



Automated extraction of weight, height, and obesity in electronic medical records are highly valid

Namneet Sandhu^{1,2}  | Alexander Krusina^{1,3} | Hude Quan^{1,2} | Robin Walker^{1,2,3} | Elliot A. Martin^{1,3} | Cathy A. Eastwood^{1,2} | Danielle A. Southern^{1,2} 

¹Centre for Health Informatics, Cumming School of Medicine, University of Calgary, Calgary, Alberta, Canada

²Department of Community Health Sciences, Cumming School of Medicine, University of Calgary, Calgary, Alberta, Canada

³Alberta Health Services, Calgary, Alberta, Canada

Correspondence

Danielle A. Southern, Department of Community Health Sciences, Centre for Health Informatics, University of Calgary, Cal Wenzel Precision Health Building, 3280 Hospital Dr NW, Calgary, AB T2N 4Z6, Canada.
Email: dasouthe@ucalgary.ca

Abstract

Objective: Coding of obesity using the International Classification of Diseases (ICD) in healthcare administrative databases is under-reported and thus unreliable for measuring prevalence or incidence. This study aimed to develop and test a rule-based algorithm for automating the detection and severity of obesity using height and weight collected in several sections of the Electronic Medical Records (EMRs).

Methods: In this cross-sectional study, 1904 inpatient charts randomly selected in three hospitals in Calgary, Canada between January and June 2015 were reviewed and linked with AllScripts Sunrise Clinical Manager EMRs. A rule-based algorithm was created which looks for patients' height and weight values recorded in EMRs. Clinical notes were split into sentences and searched for height and weight, and BMI was computed.

Results: The study cohort consisted of 1904 patients with 50.8% females and 43.3% > 64 years of age. The final model to identify obesity within EMRs resulted in a sensitivity of 92.9%, specificity of 98.4%, positive predictive value of 96.7%, negative predictive value of 96.6%, and F1 score of 94.8%.

Conclusions: This study developed a highly valid rule-based EMR algorithm that detects height and weight. This could allow large-scale analyses using obesity that were previously not possible.

KEYWORDS

algorithms, body-mass index, obesity, prevalence

1 | INTRODUCTION

Measuring the prevalence of obesity using BMI in national surveys such as the Canadian Community Health Survey (CCHS) or National Health and Nutrition Examination Survey (NHANES) relies on self-reported height and weight measures. These measures are potentially subject to respondent biases as respondents may not be aware

of their accurate height or weight. Underestimation of weight or overestimation of height by survey respondents could lead to underestimation of obesity.^{1,2} Administrative databases have previously been found to under-code certain conditions including obesity.²⁻⁴ For example, using the Discharge Abstract Database (DAD), one study found that a case definition using International Classification of Diseases (ICD) diagnostic codes E65-E68 (ICD-10) to identify obesity

This is an open access article under the terms of the [Creative Commons Attribution-NonCommercial](https://creativecommons.org/licenses/by-nc/4.0/) License, which permits use, distribution and reproduction in any medium, provided the original work is properly cited and is not used for commercial purposes.

© 2023 The Authors. Obesity Science & Practice published by World Obesity and The Obesity Society and John Wiley & Sons Ltd.

had a very low sensitivity, less than 10%.² Abstracted data in administrative databases were previously deemed not useable for obesity surveillance due to underestimation of obesity,^{2,3,5} and was found to more likely capture Class III obesity and miss those with Class I obesity.²

Electronic medical records (EMRs), which systematically collect and store patient health information (clinical and administrative) in a digital format, have been widely adopted in both acute and primary care settings in many countries.⁶⁻⁸ The use of EMR data, originally intended for clinical and administrative purposes, has expanded to disease surveillance and diverse epidemiological research over the years.⁸ The vast amount of health data generated using EMRs avails opportunities for the development and use of analytical tools for clinical decision support and improving patient outcomes.⁹⁻¹¹ Several analytical approaches for the identification of individuals with particular conditions or outcomes, also known as phenotyping, have been developed utilizing EMR data.¹² One of these techniques is rule-based phenotyping, which is based on expert-defined criteria drawn from consensus guidelines on diagnosis and treatment.^{12,13} EMRs record the patients' body weight and height for their hospital visit, which lend an opportunity for the development of an algorithm for the identification of obesity using the BMI. This algorithm could potentially have a higher validity for obesity identification in comparison to the administrative databases where identification is dependent on the documentation of obesity by healthcare providers and coding by abstractors.

