

BMJ Open

BMJ Open is committed to open peer review. As part of this commitment we make the peer review history of every article we publish publicly available.

When an article is published we post the peer reviewers' comments and the authors' responses online. We also post the versions of the paper that were used during peer review. These are the versions that the peer review comments apply to.

The versions of the paper that follow are the versions that were submitted during the peer review process. They are not the versions of record or the final published versions. They should not be cited or distributed as the published version of this manuscript.

BMJ Open is an open access journal and the full, final, typeset and author-corrected version of record of the manuscript is available on our site with no access controls, subscription charges or pay-per-view fees (<http://bmjopen.bmj.com>).

If you have any questions on BMJ Open's open peer review process please email editorial.bmjopen@bmj.com

BMJ Open

Identifying clinical features in studies using primary care electronic health record data

Journal:	<i>BMJ Open</i>
Manuscript ID	bmjopen-2017-019637
Article Type:	Research
Date Submitted by the Author:	19-Sep-2017
Complete List of Authors:	Watson, Jessica; University of Bristol, School of Social and Community Medicine Nicholson, Brian; University of Oxford, Nuffield Dept Primary Care Health Sciences Hamilton, Willie; University of Exeter Medical School, Primary Care Diagnostics Price, Sarah; University of Exeter Medical School,
Primary Subject Heading:	Health services research
Secondary Subject Heading:	General practice / Family practice, Research methods
Keywords:	Electronic Health Records, Clinical coding, PRIMARY CARE, STATISTICS & RESEARCH METHODS, EPIDEMIOLOGY

SCHOLARONE™
Manuscripts

Identifying clinical features in studies using primary care electronic health record data

Jessica Watson¹, Brian D Nicholson², Willie Hamilton³, Sarah Price³

¹Centre for Academic Primary Care, Bristol Medical School, University of Bristol, Canynge Hall, 39 Whatley Road, Bristol, BS8 2PS

²Nuffield Department of Primary Care Health Sciences, Radcliffe Primary Care Building, University of Oxford, OX2 6GG.

³University of Exeter Medical School

Correspondence to: Jessica.Watson@bristol.ac.uk

Word Count: 3,464

ABSTRACT

Objective Analysis of routinely collected Electronic Health Record (EHR) data from primary care is reliant upon the creation of codelists to define clinical features of interest. As EHR research increases, it is important to avoid waste in research through incomplete or unusable research publications. To improve scientific rigor, transparency and replicability we describe and demonstrate a standardised reproducible methodology for clinical codelist development.

Design We describe a three stage process for developing clinical codelists. First, the clear definition *a priori* of the clinical feature of interest using reliable clinical resources. Second, development of a list of potential codes using statistical software to comprehensively search all available codes. Third, a modified Delphi process to reach consensus between primary care practitioners on the most relevant codes, including the generation of an ‘uncertainty’ variable to allow sensitivity analysis. We illustrate the method by developing a codelist for shortness of breath, including modifiable syntax for commonly used statistical software.

Results

Of 78 candidate codes, 29 were excluded as inappropriate. Complete agreement was reached for 44 (90%) of the remaining codes, with partial disagreement over 5 (10%). 13,091 episodes of shortness of breath were then identified in an EHR sample of 28,216 patients diagnosed with lung cancer. Sensitivity analysis demonstrates that codes with the greatest uncertainty tend to be rarely used in clinical practice.

1
2
3 **Conclusions** Although initially time-consuming, using a rigorous and reproducible method
4
5 for codelist generation ‘future-proofs’ findings, and an auditable, modifiable syntax for
6
7 codelist generation enables sharing and replication of EHR studies. Published codelists
8
9 should be badged by quality and report the methods of codelist generation including;
10
11 definitions and justifications associated with each codelist; the syntax or search method; the
12
13 number of candidate codes identified; and the categorisation of codes after Delphi review.
14
15
16
17
18

19 **Keywords**

20
21 Electronic Health Records, Clinical Coding, Primary Health Care, Epidemiological Methods
22
23
24
25
26
27

28 **Strengths and Limitations of this study**

- 29 • This paper presents rigorous reproducible methods for codelist generation to increase
30
31 transparency and replicability in EHR studies.
- 32
33 • Clear *a priori* definition of the feature of interest ensures clinical relevance, and
34
35 enables future researchers to assess the applicability of existing codelists to future
36
37 research questions.
- 38
39 • Generation of auditable, replicable and modifiable syntax for codelists enables
40
41 replication and ‘future-proofs’ codelists.
- 42
43 • Using a Delphi approach to reach consensus on inclusion of codes allows sensitivity
44
45 analysis to explore the impact of uncertainty in coding.
- 46
47 • Using multiple clinicians in a Delphi panel reviewing codes may be unfeasible and
48
49 inefficient for studies with large numbers of codes; a compromise of using two
50
51 clinicians per feature from a panel of six offers a reasonable trade-off.
52
53
54
55
56
57
58
59
60

INTRODUCTION

Electronic Health Records (EHRs) have been used in routine primary care practice in the United Kingdom (UK) for at least 20 years.¹ EHRs are a rich resource for researchers, and are increasingly used in epidemiological and medical research resulting in over 1,500 publications since 2000, increasing from ~80 in 2005 to more than 450 in 2015/2016.

There are three well established UK primary care EHR databases: The Clinical Practice Research Datalink (CPRD) including 4.4 million currently registered patients, covering 6.9% of the UK population;² The Health Improvement Network (THIN) including 3.6 million currently registered patients giving ~5.7% coverage of the nation;³ and QResearch® including 24 million currently and previously registered patients in the UK.⁴ All three databases record coded anonymised information about patients: demographics, diagnoses, symptoms, prescriptions, immunisation history, referral information, and test results.

Linkages enable follow-up of patients beyond the primary care setting; for example, to data recorded by the Office for National Statistics (ONS), the National Cancer Registration Service (NCRS) and to Hospital Episode Statistics. Integrated primary and secondary care databases are also being developed. For example, ResearchOne includes data for over 5 million patients from General Practice, Child Health, Community Health, Out-of-Hours, Palliative Hospital, Accident and Emergency and Acute Hospital.

(<http://www.researchone.org/>).

A key stage in EHR research is identifying exposures and outcomes of interest. This apparently simple task is made more complicated by the fact that EHR clinical data is generally stored as codes, often including qualitative information, such as 'abdominal pain', 'left iliac fossa pain' and 'intermittent abdominal pain'. These separate codes need to be grouped into codelists or thesauri, with the groups containing all the codes pertaining to the variable of interest. However, the methods used to develop codelists are not standardised, and

1
2
3 are often poorly reported. They are an increasingly recognised source of bias in EHR
4
5 research, owing to both inclusion of inappropriate codes and omission of important codes. To
6
7 address this, the RECORD Statement states that '*a complete list of codes and algorithms used*
8
9 *to classify exposures, outcomes, confounders, and effect modifiers should be provided*'.⁵
10
11 Clinicalcodes.org has been developed by the University of Manchester to encourage
12
13 researchers to publish clinical codelists used in EHR research⁶ and some other Universities
14
15 are developing their own open access, citable, repositories of codelists for example the
16
17 University of Bristol⁷ and University of Cambridge.⁸ The current clinicalcodes.org repository
18
19 contains 72916 clinical codes deposited within 432 codelists (<https://clinicalcodes.org>), in the
20
21 format of a list of papers and associated codes. This repository is a necessary step forward
22
23 towards addressing transparency; however, it does not tackle the potential for bias, as it is not
24
25 sufficient to address the issues of scientific rigour and reproducibility in codelist
26
27 development.
28
29
30
31

32
33 The problem is illustrated by brief examination of codelists recently deposited on the
34
35 repository. Without a clear definition of the clinical variable a codelist is designed to
36
37 encapsulate, it is not possible to critique or evaluate it for peer review, or to decide whether it
38
39 is generalisable to other studies. For example, codelists deposited for cancer
40
41 (<https://clinicalcodes.rss.mhs.man.ac.uk/medcodes/article/50/>) do not adhere to the
42
43 standardised International Classification of Diseases definition of cancer, i.e. ICD codes C00
44
45 to C97, as 193 (~9%) of the 2254 Read codes related to carcinoma *in situ* (ICD D00 to D09).
46
47 Furthermore, 100 (~4%) codes were obsolete, or they indicated the absence of cancer or they
48
49 were completely unrelated to cancer.
50
51

52
53 This demonstrates the need to establish standardised methods for codelist development.
54

55
56 Currently recommended methods, for example Davé and Petersen⁹ and CALIBERcodelists
57
58 (<http://caliberanalysis.r-forge.r-project.org/>) need updating, not only because they omit steps
59
60

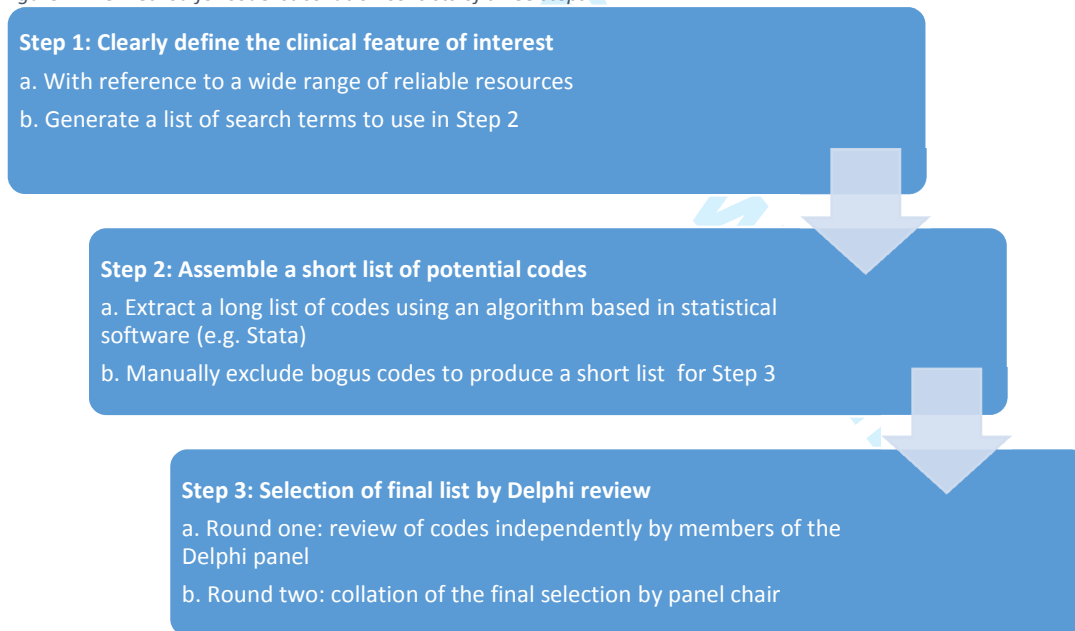
to standardise the definition of clinical terms, but also because they are based in the Read code system, which is being superseded by SNOMED CT codes (Systematized Nomenclature of Medicine -- Clinical Terms) in April 2018.

We have significant experience in EHR research, with ~40 published studies conducted in the CPRD since 2012. We have developed and refined rigorous methods for developing clinical codelists for use in CPRD studies independent of the Read code system. The aim of this paper is to report a clear, standardised, reproducible methodology, and to increase scientific rigour in conduct of EHR research. The method is illustrated using the CPRD, but applies equally well to other large EHR databases.

METHODS

Our method for collating clinical codelists involves three stages, described in Figure 1.

Figure 1 The method for codelist collation consists of three steps



Step 1: Clearly define the clinical feature of interest (symptom, disease or illness) a priori

The first step is to clearly define the clinical feature of interest and establish inclusion and exclusion criteria. This requires clinical input, particularly from GPs who are best placed to

1
2
3 understand how clinical features are coded in a primary care setting. Reliable sources of
4
5 clinical information should be used; for example:
6
7

- 8 • International Classification of Primary Care, which defines symptoms and diagnoses,
9 provides synonyms for them and, importantly, lists what should be excluded from the
10 definition¹⁰
11
- 12 • The BMJ Best Practice guidelines (<http://bestpractice.bmj.com/best-practice/welcome.html>)
13
- 14 • NICE Clinical Knowledge Summaries (<http://cks.nice.org.uk/>)
15
- 16 • ICD10 (<http://apps.who.int/classifications/icd10/browse/2016/en>) – this is less useful
17 for symptoms, as it focuses on diseases
18
- 19 • Medical Subject Headings (MeSH)
20 (https://www.nlm.nih.gov/mesh/2016/mesh_browser/MBrowser.html)
21
22
23
24
25

26 Other potential resources include patient support groups, online discussion forums, and
27 already published codelists (e.g. <https://clinicalcodes.org>). Hierarchical classifications such as
28 Read, SNOMED or ICD-10 may be useful for identifying additional search terms and
29 synonyms.
30
31
32
33
34

35
36 For some symptoms, it is necessary to tailor the definition to the context of the disease under
37 investigation. Abdominal pain is a good example, where pancreatic disease may cause pain in
38 the epigastrium and left hypochondrium, whereas disorders in the sigmoid colon generate
39 pain in the left iliac fossa.
40
41
42
43
44

45 **Step 2 – assembling list of codes that may be used to record the clinical feature**

46
47 The second stage consists of identifying all potential codes that might be used by GPs to
48 record the clinical feature of interest defined in Step 1 and collating them into a list.
49
50
51

52 This is done in several steps; we use Stata for this, but other software is possible.
53
54
55
56
57
58
59
60

1
2
3 First, using the resources listed in Step 1, an exhaustive list of synonyms for the outcome of
4
5 interest is generated. Box 1 uses the example of shortness of breath.
6
7

8 \beginbox1
9

10 **Box 1: Shortness of breath**

11 *ICPC*

- 12 • ICPC code: R02 (exclude: wheezing R03; stridor R04; hyperventilation R98)

13 *BMJ Best Practice*

- 14 • Dyspnoea, also known as shortness of breath or breathlessness, is a subjective
15 sensation of breathing discomfort (<http://bestpractice.bmj.com/best-practice/monograph/862.html>)

16 *NICE CKS*

- 17 • Breathlessness is the distressing sensation of a deficit between the body's demand for
18 breathing and the ability of the respiratory system to satisfy that demand.
19 (<http://cks.nice.org.uk/breathlessness#!backgroundsub>)
- 20 • Breathlessness can be classified by its speed of onset as:
 - 21 ○ Acute breathlessness — when it develops over minutes, hours, or days.
 - 22 ○ Chronic breathlessness — when it develops over weeks or months.

