

# Rapid Cancer Registration Dataset: data at 2nd October 2021 (CAS2110)

The National Cancer Registration and Analysis Service (NCRAS) has developed an algorithmically generated Rapid Cancer Registration Dataset (RCRD) using the standard administrative datasets which flow rapidly into NHS Digital (PHE) and are incorporated into the Cancer Analysis System (CAS) of NCRAS. The data takes the form of a series of significant events that occur to each patient as they proceed through the diagnostic and then therapeutic parts of the cancer pathway, and is available at approximately 4-5 months behind real time. The RCRD is shallower and narrower than the full NCRAS cancer registration dataset; it should be used and interpreted with reference to the caveats outlined within this document.

## Main findings

This document outlines the main features of the data to be aware of when interpreting the Rapid Cancer Registration Dataset:

- Across all cancers types included approximately 14.7% of cases are missing and 5.1% of cases are included erroneously or with incorrect cancer type or diagnosis date (when compared to 'Gold Standard' registration data for 2018 data).
- These figures vary strongly with cancer site. Broadly, more common cancers (particularly breast and prostate cancer) perform best and less common cancers (particularly bone and soft tissue and cancers of unknown primary) perform worst.
- There are more missing tumours in those aged over 70 compared to younger age groups.
- Other factors that reduce data completeness include the patient's route to diagnosis, mortality within 30 days of diagnosis, and the presence of multiple cancers.
- Usable data is available approximately 4-5 months after diagnosis or other clinical activity occurs.
- Data on cancer stage group at diagnosis is available for a number of common tumour types, although completeness is lower than that for the Gold Standard registration data. Where data is available it generally agrees with the Gold Standard stage group in 80-90% of tumours.

The dataset includes Rapid Cancer Registrations from January 2018 to the most recently available data (at the date specified in the title to this document), plus additional event data for the same period.

## Contents

Summary

Methodology

Proxy registration events (Rapid Registrations)

Data structures

Data Quality

How do the number of Rapid Registrations compare with Gold Standard Registrations?  
Comparing the matching quality of Rapid Registrations  
Sensitivity testing of matching criteria  
Counts of events over time  
Estimated completeness of Rapid Registrations and secondary datasets

Staging data in the Rapid Registrations dataset

TNM stage group 1-4  
"Early" vs "Late" stage  
Stage trends over time

Appendix 1 - List of pathway events

Appendix 2 - List of Rapid Registration fields available

Appendix 3 - Cancer groups used for matching

Appendix 4 - Alternative defining events

Appendix 5 - Counts and error tabulations

Appendix 6 - False negative errors and basis of diagnosis

## Summary

A need to make rapidly available 'proxy cancer registrations' (and associated clinical activity) for the COVID-19 period has been identified to support the public health response by NHS Digital (PHE) and other agencies, and service reorganisation by the NHS. These proxy registrations are called Rapid Registrations in contrast to the more formal detailed registration process that are used in non-clinical cancer research and the National Statistics (<https://www.gov.uk/government/statistics/cancer-registration-statistics-england-2018-final-release>).

The National Cancer Registration and Analysis Service (NCRAS) has developed a Rapid Cancer Registration Dataset (RCRD) using all standard administrative datasets which flow rapidly into PHE and are incorporated into the Cancer Analysis System (CAS) of NCRAS.

This document describes the dataset structure, creation methodology, and data quality caveats (due to the rapid automated creation process without additional data curation) behind this dataset.

These data structures and methodologies are expected to evolve over the course of the public health response to COVID-19. The data is updated monthly and is referred to by the monthly CAS snapshot upon which it is based, e.g. CAS2009 refers to the CAS snapshot from September 2020. This document is considered a 'living document' and strictly applies only to the snapshot of CAS identified in the title.

## Methodology

### Proxy registration events (Rapid Registrations)

Datasets available to PHE were surveyed for how many months in arrears that they arrive within NCRAS and are loaded in a usable format for analysis. From these datasets a selection of event types were defined similarly to those typically used for cancer pathway analysis pursued by NCRAS.

The data takes the form of a series of significant events that occur to each patient as they proceed through the diagnostic and then therapeutic parts of the cancer pathway. These events include chemotherapy cycles, radiotherapy episodes and major cancer surgery as well as events based on the Cancer Waiting Times (CWT) and Cancer Outcomes and Services Dataset (COSD) datasets. These event types are numbered in the range 1-23 in the dataset.

Some events hypothesised to be indicative of a cancer diagnosis were defined including 'Diagnosis reported in COSD' (event 51) and 'CWT estimated diagnosis date' (event 52). These are numbered in the range 50-57 in the dataset - see Appendix 1 for a full list.

The indicative events for diagnosis were explored as candidate Rapid Registration events. These candidate rapid registration events were judged as matching against a Gold Standard Registration event if it met the following two conditions:

- The difference in diagnosis dates for each event was 90 days or less.
- Both registrations fell into the same broad tumour group (as defined in Appendix 3).

Using these matching criteria False Positive errors and False Negative errors are defined as:

- **False Positive Error (FPE):** A rapid registration event has been created which does not match against a Gold Standard Registration in the comparison period.
- **False Negative Error (FNE):** There exists a Gold Standard Registration event for which no rapid registration event can be matched.

Additional filtering was applied to the candidate events and eventually event 101 was defined to minimise both false positive and false negative errors and is recommended for use by researchers as the best candidate for a rapid cancer registration. Appendix 4 briefly examines some of the alternatives examined in the development of this event definition.

## Data structures

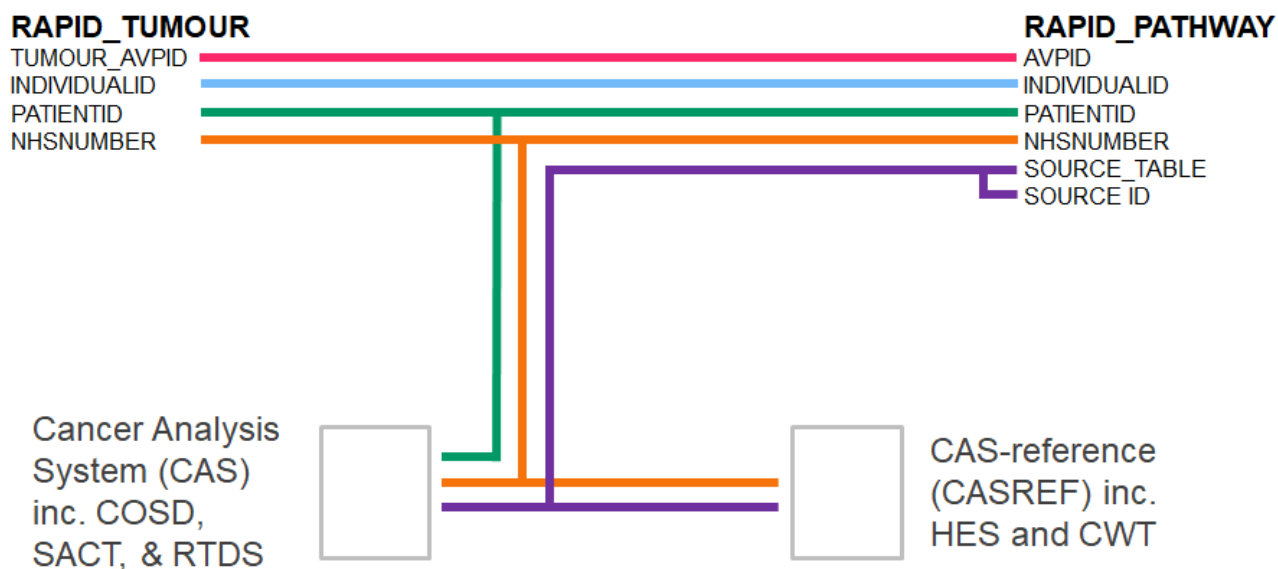
The rapid registration dataset consists of two tables:

**AT\_RAPID\_PATHWAY:** This is an event-based dataset with a number of types of event of interest defined based on the rapidly available datasets, see Appendix 1 for event definitions and properties. These are numbered in the range 1-23 for general purpose events, 50-57 for events that are candidates for combining into a rapid registration, and 101 for the final rapid registration event.

**AT\_RAPID\_TUMOUR:** This is a tumour level dataset that holds tumour and patient level data for each of the tumours defined by a rapid registration. The structure and contents of this table are presented in Appendix 3.

The rapid registration pathway and tumour table can be linked together as shown in Figure 1, and also to other datasets that are timely enough via NHSnumber.

Figure 1: Linkage diagram for the Rapid Cancer Registration Dataset



## Data Quality

### How do the number of Rapid Registrations compare with Gold Standard Registrations?

To illustrate the strengths and weaknesses of the Rapid Registrations compared to the gold standard process, registrations for tumours diagnosed during 2018 are compared in Figure 2.

For most tumour groups the counts of Rapid Registrations are significantly lower than those of standard registrations. The COSD system does not attempt to record basal cell carcinoma non-melanoma skin cancers (but they are recorded by hospital pathology systems, and thereby registered), explaining the discrepancy there. There is only one group where this situation is reversed - bone and soft tissue - for which a precise morphology is required to properly record the diagnosis. These cancers are being preferentially coded to bone and soft tissue in COSD (as the COSD standard necessitates simpler site-based coding, and this is the best choice under the circumstances) and re-coded during the gold standard registration process where more sophisticated combination of site and morphological coding is possible.

Figure 2: The number of cancer registrations by registration and tumour type, England, 2018

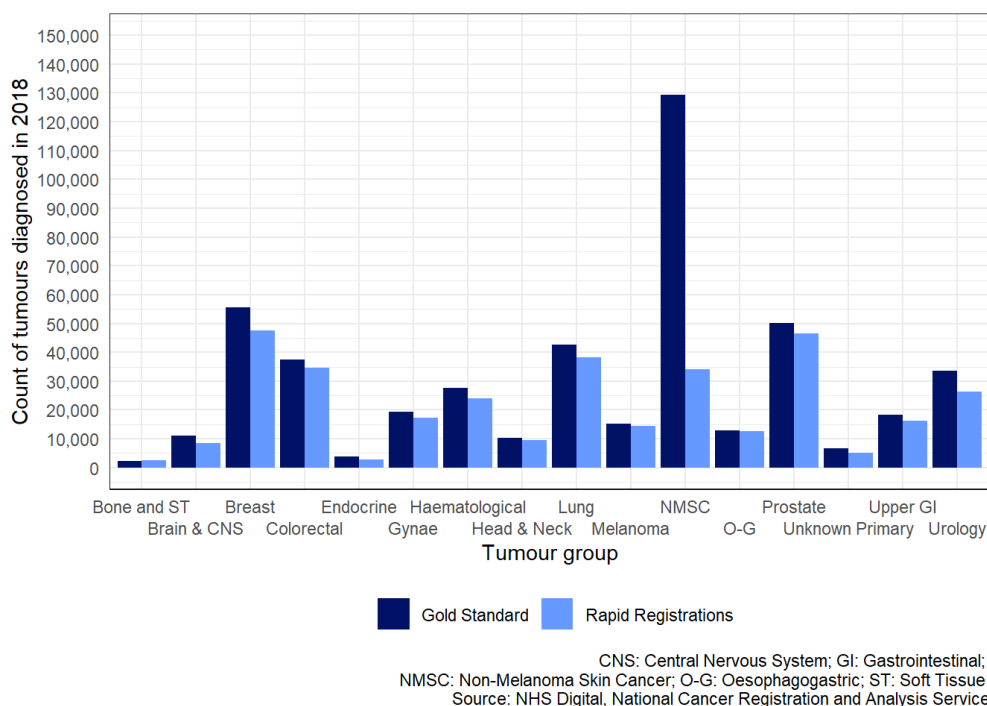
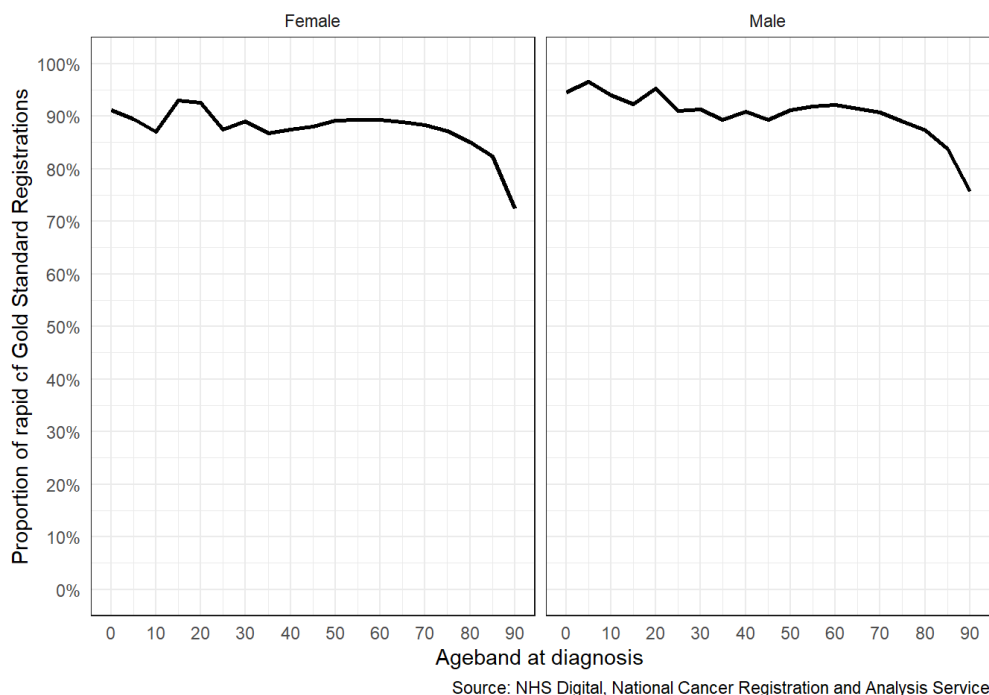


Figure 3 shows the age dependence of the ratio between Gold Standard and Rapid Registrations, Non-Melanoma Skin Cancer is excluded. The proportion of diagnoses is consistently high for both males and females until the age of 70 is reached, where it declines. This is explored further in Figure 5 below.

