



THE UNIVERSITY *of* EDINBURGH

Edinburgh Research Explorer

Linkage of primary care prescribing records and pharmacy dispensing Records in the Salford Lung Study

Citation for published version:

Tibble, H, Lay-flurrie, J, Sheikh, A, Horne, R, Mizani, MA & Tsanas, A 2020, 'Linkage of primary care prescribing records and pharmacy dispensing Records in the Salford Lung Study: application in asthma', *BMC Medical Research Methodology*, vol. 20, no. 1, 303. <https://doi.org/10.1186/s12874-020-01184-8>

Digital Object Identifier (DOI):

[10.1186/s12874-020-01184-8](https://doi.org/10.1186/s12874-020-01184-8)

Link:

[Link to publication record in Edinburgh Research Explorer](#)

Document Version:

Peer reviewed version

Published In:

BMC Medical Research Methodology

General rights

Copyright for the publications made accessible via the Edinburgh Research Explorer is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.

Take down policy

The University of Edinburgh has made every reasonable effort to ensure that Edinburgh Research Explorer content complies with UK legislation. If you believe that the public display of this file breaches copyright please contact openaccess@ed.ac.uk providing details, and we will remove access to the work immediately and investigate your claim.



1 **Linkage of Primary Care Prescribing Records and**
2 **Pharmacy Dispensing Records in the Salford Lung Study:**
3 **Application in Asthma**

4
5 **Authors:**

6 Holly Tibble ^{1,2}, James Lay-Flurrie ³, Aziz Sheikh ^{1,2,4}, Rob Horne ^{2,5}, Mehrdad A. Mizani ^{1,2},
7 Athanasios Tsanas ^{1,2} & The Salford Lung Study Team ³

8
9 **Affiliations:**

- 10 1. Usher Institute, University of Edinburgh
11 2. Asthma UK Centre for Applied Research
12 3. GlaxoSmithKline UK Ltd
13 4. Health Data Research UK
14 5. Centre for Behavioural Medicine, UCL School of Pharmacy

15
16 **Corresponding Author:**

17 Holly Tibble

18 +44 7449 053 411

19 Holly.tibble@ed.ac.uk

20 Bioquarter 9, 9 Little France Road, Edinburgh, Scotland, EH16 4UX

24 **Abstract (337 words/350 words):**

25

26 **Background:**

27 Records of medication prescriptions can be used in conjunction with pharmacy dispensing
28 records to investigate the incidence of *adherence*, which is defined as observing the treatment
29 plans agreed between a patient and their clinician. Using prescribing records alone fails to
30 identify primary non-adherence; medications not being collected from the dispensary. Using
31 dispensing records alone means that cases of conditions that resolve and/or treatments that are
32 discontinued will be unaccounted for. While using a linked prescribing and dispensing dataset
33 to measure medication non-adherence is optimal, this linkage is not routinely conducted.
34 Furthermore, without a unique common event identifier, linkage between these two datasets is
35 not straightforward.

36

37 **Methods:**

38 We undertook a secondary analysis of the Salford Lung Study dataset. A novel probabilistic
39 record linkage methodology was developed matching asthma medication pharmacy dispensing
40 records and primary care prescribing records, using semantic (meaning) and syntactic
41 (structure) harmonization, domain knowledge integration, and natural language feature
42 extraction. Cox survival analysis was conducted to assess factors associated with the time to
43 medication dispensing after the prescription was written. Finally, we used a simplified record
44 linkage algorithm in which only identical records were matched, for a naïve benchmarking to
45 compare against the results of our proposed methodology.

46

47 **Results:**

48 We matched 83% of pharmacy dispensing records to primary care prescribing records. Missing
49 data were prevalent in the dispensing records which were not matched – approximately 60%
50 for both medication strength and quantity. A naïve benchmarking approach, requiring perfect
51 matching, identified one-quarter as many matching prescribing records as our methodology.
52 Factors associated with delay (or failure) to collect the prescribed medication from a pharmacy
53 included season, quantity of medication prescribed, previous dispensing history and class of
54 medication. Our findings indicate that over 30% of prescriptions issued were not collected
55 from a dispensary (primary non-adherence).

56

57 **Conclusions:**

58 We have developed a probabilistic record linkage methodology matching a large percentage of
59 pharmacy dispensing records with primary care prescribing records for asthma medications.
60 This will allow researchers to link datasets in order to more accurately extract information
61 about asthma medication non-adherence.

62

63

64

65

66 **Background**

67

68 Medication data can be used in research to assess changes in medication prescribing trends
69 over time (1), for pharmacovigilance studies, and to investigate patients not adhering to the
70 treatment plans agreed upon with their General Practitioner (GP) (2–4). Investigating
71 medication data enables researchers to estimate the frequency, burden, and costs of non-
72 adherence (5–7), identify the most at-risk to suboptimal clinical outcomes, evaluate the
73 effectiveness of adherence interventions (8–10), and appropriately adjust for the impact of non-
74 adherence on safety and efficacy data in clinical trials (11,12).

75

76 In studies of linked (or integrated) prescribing and dispensing records, failure to collect the
77 initial asthma prescription (*primary non-adherence*) has reported incidence between 12-45%
78 (13–17), with high variance due to differences in the right censoring point. Studies across
79 multiple chronic conditions reported a pooled general primary non-adherence rate of 9-17%
80 (18–20).

81

82 In England, prescribing and dispensing of medications are recorded by separate processes.
83 After a medication prescription is issued to a patient by a GP or another authorized prescriber
84 (21), the prescription is taken to a dispensing outlet such as a community pharmacy (22). When
85 the prepared medicine is released to the patient, details relating to payment for medications are
86 recorded and managed by the NHS Business Services Authority (NHSBSA). While analysis
87 of medication adherence can be estimated using either the GP’s prescribing records or the
88 NHSBSA medication dispensing records alone, there are limitations to each approach. Without
89 linking the records together, it is not possible to ascertain whether a prescribed medication was

90 collected, or to rule out other reasons for irregularities in collection such as treatment
91 conclusion or sanctioned treatment interruptions (1,23,24).

92

93 Since 2015, NHSBSA dispensing data have included a patient identifier (NHS number) (25);
94 this is, however, not routinely linked to primary care prescribing records held by Public Health
95 England (PHE). The NHSBSA and PHE records also do not have a common unique prescribing
96 event identifier. Therefore, even with a data sharing agreement in place, matching records (one-
97 to-one) using common identifiers (known as *deterministic linkage*) is currently impossible.

98

99 Therefore, it is necessary to link records probabilistically; estimating the likelihood that two
100 records will match given the data they contain. Neither pharmacy nor primary care records are
101 written with future linkage in mind, and as such they often require substantial pre-processing.
102 The quality of the data linkage can be improved by integrating domain knowledge to identify
103 non-matching but equivalent values, for example converting between units of dose strength.

104

105 The distinction between what should be considered deterministic or probabilistic is often
106 disputed, as even complex probabilistic linkage processes can be broken down into their rule-
107 based components and both linkage types can allow for imperfect (or *fuzzy*) matching on certain
108 features (26), such as the dates of events in our case (which we would not expect to match all
109 the time). The nature of administrative data source linkage, such as with Electronic Health
110 Records, necessitates the use of fuzzy matching to overcome such prevalent qualities as
111 missing data, free-text values, non-standardised units, and generic medication substitutions
112 (resulting in different medication names). There are cases in which deterministic linkage will
113 not only reduce the overall accuracy of the linkage, but may also introduce bias (27,28).

114

115 Padmanabhan *et al.* have previously demonstrated the methodology used for linking UK health
116 datasets when the unique patient identifier (NHS number) contained missing and erroneous
117 values prohibiting deterministic linkage, including the creation of a ranking system for
118 candidate links based on the matching information between them (29).

