

# BMJ Open Sensitivity and specificity of breast cancer ICD-9-CM codes in three Italian administrative healthcare databases: a diagnostic accuracy study

Iosief Abraha,<sup>1,2</sup> Diego Serraino,<sup>3</sup> Alessandro Montedori,<sup>1</sup> Mario Fusco,<sup>4</sup> Gianni Giovannini,<sup>1</sup> Paola Casucci,<sup>5</sup> Francesco Cozzolino,<sup>1</sup> Massimiliano Orso,<sup>1</sup> Annalisa Granata,<sup>5</sup> Marcello De Giorgi,<sup>5</sup> Paolo Collarile,<sup>6</sup> Rita Chiari,<sup>7</sup> Jennifer Foglietta,<sup>7</sup> Maria Francesca Vitale,<sup>4</sup> Fabrizio Stracci,<sup>8</sup> Walter Orlandi,<sup>9</sup> Ettore Bidoli,<sup>3</sup> The D.I.V.O. Group

**To cite:** Abraha I, Serraino D, Montedori A, *et al*. Sensitivity and specificity of breast cancer ICD-9-CM codes in three Italian administrative healthcare databases: a diagnostic accuracy study. *BMJ Open* 2018;**8**:e020627. doi:10.1136/bmjopen-2017-020627

► Prepublication history for this paper is available online. To view these files, please visit the journal online (<http://dx.doi.org/10.1136/bmjopen-2017-020627>).

Received 14 November 2017  
Revised 25 April 2018  
Accepted 14 May 2018



© Author(s) (or their employer(s)) 2018. Re-use permitted under CC BY-NC. No commercial re-use. See rights and permissions. Published by BMJ.

For numbered affiliations see end of article.

## Correspondence to

Dr Alessandro Montedori;  
[amontedori@regione.umbria.it](mailto:amontedori@regione.umbria.it)

## ABSTRACT

**Objectives** To assess the accuracy of International Classification of Diseases, Ninth Revision, Clinical Modification (ICD-9-CM) codes in identifying patients diagnosed with incident carcinoma in situ and invasive breast cancer in three Italian administrative databases. **Design** A diagnostic accuracy study comparing ICD-9-CM codes for carcinoma in situ (233.0) and for invasive breast cancer (174.x) with medical chart (as a reference standard). **Case definition:** (1) presence of a primary nodular lesion in the breast and (2) cytological or histological documentation of cancer from a primary or metastatic site.

**Setting** Administrative databases from Umbria Region, Azienda Sanitaria Locale (ASL) Napoli 3 Sud (NA) and Friuli VeneziaGiulia (FVG) Region.

**Participants** Women with breast carcinoma in situ (n=246) or invasive breast cancer (n=384) diagnosed (in primary position) between 2012 and 2014.

**Outcome measures** Sensitivity and specificity for codes 233.0 and 174.x.

**Results** For invasive breast cancer the sensitivities were 98% (95% CI 93% to 99%) for Umbria, 96% (95% CI 91% to 99%) for NA and 100% (95% CI 97% to 100%) for FVG. Specificities were 90% (95% CI 82% to 95%) for Umbria, 91% (95% CI 83% to 96%) for NA and 91% (95% CI 84% to 96%) for FVG. For carcinoma in situ the sensitivities were 100% (95% CI 93% to 100%) for Umbria, 100% (95% CI 95% to 100%) for NA and 100% (95% CI 96% to 100%) for FVG. Specificities were 98% (95% CI 93% to 100%) for Umbria, 86% (95% CI 78% to 92%) for NA and 90% (95% CI 82% to 95%) for FVG.

**Conclusions** Administrative healthcare databases from Umbria, NA and FVG are accurate in identifying hospitalised news cases of carcinoma of the breast. The proposed case definition is a powerful tool to perform research on large populations of newly diagnosed patients with breast cancer.

## INTRODUCTION

The use of administrative databases is increasingly growing in various healthcare settings worldwide. These databases anonymously

## Strengths and limitations of this study

- This study is the first to have validated International Classification of Diseases, Ninth Revision, Clinical Modification (ICD-9) codes for incident breast cancer cases in three large computerised administrative databases in Italy using the same case definition.
- Case ascertainment was based on the presence of a primary nodular lesion in the breast documented by imaging and a cytological or histological documentation of cancer from a primary or metastatic site.
- This study followed recommended guidelines based on the criteria published by the Standards for Reporting of Diagnostic accuracy initiative for the accurate reporting of investigations of diagnostic studies.
- The validated ICD-9 codes for non-invasive and invasive cancer are limited to inpatient setting.
- The sample size of women with carcinoma in situ was limited due to the low prevalence of disease.

store data about residents regarding the healthcare assistance they receive including hospital admission, demographic data and disease treatment. Usually, the diagnosis of the disease is associated with a specific code from the International Classification of Diseases, Ninth Revision (ICD-9) or 10th Revision (ICD-10) edition. The ICD is designed to map health conditions to corresponding generic categories together with specific variations.<sup>1</sup> The networking of individual patient data from administrative databases and other sources such as outpatient data and prescription data allows monitoring population health status, performing outcome research<sup>2-4</sup> and exploring a wide range of significant public health questions.<sup>2</sup>

In administrative databases, while non-clinical data such as demographic or prescription data are highly accurate,<sup>5 6</sup> the accuracy of diagnoses and procedures needs to be determined.<sup>6 7</sup> Typically, the assessment of accuracy consists in confirming the reliability of information within the databases with the corresponding clinical records of patients.<sup>5</sup> To reach this goal, the content of administrative healthcare databases needs to be appropriately validated.

