

# BMJ Open

## Enhancing risk stratification for use in integrated care - A cluster analysis of high-risk patients

Journal:	<i>BMJ Open</i>
Manuscript ID	bmjopen-2016-012903
Article Type:	Research
Date Submitted by the Author:	02-Jun-2016
Complete List of Authors:	Vuik, Sabine; Imperial College London, Institute of Global Health Innovation Mayer, Erik; Imperial College London, Dept. of Biosurgery and Surgical Technology Darzi, Ara; Imperial College London, Institute of Global Health Innovation
<b>Primary Subject Heading</b>:	Health services research
Secondary Subject Heading:	Research methods, Patient-centred medicine, Evidence based practice, General practice / Family practice
Keywords:	Risk management < HEALTH SERVICES ADMINISTRATION & MANAGEMENT, Organisation of health services < HEALTH SERVICES ADMINISTRATION & MANAGEMENT, STATISTICS & RESEARCH METHODS

SCHOLARONE™  
Manuscripts

1  
2  
3 **Enhancing risk stratification for use in integrated care - A cluster**  
4  
5  
6 **analysis of high-risk patients**  
7  
8  
9

10  
11  
12 Sabine I Vuik, Erik Mayer, Ara Darzi  
13

14  
15  
16  
17  
18  
19  
20  
21  
22 Sabine I Vuik, Policy Fellow, Institute of Global Health Innovation  
23

24  
25  
26 Imperial College, 10<sup>th</sup> Floor, St Mary's Hospital, Praed Street, London, W1 2NY, UK  
27  
28

29  
30  
31  
32 Erik Mayer, Clinical Senior Lecturer, Department of Surgery  
33

34  
35  
36 Imperial College, 10<sup>th</sup> Floor, St Mary's Hospital, Praed Street, London, W1 2NY, UK  
37  
38

39  
40  
41 Ara Darzi, Professor in Surgery, Department of Surgery  
42

43  
44  
45 Imperial College, 10<sup>th</sup> Floor, St Mary's Hospital, Praed Street, London, W1 2NY, UK  
46  
47  
48  
49  
50  
51  
52  
53

54 **Correspondence to:** Sabine Vuik s.vuik@imperial.ac.uk, +44(0) 795 714 0479  
55  
56  
57  
58  
59  
60

1  
2  
3 **Keywords:** risk prediction, integrated care, care utilisation, emergency hospitalisation,  
4  
5  
6 high-risk patients  
7

8  
9 **Word count** (excl. title page, abstract, references, tables and figures): 2,456  
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

For peer review only

**ABSTRACT**

**Objective:** To show how segmentation can enhance risk stratification tools for integrated care, by providing insight into different care utilisation patterns within the high-risk population.

**Design:** A retrospective cohort study. A risk score was calculated for each person using a logistic regression, which was then used to select the top 5% high-risk individuals. This population was segmented based on utilisation of different care settings using a k-means cluster analysis. Data from 2008 to 2011 was used to create the risk score and segments, while 2012 data was used to understand the predictive abilities of the models.

**Setting and participants:** Data was collected on primary care use (CPRD) and secondary care use (HES) for a random sample of 300,000 English patients.

**Main measures:** The high-risk population was segmented based on utilisation of four different care settings: emergency acute care, elective acute care, outpatient care and GP care

**Results:** While the risk strata predicted care utilisation at a high level, within the high-risk population utilisation varied significantly. Four different groups of high-risk patients could

1  
2  
3 be identified. These four segments had distinct utilisation patterns across care settings,  
4  
5  
6 reflecting different levels and types of care needs. The 2008-2011 utilisation patterns of  
7  
8  
9 the four segments were consistent with the 2012 patterns.

10  
11  
12  
13  
14 **Discussion:** Cluster analyses revealed that the high-risk population is not homogeneous,  
15  
16  
17 as there exist four groups of patients with different needs across the care continuum.  
18  
19  
20 Since the patterns were predictive of future care use, they can be used to develop  
21  
22  
23 integrated care programmes tailored to these different groups.  
24  
25  
26  
27

28 **Conclusions:** Utilisation-based segmentation augments risk stratification by identifying  
29  
30  
31 patient groups with different care needs, around which integrated care programmes can  
32  
33  
34 be designed.  
35  
36  
37  
38  
39  
40  
41

#### 42 **STRENGTHS AND LIMITATIONS OF THIS STUDY**

- 43  
44  
45 • This study uses patient-level linked primary and secondary care administrative  
46  
47 data  
48
- 49  
50  
51 • Rather than focusing only on emergency care, this study looks at patterns of  
52  
53 utilisation across different care settings to support the development of integrated  
54  
55 care programmes  
56  
57  
58  
59  
60

- Where previous studies have focused on how to identify or manage high-risk patients, this study explores the different patient groups within the high-risk stratum
- The data used was for a random sample of English patients, and may not reflect local trends
- No data was available in linked format for other care settings, such as A&E, mental health, community and social care

## BACKGROUND

In healthcare, a small number of patients accounts for a disproportionately large share of utilisation.<sup>1 2</sup> Identifying and targeting this group can be done through risk stratification.

Risk stratification divides a population based on different levels of risk of a specific outcome, and is a core process to achieve integrated, personalised care.<sup>3-5</sup> For each stratum, a tailored care model can be developed which addresses the specific needs of the patients. Many of the interventions for high-risk patients are primary care-led integrated care programmes, like virtual wards, case management, and enhanced services and access.<sup>4 6-11</sup>

Risk stratification methods often focus on predicting emergency hospitalisations.<sup>3 12-15</sup> Unplanned hospitalisations, including readmissions, are chosen because they are costly for a health system, may indicate low quality care, and have a negative impact on patient experience.<sup>16 17</sup> As such, unplanned hospitalisations are reflective of all elements of the triple aim of healthcare – quality of care, patient experience and cost<sup>18</sup> – and can be considered a ‘triple fail event’.<sup>16</sup> Moreover, since preventing emergency hospitalisations to the acute setting requires effective primary care, they are also an important metric for integrated care.<sup>19</sup>

1  
2  
3 However, risk stratification based on emergency hospitalisations has important limitations.

4  
5  
6 Firstly, this approach only looks at one element of care. While the risk of an emergency  
7  
8 hospitalisation can be expected to correlate with overall use of emergency acute care,  
9  
10 utilisation of other care services may vary. A patient with an emergency hospitalisation  
11  
12 may be under treatment with a specialist; or regularly visit a general practitioner (GP); or  
13  
14 not access ambulatory care at all. In order to design effective integrated care programmes  
15  
16 that link up the appropriate care providers, understanding care use across all settings is  
17  
18 crucial.  
19  
20  
21  
22  
23  
24  
25  
26  
27

28 Secondly, detailed information on the characteristics of the high-risk patients, such as age,  
29  
30 morbidities and socio-economic status, is lost in the final risk score. All patients who end  
31  
32 up in the top stratum have high risk scores, but the factors driving this high score can be  
33  
34 very different. When developing interventions, these should be taken into account to  
35  
36 understand which patients are most likely to respond to different interventions.<sup>12 20</sup>  
37  
38  
39  
40  
41  
42  
43  
44

45 The aim of this study is to show how utilisation-based segmentation can enhance risk  
46  
47 stratification tools used for integrated care by, firstly, taking into account care utilisation  
48  
49 across multiple care settings and, secondly, providing insight into the characteristics of  
50  
51 different patient groups within the high-risk stratum.  
52  
53  
54  
55  
56  
57  
58  
59  
60



## METHODS

### Study design

To show how segmentation can augment risk stratification, we applied both methods to a large patient database. We first trained a risk prediction model to generate risk scores for each patient. Based on these risk scores, we identified the high-risk patient population. In this group we applied a cluster analysis to a range of different utilisation variables. The different clusters were analysed and profiled to understand the different patient types that exist within a high-risk group.

The analyses were conducted for hypothetical "historic" (2008-2011) and "future" (2012) datasets. The historic dataset reflects the information that would be available to healthcare professionals conducting risk stratification and cluster analysis at the end of 2011, while the future dataset was used to understand how accurately the models predicted actual utilisation in the following year.

### Data

A dataset covering primary and secondary care use for a random sample of 300,000 English patients was constructed from Clinical Practice Research Datalink (CPRD) and Hospital Episode Statistics (HES) data (CPRD ISAC approval under protocol 14\_211R).

1  
2  
3 Patients were eligible for inclusion if they were registered with a CPRD-participating GP  
4  
5  
6 practice during the entire study period of 2008 up to and including 2012, and if their HES  
7  
8 records could be linked to CPRD. In England, Clinical Commissioning Groups (CCG) are  
9  
10 responsible for the planning and commissioning of care for local populations. The sample  
11  
12 size in this study was set at 300,000, which is similar to the population of a CCG in the  
13  
14 75<sup>th</sup> percentile,<sup>21</sup> to reflect a typical local population in England.  
15  
16  
17  
18  
19  
20  
21  
22

23 The final dataset included patient demographics, long-term condition (LTC) diagnoses and  
24  
25 utilisation variables. We selected four high-level utilisation variables for the cluster analysis  
26  
27 of high-risk patients: inpatient emergency hospitalisations, inpatient nonemergency  
28  
29 hospitalisations, outpatient attendances and GP visits. These utilisation variables were used  
30  
31 to reflect different care settings that may be incorporated in integrated care models.  
32  
33  
34  
35  
36  
37  
38

### 39 **Risk stratification**

40  
41 We calculated our own risk prediction score, reflecting predictor variables used in PARR,  
42  
43 the Combined Predictive Model and other commonly used risk prediction algorithms. The  
44  
45 risk model was trained to predict emergency hospitalisations in 2012, using a stepwise  
46  
47 logistic regression.<sup>14 22</sup> The number of emergency hospitalisations in 2011 was included as  
48  
49 one of the predictor variables, as well as a range of other variables detailed in appendix 1.  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 The logistic regression on the training set excluded a number of diagnosis variables after  
4  
5  
6 step-wise elimination, as well as the over 75+ flag.  
7  
8  
9

10  
11 To validate the model, a split sample validation method was used. Using the random  
12  
13 sample function of SPSS,<sup>23</sup> half of the sample was defined as the training set and the  
14  
15 other half as the test set. Applying the risk model to the test set, the area under the  
16  
17 Receiver Operator Curve (ROC) was 0.75. This is in line with other models predicting  
18  
19 emergency hospitalisations, which range from 0.55 to 0.83.<sup>13 36</sup> The test population was  
20  
21 stratified into three groups, which are comprised of the top 5% highest risk patient ("High  
22  
23 risk"), the top 5-20% ("Medium risk") and the remaining 80% of the population ("Low  
24  
25 risk"), in accordance with general risk stratification practice.<sup>2 15 17</sup>  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36

### 37 **Segmentation**

38  
39 For the segmentation analysis the k-means algorithm was used to cluster the patients  
40  
41 based on their historic utilisation. This method was selected as it is efficient and produces  
42  
43 roughly similar sized segments.<sup>24</sup> Clustering solutions ranging from 2 to 8 clusters were  
44  
45 explored for the high-risk stratum. To identify the optimal number of clusters, the Pseudo-  
46  
47 F statistic was calculated for all the clustering solutions using STATA.<sup>25</sup> This statistic is  
48  
49 commonly used in healthcare clustering studies,<sup>26-30</sup> and is one of the best criteria to  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 determine the number of clusters.<sup>31</sup> It compares the between-cluster to the within-cluster  
4  
5  
6 sum-of-squares, and a large Pseudo-F statistic indicates distinct clusters.<sup>32</sup>  
7  
8  
9

## 10 11 **Analysis**

12  
13  
14 To create profiles for the segments, the utilisation variables as well as demographic  
15  
16  
17 characteristics were analysed to see if they differed significantly across segments. For the  
18  
19  
20 non-Normal utilisation and LTCs count variables, a Kruskal-Wallis test was used. For the  
21  
22  
23 continuous age and risk score variables an ANOVA test was used, and for the binary  
24  
25  
26 morbidity variables and the 2012 emergency hospitalisation flag a Chi square test. Where  
27  
28  
29 these tests found significant variation across segments, the results were then explored  
30  
31  
32 pair-wise between segments to identify which segment or segments were significantly  
33  
34  
35 different from others. For this, Mann-Whitney U tests, Student t-tests, and z-tests were  
36  
37  
38 used, respectively. To account for the multiplicity problem that occurs when performing  
39  
40  
41 multiple tests, the Bonferroni method was used to adjust the significance level.<sup>33-35</sup>  
42  
43  
44  
45  
46  
47

## 48 **RESULTS**

49  
50 The final dataset contained 298,111 people with a complete record across the variables, of  
51  
52  
53 which 149,320 observations ended up in the test set used for the analyses below. When  
54  
55  
56 the population was stratified based on risk, predictive variables such as age, long-term  
57  
58  
59  
60

conditions and historic care utilisation were all found to increase with each risk stratum (see table 1). In addition to historic utilisation, future utilisation of all care types also increased consistently with the risk strata.

*Table 1: Strata characteristics*

	<i>High risk</i>	<i>Medium risk</i>	<i>Low risk</i>	<i>Total population</i>
Number of people	7,466	22,398	119,456	<b>149,320</b>
Predicted proportion with any emergency hospitalisations in 2012	27%	9%	3%	<b>5%</b>
Actual proportion with any emergency hospitalisations in 2012	27%	11%	3%	<b>5%</b>
Age at end of study period, mean	75	65	40	<b>45</b>
Number of long-term conditions, mean	1.7	0.7	0.1	<b>0.3</b>
Number of emergency hospitalisations per year (historic), mean	0.5	0.1	0.0	<b>0.1</b>
Number of nonemergency hospitalisations per year (historic), mean	0.6	0.3	0.1	<b>0.1</b>
Number of outpatient attendances per year (historic), mean	5.8	3.0	0.8	<b>1.4</b>
Number of GP visits per year (historic), mean	15.7	9.6	3.4	<b>5.0</b>
Number of emergency hospitalisations per year (future), mean	0.4	0.1	0.0	<b>0.1</b>
Number of nonemergency hospitalisations per year (future), mean	0.5	0.4	0.1	<b>0.2</b>
Number of outpatient attendances per year (future), mean	6.1	3.4	1.0	<b>1.6</b>
Number of GP visits per year (future), mean	17.0	10.5	3.8	<b>5.5</b>

For the high-risk population, k-means cluster analyses were performed for 2- to 8-clusters and the pseudo-F statistics was obtained for each solution. A peak was observed around the 3- and 4-cluster solutions. Exploring these two sets of clusters, the 4-cluster solution included an additional, contrasting utilisation pattern and was therefore selected.