119

120 **Methods**

121

122 *Aim:*

123 The linkage of prescribing and dispensing records can enable the extraction of information
124 about adherence to prescribed medications, including the identification of uncollected
125 medications. In this study, we sought to develop a novel methodology linking primary care
126 prescribing and dispensing records without a common identifier, using heuristics and features
127 extracted from free-text fields.

128

129 The GUILD (30) and RECORD (31) guidelines for data linkage reporting were applied where
130 necessary information was not reported elsewhere (32–34)).

131

132 *Data Source:*

133 The Salford Lung Study (SLS) was a prospective, 12-month, open-label, parallel group,
134 randomised controlled trial (RCT) conducted in 74 general practice clinics in Salford and South
135 Manchester, UK (35). A total of 4,233 participants with asthma were recruited in primary care
136 settings by the healthcare professionals who provided their normal everyday care, and
137 randomly allocated to either initiate a combination fluticasone furoate/vilanterol treatment or
138 to continue their maintenance therapy (“usual care”).

139

140 Participants were at least 18 years old at the time of recruitment, with a clinical diagnosis of
141 symptomatic asthma made by a GP and had to be taking regular maintenance inhaler therapy
142 with Inhaled CorticoSteroids (ICS) either alone or in combination with a Long-Acting β_2 -
143 Agonist (LABA). The main exclusion criteria were a recent history of life-threatening asthma,
144 a history of Chronic Obstructive Pulmonary Disease (COPD), or concomitant life-threatening
145 disease (34,36). Many of the participants in the study cohort would have been excluded from
146 conventional RCTs due to their multi-morbidities (33,36), which increased the
147 representativeness of the study cohort to the target population.

148

149 The trial was registered in the National Institute of Health's database of clinical studies (32)
150 (clinicaltrials.gov identifier NCT01706198). The study was conducted in accordance with the
151 standards dictated by the National Research Ethics Service Committee North West (reference
152 12/NW/0455), as well as the International Conference on Harmonisation, Good Clinical
153 Practice, all applicable data protection requirements and the ethical principles outlined in the
154 Declaration of Helsinki 2013.

155

156 *Data Format:*

157 The dispensing data contained 225,235 records, for 4,197 unique participants, between 27th
158 November 2012 and 9th December 2016. The prescribing dataset contained 339,792 records
159 for 4,233 unique participants between 22nd November 2012 and 17th January 2017, however
160 records outside of the dispensing data period were excluded.

161

162 Both datasets contained a (common) subject ID, free text drug description, date (prescription
163 or dispensing, respectively), the dose strength, dose instructions, and a numeric quantity of

164 medication prescribed (e.g. “200 dose inhaler”). Between the two datasets, there were 8,291
165 unique (*free text*) drug descriptions.

166

167 *Inclusion and Exclusion Criteria:*

168 All unique drug descriptions, in either the prescribing or dispensing records, were searched for
169 the presence of one or more of the keywords listed in Appendix A. From here, the drug classes
170 were assigned: Short-Acting β_2 -Agonist (SABA), Long-Acting Muscarinic receptor
171 Antagonist (LAMA), LABA, theophylline, ICS, LeukoTriene Receptor Antagonist (LTRA),
172 cromoglicate, steroid, or immuno-suppressant. If only one candidate class was identified, the
173 drug class was coded according to the drug class keyword. A drug was coded as an ICS and
174 LABA combination medication (ICS+LABA) if active ingredients of both ICS and LABA
175 varieties were flagged, a SABA if a medicine containing both SABA and LAMA ingredients
176 were flagged. Medications that did not match any of the keywords in Appendix A were
177 considered to be non-asthma medications and were removed. A medication class keyword was
178 generated, containing a composite of the active ingredients, to be used in the matching
179 algorithm.

180

181 Furthermore, drug descriptions were searched for any of the exclusion keywords and brand
182 names listed in Appendix B, which signalled that a medication was being used for an indication
183 other than asthma (such as nasal spray corticosteroids for rhinitis).

184

185 *Variable Recoding:*

186 Several free text variables were recoded using custom look-up tables, to allow semantically
187 identical, but syntactically variant (such as “128mcg” vs “128 micrograms”, and other type
188 abbreviations and variations) records to be aligned. Of note, we modified the recorded

189 medication quantity to estimate the number of doses (puffs), rather than the number of units
190 (inhalers). This variable integrates domain knowledge of the number of doses per unit for each
191 medication strength combination (high potency medications are often dispensed at lower
192 volumes), calculated using the most common volumes in the data. In order to avoid candidate
193 links being ruled out as potential matches on the basis of our quantity variable modifications,
194 we included a so called ‘alias’ quantity (27), to be considered if the ‘primary’ quantity values
195 did not match. The process is summarized in Appendix C.

196

197 *Identification of Duplicates:*

198 Duplicates of prescribing and dispensing records are common due to errors in data entry (37–
199 39). Duplicate records in the data would have a strong adverse effect on the matching
200 algorithm, as it would be forced to incorrectly match distinct records in one set to duplicates in
201 the other. We identified duplicate records by searching for commonalities within the same
202 person, date (dispensing or prescribing respectively), medication brand name, and medication
203 (active ingredient) keyword, in addition to the following combinations of (modified) variables:

- 204 • Matched on quantity and dose
- 205 • Matched on dose, and the quantity was not matched due to data missingness
- 206 • Matched on quantity, and the dose was not matched due to data missingness.

207

208 *Data Linkage:*

209 The datasets of prescribing and dispensing records were merged such that a record (a *candidate*
210 *link*) was generated for each eligible (common patient identifier and medication class) pair of
211 records for matching. We note that the medication class keyword, composed of the active
212 ingredients identified, was used in the place of a brand name such that generic substitutions
213 would be identified as appropriate candidates for matching records. Pairs of records were

214 eligible if the suggested dispensing date occurred after the prescription was written, but no
215 more than six months *after* the prescription was written, at which point the prescription became
216 invalid.

217

218 Probabilistic linkage, which aims to match records based on multiple non-unique features,
219 utilizes *weights* to determine the strength of a link. These weights are numerical values
220 representing the similarity of two records, derived using domain knowledge about the
221 prevalence of dissimilarities between features in true matches.

222

223 In this linkage, a rule-based approach, based on a simplified posterior multivariate distribution
224 of clerically reviewed data and previous literature, was used to weight candidate links for
225 estimated likelihood of being a true match. Candidate links could then be ranked, and those
226 with a linkage weight (calculation detailed in Appendix D) less than 70% excluded
227 (combinations of features by match status that resulted in inclusion are listed, along with their
228 sum weights, in Appendix E).

229

230 Generic substitution for brand named medications are common (when permitted by the
231 prescriber, known as *open generic prescribing*) in asthma controller medications (15,40,41).
232 As such, brand name was assigned a lower maximum feature weight (20%) than the dose
233 strength (35%; which will vary only when one record has a missing value, or in the rare case
234 that a generic substitution requires a slightly different dosage) and quantity (35%; varying
235 when a quantity was both uncommon and missing, and was imputed with a more prevalent but
236 incorrect value). The final 10% weight corresponded to the time between the prescribing and
237 dispensing events. Prescriptions issued less than one month prior to the dispensing were
238 awarded the additional 10% weight, in line with the findings by Williams *et al.* that 95% of

239 asthma prescriptions are filled within this time window (14), however a higher weight was not
240 implemented due to the use of the time between weights in the final match selection process.
241 That is, each set of dispensing records for each person-medication combination were looped
242 through from the last to first through, as follows:

- 243 1. Identified the candidate in which the dispensing record occurs most recently after the
244 prescription was written (record with highest match weight chosen if two candidate
245 links on the same day were identified); this is a match between records,
- 246 2. Removed all other candidate links which contain the dispensing record or the
247 prescribing records relating to this match,
- 248 3. Progressed to the previous dispensing for this person-medication.