In Italy, all the Regional Health Authorities maintain large healthcare information systems containing patient data from all hospital and operative sources. These databases have the potential to address important issues in postmarketing surveillance,<sup>8 9</sup> epidemiology,<sup>10</sup> quality performance and health services research.<sup>11</sup> However, there is a concern that their considerable potential as a source of reliable healthcare information has not been realised since they have not been widely validated.<sup>12</sup>

Breast cancer is the most commonly diagnosed neoplasm in women worldwide, as well as in Italy.<sup>13</sup> Variation in the epidemiology of breast cancers,<sup>14</sup> potential heterogeneity in treatment (pharmacological or surgical) and potential clinical<sup>15 16</sup> and economic outcomes<sup>17–19</sup> can all be evaluated using validated administrative databases. Hence, assessing the accuracy of Italian administrative databases in identifying women with this oncological disease is relevant for the scientific community, the governments, as well as the industry.

As reported in our protocol,<sup>20</sup> the objective of the present study was to evaluate the accuracy of the ICD-9-Clinical Modification (CM) codes related to breast cancer in correctly identifying the respective diseases using three large Italian administrative healthcare databases.

## METHODS

### Setting and data source

#### Administrative databases

The target administrative databases were represented by the Umbria Region (890 000 residents), Local Health Unit 3 of NA (1 170 000 residents) and the Friuli Venezia Giulia (FVG) Region (1 227 000 residents). The corresponding operative units, the Regional Health Authority of Umbria (for Umbria Region), the Registro Tumori Regione Campania (for Napoli Sud Local Health Unit) and the Centro di Riferimento Oncologico Aviano (for FVG Region), conducted the same validation process.

In Italy, regional and local healthcare administrative databases routinely collect data from all patient medical records from public and private hospitals including demographics, hospital admission and discharge dates, vital statistics, the admitting hospital department, the principal diagnosis and a maximum of five secondary discharge diagnoses and the principal, and up to five secondary surgical or pharmacological treatments and diagnostic procedures. Each resident has a unique national identification code with which it is possible to link the various types of information, corresponding to each person, within the database. In Italy, healthcare is covered almost

entirely by the Italian National Health System, therefore, most residents' significant healthcare information can be found within the healthcare databases.

### Source population

The source population was represented by permanent residents aged 18 or above of Umbria Region, Local Health Unit 3 of NA and FVG Region. Any resident that has been discharged from hospital with a diagnosis of breast cancer was considered. Residents that have been hospitalised outside the regional territory of competence were excluded from analysis.

### Patient and public involvement

Patients were not directly involved. This was a retrospective study based on the consultation of medical charts. .

### Case and control selection and sampling method

In each administrative database, patients with the first occurrence of diagnosis of breast cancer between 1 January 2012 and 31 December 2014 were identified using the following ICD-9-CM codes (index test) located in primary position: (1) 233.0 for carcinoma in situ of the breast and (2) 174.x for invasive breast cancer. Estimated prevalent cases, that is, cases with the same diagnosis (ICD-9-CM codes in any position) in the 5 years (2007–2011) before the period of interest, were excluded.

For controls, within the same period, non-cases, that is, 94 female patients having in primary position a diagnosis of cancer (ICD-9 140–239) other than invasive breast cancer (ICD-9 174.0–174.9) or carcinoma in situ of the breast (ICD-9 233.0), were randomly selected.

Subsequently, for each of the above reported groups of ICD-9-CM codes, random samples of cases and non-cases were selected from each administrative database.

### Chart abstraction and case ascertainment

The medical charts of the randomly selected samples of cases and non-cases were obtained from hospitals for the validation task (considered as the reference standard). From each medical chart, the following information was retrieved: clinical chart number, hospital and ward, date of birth, sex, dates of hospital admission and discharge, signs and symptoms, any diagnostic procedures that contributed to the diagnosis of the cancer, any pharmacological or surgical interventions that were provided for the treatment of the cancer.

Within each unit, two reviewers received training on data abstraction and performed an initial consensus chart review, independently examining the same number of medical charts (n=20). The inter-rater agreement regarding the presence or absence of breast cancer among the pairs of reviewers within each unit was near perfect ( $\kappa$  statistics >0.91).

Case ascertainment of cancer within medical charts was based on (1) the presence of a primary nodular lesion in the breast documented by imaging and (2) the cytological or histological documentation of cancer from a primary or metastatic site.

Following the consensus review, data abstraction was completed independently. To ensure consistency among all the reviewers, cases with uncertainty were discussed and resolved through the involvement of an oncologist (RC).

#### Validation criteria

For non-invasive breast cancer, we considered the ICD-9-CM code 233.0 valid when there is evidence of a breast nodule documented with imaging (eg, mammography) and a histological diagnosis of ductal or lobular breast carcinoma in situ (pTis).

For invasive breast cancer, we considered the ICD-9-CM codes 174.x valid when there is evidence of a breast nodule documented with imaging (eg, mammography) and a cytological or histological diagnosis from a primary or metastatic site positive for ductal or lobular adenocarcinoma.

#### Statistical analysis

We calculated that a sample of 130 charts of cases was necessary to obtain an expected sensitivity of 80% with a 95% CI of 72% to 86% according to binomial exact calculation.<sup>21</sup> For specificity calculation, we randomly selected non-cases, that is, records without the ICD-9 codes of interest from administrative database. For controls, we calculated a sample of 94 charts of non-cases was necessary to obtain an expected specificity of 90% with a 95% CI of 83% to 95% according to binomial exact calculation.<sup>21</sup>

Sensitivity and specificity with relative 95% CIs were calculated separately for each ICD-9-CM code by constructing 2×2 tables.