Figure 3: The proportion of cancer registrations by sex, age and registration type, England, 2018 (all tumour types combined)



## Comparing the matching quality of Rapid Registrations

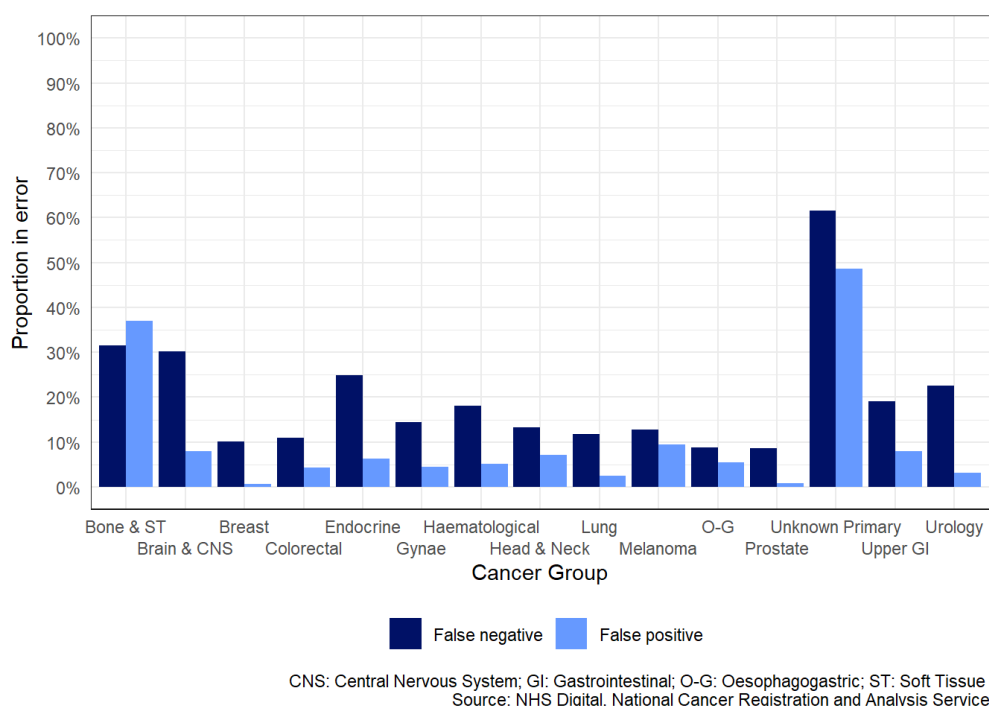
The quality of the Rapid Registrations was judged by comparing them against the gold-standard cancer registrations in the period April 2018 to September 2018. This period was chosen as available gold standard registration data was only finalised to December 2018 and a matching period of 90 days was allowed (restricting comparison to the middle six months of the twelve-month period).

Figure 4 shows the proportions of false positive and false negative events, by broad cancer type (excluding non-melanoma skin cancer), measured in the cas2110 snapshot (the tumour groups are defined in Appendix 3). A more detailed tabulation is available by tumour group and tumour site in Appendix 5.

In most tumour groups, there are more tumours missed by the rapid registrations process (false negatives) than there are falsely identified as tumours (false positives).

For breast and prostate, very few incorrect proxy registrations are made. Breast and prostate cancers are also least likely to be missing from the proxy dataset, whereas for brain and central nervous system (CNS), cancers of unknown primary, endocrine, bone and soft tissue, upper gastrointestinal and urological tumours more than 25% of cancers are missed. Bone and soft tissue tumours, which have more false positives than false negatives, are not frequently diagnosed. These tumours often require multiple pathology reports to correctly diagnose a patient and the Rapid Registrations dataset has not attempted to reconcile differences in the reported diagnoses.

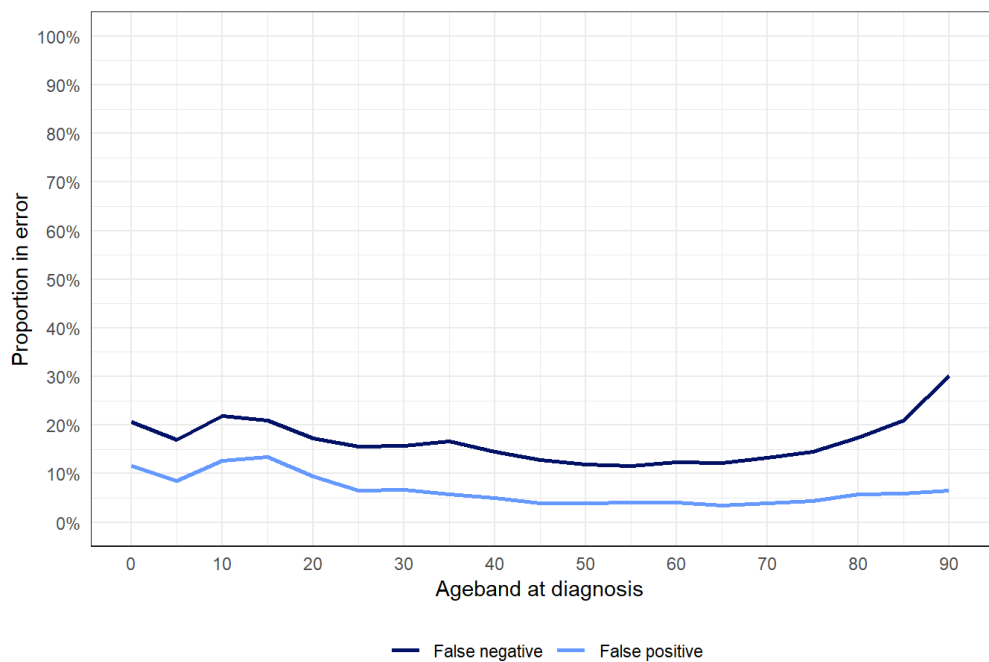
Figure 4: Types of error by tumour group



The proportion of false positive errors is fairly stable across all ages (Figure 5); the proportion of false negative errors slowly declines until age 70 when it increases significantly. The age dependence was investigated and the age-dependence of the basis of diagnosis was found to be at least partially responsible for this - see Appendix 6 for details.

The proportion of false positive cases is less sensitive to the age of the patient.

Figure 5: False negative and false positive errors by age band at diagnosis



Source: NHS Digital, National Cancer Registration and Analysis Service

The charts in Figure 6 (below) examine these patterns by tumour group. Please note that age groups for each tumour group must have a denominator of 25 patients or more or they are suppressed for reasons of statistical power.

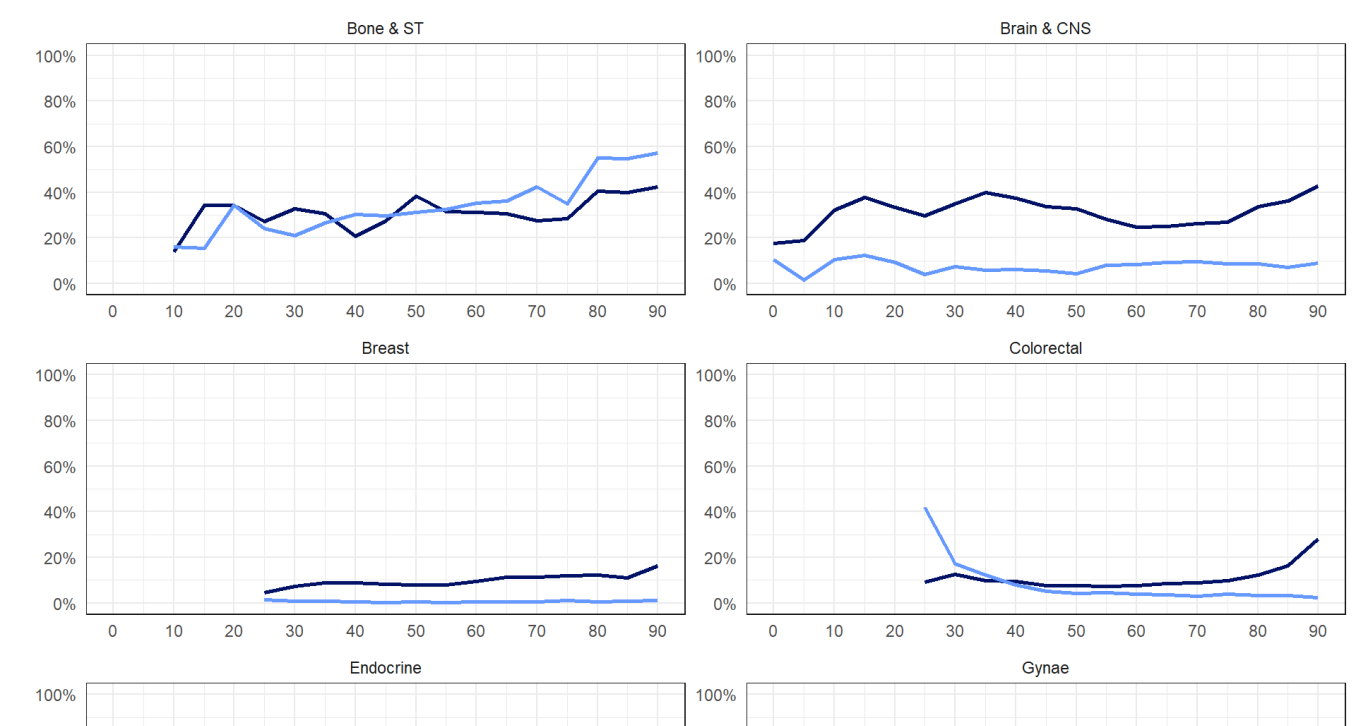
The patterns of false negative and false positive vary significantly by tumour group. Most groups have a higher proportion of false negatives than false positives at each age.

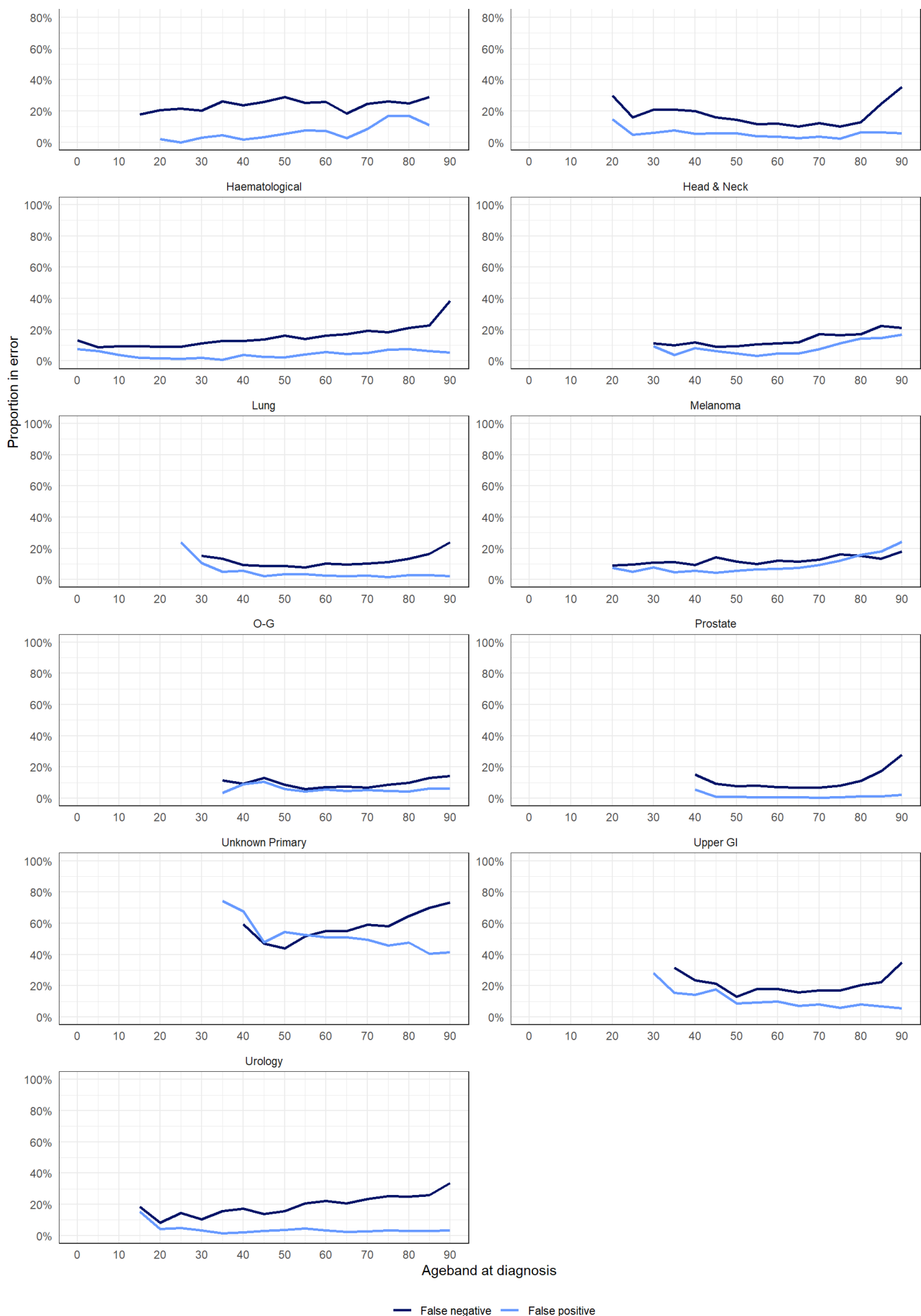
The proportion of false positives does not exhibit a trend by age for most tumour groups; the proportion rises with increasing age in the bone and soft tissue, head and neck groups and melanoma group and conversely falls with increasing age in the colorectal and unknown groups.

The proportion of false negatives rises with increasing age for all tumour groups except bone and soft tissue and endocrine. The most pronounced increases occur in the brain and central nervous system, colorectal, gynaecological, haematological, prostate, upper gastro-intestinal and unknown primary tumour groups.

The levels of both types of error are highest in tumour groups which are less likely to have solid-tissue pathology (haematological) or where survival rates are typically low. Conversely, the levels of error are lowest for tumour groups for which survival rates are typically higher.

Figure 6: False negative and false positive errors by age band at diagnosis and tumour group





CNS: Central Nervous System; GI: Gastrointestinal; O-G: Oesophagogastric; ST: Soft Tissue  
Source: NHS Digital, National Cancer Registration and Analysis Service

The variation of the false positive and false negative errors with Income deprivation quintile is shown in figure 6. While there is an overall trend visible this is likely to be due to confounding due to the variation with tumour type shown above and the known association of the incidence of many cancer types with income deprivation.

Figure 6: False negative and false positive errors by income deprivation quintile

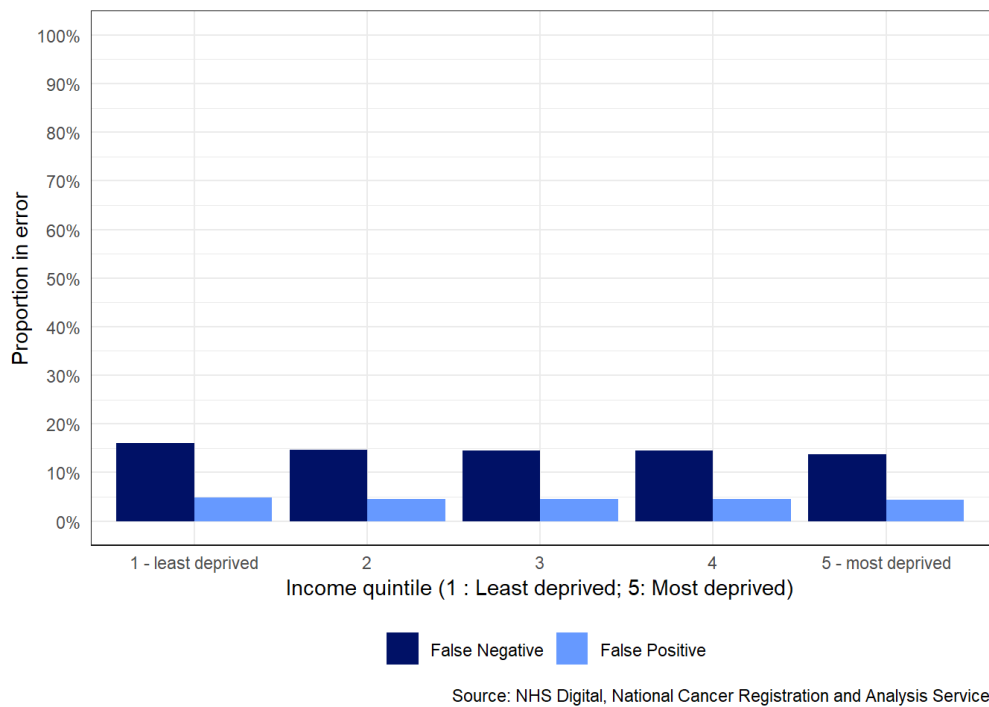


Figure 7 shows the variation of false negative and false positive errors with route to diagnosis. For false positives there is moderate variation with the lowest error rate being those cases identified through cancer screening or a two week wait referral. (These tumours are those that are likely to be captured in both the COSD dataset and the screening/Cancer Waiting Times datasets so the lower error rate is understandable.)