249 This process, illustrated in Figure 1, is also described in more detail in Appendix F.

250

251 [[Insert Figure 1 here]]

252

253 The most recent prescribing record before the dispensing was prioritised over more distant
254 records with a higher match weight, as we considered it more likely that prescription records
255 for the same person within such a short time window were for the same medication, recorded
256 differently, rather than a new treatment.

257

258 Prescriptions that did not match any dispensing record were marked as unclaimed. We also
259 noted dispensing records that were not matched (implying no corresponding prescription event)
260 to assess linkage quality.

261

262 *Statistical Analysis Plan*

263 As per the recommendations by Harron *et al.*, the characteristics of the matched and unmatched
264 records were compared in order to identify potential sources of bias (42). Specifically, the
265 missingness for each variable used in the matching was compared between matched and non-
266 matched records, factors associated with prescription collection were assessed (statistical
267 methodology described below), and the sensitivity of the algorithm parameters was tested by
268 altering certain thresholds and requirements and comparing the proportion of records that were
269 matched.

270

271 As well as estimating the incidence of primary non-adherence, we used our linked dataset to
272 analyse factors effecting the collection of prescribed medications. By comparing our results
273 to others using integrated health records (those that are linked, or linkable, inherently) we are
274 able to demonstrate the validity of our linked dataset to answer epidemiological questions about
275 high-risk individuals.

276

277 We used multivariate Cox survival analysis to assess the statistical relationship between the
278 season of the prescription, the drug class of the prescription, the number of previously
279 unclaimed prescriptions, and the strength and quantity of the medication prescribed, on the
280 time between the prescription being written and dispensed. Survival analysis calculates the
281 rates (*hazard rates*) of medications being collected at any specific time since the prescription
282 was written. Comparing the ratios (*hazard ratios*) between two levels of a factor (such as male
283 and female) allowed us to assess the difference that this factor made when everything else (age,
284 medication, etc.) remained constant. Although a prescription could be dispensed up to six
285 months after it was written, it is uncommon that their collection will be delayed for more than
286 7 days (14,15). Furthermore, a delay of beyond one month would likely result in a gap in
287 medication availability and thus be considered poor adherence. As such, we wanted to find a

288 threshold at which prescriptions could be recorded as ‘hitherto uncollected’, known as being
289 right censored. We set this threshold at the minimum number of weeks such that fewer than
290 2% of subsequently collected prescriptions would be right censored.

291

292 *Naïve Benchmarking*

293 We compared our results to those produced from a simplified algorithm in which records were
294 pseudo-deterministically matched, such that candidate links required perfect matching on
295 medication name, dose, quantity, and dose directions, without any variable recoding or removal
296 of duplicate records. The date variable, however, still allowed flexible matching as
297 medications can be dispensed up to six months following prescription.

298

299 The same iterative linkage procedure was used in the algorithm detailed previously, without
300 the inclusion of the linkage weights as a tiebreaker between candidate links on the same day.

301

302 As the dose directions were long, free-text strings, written separately by both the prescribing
303 and dispensing agents, we also repeated the benchmarking analysis, with imperfect matching
304 on the dose directions permitted.

305

306 Links identified by this process should not be considered the ground truth, or the gold
307 standard, as the algorithm will default to match records which are more distanced
308 chronologically but similar syntactically, rather than semantically similar and chronologically
309 closer record matches which are more likely to be estimated by the full algorithm. As such,
310 the matches identified between approaches will not be directly compared.

311

312 *Reporting*

313 This study has been reported in accordance with the GUILD and RECORD reporting
314 guidelines (30,31).

315

316 **Results**

317 *Data Cleaning*

318 Of the 8,291 unique drug descriptions, 928 (11%) were identified as relating to asthma
319 medications (list of keywords used in string search provided in Appendix A). Searching the
320 drug descriptions for the set of exclusion keywords led to the removal of 71 (8%) further
321 records (list and frequency of keywords in Appendix B). Removing the excluded medications
322 left 88,916 prescribing records and 64,471 dispensing records (Figure 2). Finally, duplicates
323 were removed (12,236 prescribing records and 406 dispensing records), leaving 76,680
324 prescribing records (86%) and 64,065 dispensing records (99%).

325

326 [[Insert Figure 2 here]]

327

328

329 *Matching*

330 The full join on the prescribing and dispensing records generated 265,442 candidate links for
331 linkage weight assessment (Appendix D). 62,783 candidate links were removed (23.7%) as
332 they did not fulfil the minimum linkage weight threshold, leaving 202,659 candidates to be
333 sorted through the matching algorithm. After the algorithm was applied, 53,289 candidate
334 links were confirmed as matches: 69.5% of prescribing records (n=76,680), and 83.2% of
335 dispensing records (n=64,065).

336

337 As shown in Figure 1: Diagram representing the data linkage algorithm.

338 *Figure 2: Data Linkage Flow Diagram.*

339 Figure 3, there was a substantial discrepancy between the time between the prescribing and
340 dispensing for the candidate links and the matches, with 99% of matches having less than one
341 month between prescribing and dispensing (compared to 33% of candidate links).

342

343 [[Insert Figure 3 here]]

344

345 The median percentage of prescriptions claimed by an individual was 79%, with an
346 interquartile range of 50-92% (range 0-100%). 23% of individuals claimed fewer than 50%
347 of their prescriptions.

348

349 *Quality Assurance*

350 We inspected 23,391 prescribing records (31%) and 10,776 dispensing records (17%) for
351 which a match could not be made (including those with candidate links which were not matched
352 by the matching algorithm). In the non-matched prescriptions, 9% (n=2,109/23,391) had
353 missing medication dosage, and <1% (n=87/23,391) had missing data on quantity (both
354 missing in less than <0.1%). In the non-matched *dispensing* records, however, it was 62%
355 (n=6,639/10,776) and 58% (n=6,222/10,776), respectively (both missing in 55%).

356

357 *Survival Analysis*

358 31% of prescriptions (n=23,391) were labelled as unclaimed. In claimed prescriptions
359 (n=53,289), the median time between the prescription being written and the medication being
360 dispensed was 1 day (upper-lower inter-quartiles = 0-3 days), and fewer than 5% of people
361 took longer than 1 week to claim (0.9% longer than 30 days). Considering uncollected
362 prescriptions to be right-censored at 6-months, at which point the prescription expires, the
363 median time to collection was 3 days (upper-lower inter-quartiles = 0-178 days; Figure 4).

364

365

[[Insert Figure 4 here]]

366

367 The multivariate Cox survival analysis model included 76,584 prescription records – having
368 removed 96 with missing quantity. The prescriptions were claimed in 52,186 of these records,
369 with less than 2% being collected beyond 3 weeks after the prescription was issued. As such,
370 21 days was set as our right censoring point. We found a lower hazard of claiming medications
371 in summer (June-August: 3% decrease, 95% CI = 1-6%) compared to spring (Table 1),
372 indicating that they were claimed slower in summer than in spring. There was no statistically
373 significant difference in the claiming of medications between spring and winter or spring and
374 autumn. Higher quantities (by number of doses) of prescribed medications were associated
375 with modest reduction in hazard of collecting the medication ($p < 0.001$). Finally, proportions
376 of previous prescriptions that were unclaimed (categorized into tertiles) were a strong predictor
377 – with medium vs low tertiles hazard ratio of 0.57, and high vs low of 0.20 ($p < 0.001$). Rescue
378 medication (SABA and steroids) had the highest hazard rates (1.433 and 1.839, respectively).
379 Of the controller medications, those associated with higher asthma severity (according to the
380 British Thoracic Society (BTS) treatment steps (43)), such as LAMA and LTRA medicines,
381 had higher hazards than lower severity treatments such as ICS and combination ICS+LABA
382 medications.

383

384

[[Insert Table 1 here]]

385

386 *Naïve Benchmarking*

387 There were 88,916 prescribing records and 64,471 dispensing records identified relating to an
388 asthma medication (without the removal of duplicates). Of these, 584 (0.7% of prescribing

389 records and 0.9% of dispensing records) were pseudo-deterministically linked. Even when
390 imperfect matching on dose-directions was permitted, only 15.4% of prescribing records and
391 21.2% of dispensing records could be matched (n=13,698 matches).