When missing medical charts occurred, we performed a formal sensitivity analysis based on a worst-case scenario in which the missing cases were considered as false positive and the missing-non cases were considered false negative.

## RESULTS

### Invasive breast cancer

After excluding the estimated prevalent cases from the diagnosis of invasive breast cancer in the primary position of hospital discharges, the cases were 2686, 2044 and 2107 from Umbria, NA and FVG, respectively. Subsequently, each team randomly sampled 130 cases of which the corresponding medical charts were requested for evaluation. Two and four medical charts were not available from Umbria and NA, respectively. [Figure 1](#) displays the identification of cases from the three operative units. For the non-cases, that is, female patients having a diagnosis of cancer (ICD-9 140–239) other than invasive breast cancer (ICD-9 174.0–174.9), each unit randomly selected 94 medical charts. One medical chart of non-cases from Umbria was missing.

The most common ICD-9-CM subgroup was the code 174.4, that is upper-outer quadrant breast cancer, accounting for 45% of cases for Umbria, 34% for NA and

35% for FVG. The mean age ranged between 61 and 66 years. The majority of the cases were identified in surgical departments with a percentage ranging from 84% to 94%. The types of surgical intervention were similar across the three operative units with quadrantectomy being the most reported surgical intervention. [Table 1](#) describes the basic characteristics of the incident cases across the three units. The sensitivities were 98% (95% CI 93% to 99%) for Umbria, 96% (95% CI 91% to 99%) for NA and 100% (95% CI 97% to 100%) for FVG. The specificities were 90% for Umbria, and 91% for NA and FVG. Accuracy results with their CIs are displayed in [figure 2](#).

In terms of misclassification, overall there were 28 cases that were considered false positives. The reasons for this misclassification were: histological documentation missing in the medical chart (six in Umbria, eight in NA and eight in FVG) and negative histology for invasive breast cancer (four in Umbria, one in NA and one in FVG). However, of these false positive cases with negative histology for invasive breast cancer, three were positive for breast carcinoma in situ diagnosis. Conversely, there were eight non-cases that were judged false negatives: two were possible breast cancer diagnosis and six were metastatic breast diagnoses ([table 2](#)).

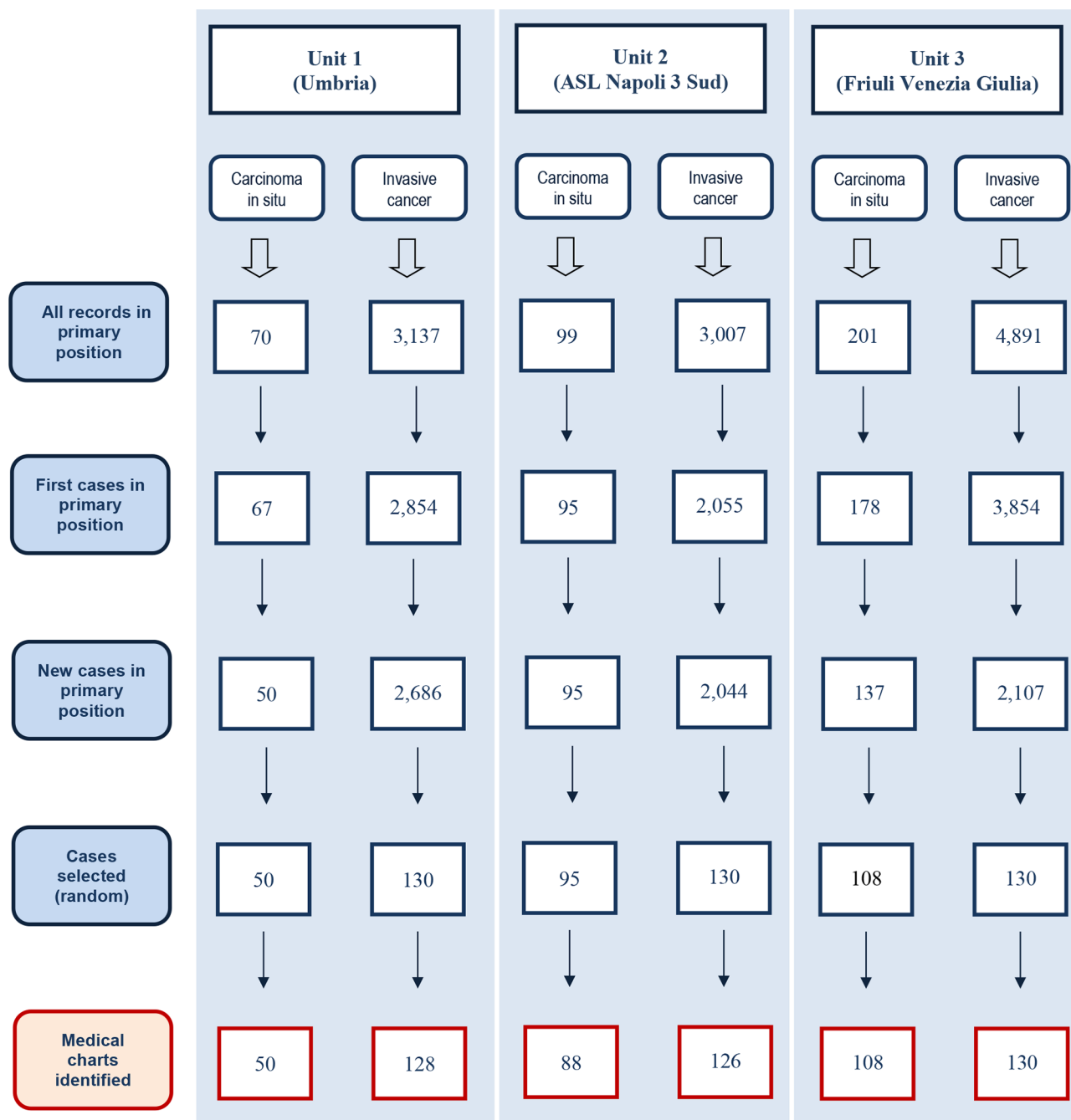
Overall, there were six missing charts: two in Umbria and four in NA. Worst-case scenario in the sensitivity analysis showed that specificity was affected marginally: it changed from 90% to 88% for Umbria and from 91% to 87% for NA. The differences between the ordinary results and the worst-case scenario analysis were not statistically significant.