In Calgary, AllScripts Sunrise Clinical Manager™ (SCM) is an inpatient EMR system that has been in use since 2006 with inpatient health data collected from more than 5.4 million people.⁶ SCM includes structured as well as free-text data including discharge summaries, clinical examination, and hospitalization progress notes from the patient-provider interactions at the five hospitals in the city.⁶ The aim of this study was to develop a rule-based algorithm to improve the phenotyping of height and weight measures for detecting obesity using several different sections of the SCM EMR data from Calgary, Alberta, Canada, and to test the performance of the algorithm by comparing with manual chart review.

2 | METHODS

2.1 | Design and study population

This is a retrospective cross-sectional study completed using EMR data for adult patient population hospitalized in Calgary between 1 January 2015, and 30 June 2015. Patients were randomly selected from three adult acute care hospitals (Foothills Medical Centre, Rockyview General Hospital, and Peter Lougheed Centre).

2.2 | Data sources

A total of 3043 patient charts were extracted from SCM EMR data and used to create and evaluate the algorithm. The free text

documents within each chart were used. There were a total of 58 unique clinical note names in the extracted EMR data, such as "discharge summary", "nursing notes", and "history and physical." A manual chart review was previously completed using the same charts which abstracted patients' demographics such as age, sex, and the presence of comorbidities, specifically diabetes, hypertension, liver disease, cancer, myocardial infarction, congestive heart failure, and cerebrovascular disease.¹⁴ This original study compared chart review and administrative data; the current study now links the EMR data from this chart review from this cohort to evaluate the rule-based algorithm.

2.3 | Final study cohort

To determine the final study cohort, charts that did not have height/weight data captured in the chart reviews or the EMRs were excluded. Figure 1 displays the process of obtaining the cohort to evaluate our algorithm. The chart review had access to paper documents that were not available for this study, resulting in the removal of 224 charts.

2.4 | Rule-based algorithm

To identify obesity in the free text of EMR, a rule-based algorithm was created. A rule-based algorithm was chosen instead of machine learning or other methods because of its simplicity. The main goal of the algorithm is to detect height and weight, clean the values if necessary, and then calculate BMI. The detection of these terms is a named-entity recognition problem, where spans of text are flagged as being the entities of interest.¹⁵ A rule-based algorithm can achieve this goal while also being easy to explain. To detect height and weight

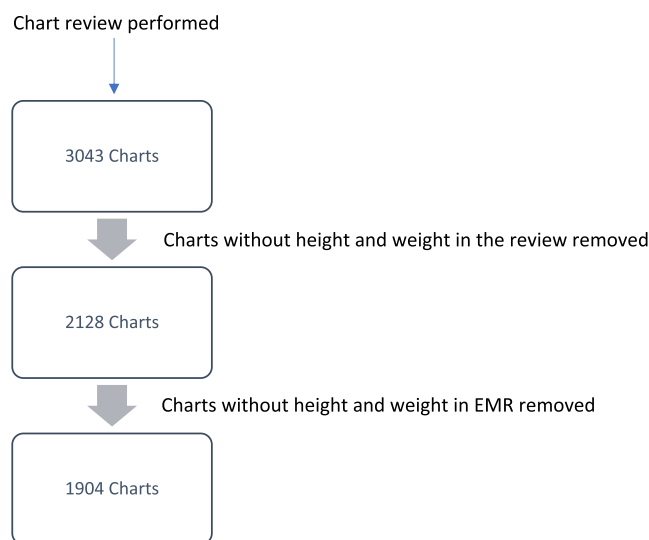


FIGURE 1 Process for obtaining obesity cohort within Electronic Medical Record data.

within EMRs, a simple spaCy pipeline was created with a customized sentencizer.¹⁶ The purpose of the pipeline is to look through all the text and define the sentence boundaries within each document. For every document within a chart, the document is fed into the pipeline to find the individual sentences. Each sentence is searched for a mention of height or weight. If one is found, regular expressions are used to find any numbers in the sentence.^{17,18}

In the EMR, most measurements were in centimeters and kilograms. If measurements were recorded in other units, they were converted to centimeters or kilograms. A common error within EMRs is that height and weight are swapped, with the height recorded where weight should be, and weight where height should be. Values found within the same document were checked and swapped if necessary. If the height value was outside five standard deviations from the mean adult height, and the weight value was inside five standard deviations, the values were swapped before computing the BMI. Five standard deviations were used as this provided the best results and is still within reasonable judgment values for a person's height. Other unrealistic values such as "height recorded 2015-02-13" were also excluded.