Most routes to diagnosis have a substantially higher false negative rate than the overall average. 'Two Week Wait' (TWW) and screening routes have a substantially lower false negative rate (and make up between them 45% of the total cohort).

Figure 7: False negative and false positive errors by route to diagnosis

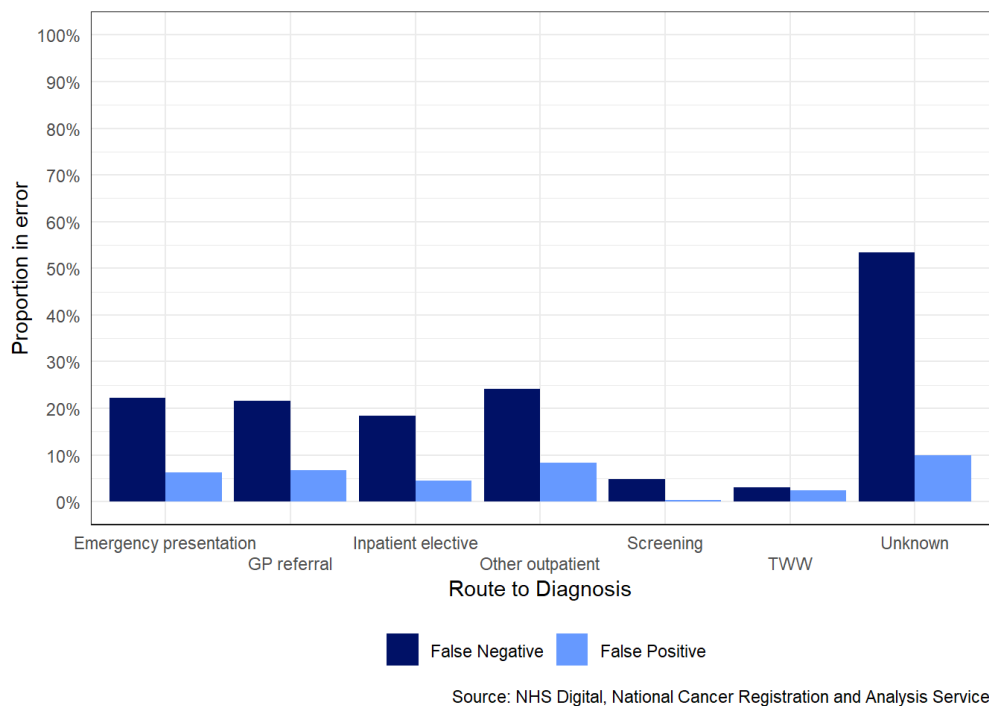
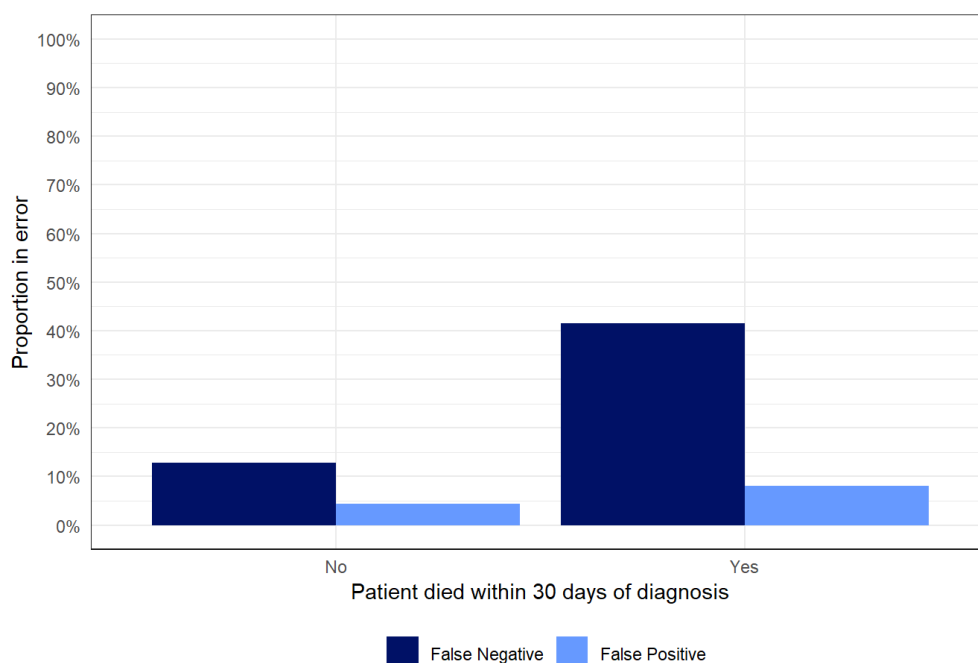


Figure 8 below shows the variation of false negative and false positive errors with whether or not the patient died within 30 days of diagnosis. The false negative error rate varies substantially between patients who die in the 30 days post-diagnosis compared to those who did, meaning that patients who die within 30 days are more likely to be missing from the dataset.

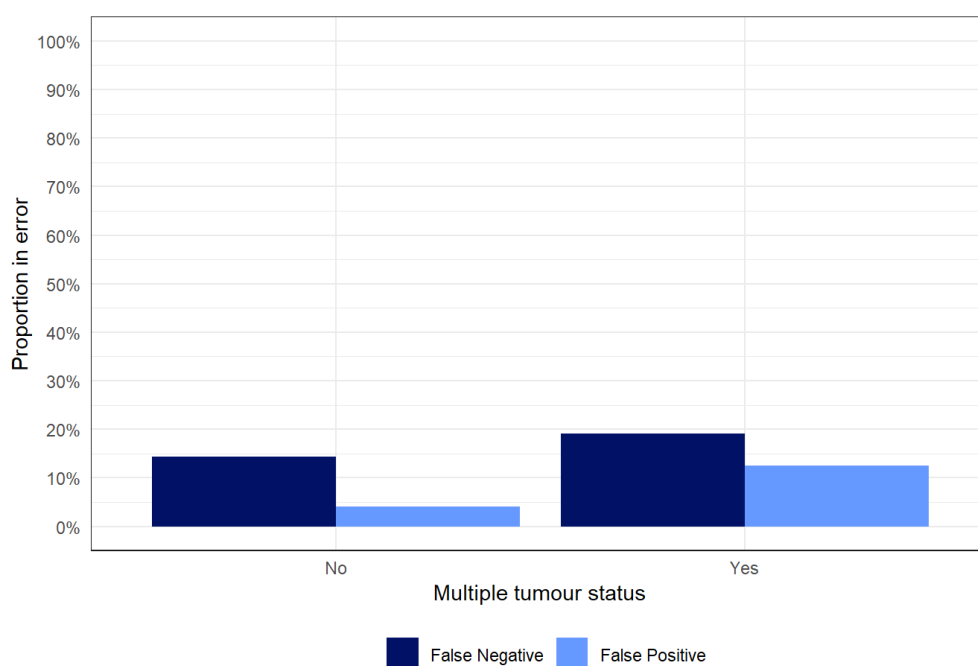
Figure 8: False negative and false positive errors by 30-day mortality



Source: NHS Digital, National Cancer Registration and Analysis Service

Figure 9 below shows the variation of false negative and false positive errors with the multiple tumour status of the patient, i.e. whether or not the patient had been diagnosed with more than one type of tumour in the period January 2018 onward. The false positive error rate varies substantially between patients with multiple tumour types and those that don't, meaning that these patients with multiple tumours are more likely to have incorrect tumour types or diagnosis dates recorded.

Figure 9: False negative and false positive errors by multiple tumour status

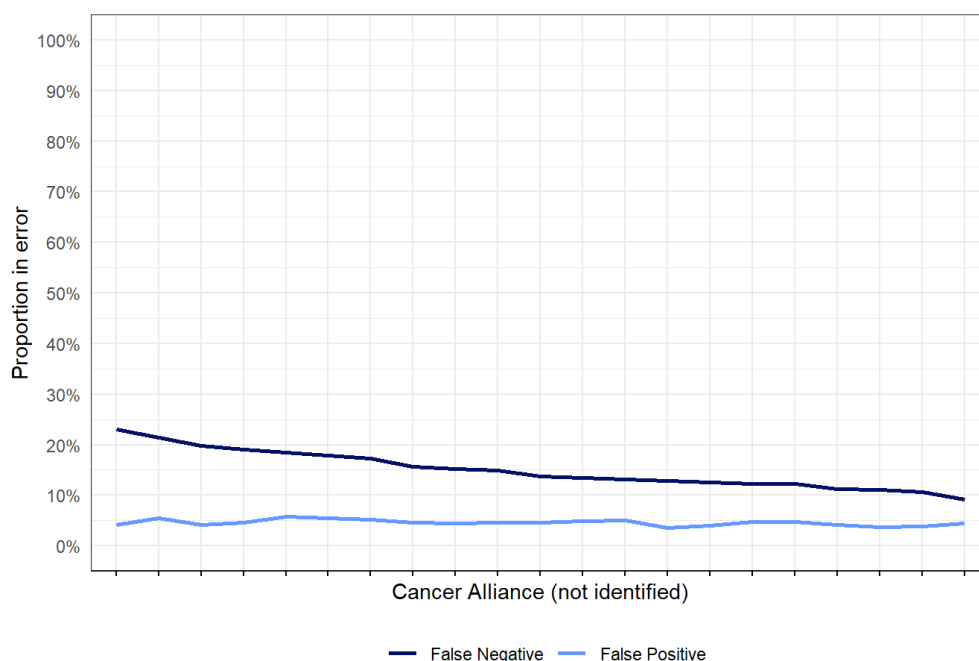


Source: NHS Digital, National Cancer Registration and Analysis Service

Figure 10 below shows the variation of false negative and false positive errors with the cancer alliance of residence of the patient at the time of diagnosis. The false negative error rate varies more in absolute terms than the false positive rate and may be driven by trust level variation (see figures 11 and 12 below).

Figure 10: False negative and false positive errors by Cancer Alliance



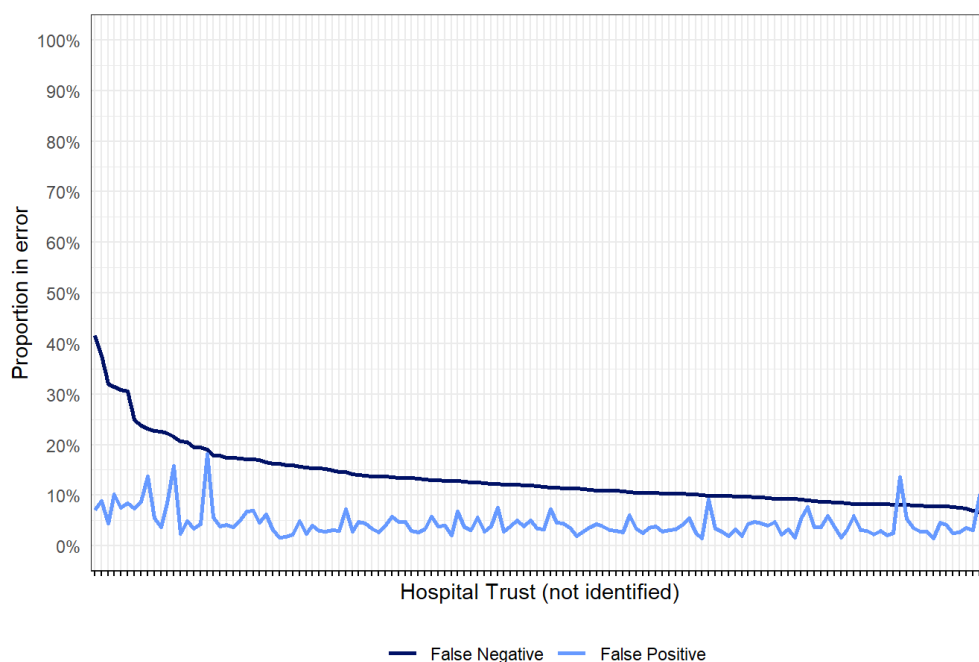


Source: NHS Digital, National Cancer Registration and Analysis Service

Figures 11 and 12 below show the variation of false negative and false positive errors with the trust that diagnosed the tumour. Figure 11 shows the error proportion and figure 12 the numerator (count) of the errors. Trusts shown are limited to NHS secondary care trusts with a denominator of at least 50 patients over the assessment period. Both figures are ordered in descending order of the false negative statistic - but note that the order is not the same in each figure.

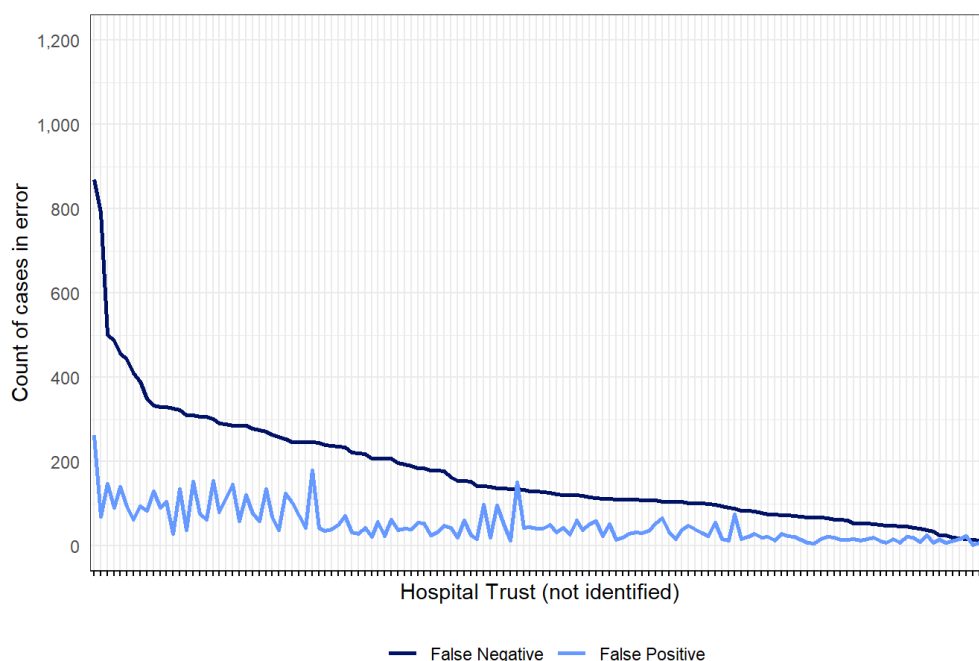
There is substantial variation in both false positive and false negative rates and counts. Some large trusts have several hundred or up to 1000 cases (over the six-month period under assessment).