### Breast carcinoma in situ

The incident cases of carcinoma in situ of the breast were 50 from Umbria, 95 from NA and 137 from FVG, from which 50, 95 and 108 were randomly selected and the corresponding medical charts were requested for assessment ([figure 1](#)). Seven charts from NA were not available. For the non-cases, that is, female patients having a diagnosis of cancer (ICD-9 140–239) other than carcinoma in situ of the breast (ICD-9 233.0), each unit randomly selected 94 medical charts. One medical chart of non-cases from Umbria was missing.

The mean age ranged between 57 (NA) and 60 years (FVG). Most of the cases were identified in surgical departments with a percentage ranging from 92% to 100%. The types of surgical intervention were similar across the three operative units with quadrantectomy being the most reported surgical intervention. [Table 1](#) describes the basic characteristics of the incident cases across the three units.

After reviewing the medical records, 100% (48/48) of the patients with carcinoma in situ of the breast from Umbria, 100% (73/73) from NA and 100% (97/97) from FVG were correctly identified by the administrative databases. The specificities were 98% for Umbria, 86% for NA and 90% for FVG. Accuracy results with their CIs are displayed in [figure 2](#).



**Figure 1** Flow chart of incident cases identification using the administrative databases and the corresponding charts identified and examined.

Table 2 describes the reasons for misclassification of cases and controls. Overall, there were 28 cases that were judged false positives. The reasons were histological documentation missing in the medical charts (eight in NA and two in FVG), and negative histology for carcinoma in situ of the breast (two in Umbria, seven in NA and nine in FVG). None of the controls resulted a false negative.

The sensitivity analysis showed that specificity for NA codes reduced to 81% (95% CI 73% to 88%) due to the seven charts of missing cases but the difference was not however statistically significant.

## DISCUSSION

We developed a case definition of breast cancer based on the presence of a primary nodular lesion in the breast documented with imaging and the cytological or histological documentation of cancer from a primary or metastatic site and the performance of the model was evaluated in terms of sensitivities and specificities for the three administrative databases. After revising the medical charts, the results showed that both codes (233.0 and 174.x) performed well in identifying new cases of hospitalised women with breast carcinoma in situ and invasive breast cancer, respectively.

**Table 1** Characteristics of patients with breast cancer who were identified in the three administrative healthcare databases

Characteristics	Unit 1 (Umbria)	Unit 2 (ASL Napoli 3 Sud)	Unit 3 (Friuli Venezia Giulia)
<b>Invasive carcinoma</b>			
Incident cases (N medical chart reviewed)	128	126	130
<b>ICD-9 code</b>			
174.0 nipple and areola	–	1 (1)	–
174.1 central portion	16 (13)	10 (8)	6 (5)
174.2 upper-inner quadrant	4 (3)	8 (6)	14 (11)
174.3 lower-inner quadrant	6 (5)	5 (4)	9 (7)
174.4 upper-outer quadrant	57 (45)	43 (34)	45 (35)
174.5 lower-outer quadrant	6 (5)	5 (4)	13 (10)
174.6 axillary tail	–	–	–
174.8 other specified sites of the female breast	–	–	–
174.9 breast female, unspecified	7 (5)	38 (30)	34 (26)
174.0 nipple and areola	32 (25)	16 (13)	9 (7)
<b>Admission to department</b>			
Medical	20 (16)	11 (9)	8 (6)
Surgical	108 (84)	115 (91)	122 (94)
<b>Age, N (%)</b>			
<40	9 (7)	6 (5)	1 (1)
40–59	40 (31)	56 (44)	45 (35)
≥60	79 (62)	64 (51)	84 (65)
<b>Instrumental diagnosis</b>			
Breast ultrasound	39 (30)	88 (70)	5 (4)
Mammography	54 (42)	60 (48)	7 (5)
CT scan (breast)	11 (8)	1 (1)	2 (2)
MRI (breast)	3 (2)	17 (14)	8 (6)
<b>Surgical procedures</b>			
Mastectomy	28 (22)	29 (23)	35 (27)
Quadrantectomy	79 (62)	73 (58)	54 (42)
<b>Histological documentation</b>			
Needle aspiration	32	34	–
Needle biopsy	27	40	5
Nodule (after surgical intervention)	115	117	112
<b>Carcinoma in situ</b>			
Incident cases (N medical chart reviewed)	50	88	108
<b>ICD-9 code</b>			
233.0	50 (100)	88 (100)	108 (100)
<b>Admission to department</b>			
Medical	–	7 (8)	–
Surgical	50 (100)	81 (92)	108 (100)
<b>Age, N (%)</b>			
<40	2 (4)	3 (3)	1 (1)
40–59	27 (54)	50 (57)	52 (48)
>60	21 (42)	35 (40)	55 (51)

Continued

Table 1 Continued

Characteristics	Unit 1 (Umbria)	Unit 2 (ASL Napoli 3 Sud)	Unit 3 (Friuli Venezia Giulia)
Instrumental diagnosis			
Breast ultrasound	3 (6)	65 (74)	30 (28)
Mammography	6 (12)	37 (42)	39 (36)
CT scan (breast)	–	–	–
MRI (breast)	2 (4)	22 (25)	4 (4)
Surgical procedures			
Mastectomy	18 (36)	13 (15)	15 (14)
Quadrantectomy	32 (64)	47 (53)	55 (52)
Histological documentation			
Needle aspiration	–	16 (18)	1 (1)
Needle biopsy	8 (16)	53 (60)	4 (4)
Nodule (after surgical intervention)	50 (100)	77 (88)	94 (87)

ICD-9, International Classification of Diseases, Ninth Revision.