23 *ICD10*

- 24 • ICD10 code: R06 – dyspnoea, orthopnoea, shortness of breath

25 *MeSH*

- 26 • MeSH: Difficult or labored breathing. Breathlessness, dyspnea

27 *Patient forums*

- 28 • Puffed, winded

29 *GP colleagues*

- 30 • Consider including ‘respiratory insufficiency’?

31 \endbox1
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

1
2
3 Second, the lookup file of all medical codes provided by the CPRD (medical.txt)¹ is opened
4
5 using Stata. This contains the alphanumeric Read code originally used by the GP to enter the
6
7 clinical information, the CPRD's proprietary 'medcode' (which is simply a numeric
8
9 equivalent of the Read code), as well as a verbal description (variable 'desc') common to
10
11 both the medcode and Read code. A variable for the clinical outcome of interest (here, 'sob'
12
13 for shortness of breath) is created and set to zero (see Box 2). Then Stata searches the verbal
14
15 description of each code, and sets 'sob' to 1 if it contains any of the synonyms. Example
16
17 syntax to replicate this process in the statistical software package R is provided in the
18
19 supplementary materials, using the lookup file of all medical codes that comes with the
20
21 CPRD browsers. Note that, in this file, the verbal description is called 'readterm' rather than
22
23 'desc'.
24
25
26
27
28
29

30 \beginbox2

```
31  
32 insheet using "medical.txt", clear  
33  
34  
35 *generate a binary variable for shortness of breath (sob) and set  
36 its value to zero  
37 g sob=0  
38  
39  
40  
41 /* search the verbal description of the Read code/medcode and change  
42 the value of variable sob from 0 to 1 if it contains words that  
43 suggest the code might be about the clinical feature of interest*/  
44 replace sob=1 if regexm(desc, "[Ss]hortness [Oo]f  
45 [Bb]reath|SHORTNESS OF BREATH")  
46 replace sob=1 if regexm(desc, "[Ss] [Oo] [Bb]|SOB")  
47 replace sob=1 if regexm(desc, "pnoea|PNOEA")  
48 replace sob=1 if regexm(desc, "pnea|PNEA")  
49 replace sob=1 if regexm(desc, "[Pp]uffed|PUFFED")  
50  
51  
52  
53  
54  
55  
56  
57
```

58 ¹ THIN and QResearch® provide equivalent files.
59
60

```
1
2
3   replace sob=1 if regexm(desc, "[Ss]hort [Oo]f [Bb]reath|SHORT OF
4   BREATH")
5
6   replace sob=1 if regexm(desc, "[Ss]hort|SHORT") & regexm(desc,
7   "[Bb]reath|BREATH")
8
9   replace sob=1 if regexm(desc, "[Ww]inded|WINDED")
10
11  replace sob=1 if regexm(desc, "[Dd]ifficult|DIFFICULT") &
12  regexm(desc, "[Bb]reath|BREATH")
13
14  replace sob=1 if regexm(desc, "[Ll]abour|LABOUR") & regexm(desc,
15  "[Bb]reath|BREATH")
16
17  replace sob=1 if regexm(desc, "[Ll]abor|LABOR") & regexm(desc,
18  "[Bb]reath|BREATH")
19
20  replace sob=1 if regexm(desc, "[Bb]reathless|BREATHLESS")
21
22  replace sob=1 if regexm(desc, "[Dd]istress|DISTRESS") & regexm(desc,
23  "[Bb]reath|BREATH")
24
25  replace sob=1 if regexm(desc, "[Dd]istress|DISTRESS") & regexm(desc,
26  "[Rr]espir|RESPIR")
27
28  replace sob=1 if regexm(desc, "[Ii]suff|INSUFF") & regexm(desc,
29  "[Bb]reath|BREATH")
30
31  replace sob=1 if regexm(desc, "[Ii]suff|INSUFF") & regexm(desc,
32  "[Rr]espir|RESPIR")
33
34
35
36  /*order the dataset so that values of variable sob==1 are all placed
37  together*/
38
39  gsort sob
40
41
42  /* Manual check for bogus codes - manually change sob==1 to sob==0
43  if the code is clearly inappropriate. */
44
45
46  edit medcode readcode desc sob
47
48
49
50  /*Retain only those codes that are specifically about sob*/
51  keep if sob==1
52
53
54  /*Retain the variables of interest*/
55  keep medcode readcode desc sob
56
57
58
59
60
```

```
1
2
3      sort medcode
4
5
6      /*Save the file as a library for sob for the Delphi process*/
7      save "sob_library.dta", replace
8
9
10     /*Export as an Excel file
11     export excel using "sob_library", replace
12
13
14
15 \endbox2
```

16
17
18 The manual check for bogus codes should err on the side of caution, only rejecting codes that
19
20 are clearly inappropriate according to predefined inclusion and exclusion criteria. Common
21
22 reasons for exclusion are that search terms can pick up bogus codes (e.g. transobturator tape
23
24 contains the letter sequence ‘sob’), or codes indicating a family history of a condition or
25
26 screening for a condition rather than presence of a condition. This generates a list of potential
27
28 codes that is then exported to Excel and reviewed manually in a Delphi-type process (Step 3).
29
30
31
32
33

34 **Step 3: Delphi review of codes**

35
36 The codelist is reviewed by one practising GP, plus at least one other GP from a panel of six,
37
38 using a modified nominal group technique¹¹. Each GP independently categorises the list,
39
40 ranking each Read code/medcode using a 3-point scale as follows:
41
42

43 1 = Definitely Include - the code accurately defines the clinical feature of interest, and GPs
44
45 would definitely use it.

46
47 2 = Uncertain – it remains unclear whether the code accurately reflects the clinical feature of
48
49 interest, or whether GPs would use it.

50
51 3 = Definitely Exclude – the code does not define the clinical feature of interest, and GPs
52
53 definitely would not use it.
54
55
56
57
58
59
60

1
2
3 Panel members are encouraged to add comments explaining their reasons for exclusion or
4
5 uncertainty, in the knowledge that these comments will be shared with an independent panel
6
7 chair who will collate all of the results.
8

9
10 Codes are retained in the final list if they are ranked ‘1=Definitely include’ by at least one of
11
12 the GPs, as this indicates sufficient evidence that the code may be used to record that clinical
13
14 feature. Codes are dropped if they are ranked as “3 = Definitely exclude”, or as “2 =
15
16 Uncertain” by *all* reviewers.
17

18
19
20 An ‘Uncertainty’ variable is also generated for retained codes, to enable sensitivity analyses
21
22 that remove codes for which any uncertainty exists about accuracy or use. The ‘uncertainty’
23
24 variable is defined as follows:
25

26
27 0 = ‘Minimal Uncertainty’, as all panel members ranked the code as ‘1=Definitely include’
28

29
30 1 = ‘Moderate Uncertainty’, at least one panel member ranked the code as ‘2=Uncertain’
31

32
33 2 = ‘Maximal Uncertainty’, at least one panel member ranked the code as ‘3=Definitely
34
35 exclude’
36

37
38 Once the codelist has been generated, a frequency check may be performed using the study’s
39
40 dataset to identify the frequency of the clinical events attributed to each clinical code. If the
41
42 Delphi process has been accurate, the most frequent events will most likely be coded as “0 =
43
44 Minimal Uncertainty”, whereas there will be fewer events for the codes ranked as “1 =
45
46 Moderate Uncertainty” or as “2 = Maximal Uncertainty”.
47
48

49 50 51 **Illustrative example using CPRD medical codes list** 52

53
54 The library of codes for shortness of breath was used to estimate the frequency of this
55
56 symptom in the year before diagnosis of lung cancer. Participants were CPRD patients aged
57
58
59
60

1
2
3 over 18 years who received an incident diagnosis of lung cancer between 1 January 2000 and
4
5 30 November 2016.

6
7 Outcome measures included the number of patients reporting shortness of breath in the year
8
9 before they were diagnosed with lung cancer, the proportion of all lung cancer patients
10
11 reporting shortness of breath, and the total number of episodes of shortness of breath.
12
13

14
15
16 In addition, a sensitivity analysis was carried out restricting the analysis to codes whose
17
18 uncertainty variable was coded 0 (=‘Minimal Uncertainty’), i.e. there was full agreement in
19
20 the Delphi process that the code should be included.
21
22

23 24 25 **RESULTS**

26
27 The codelist generated for shortness of breath is presented here to illustrate the method we
28
29 have described. The clinical resources reviewed in Step 1 (see Box 1 in Methods) indicated
30
31 that the codes used to define shortness of breath should capture evidence of ‘dyspn[o]ea’,
32
33 ‘shortness of breath’ (and its abbreviated term ‘sob’), ‘breathlessness’, ‘orthopn[o]ea’,
34
35 “‘difficult’ & ‘breathing’”, “‘labo[u]red’ & ‘breathing’”, “‘breathing’ & ‘discomfort’”,
36
37 ‘puffed’, ‘winded’, ‘respiratory distress’ and ‘respiratory insufficiency’.
38
39

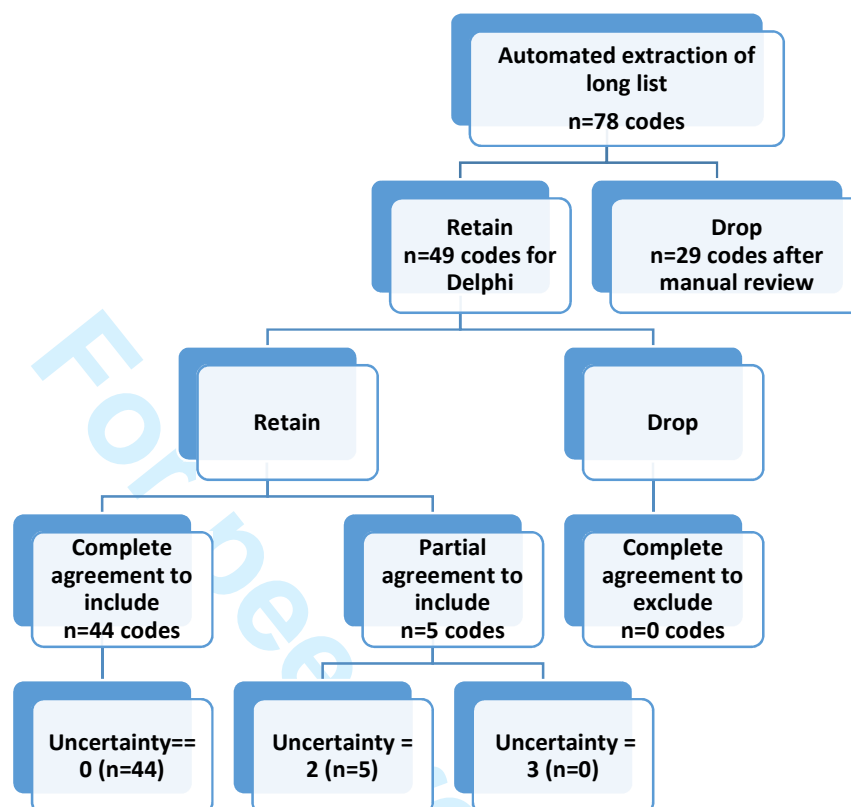
40
41 In Step 2 (Figure 1), Stata was used to produce a list of 78 possible shortness of breath codes
42
43 (for syntax see box 2 in Methods). Of the 78 potential codes, 29 were excluded because they
44
45 were clearly inappropriate (Table 1).
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Table 1 Reasons for exclusion after first round of assessment

Reason for exclusion	Number
Described ‘apnoea’ – absence of breathing – rather than breathlessness	15
Described negation of breathlessness	2
Described tachypnoea – abnormally rapid breathing – rather than breathlessness	3
Breathlessness related to pregnancy / neonate not pathology	2
Described hyperpnoea – increased rate and depth of breathing – not breathlessness	1
Description contained the string ‘sob’ but did not describe breathlessness (e.g. Removal of transobturator tape”)	6
Total	29

The remaining codes were included in Step 3, the Delphi review (Supplementary materials Table A 1). Following the Delphi process, 49 codes were included in the final library (for complete list see Supplementary materials Table A 2). There was complete agreement to include 44 of the 49 (90%) of the codes, and partial disagreement over inclusion of just 5 (10%) of codes (Figure 2). In this example, none of the codes were excluded during the Delphi process.

Figure 2 Flow chart illustrating the selection of codes



Using codelists to identify symptoms

Of 28,216 patients diagnosed with lung cancer in the study, 7,879 (28%) reported at least one episode of shortness of breath in the year before diagnosis. The total number of episodes of shortness of breath in the year before diagnosis was 13,091 (see Table 2).