Figure 11: False negative and false positive errors (proportion) by hospital trust



Source: NHS Digital, National Cancer Registration and Analysis Service

Figure 12: False negative and false positive errors (count) by hospital trust



Source: NHS Digital, National Cancer Registration and Analysis Service

## Sensitivity testing of matching criteria

In this section, the sensitivity of the Rapid Registrations dataset is illustrated for different matching criteria.

As expected, the stricter the criteria about the timing of events, more errors (both false negative and false positive) are observed. Not including a match specification on tumour type (the second line of table 1) improves both matching criteria and demonstrates that approximately 40% of false positive tumours have a cancer diagnosis of some sort when the necessity of matching by tumour group is removed.

Table 1: Proportions of false positive and negative errors under alternative matching criteria

Tumour matching	Match within N days	False Negative %	False Positive %
Broader	90	14.7%	5.1%
Broader	60	16.2%	6.6%
Broader	30	21.2%	11.9%
Broader	14	31.7%	23.6%
Broader	7	47.8%	41.7%
Broader	0	82.3%	80.0%
Narrow	90	22.3%	13.0%
None	90	13.3%	3.6%

## Counts of events over time

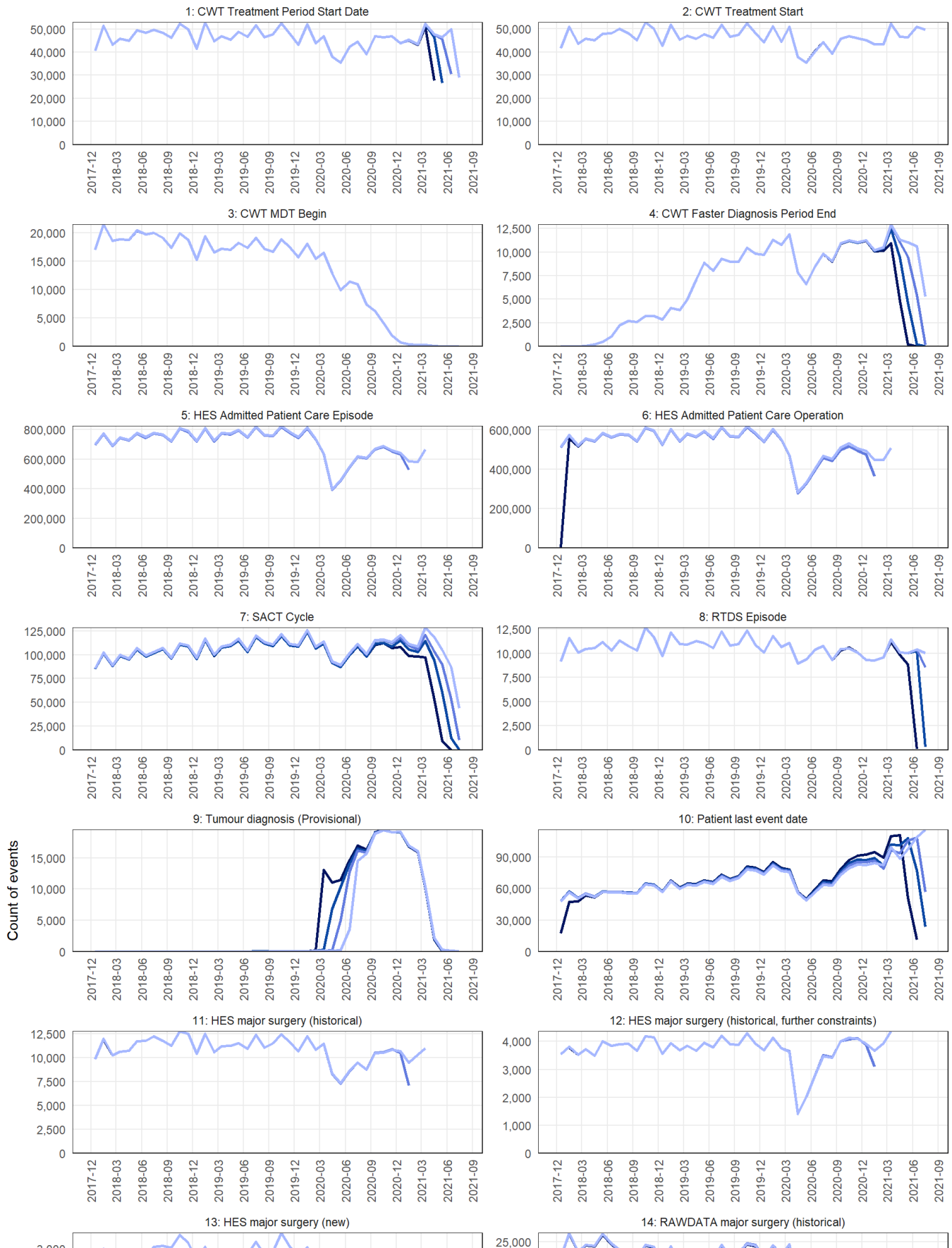
This section examines the population of events by chronological time and when they appear in successive analytical snapshots in the CAS. Figure 13 shows that most data items in the Rapid Registrations dataset are stable with respect to the snapshot month.

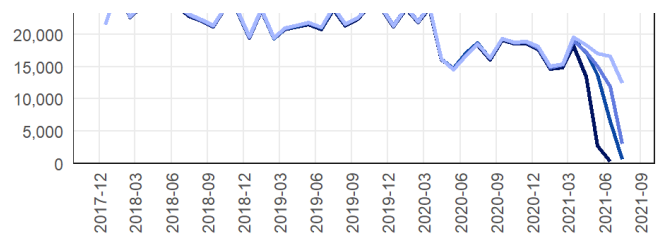
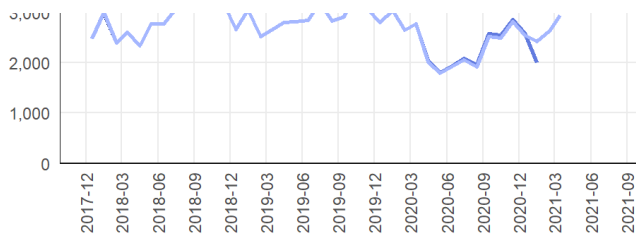
Specific comments about the events shown below are:

- Cancer Waiting Times data (events 1-4) are received based on the treatment start date, this explains the fact that for event 2 all lines lie exactly on top of each other. Other CWT events accumulate over successive snapshots where these events precede the first treatment start event.
- An issue with HES data resulting in lower than expected completeness port 2020-04-01 was resolved in cas2102, showing as increased event counts in events 5,6, 11, 12, 13 and 23.
- The definition of event 17 only includes tumour diagnoses prior to 2018, lack of data in the chart below is expected.
- Definitions of staging events may change between snapshots, this might explain higher or lower counts in one snapshot compared to others.
- The vital status shown in the event 19 is typically only assessed each January or the completion of registering each diagnosis year, explaining the large peaks in the graph.

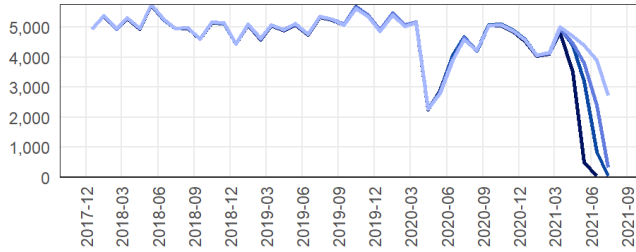
- The raw data used to populate events 21, 54, and 56 is subject to ongoing deduplication, this explains lower counts in earlier time periods for later snapshots.
- Between snapshots there is generally an increase in the Event 101-103 (Inferred diagnoses) counts, particularly for recent months as additional COSD data is submitted. However, for some earlier months there is a small decrease in these event counts. This is because the algorithm to define Events 101-103 excludes potential diagnoses where the patient has a confirmed diagnosis for the same tumour group which was more than 90 days before the potential diagnosis, to avoid double-counting the same diagnosis. These exclusions can change between snapshots due to the processing of gold standard cancer registration data, which leads to an increase in confirmed previous diagnoses. However the magnitude of this effect has been measured to be <1% of all cases in any given month.

Figure 13: Population of data items to CAS snapshot

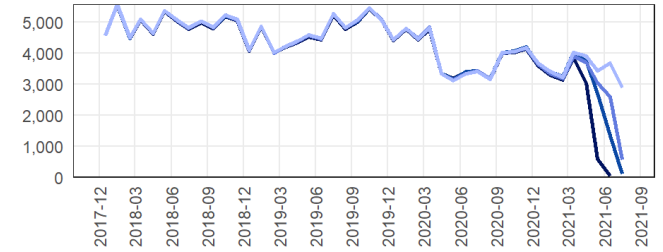




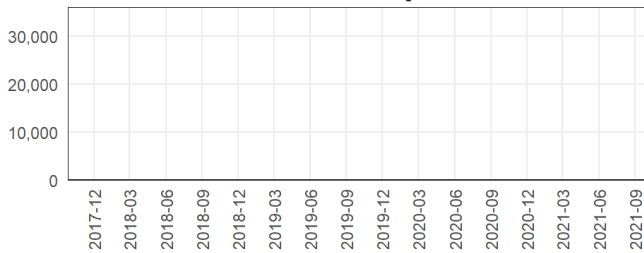
15: RAWDATA major surgery (historical, further constraints)



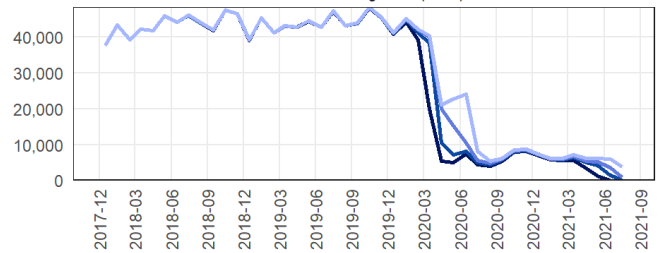
16: RAWDATA major surgery (new)



17: Prior tumour diagnosis



18: Tumour diagnosis (Final)

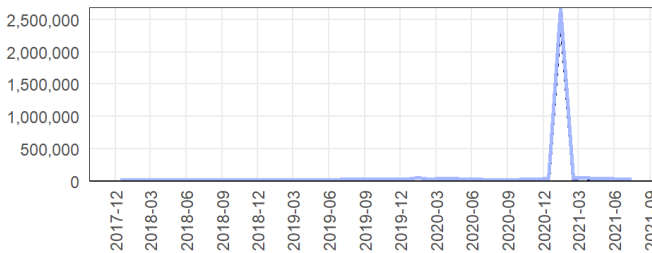


Year and Month

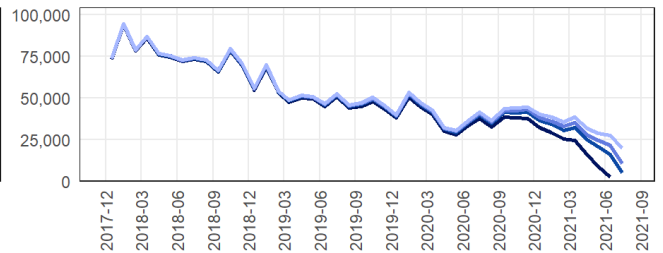
cas2107 cas2108 cas2109 cas2110

Source: NHS Digital, National Cancer Registration and Analysis Service

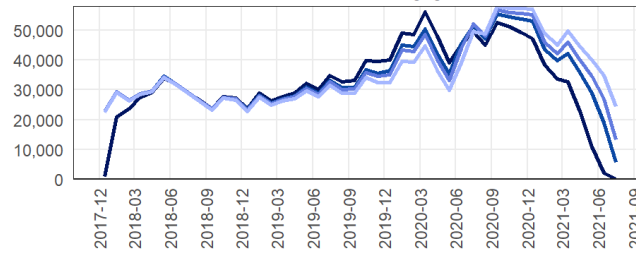
19: Patient vital status date



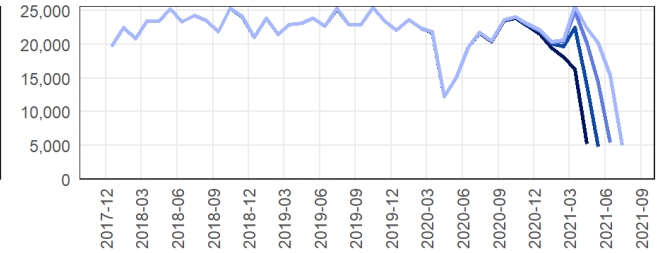
20: RAWDATA holistic needs assessment record