2.5 | Handling multiple height and weight values

For each chart, there are many documents, which can lead to more than one measurement for height or weight recorded within the chart. The patient's weight can fluctuate over their stay in the hospital, and this leads to different BMI values. If the values are on opposite sides of 30, the algorithm must have a way to predict whether obesity is present or absent. To address this, four different methods were created and tested. The first was simply taking the lowest BMI value (lowest weight, highest height). The second was taking the highest BMI value (highest weight, lowest height). The third method was taking the average of all the heights and average of all the weights and then calculating the BMI. The final method was taking the weight and height values closest to the patient's discharge (i.e., most recent) to accurately capture the state of the patient when they left the hospital.

2.6 | Statistical analysis and evaluation

The study cohort was characterized based on age, sex, and comorbidity using descriptive statistics, and the performance of the rule-based algorithm for obesity detection was tested by a previous manual chart review.¹⁴ For the performance evaluation, the study team first compared the means of heights and weights found by the algorithm to the means in the chart review. Second, the team tested whether the algorithm correctly determined whether a patient had obesity. By using the chart review as a reference standard, the study team calculated sensitivity, specificity, positive predictive value (PPV), negative predictive value (NPV), and F1 score. This comparison was conducted across different parameters including age, sex,

and number of comorbidity conditions that were identified in the chart review. The ICD codes for obesity in the DAD administrative database for the same study population were also compared to the BMI values in the chart review to assess the sensitivity, specificity, PPV, NPV and F1 scores for detecting obesity using the administrative database. The evaluation metrics were calculated on the entire chart review set that obesity labels were available for, as opposed to dividing the data into training and test sets, since a rule-based algorithm was used (as opposed to training a machine learning algorithm).

Ethics approval for this study was obtained from the Conjoint Health Research Ethics Board at the University of Calgary (REB15-0790).

3 | RESULTS

Following the process outlined in Figure 1 by excluding the charts without height and weight in the chart review and in the EMR, the final cohort created for assessing the algorithm consisted of 1904 patients with 50.8% females and 43.3% greater than 64 years of age. Baseline demographic and Charlson comorbidity characteristics of the cohort were provided by the manual chart review and are displayed in Table 1. Several patients (48%) had hypertension, 33.4% had cancer, and 17.3% had diabetes. Supplementary Table A1 shows that the characteristics of the final cohort were similar to the cohorts missing height or weight information in chart review or EMR.

Of the 1904 charts, more than one value for height was found in 1402 (73%) and more than one value for weight was found in 1518 (79%). About 3% ($n = 54$) charts had five or more heights and 16%

TABLE 1 Characteristics of the study cohort (1904).

Demographics N (%)	Final cohort N (%)
Age	
Less than 50	504 (26.5%)
50 to 64	575 (30.2%)
Greater than 64	825 (43.3%)
Sex	
Male	936 (49.2%)
Female	968 (50.8%)
Comorbidity	
Diabetes	330 (17.3%)
Hypertension	914 (48.0%)
Liver disease	138 (7.3%)
Cancer	635 (33.4%)
Myocardial infarction	49 (2.6%)
Congestive heart failure	176 (9.2%)
Cerebrovascular disease	178 (9.4%)

($n = 308$) charts had five or more weights detected. The presence of more than one value for height or weight in the patient chart resulted in the calculation of BMI using four different methods described earlier (lowest BMI, highest BMI, average BMI, closest to discharge/recent BMI). The rates for sensitivity, PPV, NPV, specificity, and F1 score for all four methods were greater than 85% (Figure 2). As can be expected, the method for taking the lowest BMI has very high specificity and PPV, but lower sensitivity. Conversely, the highest BMI method has very high sensitivity and NPV, but a lower PPV. The average and closest to discharge (i.e., recent) BMI methods are more balanced, with the closest to discharge (recent) method having a slightly higher F1 score than the average method. Therefore, further results are displayed using the closest to discharge (recent) method.