Table 2 Frequency of use of shortness of breath codes in the year before diagnosis with lung cancer

medcode	Description	Frequency	Percent	Cumulative %	Certainty variable ^a
4822	Shortness of breath	3,226	24.64	24.64	0
741	[D]Shortness of breath	1,455	11.11	35.76	0
1429	Breathlessness	1,116	8.52	44.28	0

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						
	19427	MRC Breathlessness Scale: grade 2	1,106	8.45	52.73	0
	19426	MRC Breathlessness Scale: grade 3	1,010	7.72	60.45	0
	5349	Shortness of breath symptom	816	6.23	66.68	0
	5175	Breathlessness symptom	785	6.00	72.68	0
	19430	MRC Breathlessness Scale: grade 4	764	5.84	78.51	0
	5896	Dyspnoea - symptom	437	3.34	81.85	0
	2575	Short of breath on exertion	415	3.17	85.02	0
	3092	[D]Dyspnoea	395	3.02	88.04	0
	19432	MRC Breathlessness Scale: grade 1	332	2.54	90.57	0
	6326	Breathless - moderate exertion	261	1.99	92.57	0
	19429	MRC Breathlessness Scale: grade 5	189	1.44	94.01	0
	2931	Difficulty breathing	187	1.43	95.44	0
	12474	SOBOE	166	1.27	96.71	0
	7932	Breathless - mild exertion	142	1.08	97.79	0
	735	[D]Breathlessness	66	0.50	98.30	0
	7000	O/E - dyspnoea	49	0.37	98.67	0
	57903	CLASP shortness of breath	44	0.34	99.01	0

	score				
31143	Breathless - at rest	39	0.30	99.30	0
7683	Breathless - lying flat	22	0.17	99.47	0
6434	Paroxysmal nocturnal dyspnoea	19	0.15	99.62	0
21801	Breathlessness NOS	10	0.08	99.69	0
11451	[D]Orthopnoea	9	0.07	99.76	0
9089	Orthopnoea symptom	8	0.06	99.82	0
24889	Breathless - strenuous exertion	5	0.04	99.86	0
7534	O/E - respiratory distress	4	0.03	99.89	1
18116	Nocturnal dyspnoea	3	0.02	99.92	0
2563	Adult respiratory distress syndrome	2	0.02	99.93	1
2737	Dyspnoea on exertion	2	0.02	99.95	1
24848	Respiratory distress syndrome	2	0.02	99.96	1
53771	[D]Respiratory distress	2	0.02	99.98	0
22094	Borg Breathlessness Score: 10 maximal	1	0.01	99.98	0
59860	Borg Breathlessness Score: 4 somewhat..	1	0.01	99.99	0
101843	Short of breath dressing/undressing	1	0.01	100.00	0

9297	[D]Respiratory insufficiency	0	0	100.00	1
37704	O/E - orthopnoea	0	0	100.00	0
42287	Borg Breathlessness Score: 6 severe (+)	0	0	100.00	0
57193	Borg Breathlessness Score: 3 moderate	0	0	100.00	0
57678	Adult respiratory distress syndrome	0	0	100.00	0
57759	Borg Breathlessness Score: 2 slight	0	0	100.00	0
60096	CLASP shortness of breath score	0	0	100.00	0
64049	Borg Breathlessness Score: 5 severe	0	0	100.00	0
67566	Borg Breathlessness Score: 9 very, very sev (almost maximal)	0	0	100.00	0
68707	Borg Breathlessness Score: 1 very slight	0	0	100.00	0
70061	Borg Breathlessness Score: 7 very severe	0	0	100.00	0
70818	Borg Breathlessness Score: 0.5 very, very slight	0	0	100.00	0
72334	Borg Breathlessness Score:	0	0	100.00	0

	8 very severe (+)				
Total	Total	13,091	100.00	100.00	

^aThe 'certainty variable' is coded as: 0 = 'Minimal Uncertainty' (all panel members agreed the code should be included in the list); 1 = 'Moderate Uncertainty' (at least one panel member was uncertain that the code should be included); 3 = 'Maximal Uncertainty' (at least one panel member thought the code should be excluded)

Of the 49 codes in the list for shortness of breath, 13 were never used by GPs to record this symptom (Table 2). The majority of these were related to the BORG and CLASP breathlessness scores, and one was for respiratory insufficiency, highlighted as an uncertain code in the Delphi process.

Of the 37 codes used by GPs, 12 accounted for 90% of the total number of 13,091 episodes of shortness of breath recorded. Furthermore, just 4 codes accounted for over 50% of the records (Table 2).

Sensitivity analysis

In the sensitivity analysis, the codelist was restricted to the 44 codes whose inclusion was fully agreed in the Delphi process. This resulted in the loss of just 6 patients reporting at least one episode of shortness of breath in the year before diagnosis (i.e. the number fell from 7,879 (28%) to 7,873 (28%)). The total number of episodes of shortness of breath in the year before diagnosis was 13,081, compared with 13,091 using the complete codelist (see Supplementary materials Table A 3, for complete list).

DISCUSSION

We have presented a reproducible methodology for developing clinical codelists for use when conducting EHR research. It is intended to improve scientific rigour by standardising the conduct and reporting of this generally overlooked and underreported stage of EHR research. These methods can be adapted to suit the needs of different EHR research questions. To facilitate this, we have included example syntax for two of the most widely used statistical software packages.

Reporting guidelines for observational studies aim to promote the core principles of the scientific process: discovery, transparency, and replicability.¹² For systematic reviews, where searches for eligible papers are a core part of the methods, PRISMA guidelines stipulate that eligibility criteria, information sources used, search strategy and study selection process should be reported.¹³ The process of searching for EHR codes is analogous to this. The RECORD statement requires ‘*a complete list of codes and algorithms*’; yet what is meant by ‘algorithms’ is currently open to interpretation. We suggest that if EHR studies are to be transparent and reproducible these algorithms should include: definitions associated with each codelist; the syntax or search method used; the number of candidate codes identified; and the categorisation of codes after Delphi review (see **Figure 2**). This information could either be included within the published paper, as an appendix, or via online code repositories such as clinicalcode.org.

Benefits of this methodology include: the clear *a priori* definition of the clinical feature of interest based on reliable clinical resources; use of statistical software to comprehensively search all available codes; the iterative Delphi approach to reaching consensus on the most relevant codes; the generation of an auditable, replicable and modifiable syntax for codelist generation enabling sharing and replication.

1
2
3 Clinical coding is subjective and varies between clinicians. Where Delphi panel members
4
5 differ in decisions about code inclusion or exclusion, free text comments and discussions are
6
7 important to understand these differences. In some cases, refinement of the *a priori* definition
8
9 may be required to increase concordance between reviewers. However, residual differences
10
11 are likely to persist owing to the inherent variability in clinicians' patterns of coding. This
12
13 variability can be captured by the sensitivity analysis using the 'uncertainty' variable to
14
15 explore the impact of including or excluding these codes, and by using the frequency check to
16
17 identify which codes are used most often in the dataset.
18
19

20
21 Decisions about how to manage this uncertainty will depend upon the research question, and
22
23 whether the aim is to increase sensitivity or specificity. In the example of breathlessness we
24
25 aimed to include *any code which might be used by a clinician to record this symptom*; in
26
27 other words aiming to maximise sensitivity. Codes were therefore retained if either panel
28
29 member ranked them as '*definitely include*'; as this indicates that *some* clinicians may use
30
31 this code to record this symptom.
32
33

34
35 Another option to enhance sensitivity when developing disease-specific codelists which has
36
37 been described is the use of proxy codes. For example, one study included symptoms,
38
39 referrals, tests or treatments indicative of the disease of interest, such as prescription of
40
41 disease-modifying anti-rheumatic drugs as an indicator of rheumatoid arthritis. They found
42
43 that 83.5% of 5,843 patients had at least two indicator markers before a rheumatoid arthritis
44
45 code was recorded.¹⁴ This can be applied to symptoms; for example using prescriptions of
46
47 laxatives as a proxy for constipation in a study of colorectal cancer.¹⁵
48
49

50
51 For other research questions, it may be more important to focus on specificity, aiming to
52
53 reduce the number of false positive cases by using a narrower definition, with tighter
54
55 inclusion and exclusion criteria. For these studies it may be necessary to only include codes
56
57
58
59
60

1
2
3 for which consensus exists to “definitely include”, and closer consensus may be reached
4
5 amongst Delphi participants by increasing the number of Delphi rounds or the number of
6
7 panel members. Criteria for inclusion of codes following Delphi review therefore depends on
8
9 the purpose of the codelist. Researchers should make it clear whether codelists are sensitivity
10
11 or specificity driven as this will affect the generalisability of the codelist to other studies.
12
13

14
15
16
17 Ideally a panel of six clinicians would be used to best capture the variability in coding
18
19 between clinicians; however, this is unlikely to be an efficient use of clinicians’ time for
20
21 studies with large numbers of clinical codes, so a compromise of using two clinicians from a
22
23 panel of six per clinical feature offers a reasonable trade-off. This is analogous to the methods
24
25 for systematic reviews where two independent reviewers are routinely recommended. Using
26
27 fewer than six GPs on the overall Delphi panel reduces the clinical styles incorporated and
28
29 may not capture the inherent uncertainty in coding; it is therefore important that the extent of
30
31 the clinical input into the Delphi phase of the codelist review is clearly reported.
32
33
34
35
36
37
38

39 **Comparison to existing literature**

40 Previous studies have explored the implications of using differing code lists in EHR research.
41
42 For acute stroke significant differences were found between ONS codelists and a ‘restricted’
43
44 codelist developed by a Delphi panel; with very different mortality rates and different trends
45
46 over time between these codelists.¹⁶ Another study into coding of coronary heart disease in
47
48 primary care found that limited code sets for ‘angina’ or ‘myocardial infarction’
49
50 unsurprisingly had limited sensitivity; with substantial proportions of coronary heart disease
51
52 coded by non-specific codes.¹⁷ Both these papers called for increased transparency and
53
54 increased reporting of sensitivity analysis in EHR studies. Methods for compiling medical
55
56 and drug code lists were presented by Dave and Petersen in 2009.⁹ Their process was
57
58
59
60

1
2
3 analogous to the second step described in our proposed methodology; however it omitted the
4
5 stage of defining clearly *a priori* the clinical feature of interest and the final stage of Delphi
6
7 review, which is necessary to allow uncertainty to be explored using sensitivity analysis.
8
9

10 **Future implications**

11 By April 2018 all primary care systems should have completed migration to an international
12
13 clinical terminology called SNOMED CT ([https://digital.nhs.uk/SNOMED-CT-
14
15 implementation-in-primary-care](https://digital.nhs.uk/SNOMED-CT-implementation-in-primary-care)), which does not share the same hierarchical structure as
16
17 Read codes. This means that methods of codelist generation based on Readcodes can no
18
19 longer be relied upon.⁹ Mapping SNOMED CT onto current coding systems is underway by
20
21 the major EHR providers, but will inevitably lead to a period of flux. By working
22
23 independently of these hierarchical structures, using the description of the individual codes,
24
25 we overcome these problems, allowing researchers to develop a search strategy which works
26
27 across two or more classifications.
28
29
30
31
32
33
34

35 **Conclusions**

36 We suggest that as well as publishing codelists used in EHR studies, the methods used to
37
38 generate these codelists should be reported. Collated codelists should be badged by quality
39
40 according to whether they follow recommended methods for development.
41
42

43 As EHR research increases, it is important to avoid waste in research through incomplete or
44
45 unusable research publications.¹² Although initially time-consuming, using a rigorous and
46
47 reproducible method for codelist generation ‘future-proofs’ the findings, and an auditable,
48
49 modifiable syntax for codelist generation enables sharing and replication of EHR studies.
50
51
52
53
54
55
56
57
58
59
60

Acknowledgements:

We would like to thank Daniel Dedman, Senior Research CPRD for his comments on an earlier draft of this paper and Ben Feakins, statistician at the Nuffield Department of Primary Care Health Sciences at the University of Oxford for writing the R syntax.

Competing interests:

None declared

Funding:

JW and BDN are both funded by Doctoral Research Fellowships from the National Institute for Health Research. WH is part-funded by the National Institute for Health Research (NIHR) Collaboration for Leadership in Applied Health Research and Care South West Peninsula at the Royal Devon and Exeter NHS Foundation Trust. The views expressed are those of the authors and not necessarily those of the NHS, the NIHR or the Department of Health.

Contributorship Statement

WH conceived, and SP enhanced the methods of codelist collation described in the paper. JW wrote the original outline of the paper. SP designed and performed the data analysis. JW, SP and BN developed the first draft of the paper. All authors contributed to subsequent drafts and read and approved the final manuscript.

References:

1. Benson T. Why general practitioners use computers and hospital doctors do not—Part 1: incentives. *BMJ* 2002;325(7372):1086-89. doi: 10.1136/bmj.325.7372.1086
2. Herrett E, Gallagher AM, Bhaskaran K, et al. Data Resource Profile: Clinical Practice Research Datalink (CPRD). *Int J Epidemiol* 2015;44(3):827-36. doi: 10.1093/ije/dyv098 [published Online First: 2015/06/08]
3. Lewis JD, Schinnar R, Bilker WB, et al. Validation studies of the health improvement network (THIN) database for pharmacoepidemiology research. *Pharmacoepidemiology and drug safety* 2007;16(4):393-401. doi: 10.1002/pds.1335 [published Online First: 2006/10/27]
4. Hippisley-Cox J, Stables D, Pringle M. QRESEARCH: a new general practice database for research. *Informatics in primary care* 2004;12(1):49-50. [published Online First: 2004/05/14]

- 1
- 2
- 3 5. Nicholls SG, Quach P, von Elm E, et al. The REporting of Studies Conducted Using Observational
- 4 Routinely-Collected Health Data (RECORD) Statement: Methods for Arriving at Consensus
- 5 and Developing Reporting Guidelines. *PloS one* 2015;10(5):e0125620. doi:
- 6 10.1371/journal.pone.0125620 [published Online First: 2015/05/13]
- 7
- 8 6. Springate DA, Kontopantelis E, Ashcroft DM, et al. ClinicalCodes: An Online Clinical Codes
- 9 Repository to Improve the Validity and Reproducibility of Research Using Electronic Medical
- 10 Records. *PloS one* 2014;9(6):e99825. doi: 10.1371/journal.pone.0099825
- 11
- 12 7. Payne R. CPRD codes: Severe Mental Illness, 2017.
- 13 <https://data.bris.ac.uk/data/dataset/7ymb92btsycl11ip5t9pi9zvr>
- 14
- 15 8. Cambridge Uo. CPRD @ Cambridge - codelists 2017 [Available from:
- 16 http://www.phpc.cam.ac.uk/pcu/cprd_cam/codelists/.
- 17
- 18 9. Dave S, Petersen I. Creating medical and drug code lists to identify cases in primary care
- 19 databases. *Pharmacoepidemiology and drug safety* 2009;18(8):704-7. doi: 10.1002/pds.1770
- 20 [published Online First: 2009/05/21]
- 21
- 22 10. Committee WIC. International Classification of Primary Care ICPC-2-R. Revised Se ed. Oxford:
- 23 Oxford University Press 2005.
- 24
- 25 11. Murphy MK, Black NA, Lamping DL, et al. Consensus development methods, and their use in
- 26 clinical guideline development. 1998;2(3)
- 27
- 28 12. Glasziou P, Altman DG, Bossuyt P, et al. Reducing waste from incomplete or unusable reports of
- 29 biomedical research. *Lancet* 2014;383(9913):267-76. doi: 10.1016/s0140-6736(13)62228-x
- 30 [published Online First: 2014/01/15]
- 31
- 32 13. Moher D, Liberati A, Tetzlaff J, et al. Preferred reporting items for systematic reviews and meta-
- 33 analyses: the PRISMA statement. *Bmj* 2009;339:b2535. doi: 10.1136/bmj.b2535 [published
- 34 Online First: 2009/07/23]
- 35
- 36 14. Nicholson A, Ford E, Davies KA, et al. Optimising use of electronic health records to describe the
- 37 presentation of rheumatoid arthritis in primary care: a strategy for developing code lists.
- 38 *PloS one* 2013;8(2):e54878. doi: 10.1371/journal.pone.0054878 [published Online First:
- 39 2013/03/02]
- 40
- 41 15. Hamilton W, Lancashire R, Sharp D, et al. The risk of colorectal cancer with symptoms at different
- 42 ages and between the sexes: a case-control study. *BMC medicine* 2009;7(1):17. doi:
- 43 10.1186/1741-7015-7-17
- 44
- 45 16. Gulliford MC, Charlton J, Ashworth M, et al. Selection of medical diagnostic codes for analysis of
- 46 electronic patient records. Application to stroke in a primary care database. *PloS one*
- 47 2009;4(9):e7168. doi: 10.1371/journal.pone.0007168 [published Online First: 2009/09/25]
- 48
- 49 17. Bhattarai N, Charlton J, Rudisill C, et al. Coding, recording and incidence of different forms of
- 50 coronary heart disease in primary care. *PloS one* 2012;7(1):e29776.
- 51
- 52
- 53
- 54
- 55
- 56
- 57
- 58
- 59
- 60