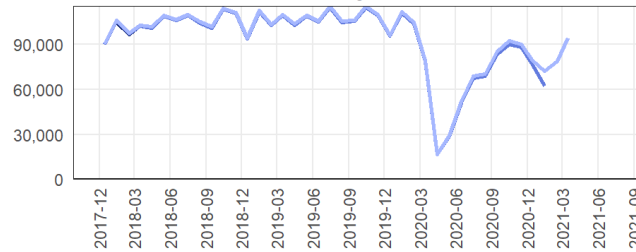
21: RAWDATA staging



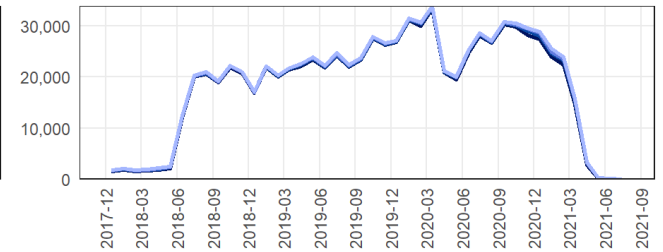
22: CWT First Seen



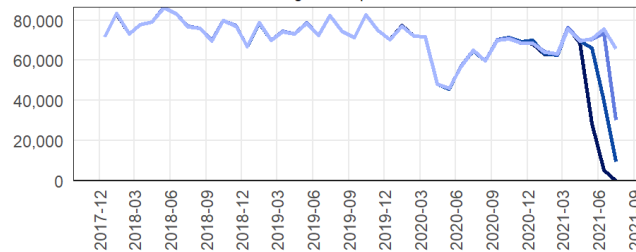
23: HES diagnostic event



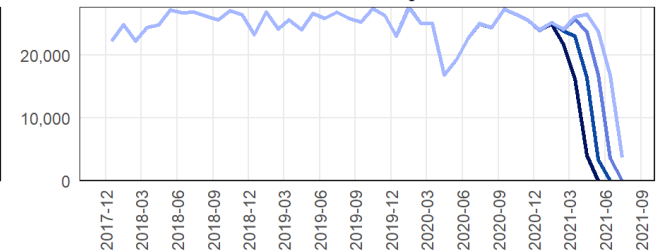
50: Skeleton Tumour creation

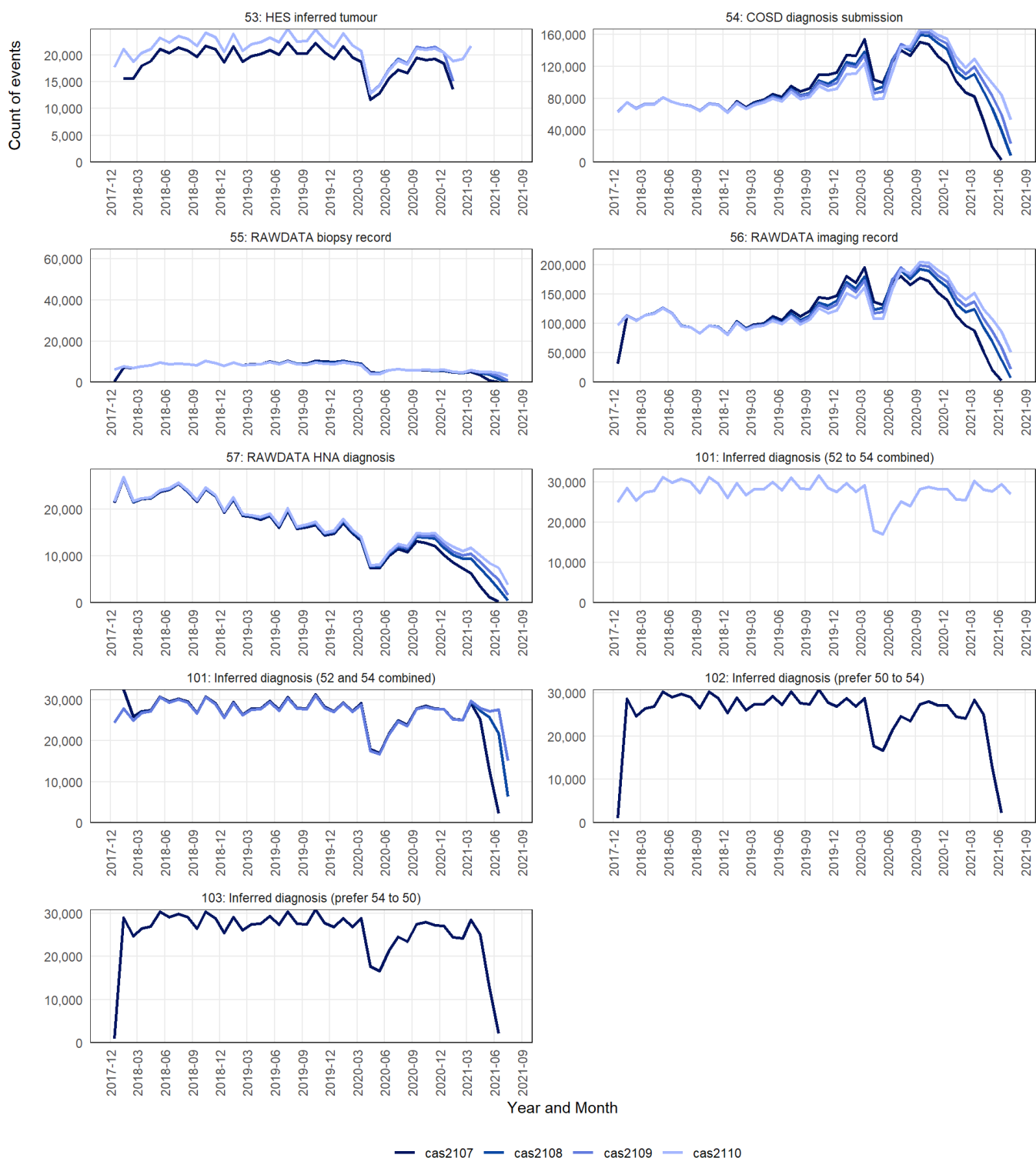


51: Diagnosis reported in COSD



52: CWT estimated diagnosis date





Source: NHS Digital, National Cancer Registration and Analysis Service

## Estimated completeness of Rapid Registrations and secondary datasets

Detailed linked rapid cancer registration, CWT, SACT and RTDS data is available at approximately a four-month lag from real time. Linked HES and raw COSD data is available at approximately 4-5 months behind real time.

Table 2 below shows data usability and completeness for Rapid Registrations and the constituent datasets. The "latest usable" column shows the 'hard limit' on data that is considered fit for analytical purposes (90% completeness), even in months prior to this though data is not necessarily considered complete and the completeness is displayed below. This should be taken into account in any use of the rapid registration data and the secondary datasets.

For the Rapid Tumour data completeness is expressed as the proportion of CCG of residence which show a cancer incidence within the normally expected range (see Table 3 below). For other datasets except CWT completeness is computed as a percentage of the number of data providers who have supplied data over those who are expected to do so.

Data completeness within the Cancer Waiting Times dataset varies at patient level with event type. Figures for the Treatment Start Date and Treatment Period Start Date are given below. Completeness of other CWT events can be estimated by inspecting Figure 13 (events 1-4).

Table 2: Rapid registration and dataset usability/completeness in cas2110

Data source	Latest usable	January 2021	February 2021	March 2021	April 2021	May 2021	June 2021	July 2021
Rapid Tumours (COSD)	July 2021	Complete	Complete	Complete	Complete	Complete	Complete	90%
HES	March 2021	Complete	Complete	95%	•	•	•	•
SACT	April 2021	96%	96%	96%	93%	•	•	•
RTDS	May 2021	98%	98%	98%	98%	98%	•	•
CWT (TSD)	July 2021	Complete	Complete	Complete	Complete	Complete	Complete	Complete
CWT (TPSD)	June 2021	Complete	Complete	Complete	Complete	Complete	98%	•

*Note:*

COSD = Cancer Outcomes and Services Dataset

TSD = Treatment Start Date

TPSD = Treatment Period Start Date

Table 3: Number of outlier CCGs in COSD dataset in cas2110

The table below shows the number of CCGs (using the April 2020 boundaries) which have 3-sigma outlier counts per month (either high or low) compared to the expectation of the fraction of the total number of new cancer registrations in England. This can be used to judge to what extent there is large scale missing data in COSD (and therefore in the Rapid Registrations in any particular month.)

Year and month	Outlier: High	Outlier: Low	In expected range	Total received
2020-01	0	1	134	135
2020-02	0	0	135	135
2020-03	0	1	134	135
2020-04	2	6	127	135
2020-05	1	3	131	135
2020-06	1	3	131	135
2020-07	0	0	135	135
2020-08	0	1	134	135
2020-09	1	0	134	135
2020-10	0	3	132	135
2020-11	0	0	135	135
2020-12	1	0	134	135
2021-01	1	0	134	135
2021-02	0	1	134	135
2021-03	0	2	133	135
2021-04	1	0	134	135
2021-05	0	0	135	135
2021-06	1	2	132	135
2021-07	1	13	121	135
2021-08	33	41	61	135

## Staging data in the Rapid Registrations dataset

### TNM stage group 1-4

The size and extent of a cancer is commonly described using the 'TNM' system (<https://www.uicc.org/resources/tnm>) for "Tumour", "Node", and "Metastases". This is often abbreviated to a number between 1 (typically a localised tumour with limited spread) to 4 (typically a tumour that has invaded or spread to distant organs). The stage at diagnosis is very strongly associated with patient outcomes.

In the current version of the Rapid Registrations dataset partial staging data is provided for a number of different cancer sites (ICD-10 codes can be found in the labels for tables 5a-k). This has been benchmarked against the gold standard cancer registry data for cas2110.

Table 4 shows the count and proportion of cases by TNM stage group for both the Rapid Registrations and the Gold Standard Registrations, for calendar year 2018. For example 32% of breast cancers are TNM stage group 1 in the Rapid Registrations, but 38% in the Gold Standard Registrations. Compared to the Gold Standard Registrations in 2018, the Rapid Registrations under report breast cancers diagnosed at stages 1 or 2; colorectal cancers diagnosed at stage 4 are under reported and prostate cancers have under reported stages 1 and 4. In all three tumour groups, there are more tumours allocated to the unknown or unstageable category. Lung cancers in the RCRD most accurately match the Gold Standard Registrations and exhibits a broadly similar stage profile from both measures.

Table 4: Summary proportions of stage at diagnosis for the Rapid Registrations and Gold Standard Registrations

Broad Cancer Group	Stage Group	Count (Rapid)	Percentage (Rapid)	Count (Gold Standard)	Percentage (Gold Standard)
Bladder	1	2318	24.5%	2860	30.3%
Bladder	2	1788	18.9%	1880	19.9%
Bladder	3	558	5.9%	884	9.4%
Bladder	4	256	2.7%	631	6.7%
Bladder	U	4525	47.9%	3190	33.8%
Breast	1	13347	31.7%	15798	37.5%
Breast	2	12799	30.4%	16181	38.4%
Breast	3	3120	7.4%	3581	8.5%
Breast	4	1103	2.6%	1846	4.4%
Breast	U	11788	28.0%	4751	11.3%
Colorectum	1	4915	15.4%	5437	17.0%
Colorectum	2	7047	22.1%	7658	24.0%
Colorectum	3	8241	25.8%	9240	29.0%
Colorectum	4	5107	16.0%	7286	22.8%
Colorectum	U	6594	20.7%	2283	7.2%
Kidney	1	2372	30.0%	3269	41.3%
Kidney	2	445	5.6%	536	6.8%
Kidney	3	1355	17.1%	1610	20.3%
Kidney	4	689	8.7%	1503	19.0%
Kidney	U	3054	38.6%	997	12.6%
Lung	1	6197	17.7%	6670	19.1%
Lung	2	2606	7.5%	2694	7.7%
Lung	3	7313	20.9%	7564	21.7%
Lung	4	14914	42.7%	16824	48.2%
Lung	U	3894	11.1%	1172	3.4%
Lymphoma	1	907	7.6%	1736	14.6%
Lymphoma	2	949	8.0%	1596	13.4%
Lymphoma	3	1196	10.1%	1969	16.6%
Lymphoma	4	2648	22.3%	4846	40.7%

Broad Cancer Group	Stage Group	Count (Rapid)	Percentage (Rapid)	Count (Gold Standard)	Percentage (Gold Standard)
Lymphoma	U	6197	52.1%	1750	14.7%
Melanoma	1	6353	48.4%	8280	63.1%
Melanoma	2	2395	18.2%	2654	20.2%
Melanoma	3	445	3.4%	1040	7.9%
Melanoma	4	183	1.4%	321	2.4%
Melanoma	U	3756	28.6%	837	6.4%
Oesophagus	1	304	3.8%	442	5.5%
Oesophagus	2	1274	15.7%	959	11.8%
Oesophagus	3	2017	24.9%	2126	26.2%
Oesophagus	4	2405	29.7%	3184	39.3%
Oesophagus	U	2106	26.0%	1395	17.2%
Ovary	1	1107	23.9%	1300	28.0%
Ovary	2	230	5.0%	269	5.8%
Ovary	3	1155	24.9%	1572	33.9%
Ovary	4	690	14.9%	1018	22.0%
Ovary	U	1455	31.4%	478	10.3%
Pancreas	1	351	4.6%	648	8.5%
Pancreas	2	624	8.2%	792	10.4%
Pancreas	3	744	9.8%	1043	13.8%
Pancreas	4	2052	27.1%	3977	52.4%
Pancreas	U	3813	50.3%	1124	14.8%
Prostate	1	11662	25.4%	16265	35.4%
Prostate	2	5490	11.9%	6545	14.2%
Prostate	3	10325	22.5%	11640	25.3%
Prostate	4	5613	12.2%	8022	17.5%
Prostate	U	12854	28.0%	3472	7.6%
Stomach	1	337	8.9%	337	8.9%
Stomach	2	599	15.8%	468	12.3%
Stomach	3	427	11.2%	710	18.7%
Stomach	4	1112	29.3%	1630	42.9%
Stomach	U	1324	34.9%	654	17.2%
Uterus	1	4594	59.1%	5315	68.4%
Uterus	2	493	6.3%	523	6.7%
Uterus	3	721	9.3%	815	10.5%
Uterus	4	493	6.3%	537	6.9%
Uterus	U	1467	18.9%	578	7.4%



In Tables 5a-m below, the distribution of the stage allocations between the Rapid Registrations and the Gold Standard Registrations are examined. The figures indicate the proportion of agreement at the 1-digit TNM stage group level, where the stage is known in the Rapid Registrations dataset. Stages 1-4 in the Rapid Registrations dataset agree with the gold standard stage variable for a high proportion.

For example, when examining the subset of Rapid Registrations breast tumours that are identified as TNM stage 1 (32%), approximately 89% of these are found to be TNM stage group 1 in the gold standard registration data, with another 11% distributed across TNM stages 2-4 and the unknown or unstageable groups.

For many but not all (e.g., late stage breast cancer), roughly 85% or more of staged cases in the Rapid Registrations table have the same stage grouping as the equivalent tumour in the standard registration data - this can be seen in the table below by inspecting the figures where the stage metrics for the Rapid Registrations and Gold Standard Registrations are the same.

Where the stage is labelled as unknown or unstageable in the rapid pathway dataset it is known for at least 70% of those cases in the gold standard data.

