

National Cancer Registration and Analysis Service's Cancer Analysis System (CAS)-SOP#1

Counting cancer cases

About the NDRS

The National Disease Registration Service (NDRS) is part of NHS Digital (NHSD). Its purpose is to collect high-quality, timely data on cancer, rare diseases and congenital anomalies to monitor changes in the health of the population.

The NDRS includes:

- the National Cancer Registration and Analysis Service (NCRAS) and
- the National Congenital Anomaly and Rare Disease Registration Service (NCARDRS)

Healthcare professionals, researchers and policy makers use data to better understand population health and disease. The data is provided by patients and collected by the NHS as part of their care and support. The NDRS uses the data to help:

- understand cancer, rare diseases, and congenital anomalies
- improve diagnosis
- plan NHS services
- improve treatment
- evaluate policy
- improve genetic counselling



National Disease Registration Service
NHS Digital (NHSD)
The Leeds Government Hub
7&8 Wellington Place
Leeds
LS1 4AP

For queries relating to this document, please contact:

NDRSenquiries@nhs.net

Improving lives with data and technology – NHS Digital support NHS staff at work, help people get the best care, and use the nation's health data to drive research and transform services.

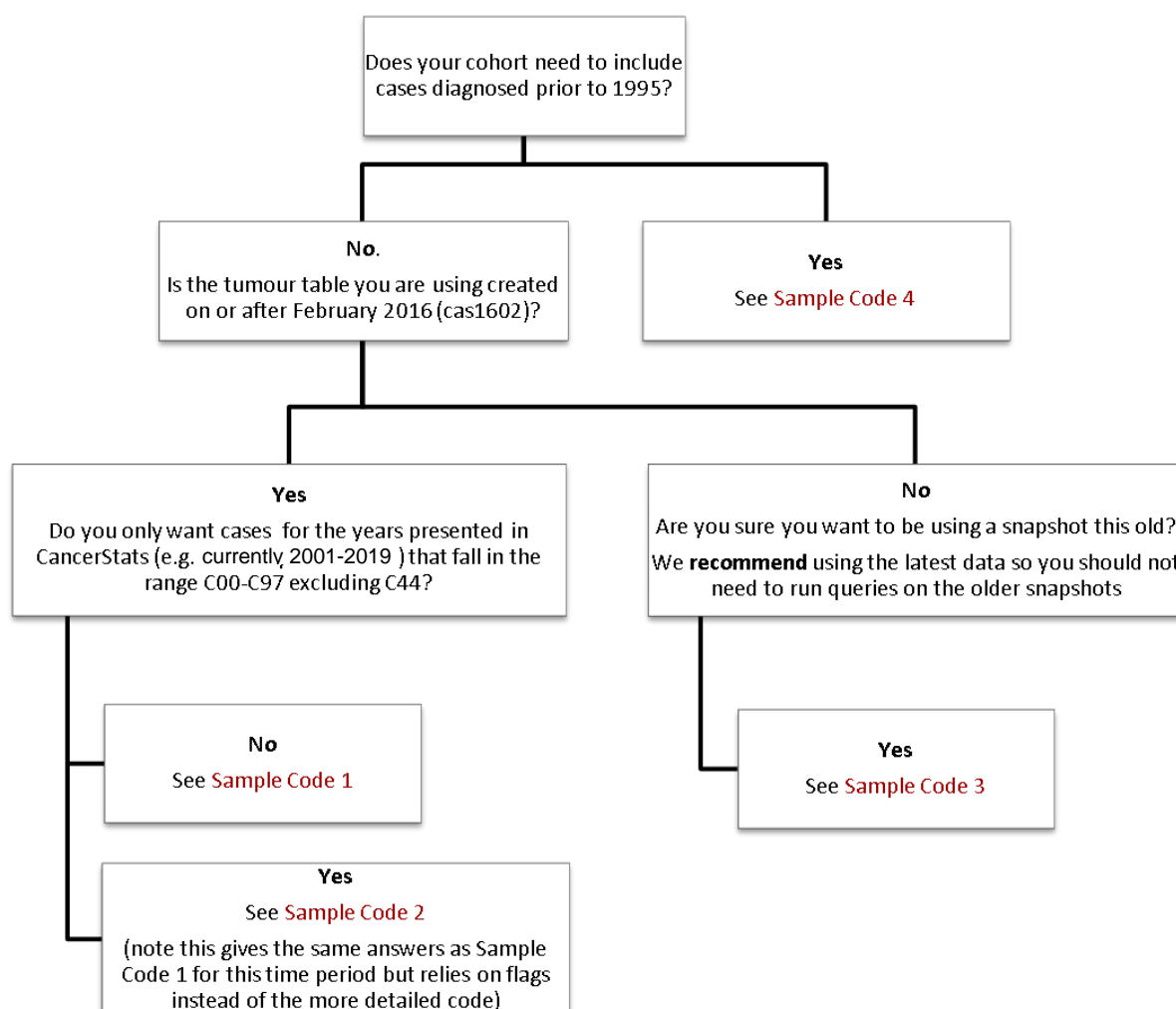


Contents

About the NDRS	1
Contents	2
Introduction	3
For cases diagnosed in or after 1995	4
Sample code 1	6
Sample code 2	7
Sample code 3	8
For cases diagnosed in or after 1995	9
Sample code 4	10

Introduction

Every year, data on over 300,000 cases of cancer are systematically registered by the National Cancer Registration and Analysis Service (NCRAS). These registrations include details on the patient, their type of cancer, how advanced it is and the treatment they receive. This Standard Operating Procedure (SOP) covers the process for counting cancer cases and extracting data on cancer incidence from the NCRAS Cancer Analysis System (CAS). It is not mandatory and depending on the nature of the request, a different approach may be adopted. It exists to outline a set of rules that can be followed to produce consistent and replicable results. The method used will depend on the diagnosis years you are interested in and the iteration of the cancer registration dataset you intend to use. The flow diagram will guide you to the relevant sample code. Note you may need to use combinations of the code depending on your project.