The cluster analysis aims to optimise the distance between groups for the clustering variables, and statistical tests confirm that historic utilisation is significantly different across segments (see table 2) In addition, non-clustering variables, including future utilisation, age, number of long-term conditions and most disease prevalence variables, also differ significantly across the clusters.

Table 2: Clusters within the high-risk population

	Cluster				ANOVA/ Kruskal-Wallis/ Chi square test
	1	2	3	4	
Number of people	1967	1807	1831	1861	
Predicted proportion with any emergency hospitalisations in 2012 (based on average risk score), %	21 ***	38 ***	20 ***	31 ***	<b>AN: &lt;0.000</b>
Actual proportion with any emergency hospitalisations in 2012, %	19 **	35 **	21 **	34 **	<b>Chi: &lt;0.000</b>
Age at end of study period, mean	79 ***	67 ***	83 ***	71 ***	<b>AN: &lt;0.000</b>
Number of long-term conditions, mean	1.8 **	2.0 **	1.4 ***	1.7 ***	<b>KW: &lt;0.000</b>
Number of emergency hospitalisations per year (historic), mean	0.1 **	0.9 ***	0.2 **	0.8 ***	<b>KW: &lt;0.000</b>
Number of nonemergency hospitalisations per year (historic), mean	1.0 ***	1.1 ***	0.1 ***	0.1 ***	<b>KW: &lt;0.000</b>
Number of outpatient attendances per year (historic), mean	7.9 ***	9.3 ***	2.5 ***	3.3 ***	<b>KW: &lt;0.000</b>
Number of GP visits per year (historic), mean	17.6 ***	16.7 ***	15.9 ***	12.5 ***	<b>KW: &lt;0.000</b>
Number of emergency hospitalisations per year (future), mean	0.3 **	0.6 **	0.3 **	0.6 **	<b>KW: &lt;0.000</b>
Number of nonemergency hospitalisations per year (future), mean	0.7 **	0.9 **	0.3 ***	0.3 ***	<b>KW: &lt;0.000</b>
Number of outpatient attendances per year (future), mean	7.7 ***	9.1 ***	3.4 ***	4.2 ***	<b>KW: &lt;0.000</b>
Number of GP visits per year (future), mean	18.5 ***	17.9 **	17.5 **	14.2 ***	<b>KW: &lt;0.000</b>
Prevalence of AMI, %	15 ***	23 ***	10 ***	19 ***	<b>Chi: &lt;0.000</b>
Prevalence of asthma, %	28 *	26	24 *	25	<b>Chi: 0.028</b>
Prevalence of cancer, %	26 ***	22 ***	8 ***	5 ***	<b>Chi: &lt;0.000</b>
Prevalence of cerebrovascular disease, %	9 **	15 **	10 **	18 **	<b>Chi: &lt;0.000</b>

Prevalence of congestive heart failure, %	8 ***	13 **	5 ***	13 **	Chi: <0.000
Prevalence of COPD, %	18 *	17 *	13 ***	18 *	Chi: <0.000
Prevalence of dementia, %	3 **	3 **	5 **	7 **	Chi: <0.000
Prevalence of diabetes, %	28 **	22 **	28 **	22 **	Chi: <0.000
Prevalence of HIV/AIDS, %	0	0	0	0	Chi: 0.39
Prevalence of learning disabilities, %	0 *	0 *	0	0	Chi: 0.032
Prevalence of liver disease, %	1	1 *	0 **	1 *	Chi: <0.000
Prevalence of mental health conditions, %	2 *	3 *	2 *	5 ***	Chi: <0.000
Prevalence of paraplegia, %	1 **	3 **	1 **	3 **	Chi: <0.000
Prevalence of peptic ulcer, %	4 *	4 *	2 **	3	Chi: <0.000
Prevalence of peripheral vascular disease, %	8 ***	11 ***	4 **	6 **	Chi: <0.000
Prevalence of renal disease, %	23 *	23 *	24 *	18 ***	Chi: <0.000
Prevalence of rheumatic disease, %	10 **	8 *	6 *	5 **	Chi: <0.000

\*\*\*: Significantly different from all 3 other clusters; \*\*: significantly different from 2 other clusters; \*: significantly different from 1 other clusters; all at  $0.05/4=0.0125$  significance level (Bonferroni adjustment)

The clusters demonstrate a great variation in future care utilisation within the high-risk stratum (see figure 1). Emergency care utilisation, which defines high-risk patients, is high for all clusters. Nevertheless, clusters 1 and 3 have emergency care utilisation rates that lie closer to the medium risk stratum than the high-risk average. Nonemergency hospitalisations and outpatient attendances for clusters 3 and 4 are at or even below the medium risk rate. GP care on the other hand is more homogenous, with the rates for each cluster close to the high-risk average.

While for each care setting there exist high and low utilisation clusters, they are not consistently the same clusters. Each cluster has a unique pattern of utilisation rates (see figure 2). Cluster 1 has high utilisation across most care types, with the exception of

1  
2  
3 emergency care. Cluster 4 has the opposite pattern, with high emergency care use but  
4  
5  
6 low utilisation of other care types. Clusters 2 and 3 have high and low utilisations across  
7  
8  
9 all settings, respectively. The differences between the clusters are strongest for historic  
10  
11  
12 care utilisation, upon which the cluster analysis is based. However, each cluster exhibits  
13  
14  
15 the same pattern of utilisation in 2012.  
16  
17  
18  
19  
20  
21

## 22 **DISCUSSION**

### 23 **Principle findings**

24  
25  
26  
27  
28 The low, medium and high risk strata broadly correlate with care utilisation. For all care  
29  
30  
31 settings, the high-risk stratum has the highest historic and future utilisation. However, this  
32  
33  
34 study shows that, within the high-risk stratum, there is significant variation in care needs  
35  
36  
37 across the care continuum. The high-risk group can be split into four segments with  
38  
39  
40 different care utilisation rates, characteristics and care priorities.  
41  
42  
43  
44

45  
46  
47 Comparing historic and future utilisation for the four clusters, similar patterns can be  
48  
49  
50 observed, indicating that cluster analysis of historic data can help predict future needs.  
51  
52  
53 However, future utilisation rates were closer to the group mean for all clusters and all care  
54  
55  
56 settings than historic rates. This can be at least partially explained by regression to the  
57  
58  
59 mean (RTM), which is known to affect care utilisation predictions.<sup>12 37 38</sup> RTM describes the  
60



1  
2  
3 phenomenon where exceptionally high or low observations tend to be followed by less  
4  
5  
6 extreme observations in repeated measurements.<sup>39</sup> This effect is compounded if subjects  
7  
8  
9 are stratified based on baseline measurements, which is the case when patients are  
10  
11  
12 clustered based on their 2008-2011 utilisation.

### 13 14 15 16 17 **Comparison to previous studies**

18  
19  
20 This study shows that, while integrated care and case management initiatives often are  
21  
22  
23 indiscriminately aimed at high-risk patients, the actual needs of these patients vary widely.  
24  
25  
26 Many studies have discussed how best to identify,<sup>13 14 40 41</sup> or care for,<sup>6 8 10 11 37 42</sup> the high-  
27  
28  
29 risk population, but few have used data analysis to better understand different types of  
30  
31  
32 high-risk patients.

33  
34  
35  
36 A major strength of this study is its reliance on data from both primary and acute care, to  
37  
38  
39 create a more comprehensive picture of care needs. While some risk prediction models,  
40  
41  
42 such as the Combined Predictive Model, include utilisation of non-acute care settings as  
43  
44  
45 predictor variables,<sup>15</sup> this detail is lost in the final risk score and the stratification. An  
46  
47  
48 utilisation-based segmentation analysis, as demonstrated in this study, can be used to  
49  
50  
51 bring out this detail.

### 52 53 54 55 56 **Limitations and future research**

1  
2  
3 While both primary and secondary care data were used in this study to understand care  
4  
5 needs across the continuum, the picture is still incomplete. No patient-level linked data  
6  
7 was available on utilisation of the A&E department, mental health, community and social  
8  
9 care, and these were therefore left out of scope. This is an important limitation, as many  
10  
11 initiatives will require integration of these settings. Future research should be done using  
12  
13 more extensive datasets where these are available.  
14  
15  
16  
17  
18  
19

20  
21  
22 Another limitation is that the population used in this study is a random sample of patients  
23  
24 in England. Local populations may see different sizes or types of segments within their  
25  
26 risk strata. Moreover, this study uses a custom risk prediction algorithm. If providers are  
27  
28 using a specific risk model, they are encouraged to replicate the analysis using their own  
29  
30 population data and risk strata.  
31  
32  
33  
34  
35  
36  
37  
38

### 39 **Implications for integrated care**

40  
41 Segmenting the high-risk stratum using cluster analysis can help tailor and target  
42  
43 integrated care programmes. For example, cluster 1 uses relatively little emergency care,  
44  
45 but has a high utilisation of nonemergency and outpatient care. Patients in this segment  
46  
47 may not be the best target for primary care-led interventions aimed at reducing  
48  
49 emergency hospitalisations, as their overall usage of emergency care is low and they may  
50  
51 already be under management of a specialist.  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3  
4  
5  
6 Cluster 2 has the highest utilisation rates, the highest risk score and the most LTCs.  
7

8  
9 Surprisingly; this segment is also the youngest of the four, with an average age of 67.  
10

11 Overall high care utilisation makes this cluster a worthwhile target for interventions aimed  
12  
13 at reducing care use. As patients in this cluster have extensive care needs across different  
14  
15 settings, they would likely benefit from care coordination and case management  
16  
17 initiatives.  
18  
19  
20  
21

22  
23  
24  
25 Cluster 3 is at 83 years the oldest segment. Despite their old age, disease prevalence  
26  
27 among the patients in this cluster is generally lower. This is reflected in their lower than  
28  
29 average care use across all settings. This segment shows that while interventions often  
30  
31 focus on elderly patients,<sup>6 37 43</sup> this population group does not necessarily have the  
32  
33 highest care usage.  
34  
35  
36  
37  
38  
39

40  
41  
42 Cluster 4 has one of the highest utilisation rates for emergency care, combined with a  
43  
44 lower use of all other care services. Even GP care, which varies little for the other clusters,  
45  
46 is below average for this group. This could indicate a lack of preventative primary care:  
47  
48 patients in this cluster have on average 1.7 LTCs, but their low usage of primary care  
49  
50 could be causing complications which require emergency care. This would make cluster 4  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 a prime target for enhances services and primary care-led interventions focused on  
4  
5  
6 preventing complications and emergency hospitalisations.  
7  
8  
9

## 10 11 CONCLUSION

12  
13  
14 This paper shows that a high risk of emergency hospitalisation is not unequivocally linked  
15  
16  
17 to high overall care needs, or a particular pattern of care use across other care settings.  
18

19  
20 While risk stratification based on emergency hospitalisation can predict general care  
21  
22  
23 utilisation rates, within the high-risk stratum there exist four very different patient types.

24  
25 Cluster analysis can enhance risk stratification by identifying groups of high-risk patients  
26  
27  
28 with unique care patterns across the care continuum, around which integrated care  
29  
30  
31 programmes can be designed.  
32  
33  
34  
35  
36  
37  
38

## 39 STATEMENTS

40  
41  
42 **Database:** This study is based on data from the Clinical Practice Research Datalink  
43  
44  
45 obtained under license from the UK Medicines and Healthcare Products Regulatory  
46  
47  
48 Agency. However, the interpretation and conclusions contained in the study are those of  
49  
50  
51 the authors alone.  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 **Data sharing:** Technical appendix available in supplementary files, statistical code available  
4  
5  
6 from the corresponding author. No additional data available.  
7  
8  
9

10  
11 **Declaration of competing interests:** All authors have completed the ICMJE uniform  
12  
13 disclosure form at [www.icmje.org/coi\\_disclosure.pdf](http://www.icmje.org/coi_disclosure.pdf) and declare: grants from Peter  
14  
15 Sowerby Foundation, during the conduct of the study; no financial relationships with any  
16  
17 organisations that might have an interest in the submitted work in the previous three  
18  
19 years; no other relationships or activities that could appear to have influenced the  
20  
21 submitted work.  
22  
23  
24  
25  
26  
27  
28  
29  
30

31 **Ethics approval:** No ethics approval was required  
32  
33  
34  
35  
36

37 **Funding:** This study was partially funded by the Soweby eHealth Forum, sponsored by the  
38  
39 Peter Sowerby Foundation. The funder had no role in the study design or analysis, or in  
40  
41 the drafting and submission of this paper. The researchers worked independent from the  
42  
43 funders.  
44  
45  
46  
47  
48  
49

50 **Contributors:** SV designed the study, created the database, analysed the data, and  
51  
52 drafted and revised the paper. She is guarantor. EM contributed to the design of the  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 study, analysed the results and revised the draft paper. AD contributed to the design of  
4  
5  
6 the study and revised the draft paper. All have approved the final version for publication.  
7  
8  
9  
10  
11  
12  
13

#### 14 **FIGURE LEGENDS**

15  
16  
17 *Figure 1: Mean future care utilisation for the risk strata - High (H), Medium (M) and Low*  
18  
19 *(L) - and the four high-risk clusters - 1, 2, 3 and 4.*  
20  
21  
22  
23  
24