Tables 5a-m: Stage comparison between Rapid Registrations and Gold Standard Registrations by cancer site

a. bladder (ICD-10 C67)

Stage Group (Gold Standard)	Stage Group (Rapid)				
	1	2	3	4	Unknown
1	84.9%	3.8%	7.7%	5.1%	17.0%
2	3.9%	72.1%	15.4%	5.9%	8.8%
3	2.6%	10.9%	65.4%	4.3%	5.6%
4	1.3%	5.0%	5.4%	80.1%	6.1%
U	7.3%	8.2%	6.1%	4.7%	62.5%

b. breast (ICD-10 C50)

Stage Group (Gold Standard)	Stage Group (Rapid)				
	1	2	3	4	Unknown
1	89.4%	4.6%	1.4%	3.3%	27.1%
2	6.4%	88.9%	10.8%	14.1%	29.3%
3	0.5%	2.7%	80.9%	5.3%	5.0%
4	0.2%	0.9%	3.0%	72.9%	6.9%
U	3.5%	2.9%	3.9%	4.4%	31.7%

c. colorectum (ICD-10 C18-C20)

Stage Group (Gold Standard)	Stage Group (Rapid)				
	1	2	3	4	Unknown
1	85.1%	2.1%	1.7%	0.7%	14.2%
2	5.7%	85.6%	5.6%	1.2%	12.5%
3	6.5%	7.5%	85.2%	4.3%	17.4%
4	0.9%	2.8%	5.8%	92.9%	27.7%
U	1.9%	2.0%	1.7%	0.9%	28.2%

d. kidney (ICD-10 C64)

Stage Group (Rapid)				
---------------------	--	--	--	--

Stage Group (Gold Standard)	1	2	Stage Group (Rapid)	4	Unknown
Stage Group (Gold Standard)	1	2	3	4	Unknown
1	91.4%	6.5%	3.0%	1.7%	33.3%
2	0.5%	78.7%	1.0%	0.7%	5.1%
3	1.7%	6.5%	86.2%	3.8%	11.4%
4	0.4%	3.4%	5.6%	92.5%	25.0%
U	5.9%	4.9%	4.2%	1.3%	25.1%

e. lung (ICD-10 C33-C34)

Stage Group (Gold Standard)	Stage Group (Rapid)				
	1	2	3	4	Unknown
1	94.0%	6.3%	1.1%	0.4%	13.9%
2	2.5%	84.9%	1.7%	0.3%	3.9%
3	1.6%	4.9%	90.9%	1.3%	12.9%
4	1.2%	3.0%	5.4%	97.6%	44.1%
U	0.7%	1.0%	0.9%	0.4%	25.1%

f. melanoma (ICD-10 C43)

Stage Group (Gold Standard)	Stage Group (Rapid)				
	1	2	3	4	Unknown
1	94.5%	1.5%	4.7%	8.2%	58.7%
2	1.9%	79.5%	9.2%	14.8%	14.9%
3	1.9%	11.8%	79.1%	14.8%	6.8%
4	0.2%	1.6%	2.5%	49.7%	4.5%
U	1.5%	5.6%	4.5%	12.6%	15.0%

g. oesophagus (ICD-10 C15)

Stage Group (Gold Standard)	Stage Group (Rapid)				
	1	2	3	4	Unknown
1	75.3%	4.5%	0.5%	0.2%	6.7%
2	7.9%	52.0%	6.7%	0.8%	5.6%
3	8.9%	33.0%	66.4%	4.2%	11.3%
4	1.6%	5.2%	21.1%	85.2%	30.3%
U	6.2%	5.3%	5.3%	9.5%	46.2%

h. ovary (ICD-10 C56-C57)

Stage Group (Rapid)

Stage Group (Gold Standard)	1	2	Stage Group (Rapid)	3	4	Unknown
Stage Group (Gold Standard)	1	2	3	4	Unknown	
1	97.5%	7.0%	0.9%	0.3%		13.3%
2	0.4%	88.3%	0.5%	0.1%		3.8%
3	0.8%	2.6%	91.9%	11.0%		28.9%
4	0.4%	0.4%	4.5%	84.3%		26.0%
U	1.0%	1.7%	2.3%	4.2%		28.0%

i. prostate (ICD-10 C61)

	Stage Group (Rapid)				
Stage Group (Gold Standard)	1	2	3	4	Unknown
1	86.4%	8.9%	3.9%	1.2%	40.8%
2	6.7%	83.8%	2.5%	0.9%	6.7%
3	4.3%	4.3%	87.1%	2.7%	13.7%
4	0.8%	0.7%	4.0%	93.4%	17.4%
U	1.9%	2.3%	2.6%	1.9%	21.4%

j. stomach (ICD-10 C16)

	Stage Group (Rapid)				
Stage Group (Gold Standard)	1	2	3	4	Unknown
1	66.2%	3.0%	NA	0.1%	7.2%
2	21.1%	45.7%	6.8%	0.8%	6.4%
3	5.6%	39.2%	66.3%	2.9%	10.6%
4	1.8%	8.5%	23.2%	94.1%	32.3%
U	5.3%	3.5%	3.7%	2.2%	43.4%

k. uterus (ICD-10 C54-C55)

	Stage Group (Rapid)				
Stage Group (Gold Standard)	1	2	3	4	Unknown
1	97.6%	11.4%	5.8%	7.3%	47.6%
2	0.6%	83.6%	1.2%	2.6%	4.2%
3	0.5%	1.8%	87.7%	7.3%	7.8%
4	0.2%	1.6%	2.4%	76.1%	8.7%
U	1.1%	1.6%	2.9%	6.7%	31.6%

l. pancreas (ICD-10 C25)

Stage Group (Rapid)

Stage Group (Gold Standard)	1	2	Stage Group (Rapid)	3	4	Unknown
Stage Group (Gold Standard)	1	2	3	4	Unknown	
1	73.5%	3.4%	1.1%	0.3%		9.3%
2	16.2%	75.0%	2.3%	0.4%		6.3%
3	4.6%	12.3%	89.2%	0.6%		7.2%
4	3.4%	6.1%	6.0%	97.9%		49.1%
U	2.3%	3.2%	1.3%	0.8%		28.1%

m. lymphoma (ICD-10 C81-86, C88)

Stage Group (Gold Standard)	Stage Group (Rapid)				
	1	2	3	4	Unknown
1	90.6%	1.2%	0.5%	0.5%	14.3%
2	0.9%	93.6%	1.2%	0.4%	10.9%
3	0.3%	1.3%	90.2%	1.5%	13.5%
4	5.8%	2.4%	7.1%	93.0%	35.9%
U	2.3%	1.6%	1.0%	4.6%	25.5%

## "Early" vs "Late" stage

Below in table 6 we repeat the above tabulations but now grouping Rapid and Gold Standard cancers into "Early" (TNM stage group 1 & 2) or "Late" (TNM stage group 3 & 4) categories. We see that 62% of breast cancers are identified as "Early" stage in the Rapid Registrations dataset compared to 76% in the Gold Standard Registration data due to the higher proportion of "Unknown" stage tumours (28% vs 10% respectively).

As with the more detailed stage data, there is a high degree of concordance between the gold standard and rapid registration stage fields if a known stage can be identified.

Table 6: Summary proportions of "Early" vs "Late" stage for Rapid Registrations and Gold Standard Registrations

Broad Cancer Group	Stage Group	Count (Rapid)	Percentage (Rapid)	Count (Gold Standard)	Percentage (Gold Standard)
Bladder	Early	4106	43.5%	4740	50.2%
Bladder	Late	814	8.6%	1515	16.0%
Bladder	Unknown	4525	47.9%	3190	33.8%
Breast	Early	26146	62.0%	31979	75.9%
Breast	Late	4223	10.0%	5427	12.9%
Breast	Unknown	11788	28.0%	4751	11.3%
Colorectum	Early	11962	37.5%	13095	41.0%
Colorectum	Late	13348	41.8%	16526	51.8%
Colorectum	Unknown	6594	20.7%	2283	7.2%
Kidney	Early	2817	35.6%	3805	48.1%
Kidney	Late	2044	25.8%	3113	39.3%
Kidney	Unknown	3054	38.6%	997	12.6%
Lung	Early	8803	25.2%	9364	26.8%
Lung	Late	22227	63.6%	24388	69.8%
Lung	Unknown	3894	11.1%	1172	3.4%

Broad Cancer Group	Stage Group	Count (Rapid)	Percentage (Rapid)	Count (Gold Standard)	Percentage (Gold Standard)
Lymphoma	Early	1856	15.6%	3332	28.0%
Lymphoma	Late	3844	32.3%	6815	57.3%
Lymphoma	Unknown	6197	52.1%	1750	14.7%
Melanoma	Early	8748	66.6%	10934	83.3%
Melanoma	Late	628	4.8%	1361	10.4%
Melanoma	Unknown	3756	28.6%	837	6.4%
Oesophagus	Early	1578	19.5%	1401	17.3%
Oesophagus	Late	4422	54.6%	5310	65.5%
Oesophagus	Unknown	2106	26.0%	1395	17.2%
Ovary	Early	1337	28.8%	1569	33.8%
Ovary	Late	1845	39.8%	2590	55.9%
Ovary	Unknown	1455	31.4%	478	10.3%
Pancreas	Early	975	12.9%	1440	19.0%
Pancreas	Late	2796	36.9%	5020	66.2%
Pancreas	Unknown	3813	50.3%	1124	14.8%
Prostate	Early	17152	37.3%	22810	49.6%
Prostate	Late	15938	34.7%	19662	42.8%
Prostate	Unknown	12854	28.0%	3472	7.6%
Stomach	Early	936	24.6%	805	21.2%
Stomach	Late	1539	40.5%	2340	61.6%
Stomach	Unknown	1324	34.9%	654	17.2%
Uterus	Early	5087	65.5%	5838	75.2%
Uterus	Late	1214	15.6%	1352	17.4%
Uterus	Unknown	1467	18.9%	578	7.4%

In Table 7a-m below the distribution of the stage allocation between the Rapid Registrations and the Gold Standard Registrations are examined, aggregated into Early and Late stage.

Tables 7a-m: "Early" vs "late" stage comparison between Rapid Registrations and Gold Standard Registrations

a. bladder (ICD-10 C67)				
Stage Category (Gold Standard)	Stage Category (Rapid)			
	Early	Late	Unknown	
Early	83.2%	19.3%	25.8%	
Late	9.1%	75.1%	11.7%	
Unknown	7.7%	5.7%	62.5%	
b. breast (ICD-10 C50)				
Stage Category (Gold Standard)	Stage Category (Rapid)			
	Early	Late	Unknown	

Stage Category (Gold Standard)	Stage Category (Rapid)		
	Early	Late	Unknown
Early	94.7%	13.6%	56.4%
Late	2.1%	82.4%	11.9%
Unknown	3.2%	4.0%	31.7%

c. colorectum (ICD-10 C18-C20)

Stage Category (Gold Standard)	Stage Category (Rapid)		
	Early	Late	Unknown
Early	88.9%	5.2%	26.7%
Late	9.1%	93.4%	45.1%
Unknown	2.0%	1.4%	28.2%

d. kidney (ICD-10 C64)

Stage Category (Gold Standard)	Stage Category (Rapid)		
	Early	Late	Unknown
Early	90.9%	3.5%	38.4%
Late	3.3%	93.3%	36.4%
Unknown	5.8%	3.2%	25.1%

e. lung (ICD-10 C33-C34)

Stage Category (Gold Standard)	Stage Category (Rapid)		
	Early	Late	Unknown
Early	94.9%	1.4%	17.8%
Late	4.3%	98.0%	57.1%
Unknown	0.7%	0.6%	25.1%

f. melanoma (ICD-10 C43)

Stage Category (Gold Standard)	Stage Category (Rapid)		
	Early	Late	Unknown
Early	92.2%	16.6%	73.6%
Late	5.2%	76.6%	11.3%
Unknown	2.6%	6.8%	15.0%

g. Oesophagus (ICD-10 C15)

Stage Category (Gold Standard)	Stage Category (Rapid)		
	Early	Late	Unknown
Early	61.6%	3.9%	12.3%
Late	32.9%	88.6%	41.5%

Stage Category (Gold Standard)	Stage Category (Rapid)		
	Early	Late	Unknown
Unknown	5.5%	7.6%	46.2%
h. ovary (ICD-10 C56-C57)			

Stage Category (Gold Standard)	Stage Category (Rapid)		
	Early	Late	Unknown
Early	97.4%	1.0%	17.0%
Late	1.5%	96.0%	54.9%
Unknown	1.1%	3.0%	28.0%

i. prostate (ICD-10 C61)

Stage Category (Gold Standard)	Stage Category (Rapid)		
	Early	Late	Unknown
Early	92.9%	4.8%	47.5%
Late	5.1%	92.8%	31.1%
Unknown	2.0%	2.3%	21.4%

j. stomach (ICD-10 C16)

Stage Category (Gold Standard)	Stage Category (Rapid)		
	Early	Late	Unknown
Early	62.6%	2.5%	13.6%
Late	33.2%	94.9%	43.0%
Unknown	4.2%	2.6%	43.4%

k. uterus (ICD-10 C54-C55)

Stage Category (Gold Standard)	Stage Category (Rapid)		
	Early	Late	Unknown
Early	97.8%	8.2%	51.9%
Late	1.0%	87.3%	16.5%
Unknown	1.2%	4.4%	31.6%

l. pancreas (ICD-10 C25)

Stage Category (Gold Standard)	Stage Category (Rapid)		
	Early	Late	Unknown
Early	82.5%	1.5%	15.6%
Late	14.7%	97.6%	56.3%
Unknown	2.9%	0.9%	28.1%

m. lymphoma (ICD-10 C81-C86, C88)

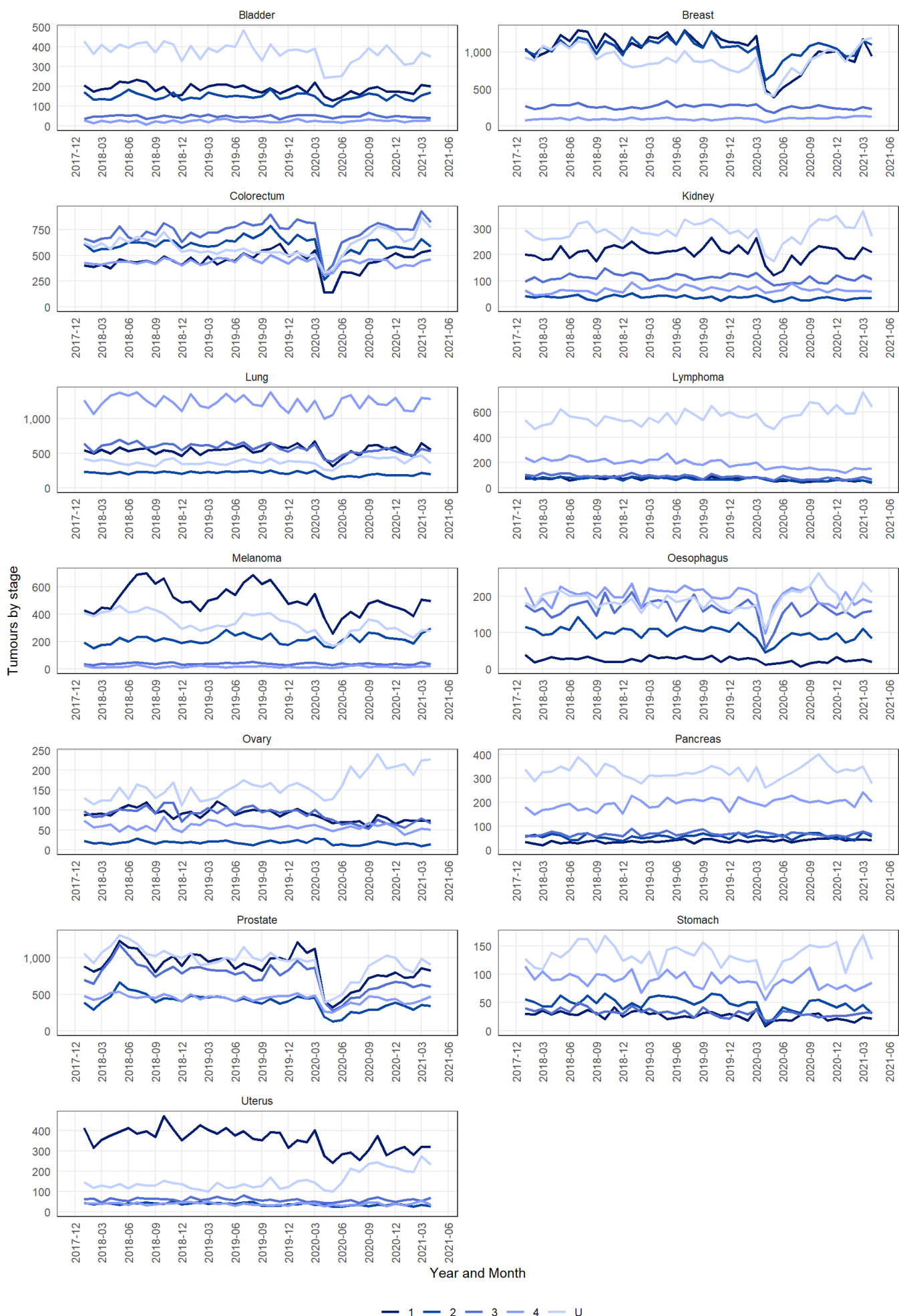
Stage Category (Gold Standard)	Stage Category (Rapid)		
	Early	Late	Unknown
Early	93.2%	1.1%	25.2%
Late	4.9%	95.4%	49.3%
Unknown	1.9%	3.5%	25.5%

### Stage trends over time

Figure 13 shows the monthly variation of the incidence count by stage at diagnosis for a number of common cancers. Allowing for variation in the number of working days in each month (which affects the overall number of tumours diagnosed per month) and for statistical fluctuation there is little evidence of any stage shift in the period displayed. The feature around May 2018 in the prostate cancer trends can be ascribed to the so called 'Turnbull-Fry effect' (<https://www.ndrs.nhs.uk/examining-the-fry-and-turnbull-effect-on-prostate-cancer-incidence-in-england/>).