Previously, researchers have assessed the accuracy of breast cancer diagnosis in administrative databases using different algorithms and in most cases using registry data as a reference standard.<sup>22</sup>

In 2008, an Italian study developed and validated an algorithm using a regional administrative database to determine incident cases of breast, lung and colorectal cancers.<sup>23</sup> The study found a sensitivity of 77% for breast cancer<sup>23</sup> and the lower value of sensitivity compared with our results can be attributed to the fact that in Baldi *et*

*al* the validity of the algorithm for each cancer site was assessed by individual matching between cases in hospital discharge and the Piedmont Cancer Registry or because authors were interested in high values of positive predictive value (PPV). Using the Surveillance, Epidemiology and End Results (SEER) database as a reference standard, Freeman *et al* developed an approach for identifying incident breast cancer cases based on a logistic regression model, which contained variables that indicate the presence of breast cancer-related diagnoses and procedures

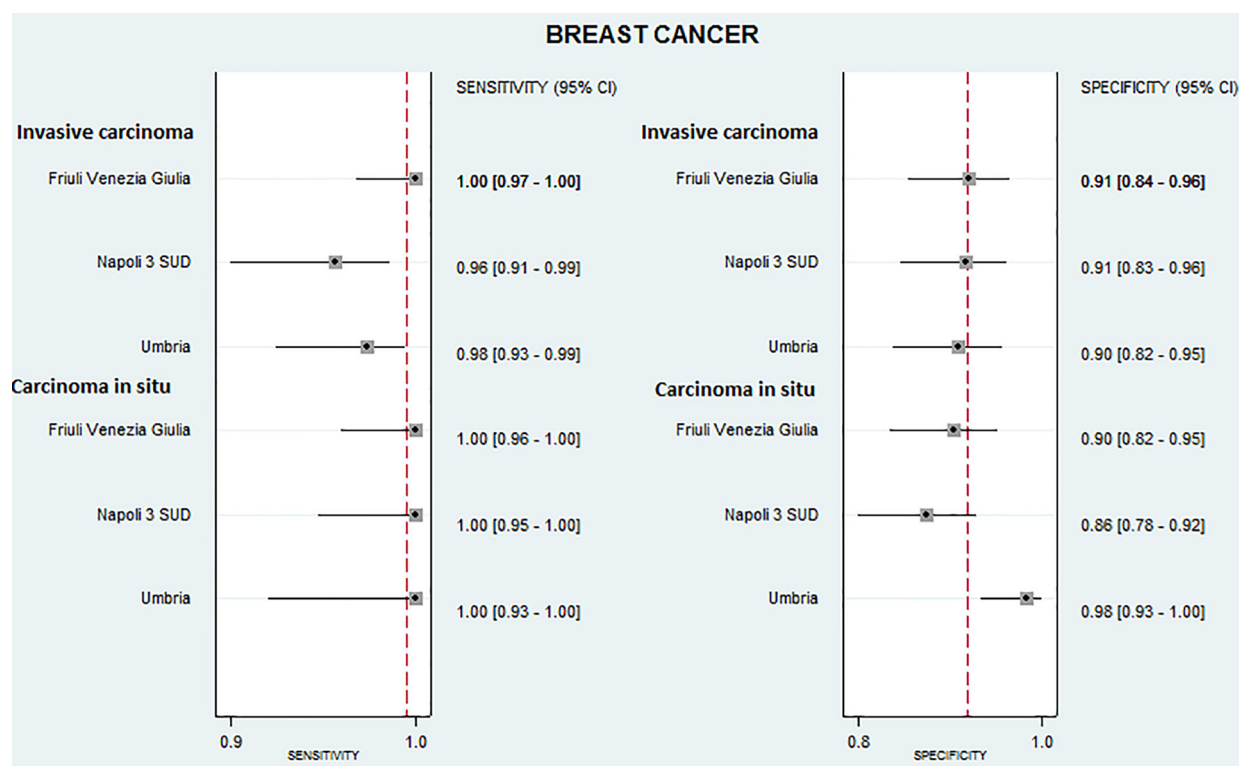


Figure 2 Sensitivity and specificity results for ICD-9-CM codes related to breast carcinoma in situ and invasive breast cancer for the three administrative databases. ICD-9-CM, International Classification of Diseases, Ninth Revision, Clinical Modification.