392

393

394 **Discussion**

395 We have developed a novel methodology matching prescribing and dispensing electronic
396 health records and demonstrated this led to matching 70% of asthma prescribing and 83% of
397 dispensing records. Fewer than 5% of prescriptions were eventually claimed after one week of
398 the issuing of the prescription. 30% of prescriptions were labelled as uncollected.

399

400 The key strength of this study is the variety of integrated mechanisms – incorporating domain
401 knowledge relating to asthma medications (such as semantic harmonization from brand name
402 to active ingredients) and rule-based natural language feature extraction and harmonization
403 (such as converting a free-text dose to a numeric value with common units).

404

405 Using a naïve benchmarking algorithm that required perfect matching between prescribing and
406 dispensing records (except for the date variable; pseudo-deterministic linkage), we were able
407 to demonstrate the superiority of our proposed methodology. In this benchmark linkage, only
408 15% of the prescribing records and 21% of dispensing records were matched, even when
409 imperfect matching on free-text dose directions was permitted. This was a result of
410 syntactically variant (different formats and value units) but semantically matching data
411 between the two sources of information.

412

413 We identified a set of records for dispensed medications (17%) for which no matching
414 prescribing record was identified. In the non-matched dispensing records, 62% had missing
415 medication strength, and 58% had missing quantity. In its current state, the algorithm will not
416 match records with high amounts of missing data even if no other match is identified.

417

418 In Appendix D, we see that 3% of matches had distinct and non-missing medication brand
419 names. This highlights that potentially brand substitutions occurring at the pharmacy need to
420 be accounted for in the matching (44). The variable with the biggest change in distribution
421 between the candidate links and the final matches was whether the medication was dispensed
422 within one month of prescribing – 33% of candidates and 99% of matches (see Figure 1:
423 Diagram representing the data linkage algorithm.

424 *Figure 2: Data Linkage Flow Diagram.*

425 Figure 3). In fact, we found that only 1% of prescriptions were claimed more than a month
426 after the prescription was written.

427

428 Our finding that 30% of prescriptions were labelled as uncollected, known as primary non-
429 adherence, was a substantially higher proportion than the 8-20% found in previous asthma
430 studies in US administrative health data studies (13–15,41,45). One might assume that
431 subsidised prescriptions, as we have in England, would result in higher primary adherence
432 rates, as a barrier to adherence has been removed. On the contrary, a recent study in Canada,
433 where prescriptions are subsidised and thus considerably more affordable than in the USA,
434 found that the fill rate for new asthma prescriptions was only 69% in adults (16). As such,
435 future work must be conducted in order to find cost-effective interventions to reduce primary
436 non-adherence in asthma.

437

438 As there is no true linkage event identifier (person-prescription), it is not possible to compare
439 our identified matches to some ground truth, a common limitation highlighted in the
440 aforementioned linkage quality assessment guidelines by Harron *et al.* (42). As the

441 benchmarking analysis allowed prescribing and dispensing date variables to differ, hence
442 pseudo-deterministic, even this does not identify ‘perfect matches’ between records. If the
443 ground truth was known, it would be possible to compare directly the matches estimated from
444 the benchmark and pseudo-deterministic analyses and evaluate how well our algorithm
445 improves the matching quality. While the ground truth may not be possible to determine in
446 challenging real-world data, even with manual review, one could also perturb data in which the
447 ground truth is known to closer approximate the real use case, and evaluate the algorithm’s
448 accuracy.

449

450 In lieu of this, we conducted quality assurance comparing features of the matched and
451 unmatched records, as recommended by Harron *et al.*’s guidelines (42). We observed that
452 prescriptions (for which the status of being non-matched might imply either medication non-
453 initiation, or not being correctly matched using the proposed algorithm) had missed medication
454 strength in fewer than 10% of records, and missing quantity in fewer than 1%. In the non-
455 matched dispensing records (which should occur only in rare emergency prescriptions and
456 indicate shortcomings in matching prescription and dispensing records), 62% had missing
457 strength and 58% had missing quantity. This indicates that one of the biggest barriers to
458 successful record linkage was poor medication dispensing record quality.

459

460 The frequency of non-matched dispensing records was our best indicator as to the quality of
461 our linkage, however we found that 95% of these records that were missing quantity (58%)
462 were also missing dose-strength. As such, reducing the weight threshold from 70% to 50%,
463 would have had a substantial effect on the pool of candidate links allowed to be used in the
464 matching algorithm. With so much missing data, however, the veracity of these matches would
465 be hard to ascertain.

466

467 The strong influence of data quality on the success of the linkage algorithm makes it difficult
468 to benchmark our results against other record linkage algorithms or even treatment initiation
469 studies in populations with linkage conducted routinely. Comparisons to algorithms derived
470 in other medication indications, such as in acute conditions such as tuberculosis, or in other
471 chronic illnesses such as mental health conditions, are even harder. Furthermore, not all
472 countries have a unique patient identifier, resulting in the use of demographic data such as
473 gender, year of birth, and postcode, to identify entries belonging to the same person (46).
474 Regardless, we find other studies have reported similar levels of inconsistency between
475 features in matched records, such as brand name, dose strength, and time between prescribing
476 and dispensing (44,47). We also observed the substantial increase in matches when variables
477 were cleaned, and recoded, and our probabilistic methodology was used in the place of a simple
478 pseudo-deterministic matching.

479

480 As with all probabilistic matching approaches, and particularly in cases such as these with
481 considerable number of missing entries and un-structured fields, it is possible that matches
482 even with high assigned weights are incorrect. Indeed, it is not likely that the matches
483 established in the benchmarking analysis are of higher accuracy than those in the primary
484 analysis, and they cannot be directly compared. In future work, this algorithm should be tested
485 in simulated data where the underlying ground truth is known for further validation, in order
486 to better determine the accuracy of the linkage. There is potential that the design of the study
487 on which this secondary analysis was conducted (a pragmatic randomised controlled trial) may
488 have influenced the linkage in some way. Validating the proposed linkage algorithm in further
489 additional randomised clinical trials would be needed to establish the generalizability of our
490 findings.

491

492 In addition to testing in other datasets, in which the true links are known and can be compared
493 to the estimated matches, further development of this study would be to test the sensitivity of
494 the model to certain parameters such as the weights for each component, the degree of influence
495 from the dates, and the minimum weight threshold. We remark that these intrinsic parameters
496 can be seen as degrees of freedom that enable data modellers to explore different levels of
497 certainty for record matching. At a higher level, these can be thought of as the equivalent free
498 parameters which need to be explored and optimised for a given dataset: for example, in
499 Support Vector Machines (SVM) one needs to optimise the penalty hyper-parameter (and
500 depending on configuration additional hyper-parameters too). Consideration must also be
501 taken to determine the acceptable limits of the false negative and positive rates, and the relative
502 importance of the two, in specific settings. For example, in adherence studies, one might
503 conservatively prefer to underestimate adherence than to overestimate it, and thus prioritise
504 lowering the false positive rate.

505

506 Additionally, accounting for how much medication supply an individual currently has, or when
507 their most recent previous prescription was issued, would allow the date component of the
508 algorithm to correspond more meaningfully to the patient's history. As previously discussed,
509 matching may also be improved by the addition of an extension allowing candidate pairs for
510 which one record had high amounts of missing data and no match was identified to be re-
511 considered.