Figure 13: Stage trends over time



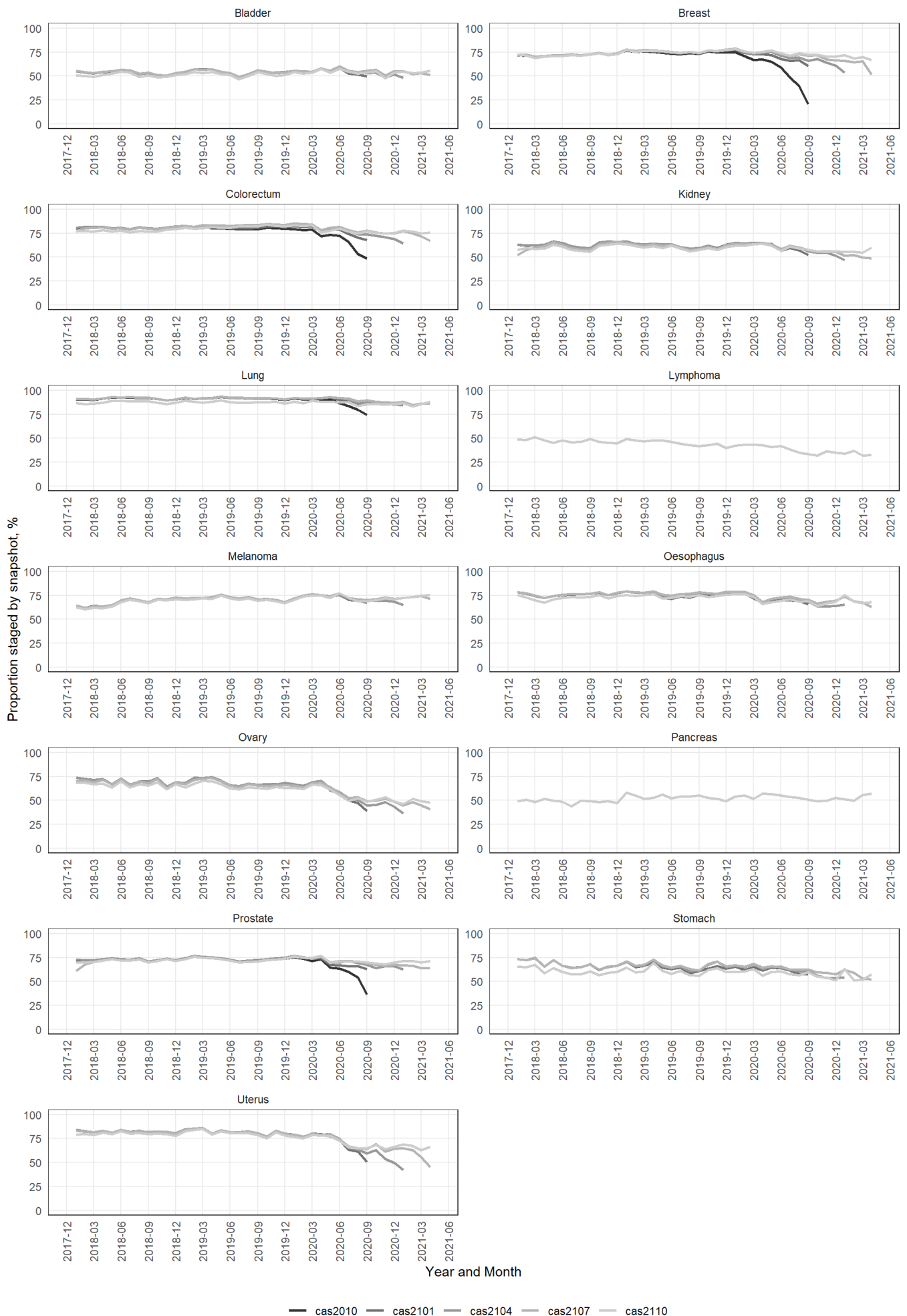


Source: NHS Digital, National Cancer Registration and Analysis Service

Stage completeness by snapshot

Figure 14 shows the completeness of stage by tumour type for one snapshot per quarter. Stage completeness continues to increase and lags behind the incidence completeness due to staging activity happening up to several months after diagnosis.

Figure 14: Stage completeness by snapshot

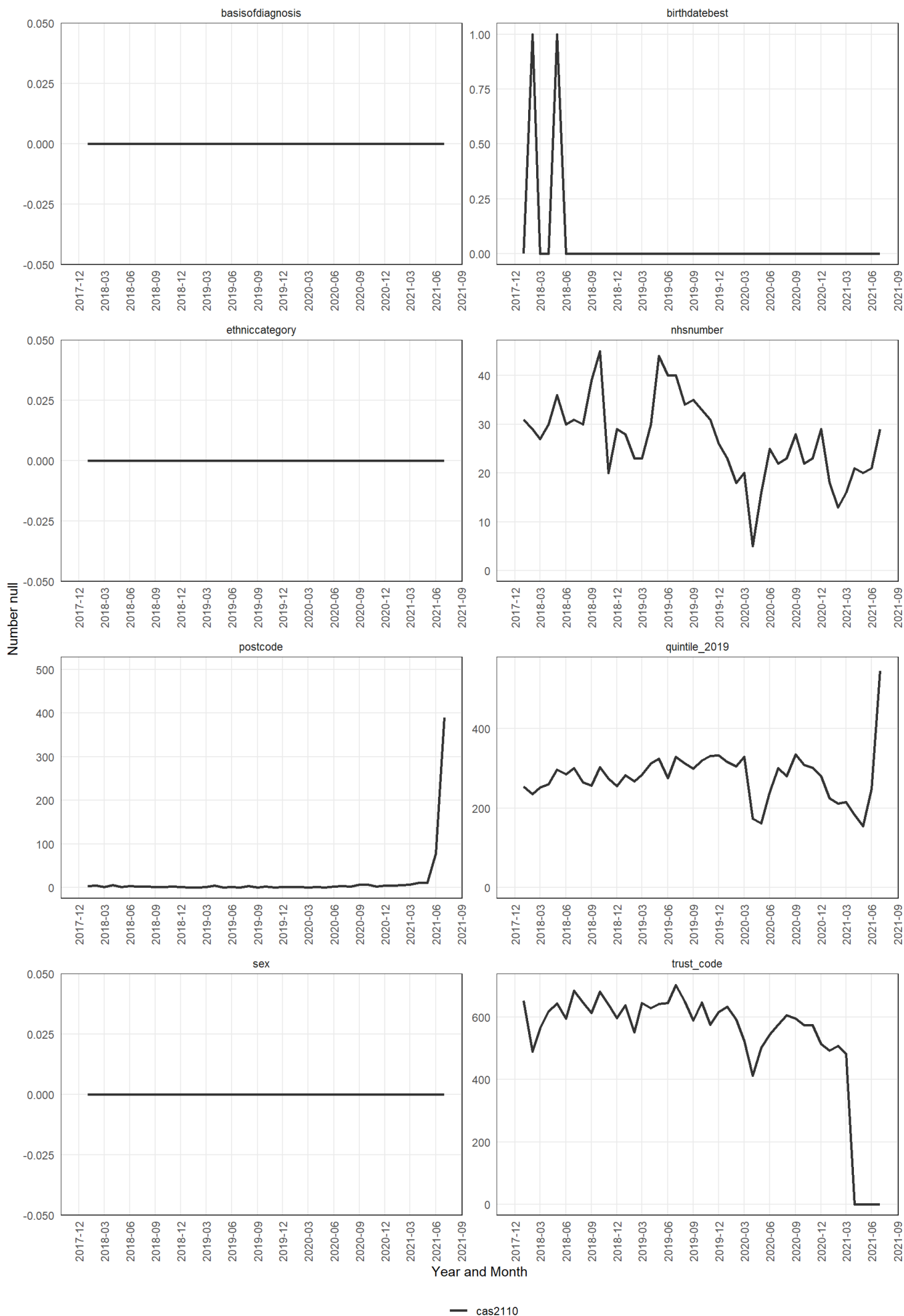


Source: NHS Digital, National Cancer Registration and Analysis Service

Counts of missing data

Figure 15 shows the count of tumours per month where the indicated data item is missing. Larger counts in the most recent months are to be expected.