The mean values for height and weight using the values closest to discharge for each chart were also calculated and compared to the means from the chart review to assess the accuracy of the algorithm (Table 2). This is important because it is possible to predict obesity correctly while incorrectly identifying height or weight. Overall, the mean height and weight from the EMR were found to be very close to the mean height and weight from the chart reviews, with the mean weight (79.5 kg) from chart review data slightly higher than the mean weight (78.0 kg) from EMR data.

The data were further analyzed to assess the severity of obesity in our study cohort and the performance metrics (sensitivity, specificity, PPV, NPV and F1) for obesity detection by age, sex, and number of comorbidity conditions (Table 3). The prevalence of obesity (i.e., BMI ≥ 30) was found to be slightly lower using the EMR algorithm (31.5%) compared to manual chart reviews (32.7%). Of those who had obesity, the majority (18.2% by EMR and 18.9% by chart review) were found to be in Class I of obesity (i.e., BMI 30 to < 35) and the least (5.0% by EMR and 5.5% by chart review) in Class III (i.e., BMI ≥ 40). The rates for sensitivity, specificity, PPV, and NPV for detecting obesity using our EMR algorithm were very high ($>90\%$) regardless of age, sex, or the number of comorbidities.

This study also found the prevalence of obesity to be higher among those who had a higher number of comorbidities. Using EMR data, the prevalence of obesity was found to be 25.7% among those with no comorbidities, 32.3% with one comorbidity and 35.2% with two or more comorbidities (Table 3). This signifies the association of obesity with other conditions as suggested by previous research. For example, the prevalence of obesity was found to be 29.0% in those with diabetes versus 17.9% in those without.²

Additionally, the ICD codes in the DAD, compared to the BMI value in the chart review, were found to have a sensitivity, specificity, PPV, NPV, and F1 scores of 8.7%, 99.8%, 96.4%, 69.2%, and 15.9%, respectively. The rule-based EMR algorithm has a comparable PPV and specificity but much higher sensitivity, NPV, and F1 score.

TABLE 2 Mean height and weight of chart review and electronic medical record.

	Mean weight (kg)		Mean height (cm)	
	Chart review	EMR	Chart review	EMR
Overall	79.5	78.0	167.8	167.9
Age				
Less than 50	80.7	78.7	170.6	170.4
50 to 64	85.5	83.1	169.1	169.5
Greater than 64	74.9	74.1	165.2	165.3
Sex				
Male	86.3	84.8	174.9	174.8
Female	72.9	71.4	160.9	161.1
Comorbidities				
0	78.7	76.7	169.6	169.5
1	79.3	78.1	167.4	167.5
2+	80.4	78.9	166.7	166.9

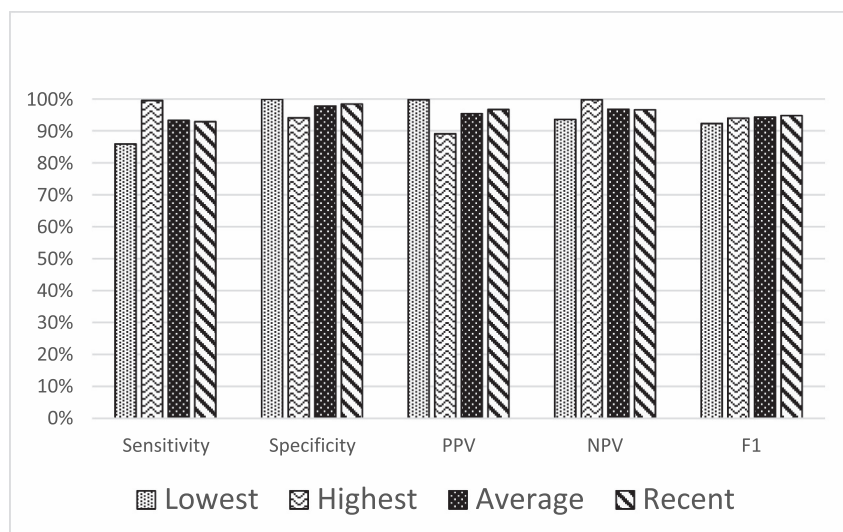


FIGURE 2 Performance metrics for obesity presence.

TABLE 3 Performance metrics of obesity detection using electronic medical records in comparison to chart review across obesity classes and patient characteristics.