Online only supplementary material

Table A 1 List of potential codes for inclusion in a list to identify patients with shortness of breath. 'sob' is the binary variable denoting shortness of breath (1: include; 2: exclude)

medcode	desc	sob	Reason for exclusion
735	[D]Breathlessness	1	
741	[D]Shortness of breath	1	
820	[D]Tachypnoea	0	Abnormally rapid breathing, not breathlessness
982	[D]Apnoea	0	Absence of breathing, not breathlessness
1429	Breathlessness	1	
2506	[D]Sleep apnoea syndrome	0	Absence of breathing, not breathlessness
2563	[D]Respiratory distress	1	
2575	Short of breath on exertion	1	
2737	Respiratory distress syndrome	1	
2931	Difficulty breathing	1	
3007	Newborn transitory tachypnoea	0	Abnormally rapid breathing, not breathlessness
3092	[D]Dyspnoea	1	
4822	Shortness of breath	1	
5175	Breathlessness symptom	1	
5349	Shortness of breath symptom	1	
5896	Dyspnoea - symptom	1	
6326	Breathless - moderate exertion	1	

1			
2			
3			
4	6434	Paroxysmal nocturnal dyspnoea	1
5			
6	7000	O/E - dyspnoea	1
7			
8	7534	O/E - respiratory distress	1
9			
10	7603	Sleep apnoea	0
11			Absence of breathing
12			
13	7683	Breathless - lying flat	1
14			
15	7932	Breathless - mild exertion	1
16			
17			
18	8148	Obstructive sleep apnoea	0
19			Absence of breathing
20			
21	9089	Orthopnoea symptom	1
22			
23	9297	[D]Respiratory insufficiency	1
24			
25	10114	O/E - tachypnoea	0
26			Abnormally rapid breathing, not
27			breathlessness
28			
29			
30	11451	[D]Orthopnoea	1
31			
32	12474	SOBOE	1
33			
34			
35	18116	Nocturnal dyspnoea	1
36			
37	19346	No breathlessness	0
38			Negation of breathlessness
39			
40	19426	MRC Breathlessness Scale: grade 3	1
41			
42	19427	MRC Breathlessness Scale: grade 2	1
43			
44	19429	MRC Breathlessness Scale: grade 5	1
45			
46			
47	19430	MRC Breathlessness Scale: grade 4	1
48			
49	19432	MRC Breathlessness Scale: grade 1	1
50			
51			
52	19576	Apnoea of newborn	0
53			Absence of breathing
54	20438	[D]Syndrome sleep apnoea	0
55			Absence of breathing
56			
57	20748	Obstructive sleep apnoea	0
58			Absence of breathing
59	21801	Breathlessness NOS	1
60			

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				
	22094	Short of breath dressing/undressing	1	
	23779	Sleep apnoea	0	Absence of breathing
	23924	Scoline apnoea	0	Absence of breathing
	24848	Adult respiratory distress syndrome	1	
	24889	Breathless - strenuous exertion	1	
	26684	Perinatal respiratory distress NOS	0	Related to pregnancy not pathology
	26871	Primary sleep apnoea of newborn	0	Absence of breathing
	31143	Breathless - at rest	1	
	31913	O/E - hyperpnoea	0	Increased depth and rate of breathing, not breathlessness
	36301	[D]Hypersomnia with sleep apnoea	0	Absence of breathing
	37667	Apnoea alarm monitoring	0	Absence of breathing
	37704	O/E - orthopnoea	1	
	42287	Borg Breathlessness Score: 6 severe (+)	1	
	48539	[D]Insomnia with sleep apnoea	0	Absence of breathing
	53771	Dyspnoea on exertion	1	
	54594	[M]Mesoblastic nephroma	0	Contains 'sob' string, but is inappropriate
	57193	Borg Breathlessness Score: 3 moderate	1	
	57283	Introduction of transobturator tape	0	Contains 'sob' string, but is inappropriate

1				
2				
3				
4	57678	Adult respiratory distress syndrome	1	
5				
6	57759	Borg Breathlessness Score: 2 slight	1	
7				
8	57903	CLASP shortness of breath score	1	
9				
10	58538	Fusobacterial necrotising tonsillitis	0	Contains 'sob' string, but is
11				inappropriate
12				
13				
14				
15	59860	Borg Breathlessness Score: 4	1	
16		somewhat severe		
17				
18				
19				
20	60096	CLASP shortness of breath score	1	
21				
22	64049	Borg Breathlessness Score: 5 severe	1	
23				
24				
25	65353	Borg Breathlessness Score: 0 none at	0	Negates breathlessness
26		all		
27				
28				
29				
30	67566	Borg Breathlessness Score: 9 very,	1	
31		very sev (almost maximal)		
32				
33				
34	68707	Borg Breathlessness Score: 1 very	1	
35		slight		
36				
37				
38				
39	70061	Borg Breathlessness Score: 7 very	1	
40		severe		
41				
42				
43				
44	70818	Borg Breathlessness Score: 0.5 very,	1	
45		very slight		
46				
47				
48				
49	72334	Borg Breathlessness Score: 8 very	1	
50		severe (+)		
51				
52				
53	72704	[X]Other respiratory distress of	0	Newborn
54		newborn		
55				
56				
57				
58				
59				
60				

1 2 3 4 5 6 7	73978	[X]Other apnoea of newborn	0	Absence of breathing, not breathlessness
8 9 10 11 12	93869	Removal of transobturator tape	0	Contains 'sob' string, but is inappropriate
13 14 15 16 17	97037	Introduction of transobturator sling	0	Contains 'sob' string, but is inappropriate
18 19 20 21 22	98965	Urine beta amino isobutyrate level	0	Contains 'sob' string, but is inappropriate
23 24	100177	Berlin questionnaire for sleep apnoea	0	Related to absence of breathing
25 26 27 28	101843	Borg Breathlessness Score: 10 maximal	1	
29 30 31 32	Total		78	

Table A 2 Delphi review of codes for shortness of breath. 1: definitely include; 2: uncertain; 3: definitely exclude

medcode	desc	Reviewer 1 Decision	Reviewer 1 comment	Reviewer 2 decision	Reviewer 2 comment
735	[D]Breathlessness	1		1	
741	[D]Shortness of breath	1		1	
1429	Breathlessness	1		1	
2563	[D]Respiratory distress	1		2	I would usually only use this term in children
2575	Short of breath on exertion	1		1	
2737	Respiratory distress syndrome	1		2	I would usually only use this term in children
2931	Difficulty breathing	1		1	
3092	[D]Dyspnoea	1		1	
4822	Shortness of breath	1		1	
5175	Breathlessness symptom	1		1	
5349	Shortness of breath symptom	1		1	
5896	Dyspnoea - symptom	1		1	
6326	Breathless - moderate exertion	1		1	
6434	Paroxysmal nocturnal dyspnoea	1		1	
7000	O/E - dyspnoea	1		1	
7534	O/E - respiratory distress	1		2	
7683	Breathless - lying flat	1		1	
7932	Breathless - mild exertion	1		1	
9089	Orthopnoea symptom	1		1	
9297	[D]Respiratory insufficiency	1		2	
11451	[D]Orthopnoea	1		1	
12474	SOBOE	1		1	

1	18116	Nocturnal dyspnoea	1		1	
2	19426	MRC Breathlessness Scale: grade 3	1		1	
3	19427	MRC Breathlessness Scale: grade 2	1		1	
4	19429	MRC Breathlessness Scale: grade 5	1		1	
5	19430	MRC Breathlessness Scale: grade 4	1		1	
6	19432	MRC Breathlessness Scale: grade 1	1		1	
7	21801	Breathlessness NOS	1		1	
8	22094	Short of breath dressing/undressing	1		1	
9	24848	Adult respiratory distress syndrome	1		2	
10	24889	Breathless - strenuous exertion	1		1	
11	31143	Breathless - at rest	1		1	
12	37704	O/E - orthopnoea	1		1	
13	42287	Borg Breathlessness Score: 6 severe (+)	1		1	I would never use this scale and not sure who would but if it is coded then it would be relevant
14	53771	Dyspnoea on exertion	1		1	
15	57193	Borg Breathlessness Score: 3 moderate	1		1	
16	57678	Adult respiratory distress syndrome	1		1	
17	57759	Borg Breathlessness Score: 2 slight	1		1	
18	57903	CLASP shortness of breath score	1		1	as above

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					

Table A 3 Sensitivity analysis: Medical codes used to record shortness of breath in a group of patients in the year before they were diagnosed with lung cancer

medcode	Description	Frequency	Percent	Cum.
4822	Shortness of breath	3,226	24.66	24.66
741	[D]Shortness of breath	1,455	11.12	35.78
1429	Breathlessness	1,116	8.53	44.32
19427	MRC Breathlessness Scale: grade 2	1,106	8.46	52.77
19426	MRC Breathlessness Scale: grade 3	1,010	7.72	60.49
5349	Shortness of breath symptom	816	6.24	66.73
5175	Breathlessness symptom	785	6.00	72.73
19430	MRC Breathlessness Scale: grade 4	764	5.84	78.57
5896	Dyspnoea - symptom	437	3.34	81.91
2575	Short of breath on exertion	415	3.17	85.09
3092	[D]Dyspnoea	395	3.02	88.10
19432	MRC Breathlessness Scale: grade 1	332	2.54	90.64
6326	Breathless - moderate exertion	261	2.00	92.64
19429	MRC Breathlessness Scale: grade 5	189	1.44	94.08
2931	Difficulty breathing	187	1.43	95.51
12474	SOBOE	166	1.27	96.78
7932	Breathless - mild exertion	142	1.09	97.87
735	[D]Breathlessness	66	0.50	98.37
7000	O/E - dyspnoea	49	0.37	98.75
57903	CLASP shortness of breath score	44	0.34	99.08

31143	Breathless - at rest	39	0.30	99.38
7683	Breathless - lying flat	22	0.17	99.55
6434	Paroxysmal nocturnal dyspnoea	19	0.15	99.69
21801	Breathlessness NOS	10	0.08	99.77
11451	[D]Orthopnoea	9	0.07	99.84
9089	Orthopnoea symptom	8	0.06	99.90
24889	Breathless - strenuous exertion	5	0.04	99.94
18116	Nocturnal dyspnoea	3	0.02	99.96
53771	Dyspnoea on exertion	2	0.02	99.98
22094	Borg Breathlessness Score: 10 maximal	1	0.01	99.98
59860	Borg Breathlessness Score: 4 somewhat..	1	0.01	99.99
101843	Short of breath dressing/undressing	1	0.01	100.00
Total	Total	13,081	100.00	

1
2
3
4 Example R syntax.
5
6

7 Credit: Ben Feakins (Benjamin.Feakins@phc.ox.ac.uk).
8

```
9 #=====#  
10 # #  
11 ##### R CODE FOR APPENDIX #####  
12 # #  
13 #=====#  
14  
15  
16  
17 ### Set Directory Objects ###  
18 browser.dir <- "/Volumes/PHC/CPRD_data/Browsers"  
19 save.dir <- "~/Desktop"  
20  
21  
22  
23 ### Read Data Into R ###  
24 setwd(browser.dir)  
25 medical <- read.delim("medical.txt", header = FALSE,  
26 na.strings = "", stringsAsFactors = FALSE, skip = 1)  
27 names(medical) <- c("medcode", "readcode",  
28 "clinicalevents", "immunisationevents",  
29 "referralevents", "testevents", "readterm",  
30 "databasebuild") # Define headers.  
31  
32  
33  
34  
35  
36 ### Define Search Terms ###  
37 search.terms.general <- c("shortness of  
38 breath|sob|pnoea|pnea|puffed|short of  
39 breath|winded|breathless")  
40 search.terms.breath <- c("breath")  
41 search.terms.breath.resid <-  
42 c("short|difficult|labour|labor|distress|insuff")  
43 search.terms.respir <- c("respir")  
44 search.terms.respir.resid <- c("insuff|distress")  
45  
46  
47  
48  
49  
50 ### Filtering ###  
51 general <- medical[grepl(search.terms.general,  
52 medical$readterm, ignore.case = TRUE), ]  
53 breath <- medical[grepl(search.terms.breath,  
54 medical$readterm, ignore.case = TRUE), ]  
55 breath <- breath[grepl(search.terms.breath.resid,  
56 breath$readterm, ignore.case = TRUE), ]  
57 respir <- medical[grepl(search.terms.respir,  
58 medical$readterm, ignore.case = TRUE), ]  
59  
60
```

```
1
2
3 respir <- respir[grepl(search.terms.respir.resid,
4 respir$readterm, ignore.case = TRUE), ]
5
6
7 ### Combine Results ###
8 medcodes <- rbind(general, breath, respir)
9
10
11 ### De-Duplicate ###
12 medcodes <- unique(medcodes)
13
14
15
16 ### Sort ###
17 medcodes <- medcodes[order(medcodes$medcode), ]
18
19
20 ### Remove Codes Not Related to Breathlessness ###
21 # Further subsetting using grepl().
22
23
24 ### Drop Useless Variables ###
25 medcodes <- medcodes[c("medcode", "readcode")]
26
27
28 ### Save a Copy ###
29 setwd(save.dir)
30 write.table(medcodes, "sob_library.txt", quote = FALSE,
31 sep = "\t", na = "", row.names = FALSE)
32
33
34
35 ### Tidying Up ###
36 rm(browser.dir, save.dir)
37 rm(search.terms.general, search.terms.breath,
38 search.terms.breath.resid, search.terms.respir,
39 search.terms.respir.resid)
40 rm(medical, general, breath, respir)
41 rm(medcodes)
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
```

BMJ Open

Identifying clinical features in primary care electronic health record studies: methods for codelist development

Journal:	<i>BMJ Open</i>
Manuscript ID	bmjopen-2017-019637.R1
Article Type:	Research
Date Submitted by the Author:	13-Oct-2017
Complete List of Authors:	Watson, Jessica; University of Bristol, School of Social and Community Medicine Nicholson, Brian; University of Oxford, Nuffield Dept Primary Care Health Sciences Hamilton, Willie; University of Exeter Medical School, Primary Care Diagnostics Price, Sarah; University of Exeter Medical School,
Primary Subject Heading:	Health services research
Secondary Subject Heading:	General practice / Family practice, Research methods
Keywords:	Electronic Health Records, Clinical coding, PRIMARY CARE, STATISTICS & RESEARCH METHODS, EPIDEMIOLOGY

SCHOLARONE™
Manuscripts

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

Identifying clinical features in primary care electronic health record studies: methods for codelist development

Jessica Watson¹, Brian D Nicholson², Willie Hamilton³, Sarah Price³

¹Centre for Academic Primary Care, Bristol Medical School, University of Bristol, Canynge Hall, 39 Whatley Road, Bristol, BS8 2PS

²Nuffield Department of Primary Care Health Sciences, Radcliffe Primary Care Building, University of Oxford, OX2 6GG.