512

513 **Conclusions**

514 The optimal dataset for measurement of medication non-adherence includes both prescribing
515 records and dispensing records, such that prescriptions that are not collected from the

516 dispensing agent and resolved/discontinued treatment regimens are accounted for. These are
517 however seldom available. We therefore developed a novel methodology that matched 83%
518 of pharmacy dispensing records to primary care prescribing records. In the 17% of dispensing
519 records for which a match could not be identified, missing information was prevalent;
520 particularly regarding the strength of the medication, and the quantity dispensed. A naïve
521 benchmarking, requiring perfect matching, identified prescribing records for only 21% of the
522 dispensing records. Although further evaluation of the quality of the data linkage is required,
523 our novel methodology enables preliminary assessment of whether patients are collecting their
524 prescribed asthma medications and can improve clinicians' understanding of patient adherence.

525 **Abbreviations**

BTS	British Thoracic Society
COPD	Chronic Obstructive Pulmonary Disease
GP	General Practitioner
ICS	Inhaled Cortico-Steroids
LABA	Long-Acting B ₂ -Agonist
LAMA	Long-Acting Muscarinic Receptor Antagonist
LTRA	Leukotriene Receptor Antagonist
NHSBSA	National Health Service Business Services Authority
PHE	Public Health England
RCT	Randomised Controlled Trial
SABA	Short-Acting B ₂ -2-Agonist
SLS	Salford Lung Study

526

527

528 **Declarations**

529

530 **Ethics approval and Consent to Participate**

531 Not Applicable

532

533 **Consent for Publication**

534 Not Applicable

535

536 **Availability of data and materials**

537 The datasets analysed during the current study are not publicly available but are available by
538 application to, and approval from, the Salford Lung Study scientific committee. Code scripts,
539 in the R language, for all components of the data cleaning, linkage, and subsequent analysis
540 will be made available in the open source GitHub website

541 (https://github.com/hollytibble/Salford-Lung-Study_Adherence-Linkage).

542

543 **Competing Interests Statement:**

544 The Salford Lung Study was funded by GlaxoSmithKline. JL-F was an employee of

545 GlaxoSmithKline during the conduct of the study, and holds shares/options in the company.

546 No other authors have any conflict pertaining to this manuscript to disclose.

547

548 **Funding:**

549 The study was supported by HT's College of Medicine and Veterinary Medicine PhD

550 (eHERC/Farr Institute) Studentship from The University of Edinburgh, and is carried out

551 with the support of the Asthma UK Centre for Applied Research [AUK-AC-2012-01].

552 MAM's Newton International Fellowship is awarded by the Academy of Medical Sciences

553 and Newton Fund. The funders had no role in study design, data collection and analysis,

554 decision to publish, or preparation of the manuscript. The Salford Lung Study in asthma

555 (HZA115150; NCT01706198) was funded by GlaxoSmithKline plc. GlaxoSmithKline

556 allows the University of Edinburgh the use of (but not access to) the study data and

557 statistician support. GlaxoSmithKline did not provide any monetary funding to this study.

558

559 **Author Contributions:**

560 HT conceived and planned the analysis. JL-F implemented the analysis scripts in the SLS data

561 platform. HT wrote the first draft, with contributions from JL-F, AS, RH, MAM, and AT. All

562 authors approved the final version and jointly take responsibility for the decision to submit this

563 manuscript to be considered for publication.

564

565 **Acknowledgements:**

566 Not applicable

567

568 References

569

- 570 1. John D, Michael W, Twigg J. Community pharmacy: an untapped patient data
571 resource. *Integr Pharm Res Pract* [Internet]. 2016 [cited 2019 Jul 24];5:19–25.
572 Available from: <http://youtu.be/IPZjCov6Obs>
- 573 2. Karter AJ, Parker MM, Moffet HH, Ahmed AT, Schmittiel JA, Selby J V. New
574 Prescription Medication Gaps: A Comprehensive Measure of Adherence to New
575 Prescriptions. *Health Serv Res* [Internet]. 2009 [cited 2018 Dec 19];44(5):1640–61.
576 Available from:
577 <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC2754552/pdf/hesr0044-1640.pdf>
- 578 3. Andrade SE, Kahler KH, Frech F, Chan KA. Methods for evaluation of medication
579 adherence and persistence using automated databases. *Pharmacoepidemiol Drug Saf*.
580 2006;15(8):565–74.
- 581 4. Hess LM, Raebel MA, Conner DA, Malone DC. Measurement of Adherence in
582 Pharmacy Administrative Databases: A Proposal for Standard Definitions and
583 Preferred Measures. *Ann Pharmacother* [Internet]. 2006 [cited 2018 Jul 1];40:1280–8.
584 Available from: www.theannals.com
- 585 5. Cutler RL, Fernandez-Llimos F, Frommer M, Benrimoj C, Garcia-Cardenas V.
586 Economic impact of medication non- adherence by disease groups: a systematic
587 review. *BMJ Open* [Internet]. 2018 [cited 2018 Apr 23];8:e016982. Available from:
588 <http://>
- 589 6. Patel AR, Campbell JR, Sadatsafavi M, Marra F, Johnston JC, Smillie K, et al. Burden
590 of non-adherence to latent tuberculosis infection drug therapy and the potential cost-
591 effectiveness of adherence interventions in Canada: A simulation study. *BMJ Open*.
592 2017 Sep 1;7.

- 593 7. Mckenzie SJ, Mclaughlin D, Clark J, Doi SAR, Mckenzie SJ, Mclaughlin AD, et al.
594 The Burden of Non-Adherence to Cardiovascular Medications Among the Aging
595 Population in Australia: A Meta-Analysis. *Drugs Aging* [Internet]. 2015 [cited 2019
596 Dec 6];32:217–25. Available from: <http://www.epigear.com>
- 597 8. Cutrona SL, Choudhry NK, Fischer MA, Servi AD, Stedman M, Liberman JN, et al.
598 Targeting cardiovascular medication adherence interventions. *J Am Pharm Assoc.*
599 2012;52(3):381–97.
- 600 9. Haberer JE, Sabin L, Amico KR, Orrell C, Galárraga O, Tsai AC, et al. Improving
601 antiretroviral therapy adherence in resource-limited settings at scale: a discussion of
602 interventions and recommendations. *J Int AIDS Soc* [Internet]. 2017 [cited 2019 Dec
603 6];20(1):21371. Available from: <http://doi.wiley.com/10.7448/IAS.20.1.21371>
- 604 10. Normansell R, Kew KM, Mathioudakis AG. Interventions to improve inhaler
605 technique for people with asthma [Internet]. *Cochrane Database of Systematic*
606 *Reviews*. 2017 Mar [cited 2020 Jan 20]. Available from:
607 <http://doi.wiley.com/10.1002/14651858.CD012286.pub2>
- 608 11. Valgimigli M, Garcia-Garcia HM, Vrijens B, Vranckx P, McFadden EP, Costa F, et al.
609 Standardized classification and framework for reporting, interpreting, and analysing
610 medication non-adherence in cardiovascular clinical trials: A consensus report from the
611 Non-adherence Academic Research Consortium (NARC). *Eur Heart J.*
612 2019;40(25):2070–85.
- 613 12. DeWorsop D, Creatura G, Bluez G, Thurnauer H, Forselius-Bielen K, Ranganathan M,
614 et al. Feasibility and success of cell-phone assisted remote observation of medication
615 adherence (CAROMA) in clinical trials. *Drug Alcohol Depend.* 2016 Jun 1;163:24–30.
- 616 13. Wu AC, Butler MG, Li L, Fung V, Kharbanda EO, Larkin EK, et al. Primary
617 Adherence to Controller Medications for Asthma Is Poor. *Ann Am Thorac Soc*