25  
26 *Figure 2: Patterns of utilisation for the four high-risk clusters – Emergency care*  
27  
28 *hospitalisations (Emg), Nonemergency hospitalisations (NonE), Outpatient attendances*  
29  
30 *(OP) and GP visits (GP) versus the high-risk population mean*  
31  
32  
33  
34  
35  
36  
37  
38  
39

#### 40 **REFERENCES**

- 41  
42 1. Zulman DM, Pal Chee C, Wagner TH, et al. Multimorbidity and healthcare utilisation  
43 among high-cost patients in the US Veterans Affairs Health Care System. *BMJ*  
44 *Open* 2015;**5**(4).  
45  
46 2. Department of Health. Supporting People with Long Term Conditions. An NHS and  
47 Social Care Model to support local innovation and integration. Leeds: Department  
48 of Health, 2005.  
49  
50 3. NHS England. Using case finding and risk stratification: A key service component for  
51 personalised care and support planning. Leeds: NHS England, 2015.  
52  
53 4. Goodwin N, Curry N. Methods for predicting risk of emergency hospitalisation:  
54 promoting self-care and integrated service responses in the home to the most  
55 vulnerable. *Int J Integr Care* 2008;**8**(5).  
56  
57  
58  
59  
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
5. Dueñas-Espín I, Vela E, Pauws S, et al. Proposals for enhanced health risk assessment and stratification in an integrated care scenario. *BMJ Open* 2016;**6**(4).
6. Roland M, Lewis R, Steventon A, et al. Case management for at-risk elderly patients in the English integrated care pilots: observational study of staff and patient experience and secondary care utilisation. *Int J Integr Care* 2012;**12**(5).
7. Lewis G. Case study: Virtual wards at Croydon Primary Care Trust. London: The King's Fund, 2006.
8. Lewis G, Bardsley M, Vaithianathan R, et al. Do 'virtual wards' reduce rates of unplanned hospital admissions, and at what cost? A research protocol using propensity matched controls. *Int J Integr Care* 2011;**11**:e079.
9. NHS England. Enhanced service specification - Avoiding unplanned admissions: proactive case finding and patient review for vulnerable people. Leeds: NHS England, 2014.
10. Wallace E, Smith SM, Fahey T, et al. Reducing emergency admissions through community based interventions. *Bmj* 2016;**352**.
11. Lewis GH, Vaithianathan R, Wright L, et al. Integrating care for high-risk patients in England using the virtual ward model: lessons in the process of care integration from three case sites. *Int J Integr Care* 2013;**13**(4).
12. Lewis G. Next Steps for Risk Stratification in the NHS. London: NHS England, 2015.
13. Billings J, Blunt I, Steventon A, et al. Development of a predictive model to identify inpatients at risk of re-admission within 30 days of discharge (PARR-30). *BMJ Open* 2012;**2**(4).
14. Billings J, Dixon J, Mijanovich T, et al. Case finding for patients at risk of readmission to hospital: development of algorithm to identify high risk patients. *Br Med J* 2006;**333**:327.
15. Wennberg D, Siegel M, Darin B, et al. Combined Predictive Model - Final report & technical documentation. London: Health Dialog, King's Fund and New York University, 2006.
16. Thomson K, Lewis G. Information Governance and Risk Stratification: Advice and Options for CCGs and GPs. London: NHS England, 2013.
17. Lewis G, Curry N, Bardsley M. Choosing a Predictive Risk Model: A guide for commissioners in England. London: Nuffield Trust, 2011.
18. Berwick DM, Nolan TW, Whittington J. The Triple Aim: Care, Health, And Cost. *Health Aff (Millwood)* 2008;**27**(3):759-69.
19. The King's Fund. Predictive risk Project - Literature Review. London: The King's Fund, 2005.
20. Lewis G, Kirkham H, Duncan I, et al. How Health Systems Could Avert 'Triple Fail' Events That Are Harmful, Are Costly, And Result In Poor Patient Satisfaction. *Health Aff (Millwood)* 2013;**32**(4):669-76.

21. NHS Health and Social Care Information Centre. Numbers of Patients Registered at a GP Practice - July 2015. Leeds: NHS Health and Social Care Information Centre, 2015.
22. Bardsley M, Billings J, Dixon J, et al. Predicting who will use intensive social care: case finding tools based on linked health and social care data. *Age Ageing* 2011;**40**:265-70.
23. IBM SPSS Statistics for Macintosh, Version 23.0 [program]. Armonk, NY: IBM Corp, 2015.
24. Han J, Kamber M. *Data Mining: Concepts and Techniques*. 1st ed. San Diego, CA: Academic Press, 2001.
25. Stata Statistical Software: Release 14. [program]. College Station, TX: StataCorp LP, 2015.
26. Armstrong JJ, Zhu M, Hirdes JP, et al. K-Means Cluster Analysis of Rehabilitation Service Users in the Home Health Care System of Ontario: Examining the Heterogeneity of a Complex Geriatric Population. *Arch Phys Med Rehabil* 2012;**93**(12):2198-205.
27. Coste J, Bouyer J, Fernandez H, et al. A population-based analytical approach to assessing patterns, determinants, and outcomes of health care with application to ectopic pregnancy. *Med Care* 2000;**38**(7):739-49.
28. Cryer PC, Saunders J, Jenkins LM, et al. Clusters within a general adult population of alcohol abstainers. *International Journal of Epidemiology* 2001;**30**(4):756-65.
29. Kendig H, Mealing N, Carr R, et al. Assessing patterns of home and community care service use and client profiles in Australia: a cluster analysis approach using linked data. *Health & Social Care in the Community* 2012;**20**(4):375-87.
30. Pud D, Ben Ami S, Cooper BA, et al. The Symptom Experience of Oncology Outpatients Has a Different Impact on Quality-of-Life Outcomes. *J Pain Symptom Manage* 2008;**35**(2):162-70.
31. Everitt BS, Landau S, Leese M, et al. *Cluster analysis*. 5th ed. Chichester: John Wiley & Sons, 2011.
32. StataCorp. Stata 14 Cluster Stop reference manual. College Station, TX: Stata Press, 2015.
33. Ng SK, Holden L, Sun J. Identifying comorbidity patterns of health conditions via cluster analysis of pairwise concordance statistics. *Stat Med* 2012;**31**(27):3393-405.
34. Chan M, Zhu MX. Investigating the health profile of Macau Chinese. *J Clin Nurs* 2008;**17**(11C):352-61.
35. Borglin G, Jakobsson U, Edberg Ak, et al. Older people in Sweden with various degrees of present quality of life: their health, social support, everyday activities and sense of coherence. *Health & Social Care in the Community* 2006;**14**(2):136-46.
36. Kansagara D, Englander H, Salanitro A, et al. Risk prediction models for hospital readmission: A systematic review. *JAMA* 2011;**306**(15):1688-98.



- 1  
2  
3 37. Roland M, Abel G. Reducing emergency admissions: are we on the right track? Br Med  
4 J 2012;**345**(e6017).  
5  
6 38. Georghiou T, Blunt I, Steventon A, et al. Predictive risk and health care: An overview.  
7 London: Nuffield Trust, 2011.  
8  
9 39. Barnett AG, van der Pols JC, Dobson AJ. Regression to the mean: what it is and how to  
10 deal with it. International Journal of Epidemiology 2005;**34**(1):215-20.  
11  
12 40. Hao S, Jin B, Shin AY, et al. Risk prediction of emergency department revisit 30 days  
13 post discharge: a prospective study. PloS one 2014;**9**(11):e112944.  
14  
15 41. Pavlou M, Ambler G, Seaman SR, et al. How to develop a more accurate risk prediction  
16 model when there are few events. Bmj 2015;**351**.  
17  
18 42. Tanio C, Chen C. Innovations At Miami Practice Show Promise For Treating High-Risk  
19 Medicare Patients. Health Aff (Millwood) 2013;**32**(6):1078-82.  
20  
21 43. Alderwick H, Ham C, Buck D. Population health systems: Going beyond integrated  
22 care. London: The King's Fund, 2015.  
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

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

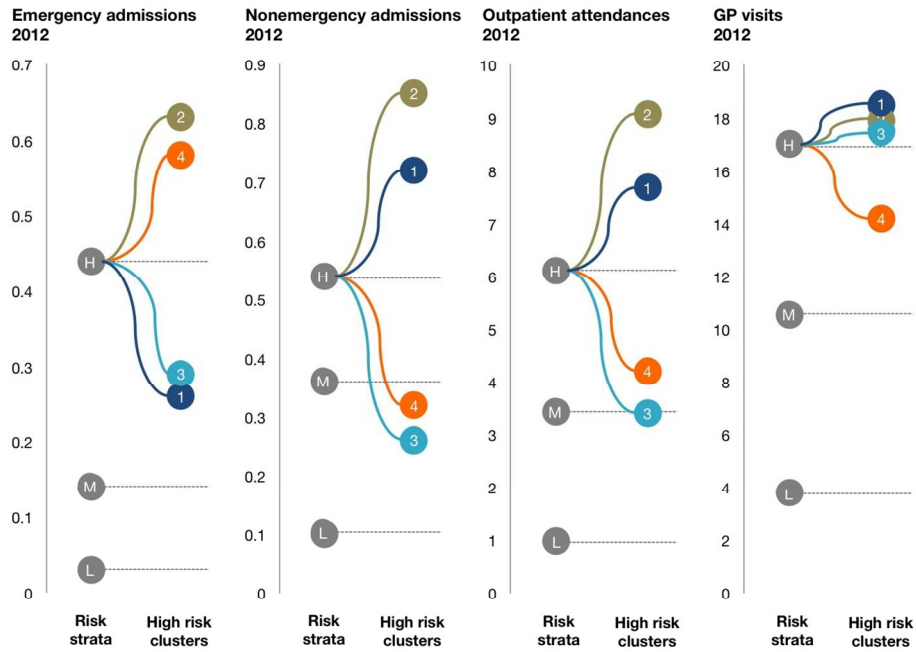
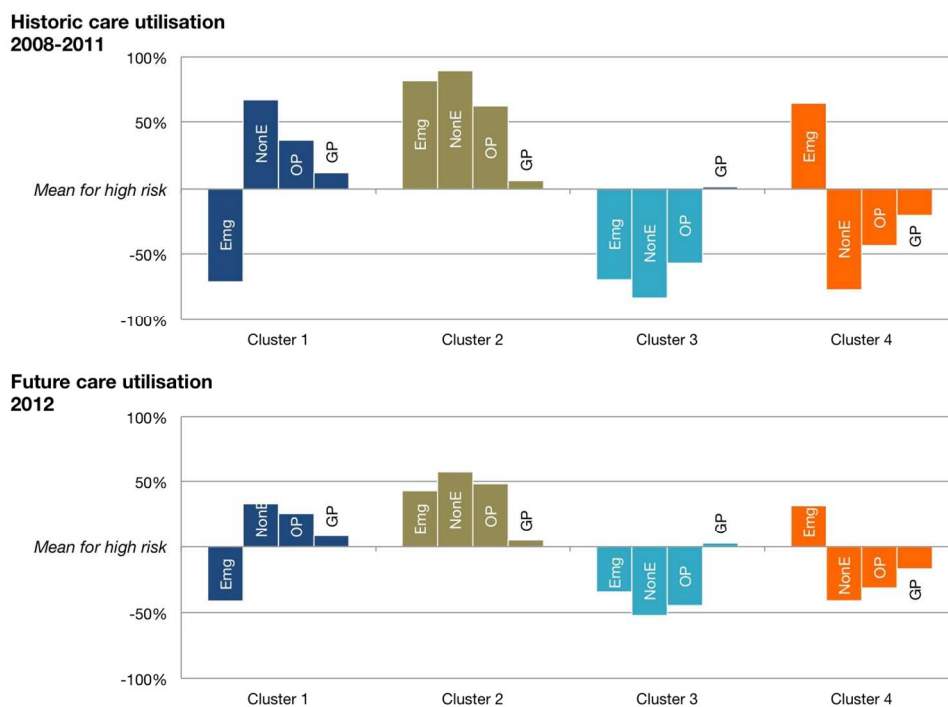


Figure 1: Mean future care utilisation for the risk strata - High (H), Medium (M) and Low (L) - and the four high-risk clusters - 1, 2, 3 and 4.

figure 1  
508x381mm (72 x 72 DPI)



Patterns of utilisation for the four high-risk clusters – Emergency care hospitalisations (Emg), Nonemergency hospitalisations (NonE), Outpatient attendances (OP) and GP visits (GP) versus the high-risk population mean

figure 2

508x381mm (72 x 72 DPI)

View Only

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

## Appendix 1

Variables included in various risk scores and variables selected for our model

Variables considered in hospital admission risk studies	Number of studies out of 30 including variable in final model <sup>1</sup>	Variable in PARR <sup>2</sup>	Variable in PARR-30 <sup>3</sup>	Variable in Combined Predictive Model <sup>4</sup> ( <i>selected variables</i> )	Included in initial model / included in final model after backwards elimination	
<b>Morbidities</b>	Medical diagnoses or comorbidity indices: 24	Cerebrovascular disease			Any diagnosis of cerebrovascular disease in 2008-2011 (in primary or secondary care)	
		Chronic obstructive pulmonary disease	Chronic pulmonary disease	COPD	<b>Any diagnosis of COPD in 2008-2011 (in primary or secondary care)</b>	
					<i>Asthma (only considered in LTC counts)</i>	Any diagnosis of Asthma in 2008-2011 (in primary or secondary care)
		Connective tissue disease/rheumatoid arthritis				<b>Any diagnosis of Rheumatic disease in 2008-2011 (in primary or secondary care)</b>
		Developmental disability				Any diagnosis of Learning disability in 2008-2011 (in primary or secondary care)
		Diabetes	Diabetes with chronic complications	<i>Diabetes (only considered in LTC counts)</i>	Any diagnosis of Diabetes in 2008-2011 (in primary or	