³University of Exeter Medical School

Correspondence to: Jessica.Watson@bristol.ac.uk

Word Count: 3,464

ABSTRACT

Objective: Analysis of routinely collected Electronic Health Record (EHR) data from primary care is reliant upon the creation of codelists to define clinical features of interest. To improve scientific rigor, transparency and replicability we describe and demonstrate a standardised reproducible methodology for clinical codelist development.

Design: We describe a three stage process for developing clinical codelists. First, the clear definition *a priori* of the clinical feature of interest using reliable clinical resources. Second, development of a list of potential codes using statistical software to comprehensively search all available codes. Third, a modified Delphi process to reach consensus between primary care practitioners on the most relevant codes, including the generation of an ‘uncertainty’ variable to allow sensitivity analysis.

Setting: These methods are illustrated by developing a codelist for shortness of breath in a primary care EHR sample, including modifiable syntax for commonly used statistical software.

Participants The codelist was used to estimate the frequency of shortness of breath in a cohort of 28,216 patients aged over 18 years who received an incident diagnosis of lung cancer between 1 January 2000 and 30 November 2016 in the Clinical Practice Research Datalink (CPRD).

Results Of 78 candidate codes, 29 were excluded as inappropriate. Complete agreement was reached for 44 (90%) of the remaining codes, with partial disagreement over 5 (10%). 13,091 episodes of shortness of breath were identified in the cohort of 28,216 patients. Sensitivity

1
2
3 analysis demonstrates that codes with the greatest uncertainty tend to be rarely used in
4
5 clinical practice.
6

7 **Conclusions** Although initially time-consuming, using a rigorous and reproducible method
8
9 for codelist generation ‘future-proofs’ findings, and an auditable, modifiable syntax for
10
11 codelist generation enables sharing and replication of EHR studies. Published codelists
12
13 should be badged by quality and report the methods of codelist generation including:
14
15 definitions and justifications associated with each codelist; the syntax or search method; the
16
17 number of candidate codes identified; and the categorisation of codes after Delphi review.
18
19
20
21
22

23 **Keywords**

24
25
26 Electronic Health Records, Clinical Coding, Primary Health Care, Epidemiological Methods
27
28
29
30
31

32 **Strengths and Limitations of this study**

- 33 • This paper presents rigorous reproducible methods for codelist generation to increase
34
35 transparency and replicability in EHR studies.
36
- 37 • Clear *a priori* definition of the feature of interest ensures clinical relevance, and
38
39 enables future researchers to assess the applicability of existing codelists to future
40
41 research questions.
42
43
- 44 • Generation of auditable, replicable and modifiable syntax for codelists enables
45
46 replication and ‘future-proofs’ codelists.
47
- 48 • Using a Delphi approach to reach consensus on inclusion of codes allows sensitivity
49
50 analysis to explore the impact of uncertainty in coding.
51
52
53
54
55
56
57
58
59
60

- Using multiple clinicians in a Delphi panel reviewing codes may be unfeasible and inefficient for studies with large numbers of codes; a compromise of using two clinicians per feature from a panel of six offers a reasonable trade-off.

INTRODUCTION

Electronic Health Records (EHRs) have been used in routine primary care practice in the United Kingdom (UK) for at least 20 years.¹ EHRs are a rich resource for researchers, and are increasingly used in epidemiological and medical research resulting in over 1,500 publications since 2000, increasing from ~80 in 2005 to more than 450 in 2015/2016.

There are three well established UK primary care EHR databases: The Clinical Practice Research Datalink (CPRD) including 4.4 million currently registered patients, covering 6.9% of the UK population;² The Health Improvement Network (THIN) including 3.6 million currently registered patients giving ~5.7% coverage of the nation;³ and QResearch® including approximately 5 million currently registered patients in the UK.⁴ All three databases record coded anonymised information about patients: demographics, diagnoses, symptoms, prescriptions, immunisation history, referral information, and test results.

Linkages enable follow-up of patients beyond the primary care setting; for example, to data recorded by the Office for National Statistics (ONS), the National Cancer Registration Service (NCRS) and to Hospital Episode Statistics. Integrated primary and secondary care databases are also being developed. For example, ResearchOne includes data for over 5 million patients from General Practice, Child Health, Community Health, Out-of-Hours, Palliative Hospital, Accident and Emergency and Acute Hospital.

(<http://www.researchone.org/>).

A key stage in EHR research is identifying exposures and outcomes of interest. This apparently simple task is made more complicated by the fact that EHR clinical data is

1
2
3 generally stored as codes, often including qualitative information, such as ‘abdominal pain’,
4
5 ‘left iliac fossa pain’ and ‘intermittent abdominal pain’. These separate codes need to be
6
7 grouped into codelists or thesauri, with the groups containing all the codes pertaining to the
8
9 variable of interest. However, the methods used to develop codelists are not standardised, and
10
11 are often poorly reported. They are an increasingly recognised source of bias in EHR
12
13 research, owing to both inclusion of inappropriate codes and omission of important codes. To
14
15 address this, the RECORD Statement states that ‘*a complete list of codes and algorithms used*
16
17 *to classify exposures, outcomes, confounders, and effect modifiers should be provided*’⁵.
18
19

20
21 Clinicalcodes.org has been developed by the University of Manchester to encourage
22
23 researchers to publish clinical codelists used in EHR research⁶ and some other Universities
24
25 are developing their own open access, citable, repositories of codelists for example the
26
27 University of Bristol⁷ and University of Cambridge.⁸ The current clinicalcodes.org repository
28
29 contains 72916 clinical codes deposited within 432 codelists (<https://clinicalcodes.org>), in the
30
31 format of a list of papers and associated codes. This repository is a necessary step forward
32
33 towards addressing transparency; however, it does not tackle the potential for bias, as it is not
34
35 sufficient to address the issues of scientific rigour and reproducibility in codelist
36
37 development.
38
39

40
41 The problem is illustrated by brief examination of codelists recently deposited on the
42
43 repository. Without a clear definition of the clinical variable a codelist is designed to
44
45 encapsulate, it is not possible to critique or evaluate it for peer review, or to decide whether it
46
47 is generalisable to other studies. For example, codelists deposited for cancer
48
49 (<https://clinicalcodes.rss.mhs.man.ac.uk/medcodes/article/50/>) do not adhere to the
50
51 standardised International Classification of Diseases definition of cancer, i.e. ICD codes C00
52
53 to C97, as 193 (~9%) of the 2254 Read codes related to carcinoma *in situ* (ICD D00 to D09).
54
55
56
57
58
59
60

1
2
3 Furthermore, 100 (~4%) codes were obsolete, or they indicated the absence of cancer or they
4
5 were completely unrelated to cancer.
6

7
8 This demonstrates the need to establish standardised methods for codelist development.
9

10 Currently recommended methods, for example Davé and Petersen⁹ and CALIBER codelists
11 (<http://caliberanalysis.r-forge.r-project.org/>), need updating. This is not only because they
12 omit steps to standardise the definition of clinical terms, but also because they are based in
13 the Read code system, which is being superseded by SNOMED CT codes (Systematized
14 Nomenclature of Medicine -- Clinical Terms) in April 2018.
15
16
17
18
19

20
21 We have significant experience in EHR research, with ~40 published studies conducted in the
22 CPRD since 2012. We have developed and refined rigorous methods for developing clinical
23 codelists for use in CPRD studies independent of the Read code system. The aim of this
24 paper is to report a clear, standardised, reproducible methodology, and to increase scientific
25 rigour in conduct of EHR research. The method is illustrated using the CPRD, but applies
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
equally well to other large EHR databases.

37 METHODS

38 Our method for collating clinical codelists involves three stages, described in Figure 1.

41 <Figure 1 here>

43 **Step 1: Clearly define the clinical feature of interest (symptom, disease or illness) a priori**

44
45
46 The first step is to clearly define the clinical feature of interest and establish inclusion and
47
48
49
50
51
52
53
54
55
56
57
58
59
60
exclusion criteria. This requires clinical input, particularly from GPs who are best placed to
understand how clinical features are coded in a primary care setting. For rare conditions,
which GPs encounter infrequently, it may also be important to get clinical input from hospital
specialist doctors. Reliable sources of clinical information should be used; for example:

- International Classification of Primary Care, which defines symptoms and diagnoses, provides synonyms for them and, importantly, lists what should be excluded from the definition¹⁰
- The BMJ Best Practice guidelines (<http://bestpractice.bmj.com/best-practice/welcome.html>)
- NICE Clinical Knowledge Summaries (<http://cks.nice.org.uk/>)
- ICD10 (<http://apps.who.int/classifications/icd10/browse/2016/en>) – this is less useful for symptoms, as it focuses on diseases
- Medical Subject Headings (MeSH) (https://www.nlm.nih.gov/mesh/2016/mesh_browser/MBrowser.html)
- NHS Digital Technology Reference data Update Distribution (TRUD): <https://isd.digital.nhs.uk/trud3/user/guest/group/0/home> Downloadable technology reference files including READ Code Browsers with cross map files.

Other potential resources include patient support groups, online discussion forums, and already published codelists (e.g. <https://clinicalcodes.org>). Hierarchical classifications such as Read, SNOMED or ICD-10 may be useful for identifying additional search terms and synonyms.

For some symptoms, it is necessary to tailor the definition to the context of the disease under investigation. Abdominal pain is a good example, where pancreatic disease may cause pain in the epigastrium and left hypochondrium, whereas disorders in the sigmoid colon generate pain in the left iliac fossa.

Step 2 – assembling list of codes that may be used to record the clinical feature

The second stage consists of identifying all potential codes that might be used by GPs to record the clinical feature of interest defined in Step 1 and collating them into a list.

This is done in several steps; we use Stata for this, but other software is possible.

1
2
3 First, using the resources listed in Step 1, an exhaustive list of synonyms for the outcome of
4
5 interest is generated. Box 1 uses the example of shortness of breath.
6
7

8 \beginbox1
9

10 **Box 1: Shortness of breath**

11 *ICPC*

- 12 • ICPC code: R02 (exclude: wheezing R03; stridor R04; hyperventilation R98)

13 *BMJ Best Practice*

- 14 • Dyspnoea, also known as shortness of breath or breathlessness, is a subjective
15 sensation of breathing discomfort ([http://bestpractice.bmj.com/best-
16 practice/monograph/862.html](http://bestpractice.bmj.com/best-practice/monograph/862.html))

17 *NICE CKS*

- 18 • Breathlessness is the distressing sensation of a deficit between the body's demand for
19 breathing and the ability of the respiratory system to satisfy that demand.
20 (<http://cks.nice.org.uk/breathlessness#!backgroundsub>)
- 21 • Breathlessness can be classified by its speed of onset as:
22
 - 23 ○ Acute breathlessness — when it develops over minutes, hours, or days.
 - 24 ○ Chronic breathlessness — when it develops over weeks or months.

25 *ICD10*

- 26 • ICD10 code: R06 – dyspnoea, orthopnoea, shortness of breath

27 *MeSH*

- 28 • MeSH: Difficult or labored breathing. Breathlessness, dyspnea

29 *Patient forums*

- 30 • Puffed, winded

31 *GP colleagues*

- 32 • Consider including ‘respiratory insufficiency’?

33 \endbox1
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

1
2
3 Second, the lookup file of all medical codes provided by the CPRD (medical.txt)¹ is opened
4
5 using Stata. This contains the alphanumeric Read code originally used by the GP to enter the
6
7 clinical information, the CPRD's proprietary 'medcode' (which is simply a numeric
8
9 equivalent of the Read code), as well as a verbal description (variable 'desc') common to
10
11 both the medcode and Read code. A variable for the clinical outcome of interest (here, 'sob'
12
13 for shortness of breath) is created and set to zero (see Box 2). Then Stata searches the verbal
14
15 description of each code, and sets 'sob' to 1 if it contains any of the synonyms. Example
16
17 syntax to replicate this process in the statistical software package R is provided in the
18
19 supplementary materials, using the lookup file of all medical codes that comes with the
20
21 CPRD browsers. Note that, in this file, the verbal description is called 'readterm' rather than
22
23 'desc'.
24
25
26
27
28
29

30 \beginbox2

```
31  
32 insheet using "medical.txt", clear  
33  
34  
35 *generate a binary variable for shortness of breath (sob) and set  
36 its value to zero  
37 generate sob=0  
38  
39  
40  
41 /* search the verbal description of the Read code/medcode and change  
42 the value of variable sob from 0 to 1 if it contains words that  
43 suggest the code might be about the clinical feature of interest*/  
44 replace sob=1 if regexm(desc, "[Ss]hortness [Oo]f  
45 [Bb]reath|SHORTNESS OF BREATH")  
46 replace sob=1 if regexm(desc, "[Ss] [Oo] [Bb]|SOB")  
47 replace sob=1 if regexm(desc, "pnoea|PNOEA")  
48 replace sob=1 if regexm(desc, "pnea|PNEA")  
49 replace sob=1 if regexm(desc, "[Pp]uffed|PUFFED")  
50  
51  
52  
53  
54  
55  
56  
57
```