For cases diagnosed in or after 1995

This SOP exists to outline the suggested exclusions that should be applied to the tumour table in CAS in order to define cancer cases for use in cancer incidence statistics.

1) People who are resident outside England

The SOP recommends only including residents in England. In most cases using the tables with the suffix '_england' will capture these as these tables are England residents only. If these tables are not available selecting records with a country code of E will work instead. In earlier datasets where the country code does not exist LSOA codes beginning with 'E' can be used instead.

2) Cancer cases that the registration officers have not finalised

Provisional cases are registered but not confirmed to be cancer until they are finalised so important details about the cancer case may be subsequently added. Therefore, this SOP recommends only using finalised cases. This may not be the best course of action in all projects and there may be good reasons to include the provisional cases in some analyses.

3) Cases that are considered to be duplicate records

The English cancer registration system holds data from the 8 former regional databases. A lot of work has been done to deduplicate these datasets but there are still duplicates in the data before 2012. The dedup_flag was developed to flag up records identified as duplicate records. Separate documentation is available for this field but briefly the flag takes account the following issues. For tumours diagnosed between 1995 and 2011, only those that can be traced in the 2013 ONS data will be counted. Cases sent late to ONS with a valid ONS ID are also included as cases.

A small quality issue with the dedup flag occurs when no ONS ID is available and some cases are potentially identified as duplicates in error. This issue applies to cancer registrations diagnosed in areas covered by West Lancashire CCG and, to a lesser extent, Eastern Cheshire, South Cheshire, and Vale Royal. This will also affect local authority data and regional data that cover this area. Cancer cases may be missing from the cohort of patients diagnosed prior to 2008. This only relates to a relatively small number of cases and so will not impact greatly on national figures, but may mean that cancer incidence is significantly underestimated in these areas.