**Table 2** Reason for misclassification of cases and controls

Type of misclassification		Unit 1 (Umbria)	Unit 2 (ASL Napoli 3 Sud)	Unit 3 (Friuli Venezia Giulia)
<b>Invasive breast cancer</b>				
False positives				
1	Histological examination missing	6	8	8
2	Negative histology	4	1	1
	(1) Carcinoma in situ	2		1
Total		10	9	9
False negative				
1	Possible breast cancer relapse	1	1	–
2	Metastatic breast cancer	2	4	–
Total		3	5	–
<b>Breast carcinoma in situ</b>				
False positives				
1	Histological examination missing	–	8	2
2	Negative histology	2	7	9
Total		2	15	11
False negative				
1	Possible carcinoma in situ	–	–	–

in three sources of claims data: hospital inpatient stays, hospital outpatient services and physician services. The Receiver operating characteristic (ROC) curve showed that the model achieved over 90% of sensitivity and specificity although a lower PPV.<sup>24</sup> Similarly, using the Medicare database, Setoguchi *et al* developed a case definition based on diagnoses and procedure codes and compared it with SEER database in Pennsylvania, USA. The authors obtained a sensitivity and specificity for identifying breast cancer cases of 83% and 99%, respectively.<sup>25</sup> Using the cancer registry data as reference standard, Kemp *et al*<sup>26</sup> evaluated Australian administrative and self-reported datasets to identify cases of invasive breast cancer and used several combinations of diagnoses and procedures obtaining the highest sensitivity and PPV (both 86%) when a flag of 'diagnosis of invasive breast cancer' was used. A systematic review of administrative databases that validated breast cancer is currently being completed and will provide a complete account of validation of administrative databases worldwide.<sup>22</sup>

### Strength and limitation

Strengths of our study include complete transparency based on prepublication of a protocol, the use of detailed and explicit eligibility criteria, and the use of duplicate and independent processes for medical chart review and abstraction following recommended guidelines based on the criteria published by the Standards for Reporting of Diagnostic accuracy (STARD) initiative for the accurate reporting of investigations of diagnostic studies.<sup>27–29</sup> In addition, we used as a required element for validation the presence of a histological or cytological documentation.

Unlike several studies that used Cancer Registries to validate breast cancer codes, in the present study medical records were reviewed directly to evaluate the accuracy of potential cases obtained from the administrative database. Generally, medical charts are considered the gold standard for the diagnosis of a disease. Cancer registries are considered also the gold standard as they produce data over decades and increasingly include data sources,<sup>30</sup> which allow complete registration of cases treated in an outpatient setting. This is relevant for the exclusion of prevalent cases and multiple primaries as well as for particular cancer sites or patients (eg, oldest old) but at a higher cost.<sup>31</sup>

Conversely, the information of administrative healthcare databases is generally limited to the information obtained from the hospital discharge register that contain the primary and secondary diagnoses, the surgical or invasive diagnostic approaches performed as well as chemotherapy or radiotherapy. When adequately validated, these databases can provide important data such length of stay and related costs,<sup>31</sup> important outcomes such as adverse events<sup>32 33</sup> as well as variations in healthcare resource utilisation.<sup>34</sup> Indeed, administrative databases are readily available for the whole of Italy, whereas cancer registry data are not. Thus, our proposed validation method using a well-defined case definition can be a good alternative in settings in which cancer registries are not available in Italy.<sup>20 35 36</sup> Our study confirms that hospital discharge data can be used for some specific purposes (eg, identification of breast cancer cases treated at a given hospital in a study on caseload). For other aims we would

recommend further refinement, even if the validity of the cancer coding is valid (eg, to provide reliable estimates of breast cancer incidence hospital discharge data should be available for many years and possibly complete anatomic pathology archives should be linked to it too).

We acknowledge specific limitations to our study. The overall number of carcinoma in situ cases identified in the three units during the period of interest was below the calculated sample size and we are unsure whether this limitation can affect the results of sensitivity and specificity for the breast carcinoma in situ ICD-9-CM code. Indeed, within the figure of the overall breast cancer diagnoses, carcinoma in situ cases diagnosed within hospitals are under-represented and any future epidemiological assessment will need to take it into account trying possibly to clarify the reason.

In addition, our assessment was limited to hospitalised patients and does not consider new cases of cancer who had the diagnoses in day hospital or day surgery. Although these cases are limited (eg, 16 carcinoma in situ cases diagnosed in day surgery or day hospital in Umbria across the 3 years and 3.8% invasive breast cancer diagnosed in day surgery or day hospital in Umbria across the 3 years), further research can be addressed the validity of ICD-9 codes in outpatient setting.

In addition, we had a higher false positive rate than false negative rate. The number of false positives is due to our stringent case ascertainment criteria, that is, the presence in the clinical chart of both imaging and histological documentation of breast cancer within the same medical chart. Twenty-two false positives cases for the invasive breast cancer and 10 false positives cases for the carcinoma in situ were due to histological documentation missing in the medical chart (table 2). This does not necessarily mean that the subjects were without the diagnosis of cancer. These cases had other elements in the medical chart such as imaging, chemotherapy or radiotherapy, that could confirm the presence of the disease. Should we have used broader case ascertainment criteria, we could have obtained a lower false positive rate. Estimates of specificity for invasive carcinoma in our three databases were high ranging from 90% to 91%. For carcinoma in situ specificity was acceptable for Napoli 3 Sud (86%) and high for FVG and Umbria (90% and 98% respectively).

In terms of generalisability, despite the success of validation processes of administrative database, any conclusion that stems from these validated database could not be generalised in other settings.

As declared in our protocol, we favoured the estimation of sensitivity and specificity rather than predictive values since predictive values are dependent on the prevalence of the disease. However, to comply with the STARD guidelines, we provide absolute numbers for true or false case and non-cases from which it is possible to obtain predictive values (table 3).