58 ¹ THIN and QResearch® provide equivalent files.
59
60

```
1
2
3   replace sob=1 if regexm(desc, "[Ss]hort [Oo]f [Bb]reath|SHORT OF
4   BREATH")
5
6   replace sob=1 if regexm(desc, "[Ss]hort|SHORT") & regexm(desc,
7   "[Bb]reath|BREATH")
8
9   replace sob=1 if regexm(desc, "[Ww]inded|WINDED")
10
11  replace sob=1 if regexm(desc, "[Dd]ifficult|DIFFICULT") &
12  regexm(desc, "[Bb]reath|BREATH")
13
14  replace sob=1 if regexm(desc, "[Ll]abour|LABOUR") & regexm(desc,
15  "[Bb]reath|BREATH")
16
17  replace sob=1 if regexm(desc, "[Ll]abor|LABOR") & regexm(desc,
18  "[Bb]reath|BREATH")
19
20  replace sob=1 if regexm(desc, "[Bb]reathless|BREATHLESS")
21
22  replace sob=1 if regexm(desc, "[Dd]istress|DISTRESS") & regexm(desc,
23  "[Bb]reath|BREATH")
24
25  replace sob=1 if regexm(desc, "[Dd]istress|DISTRESS") & regexm(desc,
26  "[Rr]espir|RESPIR")
27
28  replace sob=1 if regexm(desc, "[Ii]suff|INSUFF") & regexm(desc,
29  "[Bb]reath|BREATH")
30
31  replace sob=1 if regexm(desc, "[Ii]suff|INSUFF") & regexm(desc,
32  "[Rr]espir|RESPIR")
33
34
35
36  /*order the dataset so that values of variable sob==1 are all placed
37  together*/
38
39  gsort sob
40
41
42  /* Manual check for bogus codes - manually change sob==1 to sob==0
43  if the code is clearly inappropriate. */
44
45
46  edit medcode readcode desc sob
47
48
49
50  /*Retain only those codes that are specifically about sob*/
51  keep if sob==1
52
53
54  /*Retain the variables of interest*/
55  keep medcode readcode desc sob
56
57
58
59
60
```

```
1
2
3      sort medcode
4
5
6      /*Save the file as a library for sob for the Delphi process*/
7      save "sob_library.dta", replace
8
9
10     /*Export as an Excel file
11     export excel using "sob_library", replace
12
13
14
15 \endbox2
```

16
17
18 The manual check for bogus codes should err on the side of caution, only rejecting codes that
19
20 are clearly inappropriate according to predefined inclusion and exclusion criteria. Common
21
22 reasons for exclusion are that search terms can pick up bogus codes (e.g. transobturator tape
23
24 contains the letter sequence ‘sob’), or codes indicating a family history of a condition or
25
26 screening for a condition rather than presence of a condition.
27

28
29 The output from Step 2 is a list of potential codes that is then exported to Excel and reviewed
30
31 manually in a Delphi-type process (Step 3).
32
33

34 35 36 **Step 3: Delphi review of codes**

37
38 The codelist is reviewed by one practising GP, plus at least one other GP from a panel of six,
39
40 using a modified nominal group technique¹¹. Each GP independently categorises the list,
41
42 ranking each Read code/medcode using a 3-point scale as follows:
43
44

45
46 1 = Definitely Include - the code accurately defines the clinical feature of interest, and GPs
47
48 would definitely use it.

49
50 2 = Uncertain – it remains unclear whether the code accurately reflects the clinical feature of
51
52 interest, or whether GPs would use it.

53
54 3 = Definitely Exclude – the code does not define the clinical feature of interest, and GPs
55
56 definitely would not use it.
57
58
59
60

1
2
3 Panel members are encouraged to add comments explaining their reasons for exclusion or
4
5 uncertainty, in the knowledge that these comments will be shared with an independent panel
6
7 chair who will collate all of the results.
8

9
10 Codes are retained in the final list if they are ranked ‘1=Definitely include’ by at least one of
11
12 the GPs, as this indicates sufficient evidence that the code may be used to record that clinical
13
14 feature. Codes are dropped if they are ranked as “3 = Definitely exclude”, or as “2 =
15
16 Uncertain” by *all* reviewers.
17

18
19
20 An ‘Uncertainty’ variable is also generated for retained codes, to enable sensitivity analyses
21
22 that remove codes for which any uncertainty exists about accuracy or use. The ‘uncertainty’
23
24 variable is defined as follows:
25

26
27 0 = ‘Minimal Uncertainty’, as all panel members ranked the code as ‘1=Definitely include’
28

29
30 1 = ‘Moderate Uncertainty’, at least one panel member ranked the code as ‘2=Uncertain’
31

32
33 2 = ‘Maximal Uncertainty’, at least one panel member ranked the code as ‘3=Definitely
34
35 exclude’
36

37
38 Once the codelist has been generated, a frequency check may be performed using the study’s
39
40 dataset to identify the frequency of the clinical events attributed to each clinical code. If the
41
42 Delphi process has been accurate, the most frequent events will most likely be coded as “0 =
43
44 Minimal Uncertainty”, whereas there will be fewer events for the codes ranked as “1 =
45
46 Moderate Uncertainty” or as “2 = Maximal Uncertainty”.
47
48
49
50

51 **Illustrative example using CPRD medical codes list**

52
53
54 The library of codes for shortness of breath was used to estimate the frequency of this
55
56 symptom in the year before diagnosis of lung cancer. Participants were CPRD patients aged
57
58
59
60

1
2
3 over 18 years who received an incident diagnosis of lung cancer between 1 January 2000 and
4
5 30 November 2016.

6
7 Outcome measures included the number of patients reporting shortness of breath in the year
8
9 before they were diagnosed with lung cancer, the proportion of all lung cancer patients
10
11 reporting shortness of breath, and the total number of episodes of shortness of breath.
12
13

14
15
16 In addition, a sensitivity analysis was carried out restricting the analysis to codes whose
17
18 uncertainty variable was coded 0 (=‘Minimal Uncertainty’), i.e. there was full agreement in
19
20 the Delphi process that the code should be included.
21
22

23 24 25 **RESULTS**

26
27 The codelist generated for shortness of breath is presented here to illustrate the method we
28
29 have described. The clinical resources reviewed in Step 1 (see Box 1 in Methods) indicated
30
31 that the codes used to define shortness of breath should capture evidence of ‘dyspn[o]ea’,
32
33 ‘shortness of breath’ (and its abbreviated term ‘sob’), ‘breathlessness’, ‘orthopn[o]ea’,
34
35 “‘difficult’ & ‘breathing’”, “‘labo[u]red’ & ‘breathing’”, “‘breathing’ & ‘discomfort’”,
36
37 ‘puffed’, ‘winded’, ‘respiratory distress’ and ‘respiratory insufficiency’.
38
39

40
41 In Step 2 (Figure 1), Stata was used to produce a list of 78 possible shortness of breath codes
42
43 (for syntax see box 2 in Methods). Of the 78 potential codes, 29 were excluded because they
44
45 were clearly inappropriate (Table 1).
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Table 1 Reasons for exclusion after first round of assessment

Reason for exclusion	Number
Described ‘apnoea’ – absence of breathing – rather than breathlessness	15
Described negation of breathlessness	2
Described tachypnoea – abnormally rapid breathing – rather than breathlessness	3
Breathlessness related to pregnancy / neonate not pathology	2
Described hyperpnoea – increased rate and depth of breathing – not breathlessness	1
Description contained the string ‘sob’ but did not describe breathlessness (e.g. Removal of transobturator tape”)	6
Total	29

The remaining codes were included in Step 3, the Delphi review (Supplementary materials Table A 1). Following the Delphi process, 49 codes were included in the final library (for complete list see Supplementary materials Table A 2). There was complete agreement to include 44 of the 49 (90%) of the codes, and partial disagreement over inclusion of just 5 (10%) of codes (Figure 2). In this example, none of the codes were excluded during the Delphi process.

<Figure 2 here>

Using codelists to identify symptoms

Of 28,216 patients diagnosed with lung cancer in the study, 7,879 (28%) reported at least one episode of shortness of breath in the year before diagnosis. The total number of episodes of shortness of breath in the year before diagnosis was 13,091 (see Table 2).

Table 2 Frequency of use of shortness of breath codes in the year before diagnosis with lung cancer

medcode	Description	Frequency	Percent	Cumulative %	Certainty variable ^a
4822	Shortness of breath	3,226	24.64	24.64	0
741	[D]Shortness of breath	1,455	11.11	35.76	0
1429	Breathlessness	1,116	8.52	44.28	0
19427	MRC Breathlessness Scale: grade 2	1,106	8.45	52.73	0
19426	MRC Breathlessness Scale: grade 3	1,010	7.72	60.45	0
5349	Shortness of breath symptom	816	6.23	66.68	0
5175	Breathlessness symptom	785	6.00	72.68	0
19430	MRC Breathlessness Scale: grade 4	764	5.84	78.51	0
5896	Dyspnoea - symptom	437	3.34	81.85	0
2575	Short of breath on exertion	415	3.17	85.02	0
3092	[D]Dyspnoea	395	3.02	88.04	0
19432	MRC Breathlessness Scale: grade 1	332	2.54	90.57	0
6326	Breathless - moderate exertion	261	1.99	92.57	0
19429	MRC Breathlessness Scale: grade 5	189	1.44	94.01	0
2931	Difficulty breathing	187	1.43	95.44	0

12474	SOBOE	166	1.27	96.71	0
7932	Breathless - mild exertion	142	1.08	97.79	0
735	[D]Breathlessness	66	0.50	98.30	0
7000	O/E - dyspnoea	49	0.37	98.67	0
57903	CLASP shortness of breath score	44	0.34	99.01	0
31143	Breathless - at rest	39	0.30	99.30	0
7683	Breathless - lying flat	22	0.17	99.47	0
6434	Paroxysmal nocturnal dyspnoea	19	0.15	99.62	0
21801	Breathlessness NOS	10	0.08	99.69	0
11451	[D]Orthopnoea	9	0.07	99.76	0
9089	Orthopnoea symptom	8	0.06	99.82	0
24889	Breathless - strenuous exertion	5	0.04	99.86	0
7534	O/E - respiratory distress	4	0.03	99.89	1
18116	Nocturnal dyspnoea	3	0.02	99.92	0
2563	Adult respiratory distress syndrome	2	0.02	99.93	1
2737	Dyspnoea on exertion	2	0.02	99.95	1
24848	Respiratory distress syndrome	2	0.02	99.96	1
53771	[D]Respiratory distress	2	0.02	99.98	0
22094	Borg Breathlessness Score:	1	0.01	99.98	0

	10 maximal				
59860	Borg Breathlessness Score: 4 somewhat..	1	0.01	99.99	0
101843	Short of breath dressing/undressing	1	0.01	100.00	0
9297	[D]Respiratory insufficiency	0	0	100.00	1
37704	O/E - orthopnoea	0	0	100.00	0
42287	Borg Breathlessness Score: 6 severe (+)	0	0	100.00	0
57193	Borg Breathlessness Score: 3 moderate	0	0	100.00	0
57678	Adult respiratory distress syndrome	0	0	100.00	0
57759	Borg Breathlessness Score: 2 slight	0	0	100.00	0
60096	CLASP shortness of breath score	0	0	100.00	0
64049	Borg Breathlessness Score: 5 severe	0	0	100.00	0
67566	Borg Breathlessness Score: 9 very, very sev (almost maximal)	0	0	100.00	0
68707	Borg Breathlessness Score: 1 very slight	0	0	100.00	0

70061	Borg Breathlessness Score: 7 very severe	0	0	100.00	0
70818	Borg Breathlessness Score: 0.5 very, very slight	0	0	100.00	0
72334	Borg Breathlessness Score: 8 very severe (+)	0	0	100.00	0
Total	Total	13,091	100.00	100.00	

^aThe 'certainty variable' is coded as: 0 = 'Minimal Uncertainty' (all panel members agreed the code should be included in the list); 1 = 'Moderate Uncertainty' (at least one panel member was uncertain that the code should be included); 3 = 'Maximal Uncertainty' (at least one panel member thought the code should be excluded)

Of the 49 codes in the list for shortness of breath, 13 were never used by GPs to record this symptom (Table 2). The majority of these were related to the BORG and CLASP breathlessness scores, and one was for respiratory insufficiency, highlighted as an uncertain code in the Delphi process.

Of the 37 codes used by GPs, 12 accounted for 90% of the total number of 13,091 episodes of shortness of breath recorded. Furthermore, just 4 codes accounted for over 50% of the records (Table 2).

Sensitivity analysis

In the sensitivity analysis, the codelist was restricted to the 44 codes whose inclusion was fully agreed in the Delphi process. This resulted in the loss of just 6 patients reporting at least one episode of shortness of breath in the year before diagnosis (i.e. the number fell from

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

7,879 (28%) to 7,873 (28%)). The total number of episodes of shortness of breath in the year before diagnosis was 13,081, compared with 13,091 using the complete codelist (see Supplementary materials Table A 3, for complete list).

For peer review only

DISCUSSION

We have presented a reproducible methodology for developing clinical codelists for use when conducting EHR research. It is intended to improve scientific rigour by standardising the conduct and reporting of this generally overlooked and underreported stage of EHR research. These methods can be adapted to suit the needs of different EHR research questions. To facilitate this, we have included example syntax for two of the most widely used statistical software packages.

Reporting guidelines for observational studies aim to promote the core principles of the scientific process: discovery, transparency, and replicability.¹² For systematic reviews, where searches for eligible papers are a core part of the methods, PRISMA guidelines stipulate that eligibility criteria, information sources used, search strategy and study selection process should be reported.¹³ The process of searching for EHR codes is analogous to this. The RECORD statement requires ‘*a complete list of codes and algorithms*’; yet what is meant by ‘algorithms’ is currently open to interpretation. We suggest that if EHR studies are to be transparent and reproducible these algorithms should include: definitions associated with each codelist; the syntax or search method used; the number of candidate codes identified; and the categorisation of codes after Delphi review (see **Figure 2**). This information could either be included within the published paper, as an appendix, or via online code repositories such as clinicalcode.org.

Benefits of this methodology include: the clear *a priori* definition of the clinical feature of interest based on reliable clinical resources; use of statistical software to comprehensively search all available codes; the iterative Delphi approach to reaching consensus on the most relevant codes; the generation of an auditable, replicable and modifiable syntax for codelist generation enabling sharing and replication.

1
2
3 The way in which diagnosis is recorded in the EHR is heterogeneous, with different
4 clinicians using different codes for the same clinical features. Definitions of clinical
5 conditions also change over time, and codes are updated regularly in EHRs, often duplicating
6 pre-existing codes. As a result, decisions about inclusion or exclusion of codes will vary
7 between clinicians. Where Delphi panel members differ in decisions, free text comments and
8 discussions are important to understand these differences. In some cases, refinement of the *a*
9 *priori* definition may be required to increase concordance between reviewers. However,
10 residual differences are likely to persist owing to the inherent variability in clinicians'
11 idiosyncratic patterns of coding. This variability can be captured by the sensitivity analysis
12 using the 'uncertainty' variable to explore the impact of including or excluding these codes,
13 and by using the frequency check to identify which codes are used most often in the dataset.