4) Cases with suspected incorrect age at diagnosis

This SOP recommends including records of patients aged between 0 and 200.

5) Cases with unknown sex

Cancer cases with an unknown sex are excluded.

6) Cases where the sex is incompatible to the tumour site

For example, male patients with female reproductive cancers or female patients with testicular or prostate cancer. Sometimes this is a data quality issue, but it is also possible that the registration was based on sex at birth instead of their current sex. Due to the very small numbers and sensitivity around these cases, they are excluded from most analysis. Therefore, you should not include females with site codes in the range of C60-C63 or males with a site code in the range C51-C58.

7) Non- invasive tumours or non-melanoma skin cancers (C44)

Performance indicators and incidence trends of cancer generally focus on invasive cancers (C codes excluding C44). Non-invasive tumours (D-codes) tend to have trends over time that are affected by data quality so any analysis about these groups should be done with great care. For the purpose of this SOP, only tumours with a site code beginning with C (excluding C44 non-melanoma skin cancer) should be counted.

8) New International Statistical Classification of Diseases version 10 revision for tumours diagnosed from 2013 onwards

For diagnoses that occurred before 2013, the nine regional cancer registries each recorded cancer incidence using different coding systems. For statistical analysis and reporting purposes, these coding systems were mapped to the original version of the International Statistical Classification of Diseases 10th revision (ICD-10). Since becoming a unified cancer registration service in 2013, the International Classification of Diseases for Oncology, third edition (ICD-O3) is used to register cancer diagnoses by NCRAS. ICD-O3 can be mapped to ICD-10. The last changes that materially altered the reporting of cancer using ICD-10 were published in 2010 and implemented for some cancer diagnoses from 2012. The diagnosis year 2013, is the first complete year where this (4th) revision of ICD-10 can be used for reporting statistics on cancer diagnoses in England. The changes include updated codes primarily for non-solid tumours and the addition of new codes. Coding for revision 4 of ICD-10 is available in the CAS analytical tumour table (av_tumour_england/at_tumour_england) from 2019 (snapshot CAS2109). Revision 4 is applicable from 2013 onwards, there are four fields for tumour coding:

- SITE_ICD10R4_O2_FROM2013
- SITE_ICD10R4_O2_3CHAR_FROM2013
- SITE_ICD10_O2_PRE2013
- SITE_ICD10_O2_3CHAR_PRE2013

The pre-2013 fields should reflect ICD-10 revision 0 (R0), while the from-2013 fields are ICD10 revision 4 (R4). Both the full and 3-character versions are available. Details on the changes between code versions can be found [here](#). To align with the National Statistics most recent publications, please use the pre-2013 (R0) fields for tumours diagnosed pre 2013, and the from-2013 (R4) fields for tumours diagnosed from 2013 when using av2019 or monthly snapshots from CAS2109 onwards. This has been exemplified in sample codes 1 and 2. For previous snapshots, please use SITE_ICD10_O2 for all years.

Sample code 1

This code will count cancer cases by 3-character ICD-10 codes for diagnoses from 1995 onwards using recent snapshots (CAS2109 onwards). For previous snapshots, please use SITE_ICD10_O2 for all years.