For peer review only - <http://bmjopen.bmj.com/site/about/guidelines.xhtml>

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

For peer review only

			secondary care)
Ischaemic heart disease		<i>CAD (only considered in LTC counts)</i>	<b>Any diagnosis of Ischaemic heart disease in 2008-2011 (in primary or secondary care)</b>
Peripheral vascular disease	Peripheral vascular disease		Any diagnosis of Peripheral vascular disease in 2008-2011 (in primary or secondary care)
Renal failure	Renal disease		<b>Any diagnosis of Renal disease in 2008-2011 (in primary or secondary care)</b>
Sickle cell disease			
	Metastatic cancer with solid tumour	<i>Cancer (only considered in LTC counts)</i>	<b>Any diagnosis of Cancer in 2008-2011 (in primary or secondary care)</b>
	Other malignant cancer		
	Congestive heart failure	<i>CHF (only considered in LTC counts)</i>	Any diagnosis of Congestive heart failure in 2008-2011 (in primary or secondary care)
	Moderate/severe liver disease		Any diagnosis of Liver disease in 2008-2011 (in primary or secondary care)
	Other liver disease		
	Haemiplegia or paraplegia		Any diagnosis of Paraplegia in 2008-2011 (in primary or secondary care)

		Dementia		Any diagnosis of Dementia in 2008-2011 (in primary or secondary care)
			<i>Hypertension (only considered in LTC counts)</i>	
		Diagnostic cost groups/hierarchical condition category		
			1 LTC	<b>Flag if the sum of conditions listed is 0, 1 or 2 or more</b>
			2+ LTCS	
<b>Mental health morbidities</b>	Alcohol or substance use: 11	Alcohol related diagnosis	Psychoactive substance abuse	
	Mental illness: 9		Psychotic disorder	
			Inpatient admission with diagnosis of mental illness	<b>Any diagnosis of Mental health disorder in 2008-2011 (in primary or secondary care)</b>
			<i>Depression (only as included in LTC counts)</i>	
<b>Prior use of medical services</b>	Hospitalisations: 14	Previous admission for respiratory infection		
		Previous admission for a reference condition		
		Number of emergency admissions in previous 90, 180 and 365 days	Whether there had been a prior emergency hospital discharge in the past 30	[Combinations of] 1, 1+, 2, 2+, 3+ emergency admissions in last 30, 30 to 90, 90 to 180, 180 to

For peer review only - <http://bmjopen.bmj.com/site/about/guidelines.xhtml>

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

	days	365, 365 to 730 days	<b>Number of emergency admissions over 2008-2011</b>
Average number of episodes per spell for emergency admissions		Average number of episodes per spell for emergency admissions >=3	
	Whether the current admission was an emergency admission		
Total number of previous emergency admissions in previous three years	Number of emergency hospital discharges in the last year		
Number of non-emergency admissions in previous 365 days			
Emergency department visits: 4		A&E visits and investigations	<b>Number of non-emergency admissions over 2008-2011</b>
Clinic visits or missed visits: 3	Number of different treatment specialists seen		
		[Combinations of] 1, 1-5, 2, 3+, 6-10, 11+ out-patient specialty visits in last 30, 30 to 90, 365 to 730 days	<b>Number of outpatient visits over 2008-2011</b>
Index hospital length of stay: 4			
Other		Polypharmacy: 1-4 unique drugs in any month (last 0 to 90 days);	<b>Number of GP visits over 2008-2011 (including home</b>

For peer review only - <http://bmjopen.bmj.com/site/about/guidelines.xhtml>

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

				5-9; 10+	<i>visits)</i>
<b>Sociodemographic factors</b>	Age: 19	Age 65-74 or age 75+	Age squared	Age band (0-4, 15-39, 40-59, 5 year age bands, 85+)	<b>5-year age bands</b> & Over 75 flag
	Sex: 15	Sex		Gender	<b>Gender</b>
	Race/ ethnicity: 7	Ethnicity			
<b>Social determinants of health</b>	SES, income and employment: 5		Index of multiple deprivation band for the place of residence		<b>Townsend score (5 groups)</b>
	Insurance status: 6				
	Education: 0				
	Marital status and people in household: 4				
	Social support: 2				
	Access to care: 5				
	Discharge location: 2				
<b>Hospital specific metrics</b>	<i>Not included in review</i>	Observed:expected ratio for practice style sensitive admissions in ward of residence			
		Observed:expected ratio for rate of readmissions for hospitals of current admission	Hospital-specific variable		
<b>Illness severity</b>	Severity index: 1				
	Laboratory findings: 4				



	Other: 4	
<b>Overall health and function</b>	Functional status, ADL: 2	
	Self-rated health, QOL: 3	
	Cognitive impairment: 7	
	Visual/hearing impairment: 1	

1. Kansagara D, Englander H, Salanitro A, et al. Risk prediction models for hospital readmission: A systematic review. *JAMA* 2011;**306**(15):1688-98.
2. Billings J, Dixon J, Mijanovich T, et al. Case finding for patients at risk of readmission to hospital: development of algorithm to identify high risk patients. *Br Med J* 2006;**333**:327.
3. Billings J, Blunt I, Steventon A, et al. Development of a predictive model to identify inpatients at risk of re-admission within 30 days of discharge (PARR-30). *BMJ Open* 2012;**2**(4).
4. Wennberg D, Siegel M, Darin B, et al. Combined Predictive Model - Final report & technical documentation. London: Health Dialog, King's Fund and New York University, 2006.

# BMJ Open

## Enhancing risk stratification for use in integrated care - A cluster analysis of high-risk patients in a retrospective cohort study

Journal:	<i>BMJ Open</i>
Manuscript ID	bmjopen-2016-012903.R1
Article Type:	Research
Date Submitted by the Author:	23-Sep-2016
Complete List of Authors:	Vuik, Sabine; Imperial College London, Institute of Global Health Innovation Mayer, Erik; Imperial College London, Dept. of Biosurgery and Surgical Technology Darzi, Ara; Imperial College London, Institute of Global Health Innovation
<b>Primary Subject Heading</b>:	Health services research
Secondary Subject Heading:	Research methods, Patient-centred medicine, Evidence based practice, General practice / Family practice
Keywords:	Risk management < HEALTH SERVICES ADMINISTRATION & MANAGEMENT, Organisation of health services < HEALTH SERVICES ADMINISTRATION & MANAGEMENT, STATISTICS & RESEARCH METHODS

SCHOLARONE™  
Manuscripts

1  
2  
3 **Enhancing risk stratification for use in integrated care - A cluster**  
4  
5  
6 **analysis of high-risk patients in a retrospective cohort study**  
7  
8  
9

10  
11  
12 Sabine I Vuik, Erik Mayer, Ara Darzi  
13

14  
15  
16  
17  
18  
19  
20  
21  
22 Sabine I Vuik, Policy Fellow, Institute of Global Health Innovation  
23

24  
25  
26 Imperial College, 10<sup>th</sup> Floor, St Mary's Hospital, Praed Street, London, W1 2NY, UK  
27  
28

29  
30  
31  
32 Erik Mayer, Clinical Senior Lecturer, Department of Surgery  
33

34  
35  
36 Imperial College, 10<sup>th</sup> Floor, St Mary's Hospital, Praed Street, London, W1 2NY, UK  
37  
38

39  
40  
41 Ara Darzi, Professor in Surgery, Department of Surgery  
42

43  
44  
45 Imperial College, 10<sup>th</sup> Floor, St Mary's Hospital, Praed Street, London, W1 2NY, UK  
46  
47  
48  
49  
50  
51  
52  
53

54 **Correspondence to:** Sabine Vuik s.vuik@imperial.ac.uk, +44(0) 795 714 0479  
55  
56  
57  
58  
59  
60

1  
2  
3 **Keywords:** risk prediction, integrated care, care utilisation, emergency hospitalisation,  
4  
5  
6 high-risk patients  
7

8  
9 **Word count** (excl. title page, abstract, references, tables and figures): 2,456  
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

For peer review only

**ABSTRACT**

**Objective:** To show how segmentation can enhance risk stratification tools for integrated care, by providing insight into different care utilisation patterns within the high-risk population.

**Design:** A retrospective cohort study. A risk score was calculated for each person using a logistic regression, which was then used to select the top 5% high-risk individuals. This population was segmented based on utilisation of different care settings using a k-means cluster analysis. Data from 2008 to 2011 was used to create the risk score and segments, while 2012 data was used to understand the predictive abilities of the models.

**Setting and participants:** Data was collected from administrative datasets covering primary and secondary care for a random sample of 300,000 English patients.

**Main measures:** The high-risk population was segmented based on their utilisation of four different care settings: emergency acute care, elective acute care, outpatient care and GP care.

**Results:** While the risk strata predicted care utilisation at a high level, within the high-risk population utilisation varied significantly. Four different groups of high-risk patients could

1  
2  
3 be identified. These four segments had distinct utilisation patterns across care settings,  
4  
5  
6 reflecting different levels and types of care needs. The 2008-2011 utilisation patterns of  
7  
8  
9 the four segments were consistent with the 2012 patterns.

10  
11  
12  
13  
14 **Discussion:** Cluster analyses revealed that the high-risk population is not homogeneous,  
15  
16  
17 as there exist four groups of patients with different needs across the care continuum.  
18  
19  
20 Since the patterns were predictive of future care use, they can be used to develop  
21  
22  
23 integrated care programmes tailored to these different groups.  
24  
25  
26  
27

28 **Conclusions:** Utilisation-based segmentation augments risk stratification by identifying  
29  
30  
31 patient groups with different care needs, around which integrated care programmes can  
32  
33  
34 be designed.  
35  
36  
37  
38  
39  
40  
41

#### 42 **STRENGTHS AND LIMITATIONS OF THIS STUDY**

- 43  
44  
45 • This study uses a large dataset containing patient-level linked primary and  
46  
47  
48 secondary care administrative data
- 49  
50  
51 • Rather than focusing only on emergency care, this study looks at patterns of  
52  
53  
54 utilisation across different care settings to support the development of integrated  
55  
56  
57 care programmes  
58  
59  
60

- Where previous studies have focused on how to identify or manage high-risk patients, this study explores the different patient groups within the high-risk stratum
- The data used was for a random sample of English patients, and may not reflect local trends
- No data was available in linked format for other care settings, such as A&E, mental health, community and social care

## BACKGROUND

In healthcare, a small number of patients accounts for a disproportionately large share of utilisation.<sup>1 2</sup> Identifying and targeting this group can be done through risk stratification.

Risk stratification divides a population based on different levels of risk of a specific outcome, and is often presented as a core process to achieve integrated, personalised care.<sup>3-5</sup> For each stratum, a tailored care model can be developed which addresses the specific needs of the patients. Many of the interventions for high-risk patients are primary care-led integrated care programmes, like virtual wards, case management, and enhanced services and access.<sup>4 6-11</sup>

Risk stratification methods often focus on predicting emergency hospitalisations.<sup>3 12-15</sup> Unplanned hospitalisations, including readmissions, are chosen because they are costly for a health system, may indicate low quality care, and have a negative impact on patient experience.<sup>16 17</sup> As such, unplanned hospitalisations are reflective of all elements of the triple aim of healthcare – quality of care, patient experience and cost<sup>18</sup> – and can be considered a 'triple fail event'.<sup>16</sup> Moreover, since preventing emergency hospitalisations to the acute setting requires effective primary care, they are also an important metric for integrated care.<sup>19</sup>



1  
2  
3 However, risk stratification based on emergency hospitalisations has important limitations.

4  
5  
6 Firstly, this approach only looks at one element of care. While the risk of an emergency  
7  
8 hospitalisation can be expected to correlate with overall use of emergency acute care,  
9  
10 utilisation of other care services may vary. A patient with an emergency hospitalisation  
11  
12 may be under treatment with a specialist; or regularly visit a general practitioner (GP); or  
13  
14 not access ambulatory care at all. In order to design effective integrated care programmes  
15  
16 that link up the appropriate care providers, understanding care use across all settings is  
17  
18 crucial.  
19  
20  
21  
22  
23  
24  
25  
26  
27

28 Secondly, detailed information on the characteristics of the high-risk patients, such as age,  
29  
30 morbidities and socio-economic status, is lost in the final risk score. All patients who end  
31  
32 up in the top stratum have high risk scores, but the factors driving this high score can be  
33  
34 very different. When developing interventions, these should be taken into account to  
35  
36 understand which patients are most likely to respond to different interventions.<sup>12 20</sup>  
37  
38  
39  
40  
41  
42  
43  
44

45 The aim of this study is to show how utilisation-based segmentation can enhance risk  
46  
47 stratification tools used for integrated care by, firstly, taking into account care utilisation  
48  
49 across multiple care settings and, secondly, providing insight into the characteristics of  
50  
51 different patient groups within the high-risk stratum.  
52  
53  
54  
55  
56  
57  
58  
59  
60

## METHODS

### Study design

To show how segmentation can augment risk stratification, we applied both methods to a large patient database. We first trained a risk prediction model to generate risk scores for each patient. Based on these risk scores, we identified the high-risk patient population. In this group we applied a cluster analysis to a range of different utilisation variables. The different clusters were analysed and profiled to understand the different patient types that exist within a high-risk group.

The analyses were conducted for hypothetical "historic" (2008-2011) and "future" (2012) datasets. The historic dataset reflects the information that would be available to healthcare professionals conducting risk stratification and cluster analysis at the end of 2011, while the future dataset was used to understand how accurately the models predicted actual utilisation in the following year.

### Software

STATA (version 14)<sup>21</sup> was used to perform the cluster analyses and calculate the pseudo-F statistics. For all other analyses, including the risk prediction, SPSS (version 23)<sup>22</sup> was used.