	Chart review (%)	EMR (%)	Sensitivity (%)	Specificity (%)	PPV (%)	NPV (%)	F1 (%)
BMI							
<30 (not obese)	67.3	68.5	98.4	93.2	96.8	96.7	97.6
Class I: $30 \leq \text{BMI} < 35$	18.9	18.2	90.5	98.6	93.9	97.8	92.2
Class II: $35 \leq \text{BMI} < 40$	8.4	8.3	87.3	98.9	87.9	98.9	87.6
Class III: $\text{BMI} \geq 40$	5.5	5.0	80.8	99.3	87.5	98.9	84.0
Obesity ($\text{BMI} \geq 30$)	32.7	31.5	92.9	98.4	96.7	96.6	94.8
Age							
Less than 50	28.9	27.5	92.4	98.9	97.1	97.0	94.7
50 to 64	43.0	40.9	91.7	97.5	96.5	94.0	94.1
Greater than 64	28.2	27.6	94.5	98.7	96.6	97.9	95.5
Sex							
Male	30.7	29.9	94.1	98.5	96.4	97.4	95.2
Female	34.7	33.0	92.0	98.4	96.9	95.8	94.4
Number of comorbidities							
0	27.5	25.7	91.0	99.0	97.3	96.7	94.0
1	33.8	32.3	93.1	98.7	97.4	96.6	95.2
2 or more	35.9	35.2	93.9	97.7	95.7	96.6	94.8

4 | DISCUSSION

This study developed and validated a rule-based algorithm using EMR data for obesity detection. Overall, obesity could accurately be identified in an inpatient population using EMR data with high validity (>90% sensitivity, specificity, PPV, NPV, and F1) and across patient characteristics (age, sex, comorbidity conditions). Of the height and weight values captured multiple times during patients' hospital stay, BMI calculated using the values closest to discharge (i.e., most recent) had the highest validity.

The study findings suggest that EMR correctly captures a higher proportion of patients with obesity compared to administrative databases. EMR data had higher sensitivity (>90%) versus ICD-coded data (8.7%), while maintaining high specificity, PPV, NPV and F1 of greater than 90%. The EMR data source is valuable for obesity surveillance, health services research, as well as the evaluation of population-based public health interventions. In this study, the prevalence of obesity using the BMI closely matched that from manual chart reviews for all obesity classes, alleviating the risk of failing to detect those in earlier stages of obesity as was the case with abstracted data from administrative databases in Martin et al study.² Obesity is generally a secondary diagnosis during hospitalization, making it optional for coding per the Canadian coding standards.^{4,19} In addition, time constraints and the high volume of work faced by the coders could also impact coding. Furthermore, coders are only required to review physician or primary care provider documentation for abstracting diagnosis information. This, combined with poor documentation of obesity by physicians, leads to the inadequacy of

administrative data for obesity detection or estimation. These challenges in the administrative databases could be mitigated by an automated EMR-based algorithm, which is not dependent on the coding of obesity in a patient chart by the physician/primary care provider or the coders not abstracting obesity as secondary diagnoses. Thus, the EMR-based algorithm could overcome the potential biases associated with coding guidelines and physician/primary care provider practice of documentation. It is important to note that the Canadian Institute for Health Information (CIHI) added height and weight as new fields to the DAD as of the 2018–2019 fiscal year, which allows for the calculation of BMI for obesity.²⁰ Again, the lack of documentation of height and weight in the physician or nursing clinical notes would render the addition of these fields of low yield for obesity detection. To our knowledge, the level of completeness of the height and weight fields in the DAD is yet to be tested.

Another study conducted in an outpatient clinical setting found EMR data to significantly under-code obesity in the patient problem list compared to obesity identified using EMR-generated BMI (5.6% vs. 51.7%), where BMI was automatically calculated by EMR using the height and weight measures recorded in the patient's chart.²¹ Similar to Martin et al's study,² this study found that documentation of obesity in charts was much more likely for those who were in the higher class of obesity (i.e., $\text{BMI} > 40$), indicating a lack of recognition of obesity as a significant medical problem by the physicians.²¹