Sample code 1:

```
select
case
when diagnosisyear<2013 then SITE_ICD10_O2_3CHAR_PRE2013
else SITE_ICD10R4_O2_3CHAR_FROM2013
end as SITE_ICD10_O2_3CHAR,
diagnosisyear,
count(tumourid)
from av2019.av_tumour_england@CASREF01
where
  ctry_code ='E' -- England residents using country code
  and STATUSOFREGISTRATION ='F' -- Finalised cases
  and dedup_flag=1 -- Excluding duplicates, note quality issue under point 3 in text above
  and age between 0 and 200 -- Sensible age
  and sex in (1,2) -- Known sex
  -- Years of interest and site restrictions
  --pre 2013 using SITE_ICD10_O2_PRE2013
  and ((SITE_ICD10_O2_PRE2013 like 'C%' and SITE_ICD10_O2_3CHAR_PRE2013<>
  'C44' -- all malignant neoplasms (excl non-melanoma skin cancer) for 2013 cases
  onwards
  and ((sex = '2' and SITE_ICD10_O2_3CHAR_PRE2013 not in ('C60','C61','C62','C63'))
  or (sex = '1' and SITE_ICD10_O2_3CHAR_PRE2013 not in
  ('C51','C52','C53','C54','C55','C56','C57','C58')))) -- Sex doesn't agree with tumour site
  and (diagnosisyear >1994 and diagnosisyear <2013))
  --post 2013 using SITE_ICD10R4_O2_FROM2013
```

```
or (SITE_ICD10R4_O2_FROM2013 like 'C%' and
SITE_ICD10R4_O2_3CHAR_FROM2013<> 'C44' -- all malignant neoplasms (excl non-
melanoma skin cancer) for 2013 cases onwards
and ((sex = '2' and SITE_ICD10R4_O2_3CHAR_FROM2013 not in
('C60','C61','C62','C63'))
or (sex = '1' and SITE_ICD10R4_O2_3CHAR_FROM2013 not in
('C51','C52','C53','C54','C55','C56','C57','C58')))) -- Sex doesn't agree with tumour site
and (diagnosisyear >=2013 and diagnosisyear<=2019)))
--You can use years of interest separate to site restrictions if you just want either pre-
2013 or post-2013
group by case
when diagnosisyear<2013 then SITE_ICD10_O2_3CHAR_PRE2013
else SITE_ICD10R4_O2_3CHAR_FROM2013
end,
diagnosisyear
order by
case
when diagnosisyear<2013 then SITE_ICD10_O2_3CHAR_PRE2013
else SITE_ICD10R4_O2_3CHAR_FROM2013
end,
diagnosisyear;
```

Sample code 2

If you want to make sure your figures align with those in CancerStats i.e. currently for diagnosis years between 2001 and 2019, you can use the criteria where CASCADE_INCI_FLAG equals 1. This applies to many of the filters in Sample Code 1 so is a shorter code to do the same thing. It is also designed to give identical numbers with those in CancerStats (which used to be called Cascade) and CancerData. The flag will include both C and D site codes, so be aware of this when creating site groups. This code also includes the use of both coding methods corresponding to which year they were diagnosed.

Sample code 2:

```
select
case
when diagnosisyear<2013 then SITE_ICD10_O2_3CHAR_PRE2013
else SITE_ICD10R4_O2_3CHAR_FROM2013
end as SITE_ICD10_O2_3CHAR,
diagnosisyear,
count(tumourid)
from av2019.av_tumour_england@CASREF01
```


where cascade_inci_flag=1 -- England, finalised cases, non-duplicates, sensible age, known sex, correct sex specific cancers, diagnoses after and including 2001. Please note the cascade_inci_flag will not restrict to C codes only, some D codes will be included

-- Years of interest and site restrictions

--pre 2013 using SITE_ICD10_O2_3CHAR_PRE2013

and ((SITE_ICD10_O2_3CHAR_PRE2013 like 'C%' and

SITE_ICD10_O2_3CHAR_PRE2013<> 'C44' -- all malignant neoplasms (excl non-melanoma skin cancer) for 2013 cases onwards

and (diagnosisyear >2000 and diagnosisyear <2013))

--post 2013 using SITE_ICD10R4_O2_FROM2013

or (SITE_ICD10R4_O2_FROM2013 like 'C%' and SITE_ICD10R4_O2_FROM2013<> 'C44' -- all malignant neoplasms (excl non-melanoma skin cancer) for 2013 cases onwards

and (diagnosisyear >=2013 and diagnosisyear<=2019)))