## Data

A dataset covering primary and secondary care use for a random sample of 300,000 English patients was constructed from Clinical Practice Research Datalink (CPRD) and Hospital Episode Statistics (HES) data (CPRD ISAC approval under protocol 14\_211R). Patients were eligible for inclusion if they were registered with a CPRD-participating GP practice during the entire study period of 2008 up to and including 2012, and if their HES records could be linked to CPRD. Other than those two criteria, the sample was entirely random. The CPRD dataset is broadly representative of the age, sex and ethnicity composition of the UK population.<sup>23</sup> In England, Clinical Commissioning Groups (CCG) are responsible for the planning and commissioning of care for local populations. The sample size in this study was set at 300,000, which is similar to the population of a CCG in the 75<sup>th</sup> percentile,<sup>24</sup> to reflect a typical local population in England.

The final dataset included patient demographics, long-term condition (LTC) diagnoses and utilisation variables. We selected four high-level utilisation variables for the cluster analysis of high-risk patients: inpatient emergency hospitalisations, inpatient nonemergency hospitalisations, outpatient attendances and GP visits. These utilisation variables were used to reflect different care settings that may be incorporated in integrated care models. For the cluster analysis, the utilisation variables were log-normalised and standardised to reduce the impact of outliers and give equal weight to each variable.

## Risk stratification

We calculated our own risk prediction score, reflecting predictor variables used in Patients at Risk of Re-hospitalisation (PARR) tool, the Combined Predictive Model and other commonly used risk prediction algorithms. The risk model was trained to predict emergency hospitalisations in 2012, using a stepwise logistic regression.<sup>14 25</sup> The number of emergency hospitalisations in 2011 was included as one of the predictor variables, as well as a range of other variables used in previous risk models,<sup>13-15 26</sup> as detailed in appendix 1. The logistic regression on the training set excluded a number of diagnosis variables after step-wise elimination, as well as the 75+ flag.

To validate the model, a split sample validation method was used. Using the random sample function of SPSS, half of the sample was defined as the training set and the other half as the test set. Applying the risk model to the test set, the area under the Receiver Operator Curve (ROC) was 0.75. This is in line with other models predicting emergency hospitalisations, which range from 0.55 to 0.83.<sup>13 26</sup> The test population was stratified into three groups, which are comprised of the top 5% highest risk patient ("High risk"), the top 5-20% ("Medium risk") and the remaining 80% of the population ("Low risk"), in accordance with general risk stratification practice.<sup>2 15 17</sup>

## Segmentation

For the segmentation analysis the k-means algorithm was used to cluster the patients based on their historic utilisation. This method was selected as it is efficient and produces roughly similar sized segments.<sup>27</sup> Clustering solutions ranging from 2 to 8 clusters were explored for the high-risk stratum. To identify the optimal number of clusters, the Pseudo-F statistic was calculated for all the clustering solutions using STATA. This statistic is commonly used in healthcare clustering studies,<sup>28-32</sup> and is one of the best criteria to determine the number of clusters.<sup>33</sup> It compares the between-cluster to the within-cluster sum-of-squares, and a large Pseudo-F statistic indicates distinct clusters.<sup>34</sup> In addition, the different clustering solutions were also explored using Ward's linkage clustering and post-hoc analysis, as detailed in appendix 2. Both the k-means and Ward's clustering analyses used the Euclidian distance measure.

The clusters were evaluated based on their validity, through statistical test confirming the differences between clusters, and their stability, by comparing future care utilisation of each cluster to the historic pattern.

## Analysis

To create profiles for the segments, the utilisation variables as well as demographic characteristics were analysed to see if they differed significantly across segments. For the

non-Normal utilisation and LTCs count variables, a Kruskal-Wallis test was used. For the continuous age and risk score variables an ANOVA test was used, and for the binary morbidity variables and the 2012 emergency hospitalisation flag a Chi square test. Where these tests found significant variation across segments, the results were then explored pair-wise between segments to identify which segment or segments were significantly different from others. For this, Mann-Whitney U tests, Student t-tests, and z-tests were used, respectively. To account for the multiplicity problem that occurs when performing multiple tests, the Bonferroni method was used to adjust the significance level.<sup>35-37</sup>

## RESULTS

The final dataset contained 298,111 people with a complete record across the variables, of which 149,320 observations were allocated to the test set used for the analyses below.

When the population was stratified based on risk, predictive variables such as age, long-term conditions and historic care utilisation were all found to increase with each risk stratum (see table 1). In addition to historic utilisation, future utilisation of all care types also increased for the high-risk stratum.

*Table 1: Strata characteristics*

	<i>High risk</i>	<i>Medium risk</i>	<i>Low risk</i>	<i>Total population</i>
--	------------------	--------------------	-----------------	-------------------------

Number of people	7,466	22,398	119,456	<b>149,320</b>
Predicted proportion with any emergency hospitalisations in 2012 (based on the average risk score)	27%	9%	3%	<b>5%</b>
Actual proportion with any emergency hospitalisations in 2012	27%	11%	3%	<b>5%</b>
Age at end of study period, mean	75	65	40	<b>45</b>
Number of long-term conditions, median (Interquartile Range/IQR)	2 (1 to 2)	1 (0 to 1)	0 (0 to 0)	<b>0 (0 to 0)</b>
Number of emergency hospitalisations over 2008-2011, median (IQR)	1 (1 to 3)	0 (0 to 1)	0 (0 to 0)	<b>0 (0 to 0)</b>
Number of nonemergency hospitalisations over 2008-2011, median (IQR)	1 (0 to 3)	1 (0 to 2)	0 (0 to 0)	<b>0 (0 to 1)</b>
Number of outpatient attendances over 2008-2011, median (IQR)	16 (8 to 30)	8 (2 to 16)	1 (0 to 4)	<b>1 (0 to 6)</b>
Number of GP visits over 2008-2011 median (IQR)	55 (35 to 82)	34 (22 to 51)	10 (4 to 20)	<b>13 (6 to 27)</b>
Number of emergency hospitalisations in 2012, median (IQR)	0 (0 to 1)	0 (0 to 0)	0 (0 to 0)	<b>0 (0 to 0)</b>
Number of nonemergency hospitalisations in 2012, median (IQR)	0 (0 to 1)	0 (0 to 0)	0 (0 to 0)	<b>0 (0 to 0)</b>
Number of outpatient attendances in 2012, median (IQR)	4 (1 to 8)	1 (0 to 4)	0 (0 to 1)	<b>0 (0 to 2)</b>
Number of GP visits in 2012, median (IQR)	13 (7 to 22)	8 (5 to 14)	2 (0 to 5)	<b>3 (1 to 7)</b>

For the high-risk population, k-means cluster analyses were performed for 2- to 8-clusters and the pseudo-F statistics was obtained for each solution. A peak was observed around the 3- and 4-cluster solutions. Exploring these two sets of clusters, the 4-cluster solution included an additional, contrasting utilisation pattern and was therefore selected.

The cluster analysis aims to optimise the distance between groups for the clustering variables, and statistical tests confirm that historic utilisation is significantly different across segments (see table 2). In addition, non-clustering variables, including future utilisation,

age, number of long-term conditions and most disease prevalence variables, also differ significantly across the clusters.

Table 2: Clusters within the high-risk population

	Cluster				ANOVA/ Kruskal- Wallis/ Chi square test
	1	2	3	4	
<b>Clustering variables</b>					
Number of emergency hospitalisations over 2008-2011, median (IQR)	** 1 (0 to 1)	*** 3 (2 to 4)	1 (0 to ** 1)	3 (2 to 4)	*** KW: <0.000
Number of nonemergency hospitalisations over 2008-2011, median (IQR)	*** 3 (2 to 5)	*** 3 (2 to 5)	0 (0 to *** 1)	0 (0 to 1)	*** KW: <0.000
Number of outpatient attendances over 2008-2011, median (IQR)	24 (16 to *** 38)	29 (18 to *** 46)	7 (3 to *** 13)	10 (5 to 18)	*** KW: <0.000
Number of GP visits over 2008-2011, median (IQR)	61 (43 to *** 90)	57 (40 to *** 86)	55 (35 *** to 82)	42 (26 to *** 65)	*** KW: <0.000
<b>Post-hoc analysis of other variables</b>					
Number of people	1967	1807	1831	1861	
Predicted proportion with any emergency hospitalisations in 2012 (based on average risk score), %	21 ***	38 ***	20 ***	31 ***	AN: <0.000
Actual proportion with any emergency hospitalisations in 2012, %	19 **	35 **	21 **	34 **	Chi: <0.000
Age at end of study period, mean	79 ***	67 ***	83 ***	71 ***	AN: <0.000
Number of long-term conditions, median (IQR)	** 2 (1 to 3)	** 2 (1 to 3)	1 (1 to *** 2)	1 (1 to 2)	*** KW: <0.000
Number of emergency hospitalisations in 2012, median (IQR)	** 2 (1 to 3)	** 2 (1 to 3)	1 (1 to ** 2)	1 (1 to 2)	** KW: <0.000
Number of nonemergency hospitalisations in 2012, median (IQR)	** 0 (0 to 0)	** 0 (0 to 1)	0 (0 to *** 0)	0 (0 to 1)	*** KW: <0.000
Number of outpatient attendances in 2012, median (IQR)	*** 0 (0 to 1)	*** 0 (0 to 1)	0 (0 to *** 0)	0 (0 to 0)	*** KW: <0.000
Number of GP visits in 2012,	5 (2 to 10) ***	6 (3 to 11) **	2 (0 to **)	2 (0 to 5) ***	*** KW: <0.000



median (IQR)			4)		
Prevalence of acute myocardial infarction, %	15 ***	23 ***	10 ***	19 ***	Chi: <0.000
Prevalence of asthma, %	28 *	26	24 *	25	Chi: 0.028
Prevalence of cancer, %	26 ***	22 ***	8 ***	5 ***	Chi: <0.000
Prevalence of cerebrovascular disease, %	9 **	15 **	10 **	18 **	Chi: <0.000
Prevalence of congestive heart failure, %	8 ***	13 **	5 ***	13 **	Chi: <0.000
Prevalence of COPD, %	18 *	17 *	13 ***	18 *	Chi: <0.000
Prevalence of dementia, %	3 **	3 **	5 **	7 **	Chi: <0.000
Prevalence of diabetes, %	28 **	22 **	28 **	22 **	Chi: <0.000
Prevalence of HIV/AIDS, %	0	0	0	0	Chi: 0.39
Prevalence of learning disabilities, %	0 *	0 *	0	0	Chi: 0.032
Prevalence of liver disease, %	1	1 *	0 **	1 *	Chi: <0.000
Prevalence of mental health conditions, %	2 *	3 *	2 *	5 ***	Chi: <0.000
Prevalence of paraplegia, %	1 **	3 **	1 **	3 **	Chi: <0.000
Prevalence of peptic ulcer, %	4 *	4 *	2 **	3	Chi: <0.000
Prevalence of peripheral vascular disease, %	8 ***	11 ***	4 **	6 **	Chi: <0.000
Prevalence of renal disease, %	23 *	23 *	24 *	18 ***	Chi: <0.000
Prevalence of rheumatic disease, %	10 **	8 *	6 *	5 **	Chi: <0.000

\*\*\*: Significantly different from all 3 other clusters; \*\*: significantly different from 2 other clusters; \*: significantly different from 1 other clusters; all at  $0.05/4=0.0125$  significance level (Bonferroni adjustment)

The clusters demonstrate a great variation in future care utilisation within the high-risk stratum (see figure 1). Emergency care utilisation, which defines high-risk patients, is high for all clusters. Nevertheless, clusters 1 and 3 have emergency care utilisation rates that lie closer to the medium risk stratum than the high-risk average. Nonemergency hospitalisations and outpatient attendances for clusters 3 and 4 are at or even below the medium risk rate. GP care on the other hand is more homogenous, with the rates for each cluster close to the high-risk average.

1  
2  
3  
4  
5  
6 While for each care setting there exist high and low utilisation clusters, they are not  
7  
8 consistently the same clusters. Each cluster has a unique pattern of utilisation rates (see  
9  
10 figure 2). Cluster 1 has high utilisation across most care types, with the exception of  
11  
12 emergency care. Cluster 4 has the opposite pattern, with high emergency care use but  
13  
14 low utilisation of other care types. Clusters 2 and 3 have high and low utilisations across  
15  
16 all settings, respectively. The differences between the clusters are strongest for historic  
17  
18 care utilisation, upon which the cluster analysis is based. However, each cluster exhibits  
19  
20 the same pattern of utilisation in 2012.  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32

## 33 **DISCUSSION**

### 34 **Principle findings**

35  
36 The low, medium and high risk strata broadly correlate with care utilisation. For all care  
37  
38 settings, the high-risk stratum has the highest historic and future utilisation. However, this  
39  
40 study shows that, within the high-risk stratum, there is significant variation in care needs  
41  
42 across the care continuum. The high-risk group can be split into four segments with  
43  
44 different care utilisation rates, characteristics and care priorities.  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 Comparing historic and future utilisation for the four clusters, similar patterns can be  
4  
5  
6 observed, indicating that cluster analysis of historic data can help predict future needs.  
7  
8  
9 However, future utilisation rates were closer to the group mean for all clusters and all care  
10  
11 settings than historic rates. This can be at least partially explained by regression to the  
12  
13 mean (RTM), which is known to affect care utilisation predictions.<sup>12 38 39</sup> RTM describes the  
14  
15 phenomenon where exceptionally high or low observations tend to be followed by less  
16  
17 extreme observations in repeated measurements.<sup>40</sup> This effect is compounded if subjects  
18  
19 are stratified based on baseline measurements, which is the case when patients are  
20  
21 clustered based on their 2008-2011 utilisation.  
22  
23  
24  
25  
26  
27  
28  
29  
30

### 31 **Comparison to previous studies**

32  
33 This study shows that, while integrated care and case management initiatives often are  
34  
35 indiscriminately aimed at high-risk patients, the actual needs of these patients vary widely.  
36  
37  
38 Many studies have discussed how best to identify,<sup>13 14 41 42</sup> or care for,<sup>6 8 10 11 38 43</sup> the high-  
39  
40 risk population, but few have used data analysis to better understand different types of  
41  
42 high-risk patients.  
43  
44  
45  
46  
47  
48  
49

50  
51 A major strength of this study is its reliance on data from both primary and acute care, to  
52  
53 create a more comprehensive picture of care needs. While some risk prediction models,  
54  
55 such as the Combined Predictive Model, include utilisation of non-acute care settings as  
56  
57  
58  
59  
60

1  
2  
3 predictor variables,<sup>15</sup> this detail is lost in the final risk score and the stratification. An  
4  
5  
6 utilisation-based segmentation analysis, as demonstrated in this study, can be used to  
7  
8  
9 bring out this detail.