14
15
16
17
18
19
20
21
22
23
24
25
26
27
28 Decisions about how to manage this uncertainty will depend upon the research question, and
29 whether the aim is to increase sensitivity or specificity. In the example of breathlessness we
30 aimed to include *any code which might be used by a clinician to record this symptom*; in
31 other words aiming to maximise sensitivity. Codes were therefore retained if either panel
32 member ranked them as '*definitely include*'; as this indicates that *some* clinicians may use
33 this code to record this symptom.

34
35
36
37
38
39
40
41
42 Another option to enhance sensitivity when developing disease-specific codelists which has
43 been described is the use of proxy codes. For example, one study included symptoms,
44 referrals, tests or treatments indicative of the disease of interest, such as prescription of
45 disease-modifying anti-rheumatic drugs as an indicator of rheumatoid arthritis. They found
46 that 83.5% of 5,843 patients had at least two indicator markers before a rheumatoid arthritis
47 code was recorded.¹⁴ This can be applied to symptoms; for example using prescriptions of
48 laxatives as a proxy for constipation in a study of colorectal cancer.¹⁵

1
2
3 For other research questions, it may be more important to focus on specificity, aiming to
4 reduce the number of false positive cases by using a narrower definition, with tighter
5 inclusion and exclusion criteria. For these studies it may be necessary to only include codes
6 for which consensus exists to “definitely include”, and closer consensus may be reached
7 amongst Delphi participants by increasing the number of Delphi rounds or the number of
8 panel members. Criteria for inclusion of codes following Delphi review therefore depends on
9 the purpose of the codelist. Researchers should make it clear whether codelists are sensitivity
10 or specificity driven as this will affect the generalisability of the codelist to other studies.
11
12

13
14
15
16
17
18
19
20
21 Murphy et al suggested that a panel of at least six clinicians should be used for consensus
22 methods¹¹. This would be ideal to best capture the variability in coding between clinicians;
23 however, this is unlikely to be an efficient use of clinicians’ time for studies with large
24 numbers of clinical codes, so a compromise of using two clinicians from a panel of six per
25 clinical feature offers a reasonable trade-off. This is analogous to the methods for systematic
26 reviews where two independent reviewers are routinely recommended. Using fewer than six
27 GPs on the overall Delphi panel reduces the clinical styles incorporated and may not capture
28 the inherent uncertainty in coding; it is therefore important that the extent of the clinical input
29 into the Delphi phase of the codelist review is clearly reported.
30
31
32
33
34
35
36
37
38
39
40
41

42 These challenges are demonstrated in the example provided; although both Delphi reviewers
43 were “certain” that MRC breathlessness 1 was a code indicating shortness of breath, further
44 iterative feedback suggested that this actually indicates conditional breathlessness, being
45 defined as “*not troubled by breathlessness except on strenuous exercise*”. This emphasises
46 the importance of a transparent process of codelist development; and illustrates the fact that
47 this can be an iterative process, as Delphi reviewers, or later critics may raise issues which
48 require researchers to revisit and refine the definition or inclusion criteria to improve the
49 sensitivity and specificity of the codelist.
50
51
52
53
54
55
56
57
58
59
60

Comparison to existing literature

Previous studies have explored the implications of using differing code lists in EHR research. For acute stroke significant differences were found between ONS codelists and a 'restricted' codelist developed by a Delphi panel; with very different mortality rates and different trends over time between these codelists.¹⁶ Another study into coding of coronary heart disease in primary care found that limited code sets for 'angina' or 'myocardial infarction' unsurprisingly had limited sensitivity; with substantial proportions of coronary heart disease coded by non-specific codes.¹⁷ Both these papers called for increased transparency and increased reporting of sensitivity analysis in EHR studies. Methods for compiling medical and drug code lists were presented by Dave and Petersen in 2009.⁹ Their process was analogous to the second step described in our proposed methodology; however it omitted the stage of defining clearly *a priori* the clinical feature of interest and the final stage of Delphi review, which is necessary to allow uncertainty to be explored using sensitivity analysis.

Future implications

By April 2018 all primary care systems should have completed migration to an international clinical terminology called SNOMED CT (<https://digital.nhs.uk/SNOMED-CT-implementation-in-primary-care>), which does not share the same hierarchical structure as Read codes. This means that methods of codelist generation based on Readcodes can no longer be relied upon.⁹ Mapping SNOMED CT onto current coding systems is underway by the major EHR providers, but will inevitably lead to a period of flux. By working independently of these hierarchical structures, using the description of the individual codes, we overcome these problems, allowing researchers to develop a search strategy which works across two or more classifications. Our proposed Delphi approach to code selection aims to reduce the impact of variable coding practice between clinicians. Clinicians are rarely trained

1
2
3 in coding practice outside their individual clinical setting. An area of future development
4
5 could therefore be for standardised coding training to be delivered as part of continued
6
7 professional development.
8
9

10 **Conclusions**

11 We suggest that as well as publishing codelists used in EHR studies, the methods used to
12
13 generate these codelists should be reported. Collated codelists should be badged by quality
14
15 according to whether they follow recommended methods for development.
16
17

18
19 As EHR research increases, it is important to avoid waste in research through incomplete or
20
21 unusable research publications.¹² Although initially time-consuming, using a rigorous and
22
23 reproducible method for codelist generation ‘future-proofs’ the findings, and an auditable,
24
25 modifiable syntax for codelist generation enables sharing and replication of EHR studies.
26
27

28 **Acknowledgements:**

29
30 We would like to thank Daniel Dedman, Senior Research CPRD for his comments on an
31
32 earlier draft of this paper and Benjamin Feakins, statistician at the Nuffield Department of
33
34 Primary Care Health Sciences at the University of Oxford for writing the R syntax.
35
36
37
38
39
40
41

42 **Competing interests:**

43
44 None declared
45
46
47

48 **Funding:**

49
50 JW (DRF-2016-09-034) and BDN (DRF-2015-08-18) are both funded by Doctoral Research
51
52 Fellowships from the National Institute for Health Research Trainees Coordinating Centre.
53
54 WH is part-funded by the National Institute for Health Research (NIHR) Collaboration for
55
56 Leadership in Applied Health Research and Care South West Peninsula at the Royal Devon
57
58
59
60

and Exeter NHS Foundation Trust. The views expressed are those of the authors and not necessarily those of the NHS, the NIHR or the Department of Health.

Contributorship Statement

WH conceived, and SP enhanced the methods of codelist collation described in the paper. JW wrote the original outline of the paper. SP designed and performed the data analysis. JW, SP and BDN developed the first draft of the paper. All authors contributed to subsequent drafts and read and approved the final manuscript.

References:

1. Benson T. Why general practitioners use computers and hospital doctors do not—Part 1: incentives. *BMJ* 2002;325(7372):1086-89. doi: 10.1136/bmj.325.7372.1086
2. Herrett E, Gallagher AM, Bhaskaran K, et al. Data Resource Profile: Clinical Practice Research Datalink (CPRD). *Int J Epidemiol* 2015;44(3):827-36. doi: 10.1093/ije/dyv098 [published Online First: 2015/06/08]
3. Lewis JD, Schinnar R, Bilker WB, et al. Validation studies of the health improvement network (THIN) database for pharmacoepidemiology research. *Pharmacoepidemiology and drug safety* 2007;16(4):393-401. doi: 10.1002/pds.1335 [published Online First: 2006/10/27]
4. Hippisley-Cox J, Stables D, Pringle M. QRESEARCH: a new general practice database for research. *Informatics in primary care* 2004;12(1):49-50. [published Online First: 2004/05/14]
5. Nicholls SG, Quach P, von Elm E, et al. The REporting of Studies Conducted Using Observational Routinely-Collected Health Data (RECORD) Statement: Methods for Arriving at Consensus and Developing Reporting Guidelines. *PLoS one* 2015;10(5):e0125620. doi: 10.1371/journal.pone.0125620 [published Online First: 2015/05/13]
6. Springate DA, Kontopantelis E, Ashcroft DM, et al. ClinicalCodes: An Online Clinical Codes Repository to Improve the Validity and Reproducibility of Research Using Electronic Medical Records. *PLoS one* 2014;9(6):e99825. doi: 10.1371/journal.pone.0099825
7. Payne R. CPRD codes: Severe Mental Illness, 2017. <https://data.bris.ac.uk/data/dataset/7ymb92btsycl11ip5t9pi9zvr>
8. Cambridge Uo. CPRD @ Cambridge - codelists 2017 [Available from: http://www.phpc.cam.ac.uk/pcu/cprd_cam/codelists/].
9. Dave S, Petersen I. Creating medical and drug code lists to identify cases in primary care databases. *Pharmacoepidemiology and drug safety* 2009;18(8):704-7. doi: 10.1002/pds.1770 [published Online First: 2009/05/21]
10. Committee WIC. International Classification of Primary Care ICPC-2-R. Revised Se ed. Oxford: Oxford University Press 2005.
11. Murphy MK, Black NA, Lamping DL, et al. Consensus development methods, and their use in clinical guideline development. 1998;2(3)
12. Glasziou P, Altman DG, Bossuyt P, et al. Reducing waste from incomplete or unusable reports of biomedical research. *Lancet* 2014;383(9913):267-76. doi: 10.1016/s0140-6736(13)62228-x [published Online First: 2014/01/15]

13. Moher D, Liberati A, Tetzlaff J, et al. Preferred reporting items for systematic reviews and meta-analyses: the PRISMA statement. *Bmj* 2009;339:b2535. doi: 10.1136/bmj.b2535 [published Online First: 2009/07/23]
14. Nicholson A, Ford E, Davies KA, et al. Optimising use of electronic health records to describe the presentation of rheumatoid arthritis in primary care: a strategy for developing code lists. *PloS one* 2013;8(2):e54878. doi: 10.1371/journal.pone.0054878 [published Online First: 2013/03/02]
15. Hamilton W, Lancashire R, Sharp D, et al. The risk of colorectal cancer with symptoms at different ages and between the sexes: a case-control study. *BMC medicine* 2009;7(1):17. doi: 10.1186/1741-7015-7-17
16. Gulliford MC, Charlton J, Ashworth M, et al. Selection of medical diagnostic codes for analysis of electronic patient records. Application to stroke in a primary care database. *PloS one* 2009;4(9):e7168. doi: 10.1371/journal.pone.0007168 [published Online First: 2009/09/25]
17. Bhattarai N, Charlton J, Rudisill C, et al. Coding, recording and incidence of different forms of coronary heart disease in primary care. *PloS one* 2012;7(1):e29776.

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 legends

Figure 1 The method for codelist collation consists of three steps

Figure 2 Flow chart illustrating the selection of codes

For peer review only

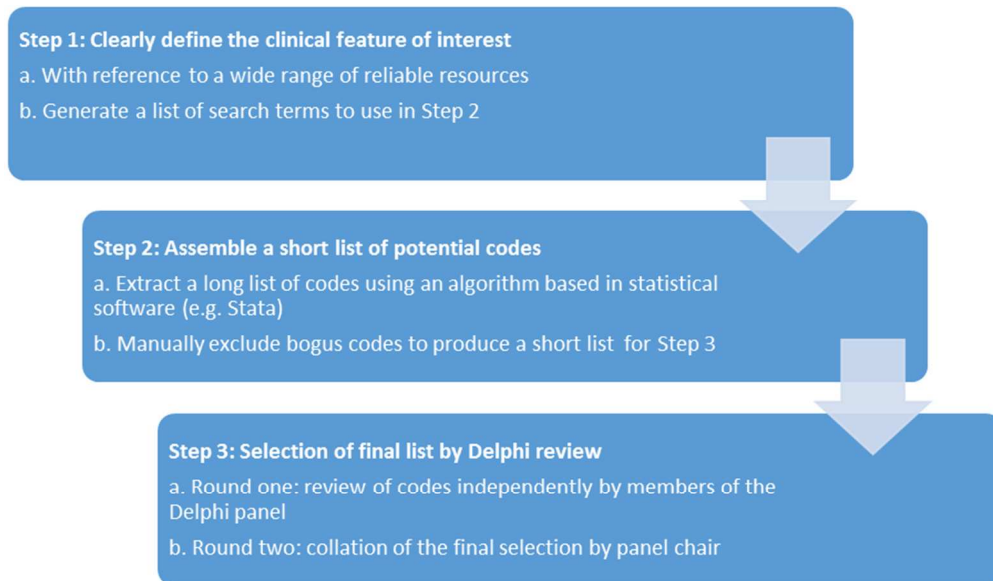


Figure 1 The method for codelist collation consists of three steps

76x44mm (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

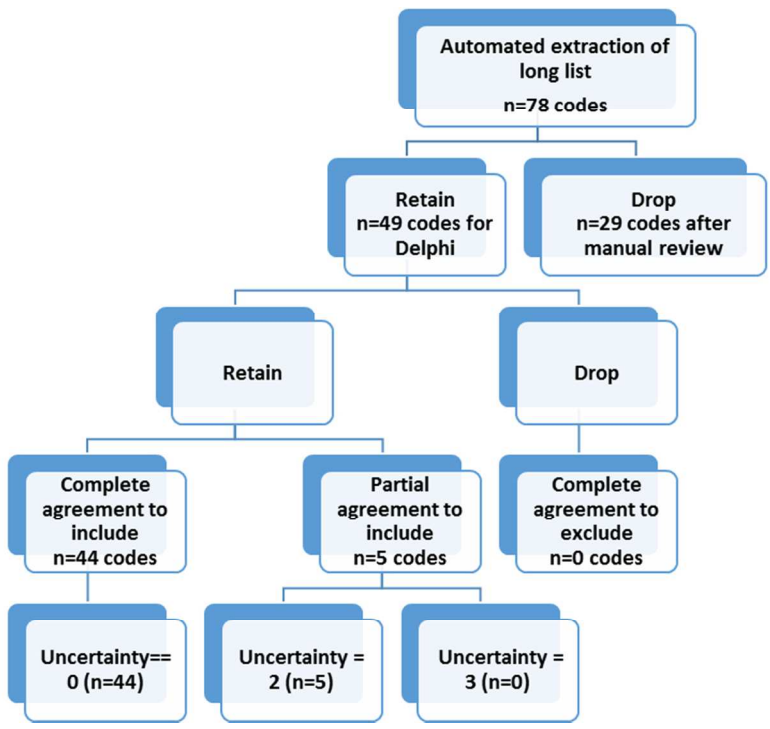


Figure 2 Flow chart illustrating the selection of codes

79x56mm (300 x 300 DPI)