- 618 [Internet]. 2015 [cited 2017 Dec 21];12(2):161–6. Available from:
619 [https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4342835/pdf/AnnalsATS.201410-](https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4342835/pdf/AnnalsATS.201410-459OC.pdf)
620 [459OC.pdf](https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4342835/pdf/AnnalsATS.201410-459OC.pdf)
- 621 14. Williams LK, Joseph CL, Peterson EL, Wells K, Wang M, Chowdhry VK, et al.
622 Patients with asthma who do not fill their inhaled corticosteroids: A study of primary
623 nonadherence. *J Allergy Clin Immunol* [Internet]. 2007 Nov [cited 2019 Sep
624 17];120(5):1153–9. Available from:
625 <https://linkinghub.elsevier.com/retrieve/pii/S0091674907015862>
- 626 15. Liberman JN, Hutchins DS, Popiel RG, Patel MH, Jan SA, Berger JE. Determinants of
627 primary nonadherence in asthma-controller and dyslipidemia pharmacotherapy. *Am J*
628 *Pharm Benefits*. 2010;2(2):111–8.
- 629 16. Blais L, Kettani FZ, Forget A, Beauchesne MF, Lemièrè C, Ducharme FM. Assessing
630 adherence to inhaled corticosteroids in asthma patients using an integrated measure
631 based on primary and secondary adherence. *Eur J Clin Pharmacol*. 2017 Jan
632 1;73(1):91–7.
- 633 17. Ducharme FM, Noya FJD, Allen-Ramey FC, Maiese EM, Gingras J, Blais L. Clinical
634 effectiveness of inhaled corticosteroids versus montelukast in children with asthma:
635 prescription patterns and patient adherence as key factors. *Curr Med Res Opin*
636 [Internet]. 2012 [cited 2019 Aug 8];28(1):111–9. Available from:
637 <https://www.tandfonline.com/action/journalInformation?journalCode=icmo20>
- 638 18. Shin J, McCombs JS, Sanchez RJ, Udall M, Deminski MC, Cheetham TC. Primary
639 nonadherence to medications in an integrated healthcare setting. *Am J Manag Care*.
640 2012 Aug 1;18(8):426–34.
- 641 19. Cheen MHH, Tan YZ, Oh LF, Wee HL, Thumboo J. Prevalence of and factors
642 associated with primary medication non-adherence in chronic disease: A systematic

- 643 review and meta-analysis. *Int J Clin Pract* [Internet]. 2019 Jun 1 [cited 2020 Mar
644 3];73(6):e13350. Available from:
645 <https://onlinelibrary.wiley.com/doi/abs/10.1111/ijcp.13350>
- 646 20. Pottegård A, Depont Christensen R, Laust Thomsen J, dePont Christensen R, Houji A,
647 Binderup Christiansen C, et al. Primary non-adherence in general practice: A Danish
648 register study. *Artic Eur J Clin Pharmacol* [Internet]. 2014 [cited 2020 Mar 3];
649 Available from: <https://www.researchgate.net/publication/261801606>
- 650 21. Duerden M, Millson D, Avery A, Smart S. *The Quality of GP Prescribing: An Inquiry*
651 *into the Quality of General Practice in England*. 2011.
- 652 22. Dispensing Doctors' Association. *All about Dispensing Practice in England: A guide*
653 *for NHS service commissioners*. 2017.
- 654 23. Feehan M, Ranker L, Durante R, Cooper DK, Jones GJ, Young DC, et al. Adherence
655 to controller asthma medications: 6-month prevalence across a US community
656 pharmacy chain. *J Clin Pharm Ther* [Internet]. 2015 Oct 1 [cited 2019 Jun
657 27];40(5):590–3. Available from: <http://doi.wiley.com/10.1111/jcpt.12316>
- 658 24. Williams AB, Amico KR, Bova C, Womack JA. A Proposal for Quality Standards for
659 Measuring Medication Adherence in Research. *AIDS Behav* [Internet]. 2013 Jan 10
660 [cited 2019 Aug 8];17(1):284–97. Available from:
661 <http://link.springer.com/10.1007/s10461-012-0172-7>
- 662 25. Henson KE, Brock R, Shand B, Coupland VH, Elliss-Brookes L, Lyratzopoulos G, et
663 al. Cohort profile: prescriptions dispensed in the community linked to the national
664 cancer registry in England. *BMJ Open* [Internet]. 2018 [cited 2019 Jul 24];8:e20980.
665 Available from: <http://bmjopen.bmj.com/>
- 666 26. Doidge JC, Harron K. Demystifying probabilistic linkage : Common myths and
667 misconceptions. *Int J Popul Data Sci*. 2018;3(1).

- 668 27. Tibble H, Law H Di, Spittal MJ, Karmel R, Borschmann R, Hail-jares K, et al. The
669 importance of including aliases in data linkage with vulnerable populations. *BMC Med*
670 *Res Methodol*. 2018;18(76).
- 671 28. Downs JM, Ford T, Stewart R, Epstein S, Shetty H, Little R, et al. An approach to
672 linking education, social care and electronic health records for children and young
673 people in South London: A linkage study of child and adolescent mental health service
674 data. *BMJ Open* [Internet]. 2019 [cited 2020 Oct 5];9(1). Available from:
675 <http://bmjopen.bmj.com/>
- 676 29. Padmanabhan S, Carty L, Cameron E, Ghosh RE, Williams R, Strongman H.
677 Approach to record linkage of primary care data from Clinical Practice Research
678 Datalink to other health-related patient data: overview and implications. *Eur J*
679 *Epidemiol* [Internet]. 2019 [cited 2020 Aug 19];34(1):91–9. Available from:
680 <https://doi.org/10.1007/s10654-018-0442-4>
- 681 30. Gilbert R, Lafferty R, Hagger-Johnson G, Harron K, Zhang LC, Smith P, et al.
682 GUILD: GUIdance for Information about Linking Data sets. *J Public Health*
683 (Bangkok). 2018;40(1):191–8.
- 684 31. Benchimol EI, Smeeth L, Guttman A, Harron K, Moher D, Petersen I, et al. The
685 REporting of studies Conducted using Observational Routinely-collected health Data
686 (RECORD) Statement. *PLOS Med* [Internet]. 2015 Oct 6 [cited 2019 Sep
687 30];12(10):e1001885. Available from:
688 <https://dx.plos.org/10.1371/journal.pmed.1001885>
- 689 32. Woodcock A, Bakerly ND, New JP, Gibson JM, Wu W, Vestbo J, et al. The Salford
690 Lung Study protocol: A pragmatic, randomised phase III real-world effectiveness trial
691 in asthma. *BMC Pulm Med* [Internet]. 2015 [cited 2019 Aug 9];15(160). Available
692 from: <https://nweh.co.uk/uploads/documents/publications/BMC-Pulmonary-Medicine->

- 693 Paper-on-SLS-Asthma-Protocol.pdf
- 694 33. Albertson T, Murin S, Sutter M, Chenoweth J. The Salford Lung Study: a pioneering
695 comparative effectiveness approach to COPD and asthma in clinical trials. *Pragmatic
696 Obs Res [Internet]*. 2017;Volume 8:175–81. Available from:
697 [https://www.dovepress.com/the-salford-lung-study-a-pioneering-comparative-
698 effectiveness-approach-peer-reviewed-article-POR](https://www.dovepress.com/the-salford-lung-study-a-pioneering-comparative-effectiveness-approach-peer-reviewed-article-POR)
- 699 34. New JP, Bakerly ND, Leather D, Woodcock A. Obtaining real-world evidence: the
700 Salford Lung Study. *Thorax [Internet]*. 2014 Dec 1 [cited 2019 Aug 9];69(12):1152–4.
701 Available from: <https://thorax.bmj.com/content/69/12/1152>
- 702 35. Woodcock A, Vestbo J, Bakerly ND, New J, Gibson JM, McCorkindale S, et al.
703 Effectiveness of fluticasone furoate plus vilanterol on asthma control in clinical
704 practice: an open-label, parallel group, randomised controlled trial. *Lancet*.
705 2017;390(10109):2247–55.
- 706 36. Bakerly ND, Woodcock A, New JP, Gibson JM, Wu W, Leather D, et al. The Salford
707 Lung Study protocol: A pragmatic, randomised phase III real-world effectiveness trial
708 in chronic obstructive pulmonary disease. *Respir Res [Internet]*. 2015;16(1):1–5.
709 Available from: <http://dx.doi.org/10.1186/s12931-015-0267-6>
- 710 37. Magid S, Forrer C, Shaha S. Duplicate Orders: An Unintended Consequence of
711 Computerized provider/physician order entry (CPOE) Implementation: Analysis and
712 Mitigation Strategies. *Appl Clin Inform [Internet]*. 2012 [cited 2019 Sep 25];3(4):377.
713 Available from: <http://www.ncbi.nlm.nih.gov/pubmed/23646085>
- 714 38. Ekedahl A, Brosius H, Jönsson J, Karlsson H, Yngvesson M. Discrepancies between
715 the electronic medical record, the prescriptions in the Swedish national prescription
716 repository and the current medication reported by patients. *Pharmacoepidemiol Drug
717 Saf*. 2011;20:1177–83.