**Table 3** Cross-tabulation of the index test (International Classification of Diseases, Ninth Revision, Clinical Modification; ICD-9-CM code) results by the results of the reference standard (medical chart)

Type of breast cancer (ICD-9-CM)	Operative unit	True Positive (TP)	False Positive (FP)	True Negative (TN)	False Negative (FN)
Invasive cancer (174.x)	Unit 1 (Umbria)	118	10	90	3
	Unit 2 (ASL Napoli 3 Sud)	117	9	89	5
	Unit 3 (Friuli Venezia Giulia)	121	9	94	0
Carcinoma in situ (233.0)	Unit 1 (Umbria)	48	2	93	0
	Unit 2 (ASL Napoli 3 Sud)	73	15	94	0
	Unit 3 (Friuli Venezia Giulia)	97	11	94	0

## CONCLUSION

In summary, the present study has demonstrated that administrative healthcare databases from Umbria, NA and FVG can be used to identify hospitalised women with newly diagnosed invasive or in situ carcinoma of the breast. The proposed case definition in the present study provides a powerful tool to perform outcome research on breast cancer based on a population of three million residents. Potential implication of this proposal includes the extension of this case definition to other Italian regional healthcare databases and the combination with other sources of data (such as prescription database) to conduct efficiently quality of care, healthcare research and pharmacoepidemiological studies that may complement randomised trials.

### Author affiliations

<sup>1</sup>Health Planning Service, Regional Health Authority of Umbria, Perugia, Italy

<sup>2</sup>Innovation and Development, Agenzia Nazionale per i Servizi Sanitari Regionali (Age.Na.S.), Rome, Italy

<sup>3</sup>Cancer Epidemiology Unit, Centro di Riferimento Oncologico Aviano, Aviano, Italy

<sup>4</sup>Registro Tumori Regione Campania, ASL Napoli 3 Sud, Brusciano, Italy

<sup>5</sup>Health ICT Service, Regional Health Authority of Umbria, Perugia, Italy

<sup>6</sup>SOC Epidemiologia Oncologica, Centro di Riferimento Oncologico Aviano, Aviano, Italy

<sup>7</sup>Dipartimento di Oncologia, Azienda Ospedaliera Perugia, Perugia, Italy

<sup>8</sup>Public Health Department, University of Perugia, Perugia, Italy

<sup>9</sup>Direzione Sanità, Regional Health Authority of Umbria, Perugia, Italy

**Collaborators** David Franchini; Giuliana Alessandrini; Roberto Cirocchi; Valerio Ciullo; Michele Gobato; Paolo Eusebi; Chiara Grisci.

**Contributors** AM, IA, MF and DS conceived the original idea of the study. IA, DS, AM, MF, EB, GG, FC, MO and WO designed the study. PCa, MDG, PCo, AG, MFV and VC identified the cohort using administrative database with the supervision of WO, EB, DS, MF, AM and FS. IA, FC, MO, AG, PCo, VC, MFV and JF undertook the data



abstraction with the supervision of AM, GG, WO, FS, MF, EB, RC and DS. IA, RC, AM and JF performed case ascertainment. IA, AM, FC, EB, MF, PE and MO performed the analysis. DS, GG, PCa, AG, MDG, PCo, RC, JF, MFV, FS, EB and WO helped in the interpretation of the data. The initial draft of the manuscript was prepared by IA, AM, FC and MO. DS, EB, GG, PCa, AG, MDG, PCo, RC, JF, MFV, FS and WO revised critically the manuscript for important intellectual content. All the authors read and approved the final manuscript. AM, MF and EB are the guarantors of the data for the respective operative units.

**Funding** This study was developed within the D.I.V.O. project (Realizzazione di un Database Interregionale Validato per l'Oncologia quale strumento di valutazione di impatto e di appropriatezza delle attività di prevenzione primaria e secondaria in ambito oncologico) supported by funding from the National Centre for Disease Prevention and Control (CCM 2014), Ministry of Health, Italy.

**Competing interests** None declared.

**Patient consent** Not required.

**Ethics approval** Ethical approval for the present study was obtained from the Ethics Committee of the Umbria Region Health Authority (CEAS).

**Provenance and peer review** Not commissioned; externally peer reviewed.

**Data sharing statement** No additional data are available.

**Open access** This is an open access article distributed in accordance with the Creative Commons Attribution Non Commercial (CC BY-NC 4.0) license, which permits others to distribute, remix, adapt, build upon this work non-commercially, and license their derivative works on different terms, provided the original work is properly cited, appropriate credit is given, any changes made indicated, and the use is non-commercial. See: <http://creativecommons.org/licenses/by-nc/4.0/>.