### 14 **Limitations and future research**

16  
17 While both primary and secondary care data were used in this study to understand care  
18  
19 needs across the continuum, the picture is still incomplete. No patient-level linked data  
20  
21 was available on utilisation of the Accident and Emergency (A&E) department, mental  
22  
23 health, community and social care, and these were therefore left out of scope. This is an  
24  
25 important limitation, as many initiatives will require integration of these settings. Future  
26  
27 research should be done using more extensive datasets where these are available.  
28  
29  
30  
31  
32

33  
34  
35  
36 Another limitation is that the population used in this study is a random sample of patients  
37  
38 in England. In this specific sample, the long-term condition prevalence was relatively low.  
39  
40 This could be attributable to the fact that conditions were identified based on coded  
41  
42 diagnoses in the administrative data rather than from disease registries, but it could also  
43  
44 be a characteristic of our sample. Local populations may see different sizes or types of  
45  
46 segments within their risk strata. Moreover, this study uses a custom risk prediction  
47  
48 algorithm. If providers are using a specific risk model, they are encouraged to replicate  
49  
50 the analysis using their own population data and risk strata.  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

### Implications for integrated care

Segmenting the high-risk stratum using cluster analysis can help tailor and target integrated care programmes. For example, cluster 1 uses relatively little emergency care, but has a high utilisation of nonemergency and outpatient care. Patients in this segment may not be the best target for primary care-led interventions aimed at reducing emergency hospitalisations, as their overall usage of emergency care is low and they may already be under management of a specialist.

Cluster 2 has the highest utilisation rates, the highest risk score and the most LTCs.

Surprisingly; this segment is also the youngest of the four, with an average age of 67.

Overall high care utilisation makes this cluster a worthwhile target for interventions aimed at reducing care use. As patients in this cluster have extensive care needs across different settings, they would likely benefit from care coordination and case management initiatives.

Cluster 3 is at 83 years the oldest segment. Despite their old age, disease prevalence among the patients in this cluster is generally lower. This is reflected in their lower than average care use across all settings. This segment shows that while interventions often

1  
2  
3 focus on elderly patients,<sup>6 38 44</sup> this population group does not necessarily have the  
4  
5  
6 highest care usage.  
7  
8  
9

10  
11 Cluster 4 has one of the highest utilisation rates for emergency care, combined with a  
12  
13 lower use of all other care services. Even GP care, which varies little for the other clusters,  
14  
15 is below average for this group. This could indicate a lack of preventative primary care:  
16  
17 patients in this cluster have on average 1.7 LTCs, but their low usage of primary care  
18  
19 could be causing complications which require emergency care. This would make cluster 4  
20  
21 a prime target for enhances services and primary care-led interventions focused on  
22  
23 preventing complications and emergency hospitalisations.  
24  
25  
26  
27  
28  
29  
30  
31  
32

33  
34 However, it is important to note that the above implications are theoretical and have not  
35  
36 been confirmed in practice. Future research is needed to translate the theoretical concepts  
37  
38 presented in this paper into actionable information, including effective interventions and  
39  
40 implementation.  
41  
42  
43  
44  
45  
46  
47

## 48 CONCLUSION

49  
50 This paper shows that a high risk of emergency hospitalisation is not unequivocally linked  
51  
52 to high overall care needs, or a particular pattern of care use across other care settings.  
53  
54

55  
56 While risk stratification based on emergency hospitalisation can predict general care  
57  
58  
59  
60

1  
2  
3 utilisation rates, within the high-risk stratum there exist four very different patient types.  
4  
5

6 Cluster analysis can enhance risk stratification by identifying groups of high-risk patients  
7

8 with unique care patterns across the care continuum, around which integrated care  
9

10 programmes can be designed.  
11  
12

### 13 14 15 16 17 18 19 20 **STATEMENTS**

21  
22 **Database:** This study is based on data from the Clinical Practice Research Datalink  
23

24 obtained under license from the UK Medicines and Healthcare Products Regulatory  
25

26 Agency. However, the interpretation and conclusions contained in the study are those of  
27

28 the authors alone.  
29  
30  
31  
32

33  
34  
35  
36 **Data sharing:** Technical appendix available in supplementary files, statistical code available  
37

38 from the corresponding author. No additional data available.  
39  
40  
41  
42  
43  
44

45 **Declaration of competing interests:** All authors have completed the ICMJE uniform  
46

47 disclosure form at [www.icmje.org/coi\\_disclosure.pdf](http://www.icmje.org/coi_disclosure.pdf) and declare: grants from Peter  
48

49 Sowerby Foundation, during the conduct of the study; no financial relationships with any  
50

51 organisations that might have an interest in the submitted work in the previous three  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 years; no other relationships or activities that could appear to have influenced the  
4  
5  
6 submitted work.  
7  
8  
9

10  
11 **Ethics approval:** No ethics approval was required  
12  
13

14  
15  
16  
17 **Funding:** This study was partially funded by the Sowerby eHealth Forum, sponsored by  
18  
19 the Peter Sowerby Foundation. The funder had no role in the study design or analysis, or  
20  
21 in the drafting and submission of this paper. The researchers worked independent from  
22  
23 the funders.  
24  
25  
26  
27  
28  
29  
30

31 **Contributors:** SV designed the study, created the database, analysed the data, and  
32  
33 drafted and revised the paper. She is guarantor. EM contributed to the design of the  
34  
35 study, analysed the results and revised the draft paper. AD contributed to the design of  
36  
37 the study and revised the draft paper. All have approved the final version for publication.  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47

#### 48 **FIGURE LEGENDS**

49

50 *Figure 1: Mean future care utilisation for the risk strata - High (H), Medium (M) and Low*  
51  
52 *(L) - and the four high-risk clusters - 1, 2, 3 and 4.*  
53  
54  
55  
56  
57  
58  
59  
60



Figure 2: Patterns of utilisation for the four high-risk clusters – Emergency care

hospitalisations (Emg), Nonemergency hospitalisations (NonE), Outpatient attendances

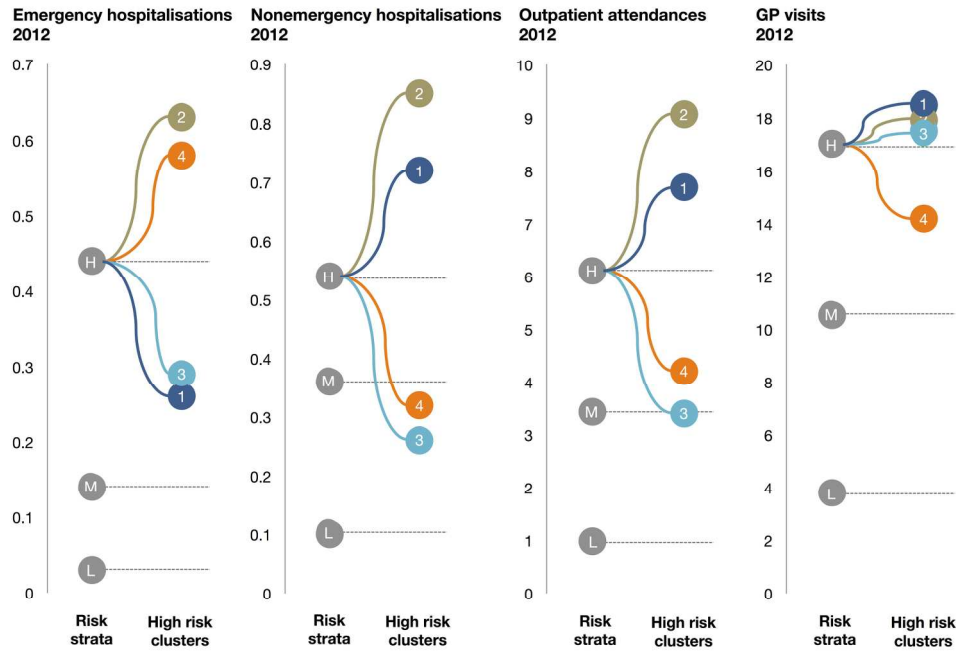
(OP) and GP visits (GP) versus the high-risk population mean

## REFERENCES

1. Zulman DM, Pal Chee C, Wagner TH, et al. Multimorbidity and healthcare utilisation among high-cost patients in the US Veterans Affairs Health Care System. *BMJ Open* 2015;**5**(4).
2. Department of Health. Supporting People with Long Term Conditions. An NHS and Social Care Model to support local innovation and integration. Leeds: Department of Health, 2005.
3. NHS England. Using case finding and risk stratification: A key service component for personalised care and support planning. Leeds: NHS England, 2015.
4. Goodwin N, Curry N. Methods for predicting risk of emergency hospitalisation: promoting self-care and integrated service responses in the home to the most vulnerable. *Int J Integr Care* 2008;**8**(5).
5. Dueñas-Espín I, Vela E, Pauws S, et al. Proposals for enhanced health risk assessment and stratification in an integrated care scenario. *BMJ Open* 2016;**6**(4).
6. Roland M, Lewis R, Steventon A, et al. Case management for at-risk elderly patients in the English integrated care pilots: observational study of staff and patient experience and secondary care utilisation. *Int J Integr Care* 2012;**12**(5).
7. Lewis G. Case study: Virtual wards at Croydon Primary Care Trust. London: The King's Fund, 2006.
8. Lewis G, Bardsley M, Vaithianathan R, et al. Do 'virtual wards' reduce rates of unplanned hospital admissions, and at what cost? A research protocol using propensity matched controls. *Int J Integr Care* 2011;**11**:e079.
9. NHS England. Enhanced service specification: Avoiding unplanned admissions: proactive case finding and patient review for vulnerable people. Leeds: NHS England, 2014.
10. Wallace E, Smith SM, Fahey T, et al. Reducing emergency admissions through community based interventions. *Br Med J* 2016;**352**.
11. Lewis GH, Vaithianathan R, Wright L, et al. Integrating care for high-risk patients in England using the virtual ward model: lessons in the process of care integration from three case sites. *Int J Integr Care* 2013;**13**(4).
12. Lewis G. Next Steps for Risk Stratification in the NHS. London: NHS England, 2015.
13. Billings J, Blunt I, Steventon A, et al. Development of a predictive model to identify inpatients at risk of re-admission within 30 days of discharge (PARR-30). *BMJ Open* 2012;**2**(4).
14. Billings J, Dixon J, Mijanovich T, et al. Case finding for patients at risk of readmission to hospital: development of algorithm to identify high risk patients. *Br Med J* 2006;**333**:327.
15. Wennberg D, Siegel M, Darin B, et al. Combined Predictive Model - Final report & technical documentation. London: Health Dialog, King's Fund and New York University, 2006.
16. Thomson K, Lewis G. Information Governance and Risk Stratification: Advice and Options for CCGs and GPs. London: NHS England, 2013.
17. Lewis G, Curry N, Bardsley M. Choosing a Predictive Risk Model: A guide for commissioners in England. London: Nuffield Trust, 2011.
18. Berwick DM, Nolan TW, Whittington J. The Triple Aim: Care, Health, And Cost. *Health Aff (Millwood)* 2008;**27**(3):759-69.
19. The King's Fund. Predictive Risk Project - Literature Review. London: The King's Fund, 2005.
20. Lewis G, Kirkham H, Duncan I, et al. How Health Systems Could Avert 'Triple Fail' Events That Are Harmful, Are Costly, And Result In Poor Patient Satisfaction. *Health Aff (Millwood)* 2013;**32**(4):669-76.
21. Stata Statistical Software: Release 14. [program]. College Station, TX: StataCorp LP, 2015.
22. IBM SPSS Statistics for Macintosh, Version 23.0 [program]. Armonk, NY: IBM Corp, 2015.
23. Herrett E, Gallagher AM, Bhaskaran K, et al. Data Resource Profile: Clinical Practice Research Datalink (CPRD). *International Journal of Epidemiology* 2015;**44**(3):827-36.