--You can use years of interest separate to site restrictions if you just want either pre-2013 or post-2013

group by case

when diagnosisyear<2013 then SITE_ICD10_O2_3CHAR_PRE2013

else SITE_ICD10R4_O2_3CHAR_FROM2013

end, diagnosisyear

order by case

when diagnosisyear<2013 then SITE_ICD10_O2_3CHAR_PRE2013

else SITE_ICD10R4_O2_3CHAR_FROM2013

end, diagnosisyear;

Sample code 3

Before the dedup_flag was included in the tumour table, we needed more complex code to identify and remove duplicates. This process used the ONS dataset to help with the deduplication. We do not recommend using this code unless absolutely necessary, but it is included for completeness. It can be used on the CAS snapshots including CAS1502 or AV2013. Note some specific IDs in the code are only available internally.

Sample code 3:

with

tidycanregcodes as(

select decode(av.centre, '0402', '0401', '0403', '0401', '0404', '0401', av.centre) as

canreg , substr(av.onsid, 5, 11) as canregno , tumourid

from av2013.av_tumour av where av.centre is not null and substr(av.onsid, 5, 11) is not null)

```

, findpairs as(select canreg, canregno, count(*) as paircount from tidycanregcodes
group by canreg, canregno
)
, dupflags as(select canreg, canregno, case when paircount =1 then 0 else 1 end as
dupflag from findpairs
)
, table1 as (select * from (
(
select tumourid, diagnosisyear, substr(site_ICD10_O2,1,3) as site3
from av2013.av_tumour T INNER JOIN ONS2012.ONSINCIDENCE@CASREF01 N ON
DECODE(T.CENTRE, '0402', '0401', '0403', '0401','0404', '0401', T.CENTRE) =
N.CANREG AND SUBSTR (T.ONSID, 5, 11) = N.CANREGNO
left outer join dupflags d on d.canreg = N.canreg and d.canregno = N.canregno where
diagnosisyear>1994 and diagnosisyear<2012 --1995-2011 cases
and SUBSTR(T.LSOA11_CODE, 1, 1) ='E' and substr(site_ICD10_O2,1,1)= 'C' and
substr(site_ICD10_O2,1,3)<> 'C44' and not(T.tumourid between [specific ID 1] and
[specific ID 2] and substr(T.onsid, 1, 4) = '[Specific ONS ID]' and dupflag = 1) and
STATUSOFREGISTRATION='F') union
(select tumourid, diagnosisyear, substr(site_ICD10_O2,1,3) as site3 from
av2013.av_tumour T
where diagnosisyear>2011 and diagnosisyear<2014 --2012-2013 cases
and SUBSTR(T.LSOA11_CODE, 1, 1) ='E'
and substr(site_ICD10_O2,1,1)= 'C' and substr(site_ICD10_O2,1,3)<> 'C44' and
STATUSOFREGISTRATION ='F'))
,
table2 as (select diagnosisyear, site3, count(tumourid) from table1 group by
diagnosisyear, site3 order by diagnosisyear, site3)

select * from table2;

```

For cases diagnosed before 1995

Due to historical duplicates on CAS and the dedup_flag only being available for cases back to 1995, it is necessary to use the Office for National Statistics incidence data to count cases between 1971 and 1994. The ONS dataset is stored in the ONS1971_1994 schema in CASREF01. The numbers of cases produced by the code in this section should be the same as the ONS publication covering the same period.