## REFERENCES

- World Health Organization. *International statistical classification of diseases and health related problems*. 10th edn. Geneva: WHO, 1992.
- Jutte DP, Roos LL, Brownell MD. Administrative record linkage as a tool for public health research. *Annu Rev Public Health* 2011;32:91–108.
- Jick SS. Fresh evidence confirms links between newer contraceptive pills and higher risk of venous thromboembolism. *BMJ* 2015;350:h2422.
- García Rodríguez LA, Pérez-Gutthann S, Jick S. The UK General practice research database. *Pharmacoepidemiology: John Wiley & Sons, Ltd* 2002:375–85.
- Rawson NSB, Shatin D. Assessing the validity of diagnostic data in large administrative healthcare utilization databases. In: Hartzema A, Tilson H, Chan K, eds. *Pharmacoepidemiology and therapeutic risk management*: Harvey Whitney Books, 2008.
- West SL, Ritchey ME, Poole C. Validity of pharmacoepidemiologic drug and diagnosis data. *Pharmacoepidemiology: Wiley-Blackwell* 2012:757–94.
- Campbell SE, Campbell MK, Grimshaw JM, et al. A systematic review of discharge coding accuracy. *J Public Health Med* 2001;23:205–11.
- Traversa G, Bianchi C, Da Cas R, et al. Cohort study of hepatotoxicity associated with nimesulide and other non-steroidal anti-inflammatory drugs. *BMJ* 2003;327:18–22.
- Trifirò G, Patadia V, Schuemie MJ, et al. EU-ADR healthcare database network vs. spontaneous reporting system database: preliminary comparison of signal detection. *Stud Health Technol Inform* 2011;166:25–30.
- Gini R, Francesconi P, Mazzaglia G, et al. Chronic disease prevalence from Italian administrative databases in the VALORE project: a validation through comparison of population estimates with general practice databases and national survey. *BMC Public Health* 2013;13:15.
- Colais P, Pinnarelli L, Fusco D, et al. The impact of a pay-for-performance system on timing to hip fracture surgery: experience from the Lazio Region (Italy). *BMC Health Serv Res* 2013;13:393.
- Abraha I, Montedori A, Eusebi P, et al. The current state of validation of administrative healthcare databases in Italy: a systematic review. *Pharmacoepidemiology and Drug Safety* 2012;21:400–00.
- Ferlay J, Soerjomataram I, Dikshit R, et al. Cancer incidence and mortality worldwide: sources, methods and major patterns in GLOBOCAN 2012. *Int J Cancer* 2015;136:E359–86.
- Yuan Y, Li M, Yang J, et al. Using administrative data to estimate time to breast cancer diagnosis and percent of screen-detected breast cancers – a validation study in Alberta, Canada. *Eur J Cancer Care* 2015;24:367–75.
- Lambertini M, Anserini P, Fontana V, et al. The PREgnancy and FERtility (PREFER) study: an Italian multicenter prospective cohort study on fertility preservation and pregnancy issues in young breast cancer patients. *BMC Cancer* 2017;17:346.
- De Placido S, De Angelis C, Giuliano M, et al. Imaging tests in staging and surveillance of non-metastatic breast cancer: changes in routine clinical practice and cost implications. *Br J Cancer* 2017;116:821–7.
- Dehal A, Abbas A, Johna S. Comorbidity and outcomes after surgery among women with breast cancer: analysis of nationwide in-patient sample database. *Breast Cancer Res Treat* 2013;139:469–76.
- Konski A. Clinical and economic outcomes analyses of women developing breast cancer in a managed care organization. *Am J Clin Oncol* 2005;28:51–7.
- Mittmann N, Liu N, Porter J, et al. Utilization and costs of home care for patients with colorectal cancer: a population-based study. *CMAJ Open* 2014;2:E11–17.
- Abraha I, Serraino D, Giovannini G, et al. Validity of ICD-9-CM codes for breast, lung and colorectal cancers in three Italian administrative healthcare databases: a diagnostic accuracy study protocol. *BMJ Open* 2016;6:e010547.
- Wilson EB. Probable inference, the law of succession, and statistical inference. *J Am Stat Assoc* 1927;22:209–12.
- Abraha I, Giovannini G, Serraino D, et al. Validity of breast, lung and colorectal cancer diagnoses in administrative databases: a systematic review protocol. *BMJ Open* 2016;6:e010409.
- Baldi I, Vicari P, Di Cuozzo D, et al. A high positive predictive value algorithm using hospital administrative data identified incident cancer cases. *J Clin Epidemiol* 2008;61:373–9.
- Freeman JL, Zhang D, Freeman DH, et al. An approach to identifying incident breast cancer cases using Medicare claims data. *J Clin Epidemiol* 2000;53:605–14.
- Setoguchi S, Solomon DH, Glynn RJ, et al. Agreement of diagnosis and its date for hematologic malignancies and solid tumors between medicare claims and cancer registry data. *Cancer Causes Control* 2007;18:561–9.
- Kemp A, Preen DB, Saunders C, et al. Ascertaining invasive breast cancer cases; the validity of administrative and self-reported data sources in Australia. *BMC Med Res Methodol* 2013;13:17.
- Benchimol EI, Manuel DG, To T, et al. Development and use of reporting guidelines for assessing the quality of validation studies of health administrative data. *J Clin Epidemiol* 2011;64:821–9.
- De Coster C, Quan H, Finlayson A, et al. Identifying priorities in methodological research using ICD-9-CM and ICD-10 administrative data: report from an international consortium. *BMC Health Serv Res* 2006;6:77.
- Bossuyt PM, Reitsma JB, Bruns DE, et al. Towards complete and accurate reporting of studies of diagnostic accuracy: The STARD Initiative. *Ann Intern Med* 2003;138:40–4.
- Gliklich RE, Dreyer NA, Leavy MB. Registries for evaluating patient outcomes: a user's guide. 2014.
- Rice HE, Englum BR, Gulack BC, et al. Use of patient registries and administrative datasets for the study of pediatric cancer. *Pediatr Blood Cancer* 2015;62:1495–500.
- Rodrigo-Rincon I, Martin-Vizcaino MP, Tirapu-Leon B, et al. Validity of the clinical and administrative databases in detecting post-operative adverse events. *Int J Qual Health Care* 2015;27:267–75.
- Rashid N, Koh H, Baca H, et al. Economic burden related to chemotherapy-related adverse events in patients with metastatic breast cancer in an integrated health care system. *Breast Cancer: Targets and Therapy* 2016;8:173–81.
- Mahar AL, Coburn NG, Kagedan DJ, et al. Regional variation in the management of metastatic gastric cancer in Ontario. *Curr Oncol* 2016;23:250–7.
- Orso M, Serraino D, Abraha I, et al. Validating malignant melanoma ICD-9-CM codes in Umbria, ASL Napoli 3 Sud and Friuli Venezia Giulia administrative healthcare databases: a diagnostic accuracy study. *BMJ Open* 2018;8:e020631.
- Montedori A, Bidoli E, Serraino D, et al. Accuracy of lung cancer ICD-9-CM codes in Umbria, Napoli 3 Sud and Friuli Venezia Giulia administrative healthcare databases: a diagnostic accuracy study. *BMJ Open* 2018;8:e020628.