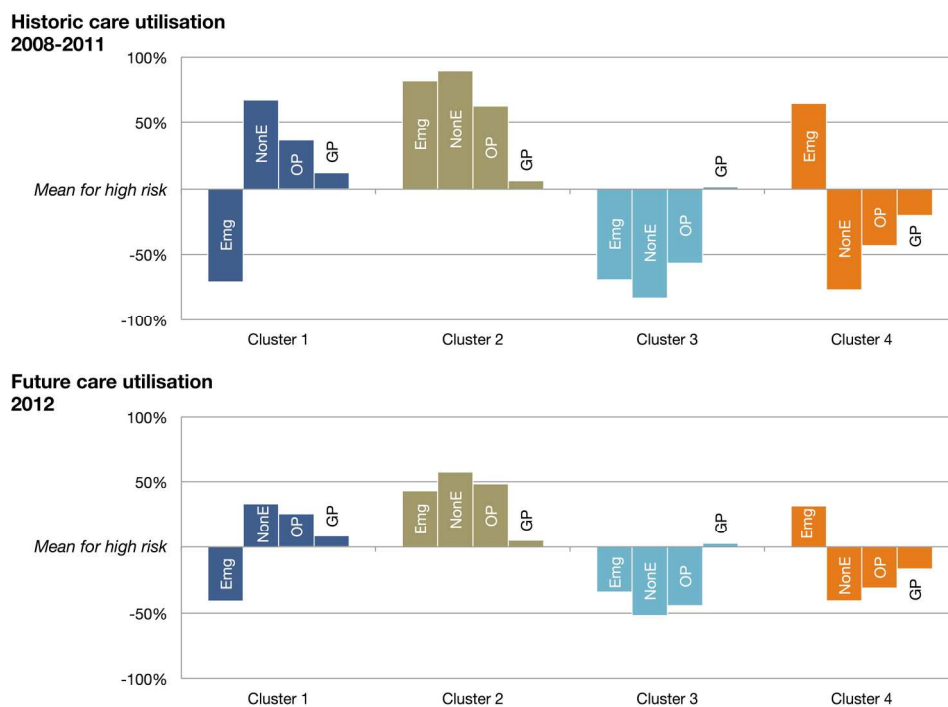
24. NHS Health and Social Care Information Centre. Numbers of Patients Registered at a GP Practice - July 2015. Leeds: NHS Health and Social Care Information Centre, 2015.
25. Bardsley M, Billings J, Dixon J, et al. Predicting who will use intensive social care: case finding tools based on linked health and social care data. *Age Ageing* 2011;**40**:265-70.
26. Kansagara D, Englander H, Salanitro A, et al. Risk prediction models for hospital readmission: A systematic review. *JAMA* 2011;**306**(15):1688-98.
27. Han J, Kamber M. *Data Mining: Concepts and Techniques*. 1st ed. San Diego, CA: Academic Press, 2001.
28. Armstrong JJ, Zhu M, Hirdes JP, et al. K-Means Cluster Analysis of Rehabilitation Service Users in the Home Health Care System of Ontario: Examining the Heterogeneity of a Complex Geriatric Population. *Arch Phys Med Rehabil* 2012;**93**(12):2198-205.
29. Coste J, Bouyer J, Fernandez H, et al. A population-based analytical approach to assessing patterns, determinants, and outcomes of health care with application to ectopic pregnancy. *Med Care* 2000;**38**(7):739-49.
30. Cryer PC, Saunders J, Jenkins LM, et al. Clusters within a general adult population of alcohol abstainers. *International Journal of Epidemiology* 2001;**30**(4):756-65.
31. Kendig H, Mealing N, Carr R, et al. Assessing patterns of home and community care service use and client profiles in Australia: a cluster analysis approach using linked data. *Health & Social Care in the Community* 2012;**20**(4):375-87.
32. Pud D, Ben Ami S, Cooper BA, et al. The Symptom Experience of Oncology Outpatients Has a Different Impact on Quality-of-Life Outcomes. *J Pain Symptom Manage* 2008;**35**(2):162-70.
33. Everitt BS, Landau S, Leese M, et al. *Cluster analysis*. 5th ed. Chichester: John Wiley & Sons, 2011.
34. StataCorp. *Stata 14 Cluster Stop reference manual*. College Station, TX: Stata Press, 2015.
35. Ng SK, Holden L, Sun J. Identifying comorbidity patterns of health conditions via cluster analysis of pairwise concordance statistics. *Stat Med* 2012;**31**(27):3393-405.
36. Chan M, Zhu MX. Investigating the health profile of Macau Chinese. *J Clin Nurs* 2008;**17**(11C):352-61.
37. Borglin G, Jakobsson U, Edberg Ak, et al. Older people in Sweden with various degrees of present quality of life: their health, social support, everyday activities and sense of coherence. *Health & Social Care in the Community* 2006;**14**(2):136-46.
38. Roland M, Abel G. Reducing emergency admissions: are we on the right track? *Br Med J* 2012;**345**(e6017).
39. Georghiou T, Blunt I, Steventon A, et al. Predictive risk and health care: An overview. London: Nuffield Trust, 2011.
40. Barnett AG, van der Pols JC, Dobson AJ. Regression to the mean: what it is and how to deal with it. *International Journal of Epidemiology* 2005;**34**(1):215-20.
41. Hao S, Jin B, Shin AY, et al. Risk prediction of emergency department revisit 30 days post discharge: a prospective study. *PloS one* 2014;**9**(11):e112944.
42. Pavlou M, Ambler G, Seaman SR, et al. How to develop a more accurate risk prediction model when there are few events. *Bmj* 2015;**351**.
43. Tanio C, Chen C. Innovations At Miami Practice Show Promise For Treating High-Risk Medicare Patients. *Health Aff (Millwood)* 2013;**32**(6):1078-82.
44. Alderwick H, Ham C, Buck D. *Population health systems: Going beyond integrated care*. London: The King's Fund, 2015.

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



Mean future care utilisation for the risk strata - High (H), Medium (M) and Low (L) - and the four high-risk clusters - 1, 2, 3 and 4.

Figure 1  
190x142mm (300 x 300 DPI)



Patterns of utilisation for the four high-risk clusters – Emergency care hospitalisations (Emg), Nonemergency hospitalisations (NonE), Outpatient attendances (OP) and GP visits (GP) versus the high-risk population mean

Figure 2  
190x142mm (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

## Appendix 1

Variables included in various risk scores and variables selected for our model

Variables considered in hospital admission risk studies	Number of studies out of 30 including variable in final model <sup>1</sup>	Variable in PARR <sup>2</sup>	Variable in PARR-30 <sup>3</sup>	Variable in Combined Predictive Model <sup>4</sup> (selected variables)	Included in initial model / included in final model after backwards elimination
<b>Morbidities</b>	Medical diagnoses or comorbidity indices: 24	Cerebrovascular disease			Any diagnosis of cerebrovascular disease in 2008-2011 (in primary or secondary care)
		Chronic obstructive pulmonary disease	Chronic pulmonary disease	COPD	<b>Any diagnosis of COPD in 2008-2011 (in primary or secondary care)</b>
				Asthma (only considered in LTC counts)	Any diagnosis of Asthma in 2008-2011 (in primary or secondary care)
		Connective tissue disease/rheumatoid arthritis			<b>Any diagnosis of Rheumatic disease in 2008-2011 (in primary or secondary care)</b>
		Developmental disability			Any diagnosis of Learning disability in 2008-2011 (in primary or secondary care)
		Diabetes	Diabetes with chronic complications	Diabetes (only considered in LTC counts)	Any diagnosis of Diabetes in 2008-2011 (in primary or secondary care)
		Ischaemic heart disease		CAD (only considered in LTC counts)	<b>Any diagnosis of Ischaemic heart disease in 2008-2011 (in primary or secondary care)</b>
		Peripheral vascular disease	Peripheral vascular disease		Any diagnosis of Peripheral vascular disease in 2008-2011 (in primary or secondary care)
	Renal failure	Renal disease		<b>Any diagnosis of Renal disease in 2008-2011 (in primary or secondary care)</b>	

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

		Sickle cell disease			
		Metastatic cancer with solid tumour	<i>Cancer (only considered in LTC counts)</i>		<b>Any diagnosis of Cancer in 2008-2011 (in primary or secondary care)</b>
		Other malignant cancer			
		Congestive heart failure	<i>CHF (only considered in LTC counts)</i>		Any diagnosis of Congestive heart failure in 2008-2011 (in primary or secondary care)
		Moderate/severe liver disease			Any diagnosis of Liver disease in 2008-2011 (in primary or secondary care)
		Other liver disease			
		Haemiplegia or paraplegia			Any diagnosis of Paraplegia in 2008-2011 (in primary or secondary care)
		Dementia			Any diagnosis of Dementia in 2008-2011 (in primary or secondary care)
			<i>Hypertension (only considered in LTC counts)</i>		
		Diagnostic cost groups/hierarchical condition category			
			1 LTC		<b>Flag if the sum of conditions listed is 0, 1 or 2 or more</b>
			2+ LTCS		
<b>Mental health morbidities</b>	Alcohol or substance use: 11	Alcohol related diagnosis		Psychoactive substance abuse	
	Mental illness: 9			Psychotic disorder	
				Inpatient admission with diagnosis of mental illness	<b>Any diagnosis of Mental health disorder in 2008-2011 (in primary or secondary care)</b>
				<i>Depression (only as included in LTC counts)</i>	
<b>Prior use of medical services</b>	Hospitalisations: 14	Previous admission for respiratory infection			
		Previous admission for a reference condition			
		Number of emergency admissions in previous 90, 180	Whether there had been a prior emergency hospital	[Combinations of] 1+, 2, 2+, 3+ emergency admissions in last 30,	

bmjopen-2016-012903 on December 20, 2016. Downloaded from http://bmjopen.bmj.com/ on April 20, 2014 by guest. Protected by copyright.

	and 365 days	discharge in the past 30 days	30 to 90, 90 to 180, 180 to 365, 365 to 730 days	<b>admissions in 2011 &amp; Number of emergency admissions over 2008-2011</b>
	Average number of episodes per spell for emergency admissions		Average number of episodes per spell for emergency admissions >=3	
		Whether the current admission was an emergency admission		
	Total number of previous emergency admissions in previous three years	Number of emergency hospital discharges in the last year		
	Number of non-emergency admissions in previous 365 days			<b>Number of non-emergency admissions over 2008-2011</b>
	Emergency department visits: 4		A&E visits and investigations	
	Clinic visits or missed visits: 3	Number of different treatment specialists seen		
			[Combinations of] 1-5, 2, 3+, 6-10, 11+ out-patient specialty visits in last 30, 30 to 90, 90 to 365 to 730 days	<b>Number of outpatient visits over 2008-2011</b>
	Index hospital length of stay: 4			
	Other		Polypharmacy: 1-4 unique drugs in any month (last 30 to 90 days); 5-9; 10+	<b>Number of GP visits over 2008-2011 (including home visits)</b>
<b>Sociodemographic factors</b>	Age: 19	Age 65-74 or age 75+	Age squared	<b>5-year age bands &amp; Over 75 flag</b>
	Sex: 15	Sex	Gender	<b>Gender</b>
	Race/ ethnicity: 7	Ethnicity		
<b>Social determinants of health</b>	SES, income and employment: 5		Index of multiple deprivation band for the place of residence	<b>Townsend score (5 groups)</b>
	Insurance status: 6			

bmjopen-2016-012903 on 19 December 2016. Downloaded from <http://bmjopen.bmj.com/> on April 20, 2024 by guest. Protected by copyright.

	Education: 0		
	Marital status and people in household: 4		
	Social support: 2		
	Access to care: 5		
	Discharge location: 2		
<b>Hospital specific metrics</b>	<i>Not included in review</i>	Observed:expected ratio for practice style sensitive admissions in ward of residence	
		Observed:expected ratio for rate of readmissions for hospitals of current admission	Hospital-specific variable
<b>Illness severity</b>	Severity index: 1		
	Laboratory findings: 4		
	Other: 4		
<b>Overall health and function</b>	Functional status, ADL: 2		
	Self-rated health, QOL: 3		
	Cognitive impairment: 7		
	Visual/hearing impairment: 1		

1. Kansagara D, Englander H, Salanitro A, et al. Risk prediction models for hospital readmission: A systematic review. JAMA 2011;**306**(15):1688-98.
2. Billings J, Dixon J, Mijanovich T, et al. Case finding for patients at risk of readmission to hospital: development of an algorithm to identify high risk patients. Br Med J 2006;**333**:327.
3. Billings J, Blunt I, Steventon A, et al. Development of a predictive model to identify inpatients at risk of re-admission within 30 days of discharge (PARR-30). BMJ Open 2012;**2**(4).
4. Wennberg D, Siegel M, Darin B, et al. Combined Predictive Model - Final report & technical documentation. London: Health Dialog, King's Fund and New York University, 2006.



## Supplementary file

### DECIDING ON THE NUMBER OF CLUSTERS - METHODOLOGICAL EXPLORATION

To segment the high-risk population, a k-means method was used. This method is efficient even for large sample sizes and produces roughly similar sized segments.<sup>1</sup> However, this method also require the number of clusters (k) to be specified before the analysis, rather than deducing it from the results afterwards. Therefore, a number of steps were taking to identify the optimal number of clusters for this population.

#### PSEUDO-F STATISTIC

The main method for determining the number of clusters was the Pseudo-F statistic.<sup>2</sup> This statistic is commonly used in healthcare clustering studies,<sup>3-7</sup> and has been identified as one of the best criteria to determine the number of clusters.<sup>8</sup> It compares the between-cluster to the within-cluster sum-of-squares, and a large Pseudo-F statistic indicates distinct clusters.<sup>9</sup>

The k-means analysis was run for 2 to 8 clusters, and the Pseudo-F statistic was calculated for each solution (see table 1). A peak could be observed around the 3- and 4-cluster solutions.