- 718 39. Burden AM, Paterson JM, Gruneir A, Cadarette SM. Adherence to osteoporosis
719 pharmacotherapy is underestimated using days supply values in electronic pharmacy
720 claims data. *Pharmacoepidemiol Drug Saf.* 2015;24:67–74.
- 721 40. Duerden MG, Hughes DA. Generic and therapeutic substitutions in the UK: are they a
722 good thing? *Br J Clin Pharmacol* [Internet]. 2010 Sep 1 [cited 2019 Sep
723 25];70(3):335–41. Available from: [http://doi.wiley.com/10.1111/j.1365-](http://doi.wiley.com/10.1111/j.1365-2125.2010.03718.x)
724 [2125.2010.03718.x](http://doi.wiley.com/10.1111/j.1365-2125.2010.03718.x)
- 725 41. Fischer MA, Stedman MR, Lii J, Vogeli C, Shrank WH, Brookhart MA, et al. Primary
726 medication non-adherence: analysis of 195,930 electronic prescriptions. *J Gen Intern
727 Med* [Internet]. 2010 Apr [cited 2017 Dec 21];25(4):284–90. Available from:
728 <http://www.ncbi.nlm.nih.gov/pubmed/20131023>
- 729 42. Harron KL, Doidge JC, Knight HE, Gilbert RE, Goldstein H, Cromwell DA, et al. A
730 guide to evaluating linkage quality for the analysis of linked data. *Int J Epidemiol.*
731 2017;46(5):1699–710.
- 732 43. British Thoracic Society. British Guideline on the Management of Asthma: Quick
733 Reference Guide [Internet]. Scottish Intercollegiate Guidelines Network. 2016.
734 Available from: <http://www.ncbi.nlm.nih.gov/pubmed/19209371>
- 735 44. Hoopes M, Angier H, Raynor LA, Suchocki A, Muench J, Marino M, et al.
736 Development of an algorithm to link electronic health record prescriptions with
737 pharmacy dispense claims. *J Am Med Informatics Assoc* [Internet]. 2018 Oct 1 [cited
738 2019 Jul 24];25(10):1322–30. Available from:
739 <https://academic.oup.com/jamia/article/25/10/1322/5065692>
- 740 45. Berger Z, Kimbrough W, Gillespie C, Boscarino JA, Wood GC, Qian Z, et al. Lower
741 copay and oral administration: Predictors of first-fill adherence to new asthma
742 prescriptions. *Am Heal Drug Benefits.* 2009;2(4):174–9.

- 743 46. Florentinus SR, Souverein PC, Griens FA, Groenewegen PP, Leufkens HG, Heerdink
744 ER. Community pharmacy dispensing data to prescribing data of general practitioners.
745 BMC Med Informatics Decis Mak Link [Internet]. 2006 [cited 2019 Jul 24];6(18).
746 Available from: <http://www.biomedcentral.com/1472-6947/6/18>
- 747 47. Johannesdottir SA, Lund M, Jens M, Hansen G, Lash TL, Pedersen L, et al.
748 Correspondence between general practitioner-reported medication use and timing of
749 prescription dispensation. Clin Epidemiol [Internet]. 2012 [cited 2019 Jul 24];4:13–8.
750 Available from: <http://dx.doi.org/10.2147/CLEP.S26958>
- 751
- 752

753 *Table 1: Cox Proportional hazards model risk factors associated with time to collecting a*
 754 *prescribed medication.*

	Hazard Ratio (95% Confidence Interval)	Statistical significance (p-value)
Season		
Spring	{reference}	
Summer	0.967 (0.944 - 0.991)	0.008 *
Autumn	0.981 (0.958 – 1.005)	0.123
Winter	1.003 (0.979 – 1.028)	0.791
Drug Class		
SABA	1.433 (1.387 - 1.479)	<0.001 *
LABA	0.938 (0.890 - 0.990)	0.019 *
ICS	{reference}	
ICS+LABA	1.067 (1.033 – 1.102)	<0.001 *
Cromogliclate	0.778 (0.389 – 1.558)	0.479
Immuno-suppressants	1.244 (1.100 - 1.408)	<0.001 *
LAMA	1.349 (1.161 - 1.567)	<0.001 *
LTRA	1.350 (1.289 – 1.414)	<0.001 *
Theophylline	1.040 (0.897 - 1.205)	0.604
Oral steroids	1.839 (1.743 - 1.940)	<0.001 *
Previously unclaimed medications		
Low tertile	{reference}	
Mid tertile	0.565 (0.553 – 0.577)	<0.001 *
High tertile	0.198 (0.193 – 0.204)	<0.001 *
Quantity of doses prescribed	1.000 ** (1.000 – 1.000)	<0.001 *

755 Statistically significant variables (using a threshold of $p=0.05$) are denoted by a star (*).

756 ** Coefficient 0.9999 to four decimal places, and therefore lower than the reference value

757

758

759 **Figure Legends:**

760

761 *Figure 1: Diagram representing the data linkage algorithm.*

762 *Figure 2: Data Linkage Flow Diagram.*

763 *Figure 3: Distributions of linkage weight points per variable, for candidates and final matches.*

764 *Figure 4: Kaplan-Meier of the time to collecting prescriptions, censored at three weeks.*

765 APPENDIX A: String Search Keywords by Medication and Drug Class Keyword Categories.

Drug Class Keyword	Medication Keyword	String Search Keywords
SABA	SALBUTAMOL	"SALBUTAMOL", "ALBUTEROL", "VENTOLIN", "AIROMIR", "SALAMOL", "AIRSALB", "SALAPIN", "VENTMAX", "ASMASAL", "ESI- BREATHE", "SALBULIN", "SALIPRANEB", "IPRAMOL", "COMBIVENT"
SABA	BAMBUTEROL	"BAMBUTEROL", "BAMBEC"
LABA	FORMOTEROL	"FORMOTEROL", "FORADIL", "FOSTAIR", "SYMBICORT", "FLUTIFORM", "SPIROMAX", "OXIS", "ATIMOS"
LABA	SALMETEROL	"SALMETEROL", "NEOVENT", "SEREVENT", "SERETIDE", "SIRDUPLA", "AIRFLUSAL"
LABA	TERBUTALINE	"TERBUTALINE", "BRICANYL"
LABA	TIOTROPIUM	"TIOTROPIUM", "SPIRIVA"
LABA	VILANTEROL	"VILANTEROL", "RELVAR", "VILENTEROL"
LAMA	GLYCOPYRRONIUM BROMIDE	"SEEBRI"
LAMA	IPRATROPIUM	"IPRATROPIUM", "ATROVENT", "RESPONTIN", "IPRAVENT", "SALIPRANEB",