Direct implementation of this algorithm in an EMR system has the potential to enhance patient care by aiding physicians by recognizing the presence of obesity at the point-of-care leading to

patient counseling and referral for interventions targeting obesity and obesity-related complications. A study by Roth et al (2014) further explored the secondary use of EMRs by linking EMR-derived BMI and community data to study community-level factors, such as farmers markets, grocery stores, and level of education, associated with overweight and obesity.²² Their study suggested the value of integrating EMR and community data in enabling physicians to factor in the patient environment into health recommendations at the point-of-care.²² In Alberta, Alberta Health Services (AHS) has launched a provincial EMR initiative, known as Connect Care, built by Epic Systems Corporation. The deployment of Connect Care offers opportunities for the implementation of EMR algorithms across the province within inpatient and outpatient healthcare service centres. To explore the opportunity of automating obesity detection by implementing this algorithm in the Connect Care EMR, further discussions will be required with the implementation and authoritative teams of Connect Care and AHS.

Some of the challenges faced in developing this algorithm that affected the results included human error and inconsistency in documentation leading to poor data captured in the EMR. Examples of poor data encountered included recording of height as 108 instead of 180, swapping of height and weight values, recording of weight value for both height and weight, and use of incorrect units (such as, height recorded as 66 cm instead of 66 inches). These challenges were addressed by creating rules as outlined in the methodology section above. However, there were other scenarios which were more difficult to address by creating rules. An example scenario is “the patient is a small person of 152 cm”, where it is easy for a human to infer the height correctly but more challenging for an algorithm to detect. Thus, not every possible situation was addressed due to the increasingly large number of rules that would be required for a very small performance increase. Furthermore, this study found that 11% ($n = 224$) of medical charts were missing height or weight in the EMR. This makes it impossible to use an algorithm that uses height and weight to detect obesity. Additionally, further research is needed to assess the use of this EMR algorithm for capturing patient BMI trajectory over time as a potential method for long-term follow-up.

This study also has some limitations that are important to note. First, the study team used inpatient documentation only and is aware that obesity is largely managed in outpatient settings. However, this study was aimed at developing EMR-based obesity case identification in order to overcome under-coding issues in ICD administrative databases, such as the DAD. Second, this study was performed using data from one large urban city and the practice of documentation in the EMRs may vary across geographic areas. Thus, further studies are needed to conduct external validation of our rule-based algorithm using data from other regions and jurisdictions, including community-based EMRs.

This study found EMR data from Calgary, Alberta to be suitable for developing an algorithm for detecting obesity within an inpatient setting. An EMR algorithm was developed and evaluated that detects height and weight and was found to be highly valid. This could allow for large-scale secondary data studies with BMI as either a predictor

or outcome. The study team suggests that the validity of this algorithm should be tested in other regions of Canada using different EMRs used across the country.

AUTHOR CONTRIBUTIONS

Hude Quan conceived the project. Alexander Krusina conducted data analysis and wrote the manuscript. Namneet Sandhu contributed to the writing of the manuscript. Robin Walker reviewed and edited the manuscript. Elliot A. Martin provided guidance to data analysis and edited the manuscript. Cathy A. Eastwood reviewed and edited the manuscript. Danielle A. Southern provided guidance to data analysis and edited the manuscript. All authors have read and approved the final manuscript.

ACKNOWLEDGMENTS

There are no further contributions to acknowledge. No specific funding agency or grant supported this project.

CONFLICT OF INTEREST STATEMENT

The authors declare that they have no competing interests.