Table 1: Pseudo-F statistics for 2- to 8-cluster solutions

2 clusters	2249
3 clusters	2745
4 clusters	2662
5 clusters	2374
6 clusters	2267
7 clusters	2131
8 clusters	2041

#### WARD'S LINKAGE

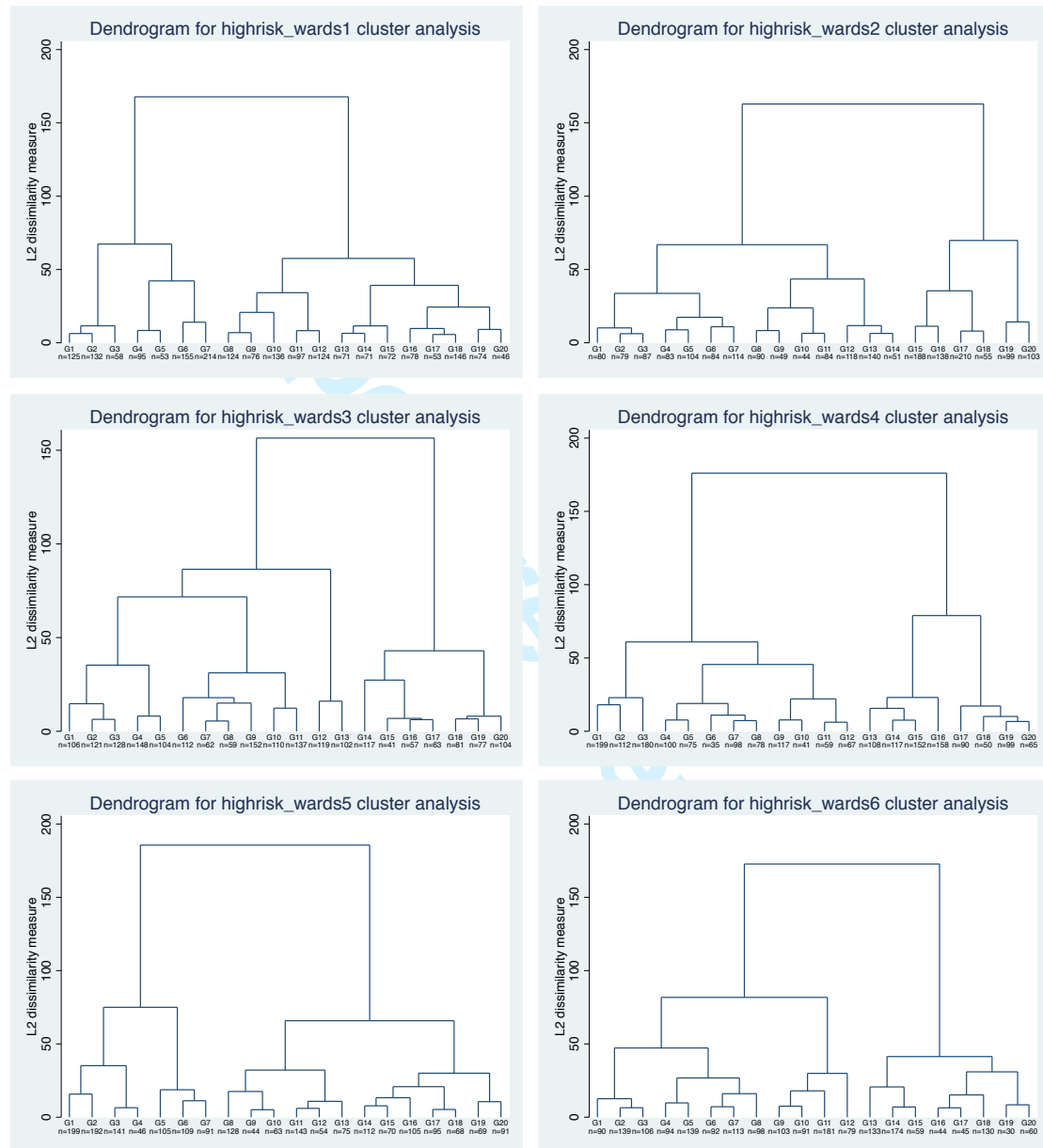
K-means is a non-hierarchical clustering method. Hierarchical methods, including the popular Ward's method, do not require k to be specified before the analysis.<sup>8</sup> Hierarchical clustering can be used to gain more insight into the data's structure. By displaying the results as a dendrogram (a tree-like plot detailing each hierarchical step in the model) different clustering solutions can be visually explored.<sup>10,11</sup> Indeed, many studies combine hierarchical clustering with k-means in a two-stepped approach.<sup>12-15</sup>

However, hierarchical methods present some limitations. The approach is computational intensive and struggles to handle large datasets with more than a thousand observations.<sup>10,11</sup> In addition, hierarchical clustering based on Ward's method can be sensitive to outliers.<sup>8</sup>

The high-risk population in the test sample, consisting of 7,433 people, was too large to include in its entirety in a hierarchical cluster analysis. Therefore, three unique, random samples of 2,000 people we used. After reshuffling the data, another three 2,000 people samples were taken and clustered. These results were then analysed through dendrograms

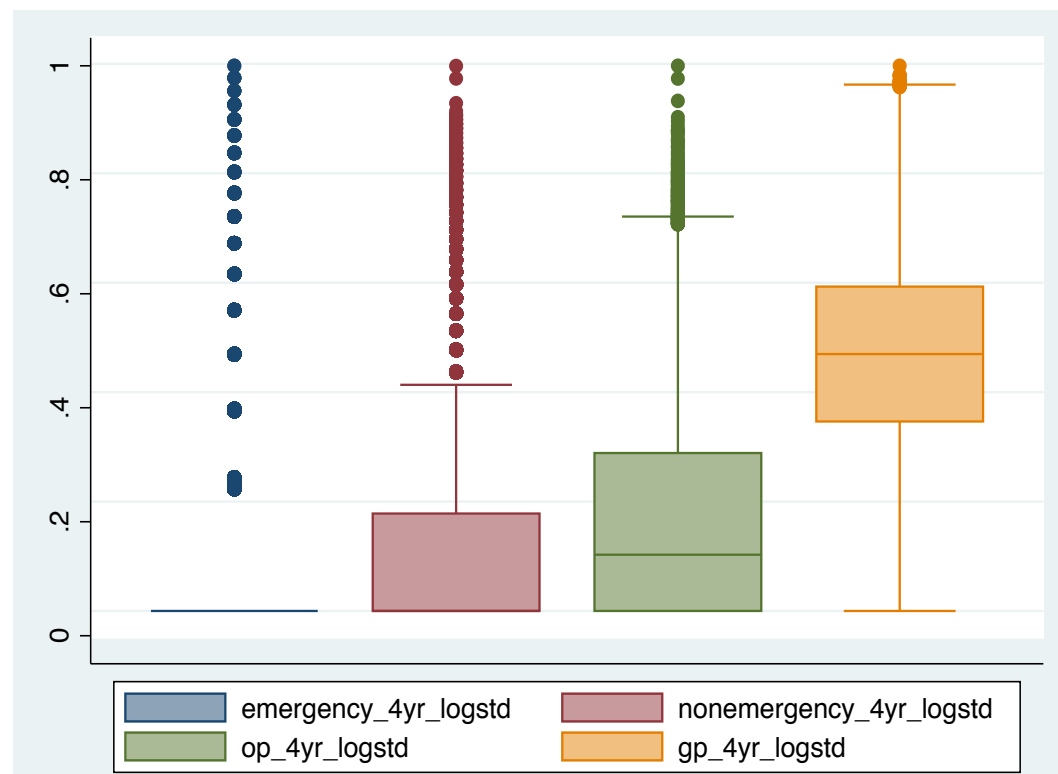
(see figure 1). All samples favoured a two-cluster solution, reflecting high- and low-utilisation groups, with the next split being further down the graph. The samples showed different results regarding the next best split. Sample 1,2, 3 and 5 can be interpreted as indicating the existence of four distinct clusters. Sample 4 favoured five clusters, and sample 6 could be interpreted as three or five clusters. Overall, the differences at this level are small.

Figure 1: Dendrograms for the six 2,000 people samples clustered using Ward's linkage



One of the reasons the results are different across the samples is the impact of outliers, which Ward's method is sensitive to.<sup>8</sup> Despite the log-normalisation of the clustering variables, there still exist a large number of outliers (see figure 2). Especially in the smaller samples used for the clustering, these outliers could have changed the resulting clusters.

Figure 2: Box plots of the standardised, log-normalised clustering variables



### POST-HOC ANALYSIS

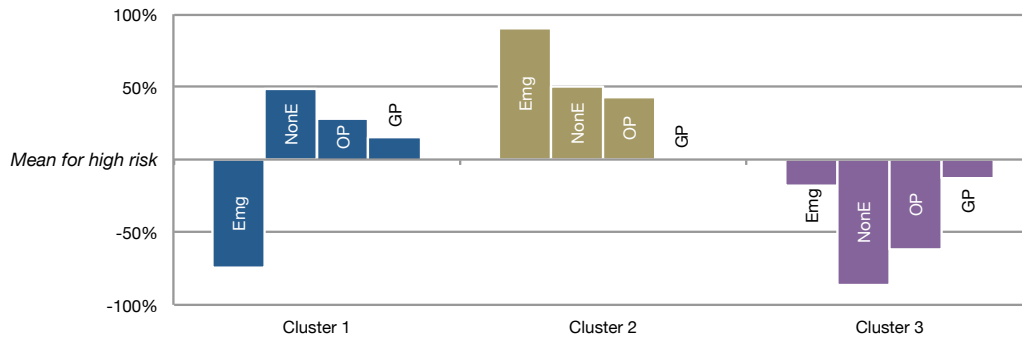
It is important to keep in mind that in cluster analysis, there is no absolute 'right' answer<sup>16</sup> - it all depends on the purpose of the clustering. Some other aspects to consider in evaluating the number of segments are, for example, interpretability, actionability and ease of use.<sup>11</sup>

The cluster means of the 3- and 4-cluster solutions were compared to review the practical usefulness of the resulting population groups (see figure 3). Both solutions found clusters of people with high utilisation but low emergency care use (clusters one), and people with overall high utilisations (clusters two). As the third group, the 3-cluster solution identified people with low overall utilisation but average emergency care use. However, the 4-cluster solution split this final cluster into two very distinct groups: people with overall low utilisation, and people with low utilisation but high emergency care use.

Figure 3: Practical comparison of the 3- and 4-cluster solutions

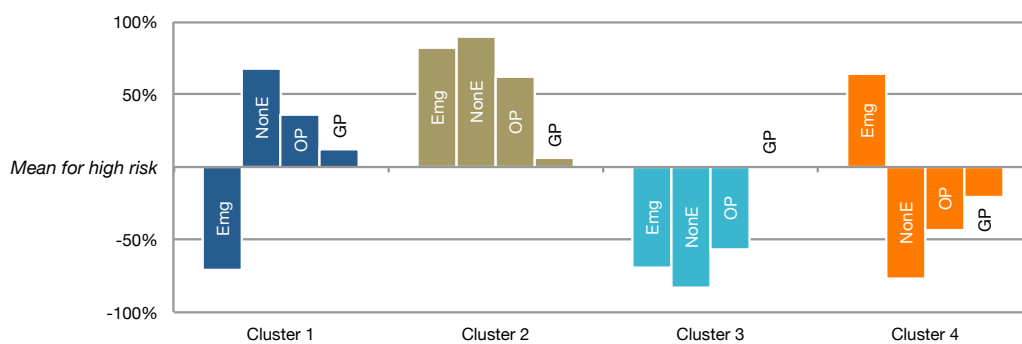
**3-cluster solution**

Care utilisation 2008-2011



**4-cluster solution**

Care utilisation 2008-2011



Considering the relevance of emergency care use for risk stratification, the difference between clusters three and four are important to the interpretability of the results. In terms of actionability, differentiating between these two groups allows tailored initiatives to be developed that target those people with low care utilisation but high emergency admissions. Taking this, and the previously described analyses into account, the 4-cluster solution was ultimately selected.

## References

1. Han J, Kamber M. *Data Mining: Concepts and Techniques*. 1st ed. San Diego, CA: Academic Press, 2001.
2. Stata Statistical Software: Release 14. [program]. College Station, TX: StataCorp LP, 2015.
3. Armstrong JJ, Zhu M, Hirdes JP, et al. K-Means Cluster Analysis of Rehabilitation Service Users in the Home Health Care System of Ontario: Examining the Heterogeneity of a Complex Geriatric Population. *Arch Phys Med Rehabil* 2012;**93**(12):2198-205.
4. Coste J, Bouyer J, Fernandez H, et al. A population-based analytical approach to assessing patterns, determinants, and outcomes of health care with application to ectopic pregnancy. *Med Care* 2000;**38**(7):739-49.
5. Cryer PC, Saunders J, Jenkins LM, et al. Clusters within a general adult population of alcohol abstainers. *International Journal of Epidemiology* 2001;**30**(4):756-65.
6. Kendig H, Mealing N, Carr R, et al. Assessing patterns of home and community care service use and client profiles in Australia: a cluster analysis approach using linked data. *Health & Social Care in the Community* 2012;**20**(4):375-87.
7. Pud D, Ben Ami S, Cooper BA, et al. The Symptom Experience of Oncology Outpatients Has a Different Impact on Quality-of-Life Outcomes. *J Pain Symptom Manage* 2008;**35**(2):162-70.
8. Everitt BS, Landau S, Leese M, et al. *Cluster analysis*. 5th ed. Chichester: John Wiley & Sons, 2011.
9. StataCorp. Stata 14 Cluster Stop reference manual. College Station, TX: Stata Press, 2015.
10. IBM Corp. IBM SPSS Statistics Base 22. Chicago, IL: IBM Software Group, 2013.
11. Tsipis K, Chorianopoulos A. *Data Mining Techniques in CRM: Inside Customer Segmentation*. Chichester: John Wiley & Sons, 2009.
12. Lega F, Mengoni A. Profiling the different needs and expectations of patients for population-based medicine: a case study using segmentation analysis. *BMC Health Serv Res* 2012;**12**(1):473.
13. Liu C-Y, Liu J-S. Mining the optimal clustering of people's characteristics of health care choices. *Expert Systems with Applications* 2011;**38**(3):1400-04.
14. Nagel GC, Schmidt S, Strauss BM, et al. Quality of life in breast cancer patients: a cluster analytic approach. *Breast Cancer Res Treat* 2001;**68**(1):75-87.
15. Leijon O, Härenstam A, Waldenström K, et al. Target groups for prevention of neck/shoulder and low back disorders: an exploratory cluster analysis of working and living conditions. *Work* 2006;**27**(2):189.
16. Han J, Kamber M, Pei J. *Data Mining: Concepts and Techniques*. 3rd ed. Waltham, MA: Morgan Kaufmann, 2011.