroups and interactions, and sensitivity analyses
		P9-11
Discussion		
Key results	18	Summarise key results with reference to study objectives
		P11-12
Limitations	19	Discuss limitations of the study, taking into account sources of potential bias or imprecision. Discuss both direction and magnitude of any potential bias
		P15-16
Interpretation	20	Give a cautious overall interpretation of results considering objectives, limitations, multiplicity of analyses, results from similar studies, and other relevant evidence
		P16
Generalisability	21	Discuss the generalisability (external validity) of the study results
		P15
Other information		
Funding	22	Give the source of funding and the role of the funders for the present study and, if

applicable, for the original study on which the present article is based

P17

*Give information separately for exposed and unexposed groups.

Note: An Explanation and Elaboration article discusses each checklist item and gives methodological background and published examples of transparent reporting. The STROBE checklist is best used in conjunction with this article (freely available on the Web sites of PLoS Medicine at <http://www.plosmedicine.org/>, Annals of Internal Medicine at <http://www.annals.org/>, and Epidemiology at <http://www.epidem.com/>). Information on the STROBE Initiative is available at <http://www.strobe-statement.org>.

For peer review only