		"IPRAMOL", "COMBIVENT"
THEOPHYLLINE	THEOPHYLLINE	"THEOPHYLLINE", "NEULIN", "SLO- PHYLLIN", "UNIPHYLLIN"
THEOPHYLLINE	AMINOPHYLLINE	"AMINOPHYLLINE", "PHYLLOCONTIN"
ICS	BECLOMETASONE	"BECLOMETASONE", "ASMABEC", "BECODISKS", "CLENIL", "QVAR", "FOSTAIR"
ICS	CICLESONIDE	"CICLESONIDE", "ALVESCO"
ICS	BUDESONIDE	"BUDESONIDE", "BUDELIN", "PULMICORT", "SYMBICORT", "SPIROMAX"
ICS	FLUTICASONE	"FLUTICASONE", "FLIXOTIDE", "FLUTIFORM", "SERETIDE", "SIRDUPLA", "AIRFLUSAL", "RELVAR"
ICS	MOMETASONE	"MOMETASONE", "TWISTHALER", "ASMANEX"
LTRA	MONTELUKAST	"MONTELUKAST", "SINGULAIR"
LTRA	ZAFIRLUKAST	"ZAFIRLUKAST", "ACCOLATE"
LTRA	ZILEUTON	"ZILEUTON", "ZYFLO"
CROMOGLICATE	NEDOCROMIL	"NEDOCROMIL", "TILADE"
CROMOGLICATE	CROMOGLICATE	"CROMOGLICATE", "CROMOGLYCATE", "INTAL"
STEROID	OMALIZUMAB	"OMALIZUMAB", "XOLAIR"
STEROID	PREDNISOLONE	"PREDNISOLONE"

IMMUNO-SUPPRESSANT	METHOTREXATE	"METHOTREXATE", "MAXTREX", "METOJECT", "METHOFILL", "NORDIMET", "ZLATAL"
IMMUNO-SUPPRESSANT	CICLOSPORIN	"CICLOSPORIN", "CAPIMUNE", "CAPSORIN", "DEXIMUNE", "NEORAL", "SANDIMMUN"
IMMUNO-SUPPRESSANT	AZATHIOPRINE	"AZATHIOPRINE", "IMURAN"

766 String search keywords may appear under multiple medication and drug class keyword
767 categories, if they contain more than one active ingredient, such as combination ICS LABA
768 medications.

769 Bold string search keywords indicate brand names

770

771

772

773

774 APPENDIX B: EXCLUSION KEYWORDS AND FREQUENCY

775

Exclusion Keyword	Unique Drug Descriptions (N=928)
NASAL	39
NOSE	1
NOSTRIL	0
NASULE	0
HAYFEVER	0
EYE	11
EAR	0
DROP	16
TONGUE	0
FOAM	2
ENEMA	1
RECTAL	0
GASTRO *	1
MODIFIED *	0
CREAM	4
APPLY	0
SKIN	0
ULCER	0
OINTMENT	6
PATCH	0
CAPSULE**	2
SACHET	0
SPRAY	33
AZELASTINE	4
NASONEX	0
FLIXONASE	0
ANORA ELLIPTA	0
SUMATRIPTAN	0
AVAMYS	0
RHINOCORT	0
NASOBEC	0
NASOFAN	0
TOTAL EXCLUDED	71 (7.7%)

776 * Excluding medications of drug class “steroid” or “theophylline”

777 ** Excluding medications of drug class “steroid”, “theophylline”, “tiotropium” or
778 “glycopyrronium bromide”

779

780

781

782 APPENDIX C: Variable Recoding

783

784 *Quantity Recoding:*

785 Quantities with values of over 28 were assumed to be the number of doses, rather than the
786 number of units/inhalers. The most common recorded number of dose quantity was imputed
787 as the most commonly occurring number of doses per unit (as the most common number of
788 units prescribed is one) for that medication class. If the quantity was recorded in doses, this
789 was set as the primary dose quantity, with the second most commonly occurring dose quantity
790 as the alias value. If the quantity was recorded in units, the number of units multiplied by the
791 most commonly occurring dose quantity was imputed as the primary value, and the second
792 most likely as the alias.

793

794 *Dose Strength Recoding:*

795 All dose strengths were converted into upper case, spaces were removed, and the following
796 string substitutions were made:

- 797 • "MICROGRAMS" replaced with "MCG",
- 798 • "MICROGRAM" replaced with "MCG",
- 799 • "MICROG" replaced with "MCG",
- 800 • "UNITS" replaced with "U"

801 Strings were then searched for the first pattern of "0.5", "500", "400", "320", "200", "184",
802 "160", "125", "100", "92", "80", "50", "25", "20", "10", "5", "4", "2", or "1", followed by any
803 of "MG", "MCG" or "/". ICS+LABA medications often recorded as X/X dose, in which the
804 larger number relates to the ICS and the lower to the LABA. Some records listed the

805 ICS+LABA combination medicines as ICS/LABA dose, and some as LABA/ICS dose; as
806 such, the possible patterns were searched in order of size, rather than position in string.
807

Factor	Criteria	Points	Factor Range	% of candidates	% of matches
Brand Name *	Both records had non-missing, and distinct, brand names	0	0-20	6.3%	2.8%
	One or both of the records had a missing brand name	10		0%	0%
	Both records had non-missing, and matching, brand names	20		93.7%	97.2%
(Modified) Dose Strength	Both records had non-missing, and distinct, dose strengths	0	0-35	4.8%	0%
	One or both of the records had a missing dose strength	10		18.1%	9.0%
	Both records had non-missing, and matching, dose strengths	35		77.2%	91.0%
(Modified) Medication Quantity	Both records had non-missing, and distinct, primary and alias dose quantities	0	0-35	4.2%	0%
	One or both of the records had a missing primary quantity value, indicating that no value was observed or could be imputed	10		9.8%	<0.1%
	Both records had non-missing, and distinct, primary dose quantities, but the alias of one record matched to the primary of the other	15		4.9%	1.5%
	Both records had non-missing, and matching, primary dose quantities	35		81.1%	98.5%
Date difference	Dispensing occurred more than one month after prescription (but less than six months)	0	0-10	67.2%	1.3%
	Dispensing occurred within one month of prescription	10		32.8%	98.7%

809 * If a generic medication was used, the brand name was listed as 'generic'

810 APPENDIX E: INCLUDED FEATURE WEIGHT COMBINATIONS

WEIGHT	BRAND NAME	DOSE STRENGTH	QUANTITY	DATES
100	Non-missing and matching	Non-missing and matching	Non-missing and matching	Less than one-month delay
90	One or more missing	Non-missing and matching	Non-missing and matching	Less than one-month delay
	Non-missing and matching			More than one-month delay
80	Non-missing and distinct	Non-missing and matching	Non-missing and matching	Less than one-month delay
	Non-missing and matching		Primary/alias match	
	One or more missing		Non-missing and matching	More than one-month delay
75	Non-missing and matching	One or more missing	Non-missing and matching	Less than one-month delay
		Non-missing and matching	One or more missing	
70	Non-missing and distinct	Non-missing and matching	Non-missing and matching	More than one-month delay
	Non-missing and matching		Primary/alias match	Less than one-month delay
	One or more missing			

811

812 APPENDIX F: LINKAGE ALGORITHM DESCRIPTION

813

814 The matching algorithm iteratively searches through dispensing records, finding the closest
815 matching prescription record and subsequently removing it from future iterations, for each
816 person and medication class keyword. The medication class keyword is generated by
817 identifying the key active ingredients in a medication that are common between both generic
818 and brand name equivalents, using a domain-knowledge look-up table.

819

820 Starting with the first dispensing record, all candidate prescription record links (linkage weight
821 over the threshold and prescription date up to a maximum of six months prior to dispensing)
822 are identified. The most recently prescribed candidate link for the dispensing is selected as the
823 most likely match, using highest linkage weights to break ties, and the non-selected candidate
824 links for both the matched dispensing record and the matched prescription record are excluded
825 from future iterations. The process repeats until every dispensing record has been considered,
826 although it is possible that no candidate links will be available for some dispensing records at
827 later iterations if all initial prescription candidates have been successfully matched to other
828 dispensing records.