Review only

Online only supplementary material

Table A 1 List of potential codes for inclusion in a list to identify patients with shortness of breath. 'sob' is the binary variable denoting shortness of breath (1: include; 0: exclude)

medcode	desc	sob	Reason for exclusion
735	[D]Breathlessness	1	
741	[D]Shortness of breath	1	
820	[D]Tachypnoea	0	Abnormally rapid breathing, not breathlessness
982	[D]Apnoea	0	Absence of breathing, not breathlessness
1429	Breathlessness	1	
2506	[D]Sleep apnoea syndrome	0	Absence of breathing, not breathlessness
2563	[D]Respiratory distress	1	
2575	Short of breath on exertion	1	
2737	Respiratory distress syndrome	1	
2931	Difficulty breathing	1	
3007	Newborn transitory tachypnoea	0	Abnormally rapid breathing, not breathlessness
3092	[D]Dyspnoea	1	
4822	Shortness of breath	1	
5175	Breathlessness symptom	1	
5349	Shortness of breath symptom	1	
5896	Dyspnoea - symptom	1	
6326	Breathless - moderate exertion	1	

1			
2			
3			
4	6434	Paroxysmal nocturnal dyspnoea	1
5			
6	7000	O/E - dyspnoea	1
7			
8	7534	O/E - respiratory distress	1
9			
10	7603	Sleep apnoea	0
11			Absence of breathing
12			
13	7683	Breathless - lying flat	1
14			
15	7932	Breathless - mild exertion	1
16			
17			
18	8148	Obstructive sleep apnoea	0
19			Absence of breathing
20			
21	9089	Orthopnoea symptom	1
22			
23	9297	[D]Respiratory insufficiency	1
24			
25	10114	O/E - tachypnoea	0
26			Abnormally rapid breathing, not
27			breathlessness
28			
29			
30	11451	[D]Orthopnoea	1
31			
32	12474	SOBOE	1
33			
34			
35	18116	Nocturnal dyspnoea	1
36			
37	19346	No breathlessness	0
38			Negation of breathlessness
39			
40	19426	MRC Breathlessness Scale: grade 3	1
41			
42	19427	MRC Breathlessness Scale: grade 2	1
43			
44	19429	MRC Breathlessness Scale: grade 5	1
45			
46			
47	19430	MRC Breathlessness Scale: grade 4	1
48			
49	19432	MRC Breathlessness Scale: grade 1	1
50			
51			
52	19576	Apnoea of newborn	0
53			Absence of breathing
54	20438	[D]Syndrome sleep apnoea	0
55			Absence of breathing
56			
57	20748	Obstructive sleep apnoea	0
58			Absence of breathing
59	21801	Breathlessness NOS	1
60			

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				
	22094	Short of breath dressing/undressing	1	
	23779	Sleep apnoea	0	Absence of breathing
	23924	Scoline apnoea	0	Absence of breathing
	24848	Adult respiratory distress syndrome	1	
	24889	Breathless - strenuous exertion	1	
	26684	Perinatal respiratory distress NOS	0	Related to pregnancy not pathology
	26871	Primary sleep apnoea of newborn	0	Absence of breathing
	31143	Breathless - at rest	1	
	31913	O/E - hyperpnoea	0	Increased depth and rate of breathing, not breathlessness
	36301	[D]Hypersomnia with sleep apnoea	0	Absence of breathing
	37667	Apnoea alarm monitoring	0	Absence of breathing
	37704	O/E - orthopnoea	1	
	42287	Borg Breathlessness Score: 6 severe (+)	1	
	48539	[D]Insomnia with sleep apnoea	0	Absence of breathing
	53771	Dyspnoea on exertion	1	
	54594	[M]Mesoblastic nephroma	0	Contains 'sob' string, but is inappropriate
	57193	Borg Breathlessness Score: 3 moderate	1	
	57283	Introduction of transobturator tape	0	Contains 'sob' string, but is inappropriate

1				
2				
3				
4	57678	Adult respiratory distress syndrome	1	
5				
6	57759	Borg Breathlessness Score: 2 slight	1	
7				
8	57903	CLASP shortness of breath score	1	
9				
10	58538	Fusobacterial necrotising tonsillitis	0	Contains 'sob' string, but is
11				inappropriate
12				
13				
14				
15	59860	Borg Breathlessness Score: 4	1	
16		somewhat severe		
17				
18				
19				
20	60096	CLASP shortness of breath score	1	
21				
22	64049	Borg Breathlessness Score: 5 severe	1	
23				
24				
25	65353	Borg Breathlessness Score: 0 none at	0	Negates breathlessness
26		all		
27				
28				
29				
30	67566	Borg Breathlessness Score: 9 very,	1	
31		very sev (almost maximal)		
32				
33				
34	68707	Borg Breathlessness Score: 1 very	1	
35		slight		
36				
37				
38				
39	70061	Borg Breathlessness Score: 7 very	1	
40		severe		
41				
42				
43				
44	70818	Borg Breathlessness Score: 0.5 very,	1	
45		very slight		
46				
47				
48				
49	72334	Borg Breathlessness Score: 8 very	1	
50		severe (+)		
51				
52				
53	72704	[X]Other respiratory distress of	0	Newborn
54		newborn		
55				
56				
57				
58				
59				
60				

1 2 3 4 5 6 7	73978	[X]Other apnoea of newborn	0	Absence of breathing, not breathlessness
8 9 10 11 12	93869	Removal of transobturator tape	0	Contains 'sob' string, but is inappropriate
13 14 15 16 17	97037	Introduction of transobturator sling	0	Contains 'sob' string, but is inappropriate
18 19 20 21 22	98965	Urine beta amino isobutyrate level	0	Contains 'sob' string, but is inappropriate
23 24	100177	Berlin questionnaire for sleep apnoea	0	Related to absence of breathing
25 26 27 28	101843	Borg Breathlessness Score: 10 maximal	1	
29 30 31 32	Total		78	

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

Table A 2 Delphi review of codes for shortness of breath. 1: definitely include; 2: uncertain; 3: definitely exclude

medcode	desc	Reviewer 1 Decision	Reviewer 1 comment	Reviewer 2 decision	Reviewer 2 comment
735	[D]Breathlessness	1		1	
741	[D]Shortness of breath	1		1	
1429	Breathlessness	1		1	
2563	[D]Respiratory distress	1		2	I would usually only use this term in children
2575	Short of breath on exertion	1		1	
2737	Respiratory distress syndrome	1		2	I would usually only use this term in children
2931	Difficulty breathing	1		1	
3092	[D]Dyspnoea	1		1	
4822	Shortness of breath	1		1	
5175	Breathlessness symptom	1		1	
5349	Shortness of breath symptom	1		1	
5896	Dyspnoea - symptom	1		1	
6326	Breathless - moderate exertion	1		1	
6434	Paroxysmal nocturnal dyspnoea	1		1	
7000	O/E - dyspnoea	1		1	
7534	O/E - respiratory distress	1		2	
7683	Breathless - lying flat	1		1	
7932	Breathless - mild exertion	1		1	
9089	Orthopnoea symptom	1		1	
9297	[D]Respiratory insufficiency	1		2	

1					
2					
3					
4	11451	[D]Orthopnoea	1		1
5	12474	SOBOE	1		1
6	18116	Nocturnal dyspnoea	1		1
7	19426	MRC Breathlessness Scale: grade 3	1		1
8	19427	MRC Breathlessness Scale: grade 2	1		1
9	19429	MRC Breathlessness Scale: grade 5	1		1
10	19430	MRC Breathlessness Scale: grade 4	1		1
11	19432	MRC Breathlessness Scale: grade 1	1		1
12	21801	Breathlessness NOS	1		1
13	22094	Short of breath dressing/undressing	1		1
14	24848	Adult respiratory distress syndrome	1		2
15	24889	Breathless - strenuous exertion	1		1
16	31143	Breathless - at rest	1		1
17	37704	O/E - orthopnoea	1		1
18	42287	Borg Breathlessness Score: 6 severe (+)	1		1
19	53771	Dyspnoea on exertion	1		1
20	57193	Borg Breathlessness Score: 3 moderate	1		1
21	57678	Adult respiratory distress syndrome	1		1
22	57759	Borg Breathlessness Score: 2 slight	1		1
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					

1					
2					
3					
4	57903	CLASP shortness of breath score	1		1 as above
5					
6	59860	Borg Breathlessness Score: 4 somewhat severe	1		1
7					
8	60096	CLASP shortness of breath score	1		1
9					
10	64049	Borg Breathlessness Score: 5 severe	1		1
11					
12	67566	Borg Breathlessness Score: 9 very, very sev (almost maximal)	1		1
13					
14	68707	Borg Breathlessness Score: 1 very slight	1		1
15					
16	70061	Borg Breathlessness Score: 7 very severe	1		1
17					
18	70818	Borg Breathlessness Score: 0.5 very, very slight	1		1
19					
20	72334	Borg Breathlessness Score: 8 very severe (+)	1		1
21					
22	101843	Borg Breathlessness Score: 10 maximal	1		1
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					

Table A 3 Sensitivity analysis: Medical codes used to record shortness of breath in a group of patients in the year before they were diagnosed with lung cancer

medcode	Description	Frequency	Percent	Cum.
4822	Shortness of breath	3,226	24.66	24.66
741	[D]Shortness of breath	1,455	11.12	35.78
1429	Breathlessness	1,116	8.53	44.32
19427	MRC Breathlessness Scale: grade 2	1,106	8.46	52.77
19426	MRC Breathlessness Scale: grade 3	1,010	7.72	60.49
5349	Shortness of breath symptom	816	6.24	66.73
5175	Breathlessness symptom	785	6.00	72.73
19430	MRC Breathlessness Scale: grade 4	764	5.84	78.57
5896	Dyspnoea - symptom	437	3.34	81.91
2575	Short of breath on exertion	415	3.17	85.09
3092	[D]Dyspnoea	395	3.02	88.10
19432	MRC Breathlessness Scale: grade 1	332	2.54	90.64
6326	Breathless - moderate exertion	261	2.00	92.64
19429	MRC Breathlessness Scale: grade 5	189	1.44	94.08
2931	Difficulty breathing	187	1.43	95.51
12474	SOBOE	166	1.27	96.78
7932	Breathless - mild exertion	142	1.09	97.87
735	[D]Breathlessness	66	0.50	98.37
7000	O/E - dyspnoea	49	0.37	98.75
57903	CLASP shortness of breath score	44	0.34	99.08
31143	Breathless - at rest	39	0.30	99.38

7683	Breathless - lying flat	22	0.17	99.55
6434	Paroxysmal nocturnal dyspnoea	19	0.15	99.69
21801	Breathlessness NOS	10	0.08	99.77
11451	[D]Orthopnoea	9	0.07	99.84
9089	Orthopnoea symptom	8	0.06	99.90
24889	Breathless - strenuous exertion	5	0.04	99.94
18116	Nocturnal dyspnoea	3	0.02	99.96
53771	Dyspnoea on exertion	2	0.02	99.98
22094	Borg Breathlessness Score: 10 maximal	1	0.01	99.98
59860	Borg Breathlessness Score: 4 somewhat..	1	0.01	99.99
101843	Short of breath dressing/undressing	1	0.01	100.00
Total	Total	13,081	100.00	

1
2
3
4 Example R syntax.
5
6

7 Credit: Ben Feakins (Benjamin.Feakins@phc.ox.ac.uk).
8

```
9 #=====#  
10 # #  
11 ##### R CODE FOR APPENDIX #####  
12 # #  
13 #=====#  
14  
15  
16  
17 ### Set Directory Objects ###  
18 browser.dir <- "/Volumes/PHC/CPRD_data/Browsers"  
19 save.dir <- "~/Desktop"  
20  
21  
22  
23 ### Read Data Into R ###  
24 setwd(browser.dir)  
25 medical <- read.delim("medical.txt", header = FALSE,  
26 na.strings = "", stringsAsFactors = FALSE, skip = 1)  
27 names(medical) <- c("medcode", "readcode",  
28 "clinicalevents", "immunisationevents",  
29 "referralevents", "testevents", "readterm",  
30 "databasebuild") # Define headers.  
31  
32  
33  
34  
35  
36 ### Define Search Terms ###  
37 search.terms.general <- c("shortness of  
38 breath|sob|pnoea|pnea|puffed|short of  
39 breath|winded|breathless")  
40 search.terms.breath <- c("breath")  
41 search.terms.breath.resid <-  
42 c("short|difficult|labour|labor|distress|insuff")  
43 search.terms.respir <- c("respir")  
44 search.terms.respir.resid <- c("insuff|distress")  
45  
46  
47  
48  
49  
50 ### Filtering ###  
51 general <- medical[grepl(search.terms.general,  
52 medical$readterm, ignore.case = TRUE), ]  
53 breath <- medical[grepl(search.terms.breath,  
54 medical$readterm, ignore.case = TRUE), ]  
55 breath <- breath[grepl(search.terms.breath.resid,  
56 breath$readterm, ignore.case = TRUE), ]  
57 respir <- medical[grepl(search.terms.respir,  
58 medical$readterm, ignore.case = TRUE), ]  
59  
60
```



```
1
2
3 respir <- respir[grepl(search.terms.respir.resid,
4 respir$readterm, ignore.case = TRUE), ]
5
6
7 ### Combine Results ###
8 medcodes <- rbind(general, breath, respir)
9
10
11 ### De-Duplicate ###
12 medcodes <- unique(medcodes)
13
14
15
16 ### Sort ###
17 medcodes <- medcodes[order(medcodes$medcode), ]
18
19
20 ### Remove Codes Not Related to Breathlessness ###
21 # Further subsetting using grepl().
22
23
24 ### Drop Useless Variables ###
25 medcodes <- medcodes[c("medcode", "readcode")]
26
27
28 ### Save a Copy ###
29 setwd(save.dir)
30 write.table(medcodes, "sob_library.txt", quote = FALSE,
31 sep = "\t", na = "", row.names = FALSE)
32
33
34
35 ### Tidying Up ###
36 rm(browser.dir, save.dir)
37 rm(search.terms.general, search.terms.breath,
38 search.terms.breath.resid, search.terms.respir,
39 search.terms.respir.resid)
40 rm(medical, general, breath, respir)
41 rm(medcodes)
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
```