Things to include/exclude:

- 1) The dataset includes 710 cases registered in Wales, identified by a specific canreg number, but these should be retained to make sure the numbers agree to previous ONS publications using this data.
- 2) There are several filters applied to the more recent data that do not need to be applied to the ONS data. For example, ONS data only includes records with a known sex and duplicates have already been removed. Cases where the sex does not agree with the tumour site have also been accounted for. For information ICD8/9 codes 185-187 are specific to men only and 179-184 are specific to women only (ICD8 174 was only split between men and women in the 1979 version of ICD9).
- 3) Cases are coded in ICD8 for 1971-1978 and in ICD9 for 1979-1994. For ICD8 invasive cancers include 140-207 (excluding 173) and ICD9 includes codes 140-208 (excluding 173). To compare the site distribution with more recent data mapping the ICD 8/9 codes to ICD10 / ICD10-02 will be necessary. Similarly for morphology (type5) as this is coded in MOTNAC in 1971-1989 and ICD00 in 1990-1994. If you are interested in both invasive and in-situ cancer registrations you will need to include the ICD range of 140 to 239 e.g. `substr(site4,1,3) > '139'` and `substr(site4,1,3) < '240'`.
- 4) Age at diagnosis needs to be derived from date of birth (DOB) and date of diagnosis (DIAGDATE). There are 2 dates of birth (DOB1 & DOB2) in the dataset. This is because it might be known that a patient died in a particular month ie April but not exactly when. So DOB1 will be 1 April and DOB2 will be the 30 April. Therefore, the mid-point between these 2 DOB should be used to derive age at diagnosis. Only patients with ages between 0 and 200 at diagnosis should be included.

Sample code 4

This code will allow you to count cancer cases diagnosed from 1971 to 1994 using the ONS incidence data.

Sample code 4:

```
select extract(year from DIAGDATE) as diagyear, substr(site4,1,3) as site, count(*) from
ONS1971_1994.ONSINCIDENCE
where Trunc(trunc(months_between(diagdate, (dob1+(dob2-dob1)/2)))/12)>0 and
Trunc(trunc(months_between(diagdate, (dob1+(dob2-dob1)/2)))/12)<=200 and
substr(site4,1,3) != '173'
and (((substr(site4,1,3) < '208' and (substr(site4,1,3)>'139' and to_char(DIAGDATE,
'yyyy') < '1979'))
or (substr(site4,1,3) < '209' and (substr(site4,1,3)>'139' and to_char(DIAGDATE,
'yyyy') > '1978'))))
```

group by extract(year from DIAGDATE), substr(site4,1,3), sex

For information:

1) Deriving age at diagnosis in the ONS dataset:

Trunc(trunc(months_between(diagdate, (dob1+(dob2-dob1)/2)))/12) as diage

2) Deriving age group:

case when

Trunc(trunc(months_between(diagdate, (dob1+(dob2-dob1)/2)))/12) between 0 and 39 then '<39'

when Trunc(trunc(months_between(diagdate, (dob1+(dob2-dob1)/2)))/12) between 40 and 44 then '4044'

when Trunc(trunc(months_between(diagdate, (dob1+(dob2-dob1)/2)))/12) between 45 and 49 then '4549'

when Trunc(trunc(months_between(diagdate, (dob1+(dob2-dob1)/2)))/12) between 50 and 54 then '5054'

when Trunc(trunc(months_between(diagdate, (dob1+(dob2-dob1)/2)))/12) between 55 and 59 then '5559'

when Trunc(trunc(months_between(diagdate, (dob1+(dob2-dob1)/2)))/12) between 60 and 64 then '6064'

when Trunc(trunc(months_between(diagdate, (dob1+(dob2-dob1)/2)))/12) between 65 and 69 then '6569'

when Trunc(trunc(months_between(diagdate, (dob1+(dob2-dob1)/2)))/12) between 70 and 74 then '7074'

when Trunc(trunc(months_between(diagdate, (dob1+(dob2-dob1)/2)))/12) between 75 and 79 then '7579'

when Trunc(trunc(months_between(diagdate, (dob1+(dob2-dob1)/2)))/12) between 80 and 84 then '8084' else '85+' end age,

Note: You may need to extend these age groups to 90+ depending on your project.

3) Deriving date of death:

(dod1+(dod2-dod1)/2) as dod