ORCID

Namneet Sandhu  <https://orcid.org/0000-0003-3870-0320>

Danielle A. Southern  <https://orcid.org/0000-0002-0006-0033>

REFERENCES

1. Statistics Canada: Health Fact Sheets Overweight and obese adults. 2018. <https://www150.statcan.gc.ca/n1/pub/82-625-x/2019001/article/00005-eng.htm>
2. Martin BJ, Chen G, Graham M, Quan H. Coding of obesity in administrative hospital discharge abstract data: accuracy and impact for future research studies. *BMC Health Serv Res.* 2014;14(1):70. <https://doi.org/10.1186/1472-6963-14-70>
3. Quan H, Li B, Duncan Saunders L, et al. Assessing validity of ICD-9-CM and ICD-10 administrative data in recording clinical conditions in a unique dually coded database. *Health Serv Res.* 2008;43(4):1424-1441. <https://doi.org/10.1111/j.1475-6773.2007.00822.x>
4. Peng M, Southern DA, Williamson T, Quan H. Under-coding of secondary conditions in coded hospital health data: impact of co-existing conditions, death status and number of codes in a record. *Health Inf J.* 2017;23(4):260-267. <https://doi.org/10.1177/1460458216647089>
5. Woo JG, Zeller MH, Wilson K, Inge T. Obesity identified by discharge ICD-9 codes underestimates the true prevalence of obesity in hospitalized children. *J Pediatr.* 2009;154(3):327-331. <https://doi.org/10.1016/j.jpeds.2008.09.022>
6. Lee S, Doktorchik C, Martin EA, et al. Electronic medical record-based case phenotyping for the Charlson conditions: scoping review. *JMIR Med Inf.* 2021;9(2):e23934. <https://doi.org/10.2196/23934>
7. Chang F, Gupta N. Progress in electronic medical record adoption in Canada. *Can Fam Physician.* 2015;61(12):1076-1084.
8. Casey JA, Schwartz BS, Stewart WF, Adler NE. Using electronic health records for population health research: a review of methods and applications. *Annu Rev Publ Health.* 2016;37(1):61-81. <https://doi.org/10.1146/annurev-publhealth-032315-021353>
9. Kapoor A, Kim J, Zeng X, Harris ST, Anderson A. Weighing the odds: assessing underdiagnosis of adult obesity via electronic medical record problem list omissions. *Digit Health.* 2020;6:2055207620918715. <https://doi.org/10.1177/2055207620918715>

10. Castaneda C, Nalley K, Mannion C, et al. Clinical decision support systems for improving diagnostic accuracy and achieving precision medicine. *J Clin Bioinf.* 2015;5(1):4. <https://doi.org/10.1186/s13336-015-0019-3>
11. Kohn MS, Sun J, Knoop S, et al. IBM's health analytics and clinical decision support. *Yearb Med Inf.* 2014;9(1):154-162. <https://doi.org/10.15265/IY-2014-0002>
12. Banda JM, Seneviratne M, Hernandez-Boussard T, Shah NH. Advances in electronic phenotyping: from rule-based definitions to machine learning models. *Annu Rev Biomed Data Sci.* 2018;1:53-68. <https://doi.org/10.1146/annurev-biodatasci-080917-013315>
13. Shivade C, Raghavan P, Fosler-Lussier E, et al. A review of approaches to identifying patient phenotype cohorts using electronic health records. *J Am Med Inf Assoc.* 2014;21(2):221-230. <https://doi.org/10.1136/amiajnl-2013-001935>
14. Wiebe N, Quan H, Southern DA, Doktorchik C, Eastwood C. Describing agreement in the main condition coding field using Canadian ICD-11 inpatient data. *Int J Popul Data Sci.* 2021;6(1):1397. <https://doi.org/10.23889/ijpds.v6i1.1397>
15. Jurafsky D, Martin JH. *Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition.* 2nd ed. Pearson Prentice Hall; 2009.988
16. spaCy: Industrial-Strength Natural Language Processing. <https://spacy.io/>
17. Regular expression: From Wikipedia, the free encyclopedia. https://en.wikipedia.org/wiki/Regular_expression
18. Stubblebine T. *Regular Expression Pocket Reference.* O'Reilly Media, Inc.; 2003.
19. Canadian Coding Standards for Version 2022 ICD-10-CA and CCI. 2022.
20. Data Quality Documentation: Discharge Abstract Database.
21. Mattar A, Carlston D, Sariol G, et al. The prevalence of obesity documentation in primary care electronic medical records: are we acknowledging the problem? *Appl Clin Inf.* 2017;8(1):67-79. <https://doi.org/10.4338/ACI-2016-07-RA-0115>
22. Roth C, Foraker RE, Payne PRO, Embi PJ. Community-level determinants of obesity: harnessing the power of electronic health records for retrospective data analysis. *BMC Med Inf Decis Making.* 2014;14(1):36. <https://doi.org/10.1186/1472-6947-14-36>

SUPPORTING INFORMATION

Additional supporting information can be found online in the Supporting Information section at the end of this article.

How to cite this article: Sandhu N, Krusina A, Quan H, et al. Automated extraction of weight, height, and obesity in electronic medical records are highly valid. *Obes Sci Pract.* 2024:e705. <https://doi.org/10.1002/osp4.705>