Figure 15: Counts of missing data

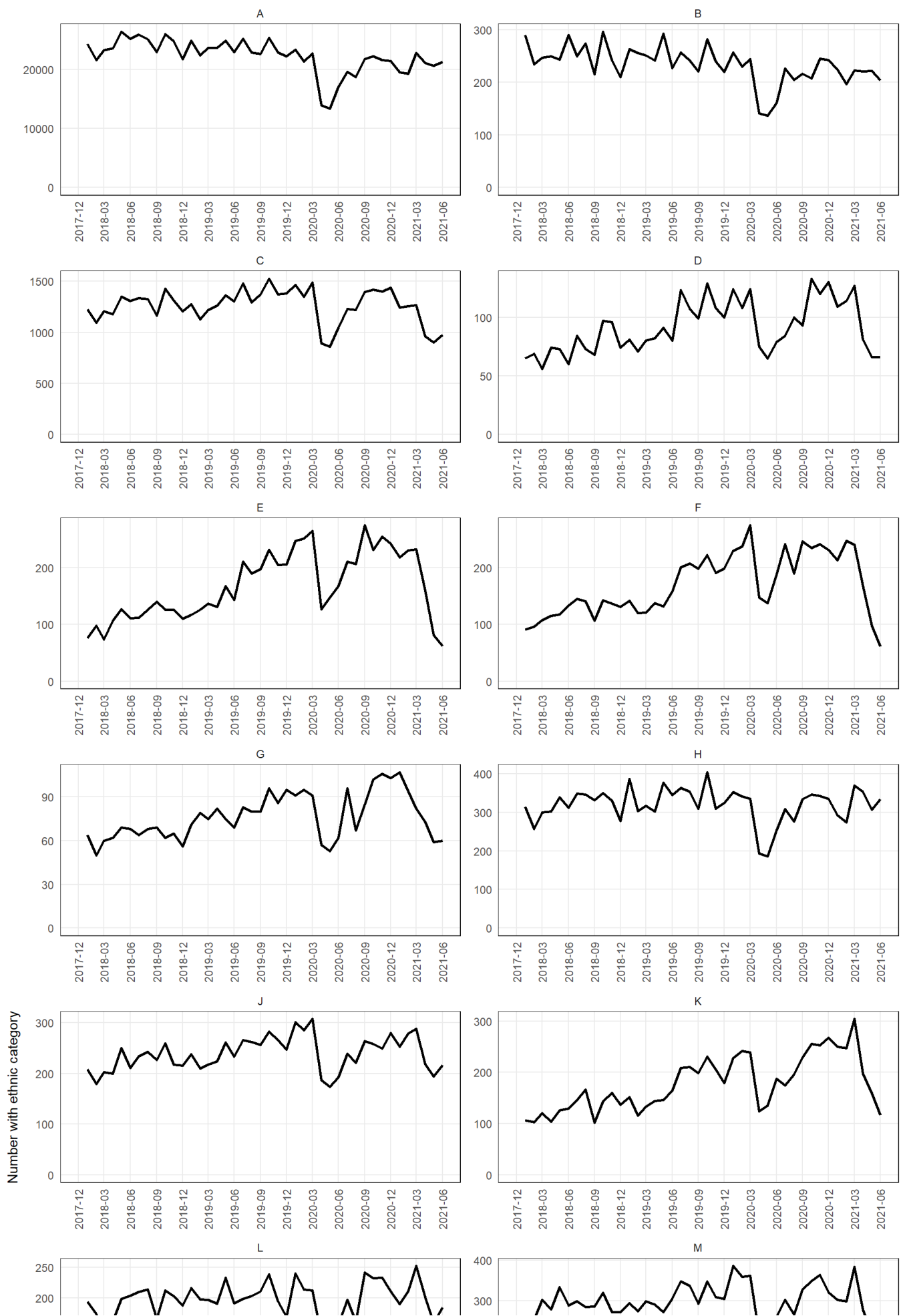


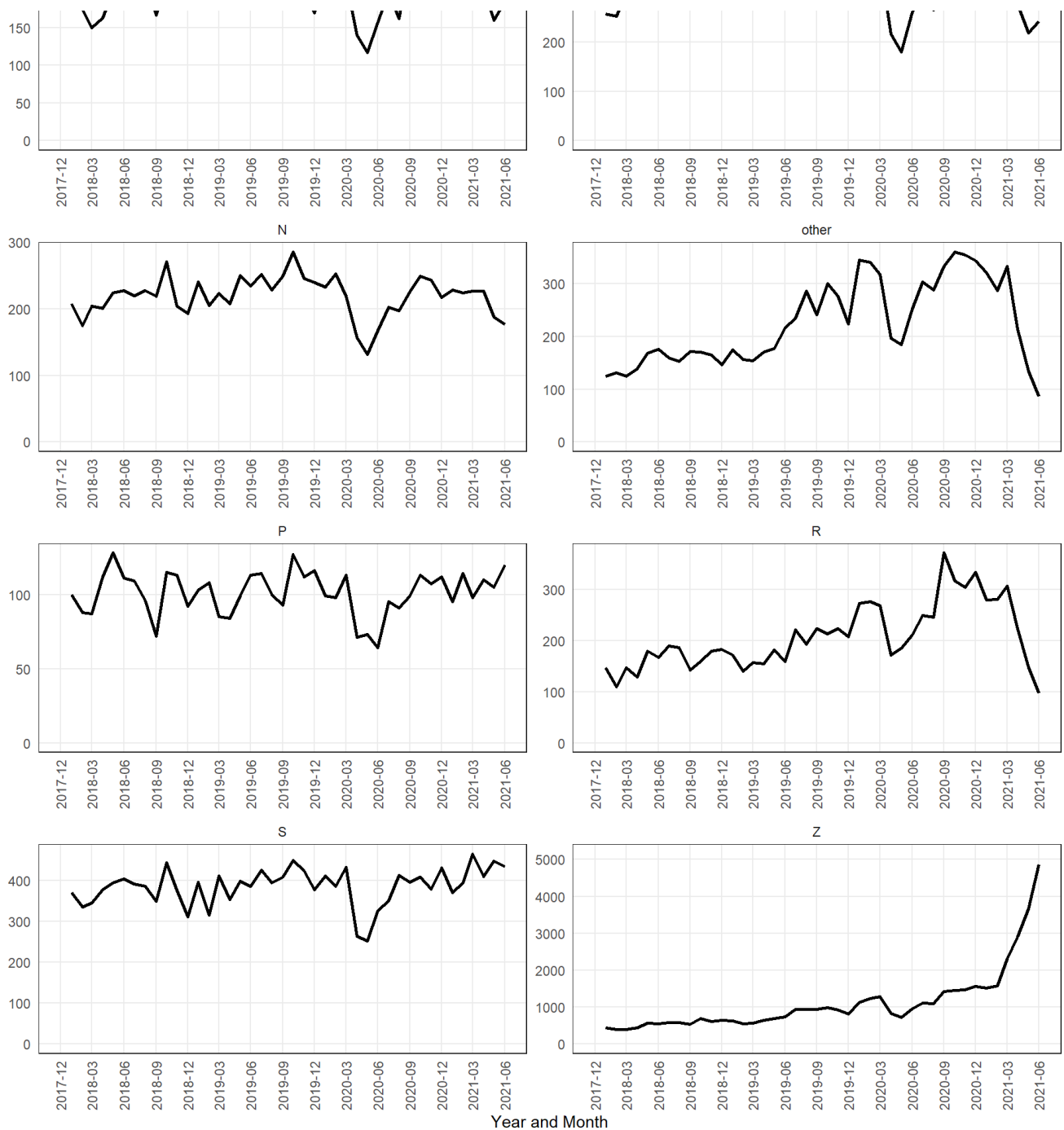
Source: NHS Digital, National Cancer Registration and Analysis Service

Ethnicity completeness

Figure 16 shows the count of tumours per month where the indicated data item is missing. Larger counts in the most recent months are to be expected.

Figure 16: Ethnicity completeness





Source: NHS Digital, National Cancer Registration and Analysis Service

## Tumour source

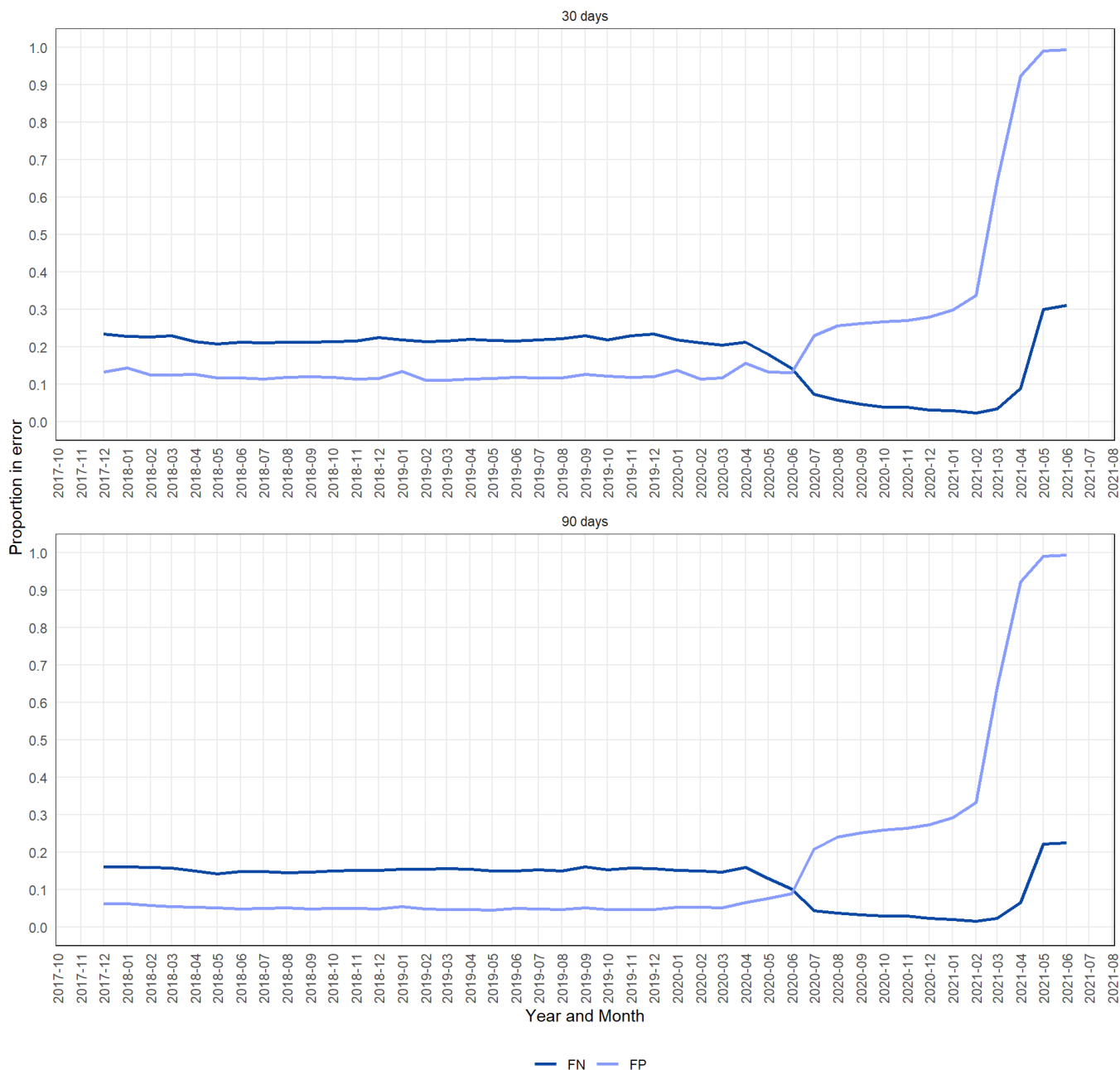
Figure 17 shows the proportion of tumours created by the source of the diagnosis - i.e., which dataset was used to create them, by month

Figure 17: Tumour source dataset

## False positive and falsenegative proportion by month

Figure 18 shows the False Negative and False Positive error proportions by month for the broader matching criteria and a matching period of 90 and 30 days.

Figure 18: Monthly False Positive and False Negative proportions



Source: NHS Digital, National Cancer Registration and Analysis Service

## Appendix 1 - List of pathway events

Table A1: AT\_RAPID\_PATHWAY: event list

EVENT_TYPE	EVENT_DESC	EVENT_PROPERTY_1	EVENT_PROPERTY_2	EVENT_PROPERTY_3	EVENT_DATE	Linkage
1	CWT Treatment Period Start Date	CWT First Treatment Flag	CWT SITE_ICD10	CWT Cancer Treatment Event Type	Treat period start	NHSNUMBER
2	CWT Treatment Start	CWT Treatment Modality	CWT Cancer Treatment Event type		Treatment start date	NHSNUMBER
3	CWT MDT Begin	CWT MDT Cancer Care Plan discussed indicator			MDT date	NHSNUMBER
4	CWT Faster Diagnosis Period End	(null)	Faster Diagnosis Period site		Faster Diagnosis Period end date	NHSNUMBER



EVENT_TYPE	EVENT_DESC	EVENT_PROPERTY_1	EVENT_PROPERTY_2	EVENT_PROPERTY_3	EVENT_DATE	Linkage
5	HES Admitted Patient Care Episode	Treatment speciality	All ICD-10 codes (for episode)	All OPCS-4 codes (for episode)	Episode Start date - Episode end date	NHSNUMBER
6	HES Admitted Patient Care Operation	OPCS codes (for date) in POS order	ICD-10 codes (for episode)		Operation date	NHSNUMBER
7	SACT Cycle	Benchmark group	Cycle number	Treatment intent	Cycle start date	PATIENTID
8	RTDS Episode	Radiotherapy intent	ICD-10 diagnosis code		Episode treatment start date	PATIENTID
9	Tumour diagnosis (Provisional)	Statusofregistration	ICD-10 diagnosis code	Stage_best	Diagnosisdatebest	PATIENTID
10	Patient last event date	Vitalstatus			Dateofvitalstatus1 (start of range)	PATIENTID
11	HES major surgery (historical)	OPCS-4 code	ICD-10 diagnosis code	Further notes/constraints	Operation date	NHSNUMBER
12	HES major surgery (historical, further constraints)	OPCS-4 code	ICD-10 diagnosis code	Further notes/constraints	Operation date	NHSNUMBER
13	HES major surgery (new)	OPCS-4 code	ICD-10 diagnosis code	Further notes/constraints	Operation date	NHSNUMBER
14	RAWDATA major surgery (historical)	OPCS-4 code	ICD-10 diagnosis code	Further notes/constraints	Operation date	PATIENTID
15	RAWDATA major surgery (historical, further constraints)	OPCS-4 code	ICD-10 diagnosis code	Further notes/constraints	Operation date	PATIENTID
16	RAWDATA major surgery (new)	OPCS-4 code	ICD-10 diagnosis code	Further notes/constraints	Operation date	PATIENTID
17	Prior tumour diagnosis	Statusofregistration	ICD-10 diagnosis code	Stage_best	Diagnosisdatebest	PATIENTID
18	Tumour diagnosis (Final)	Statusofregistration	ICD-10 diagnosis code	Stage_best	Diagnosisdatebest	PATIENTID
19	Patient vital status date	Vitalstatus	ICD-10 underlying cause of death		Vitalstatusdate	PATIENTID
20	RAWDATA holistic needs assessment record	HNA point of pathway **	Primary diagnosis	Laterality	Date of HNA	PATIENTID
21	RAWDATA staging	Inferred best stage	ICD-10 diagnosis code	TNM components	Collected stage date	PATIENTID
22	CWT First Seen	REF_SOURCE	Categorisation of TWW, screening and consultant upgrade cases, where relevant	Suspected cancer referral type	Date first seen	NHSNUMBER

EVENT_TYPE	EVENT_DESC	EVENT_PROPERTY_1	EVENT_PROPERTY_2	EVENT_PROPERTY_3	EVENT_DATE	Linkage
23	HES diagnostic event	OPCS-4 code	Description	BX/LD	Operation date	NHSNUMBER
50	Skeleton Tumour creation	E_base_record type	ICD-10 diagnosis code		Diagnosisdate	PATIENTID
51	Diagnosis reported in COSD	Number of times reported	ICD-10 diagnosis code	E_base_record type	Diagnosisdate	NHSNUMBER
52	CWT estimated diagnosis date	CWT First Treatment Flag	CWT SITE_ICD10	CWT Cancer Treatment Event Type	Adjusted treat period start	NHSNUMBER
53	HES inferred tumour	HES cancer group	ICD-10 diagnosis code		Episode start date	NHSNUMBER
54	COSD diagnosis submission	E_base_record primary diagnoses	ICD-10 diagnosis code (submission)		Diagnosis date (submission)	PATIENTID
55	RAWDATA biopsy record	Laterality	ICD-10 diagnosis code		Collected date/authorised date	PATIENTID
56	RAWDATA imaging record	Laterality	ICD-10 diagnosis code	Procedure_date - diagdate	Diagdate	PATIENTID
57	RAWDATA HNA diagnosis	Laterality	Primary diagnosis (ICD-10)		Diagdate	PATIENTID
101	Inferred diagnosis (54 only)	Event_property_1	ICD-10 diagnosis code	Cancer group	First recorded date	PATIENTID

\*: [https://www.datadictionary.nhs.uk/data\\_dictionary/attributes/p/prev/primary\\_cancer\\_site\\_for\\_cancer\\_faster\\_diagnosis\\_pathway\\_de.asp?shownav=0](https://www.datadictionary.nhs.uk/data_dictionary/attributes/p/prev/primary_cancer_site_for_cancer_faster_diagnosis_pathway_de.asp?shownav=0)  
 (https://www.datadictionary.nhs.uk/data\_dictionary/attributes/p/prev/primary\_cancer\_site\_for\_cancer\_faster\_diagnosis\_pathway\_de.asp?shownav=0)

\*\* : [https://www.datadictionary.nhs.uk/data\\_dictionary/attributes/h/ho/holistic\\_needs\\_assessment\\_point\\_of\\_pathway\\_for\\_cancer\\_de.asp?shownav=0](https://www.datadictionary.nhs.uk/data_dictionary/attributes/h/ho/holistic_needs_assessment_point_of_pathway_for_cancer_de.asp?shownav=0)  
 (https://www.datadictionary.nhs.uk/data\_dictionary/attributes/h/ho/holistic\_needs\_assessment\_point\_of\_pathway\_for\_cancer\_de.asp?shownav=0)

## Appendix 2 - List of Rapid Registration fields available

Table A2: AT\_RAPID\_TUMOUR: field list

COLUMN_NAME	DATA_TYPE	Notes
INDIVIDUALID	NUMBER(11,0)	Matches AT_RAPID_PATHWAY for each event with event_type=101
PATIENTID	NUMBER(19,0)	Matches AT_RAPID_PATHWAY for each event with event_type=101
NHSNUMBER	VARCHAR2(12 BYTE)	Matches AT_RAPID_PATHWAY for each event with event_type=101
TUMOUR_AVPID	NUMBER	Matches AT_RAPID_PATHWAY for each event with event_type=101
DIAGNOSISDATE	DATE	Matches AT_RAPID_PATHWAY for each event with event_type=101
TUMOUR_SITE	VARCHAR2(255 BYTE)	Matches AT_RAPID_PATHWAY for each event with event_type=101 (event_property_2)
BIRTHDATEBEST	DATE	Taken from Encore

COLUMN_NAME	DATA_TYPE	Notes
SEX	VARCHAR2(255 BYTE)	Taken from Encore
POSTCODE	VARCHAR2(255 BYTE)	Taken from Encore
SURNAME	VARCHAR2(64 BYTE)	Taken from Encore
FORENAME	VARCHAR2(64 BYTE)	Taken from Encore
STAGE	VARCHAR2(255 BYTE)	Defined for selected cancer sites
ETHNICITY	VARCHAR2(255 BYTE)	Taken from Encore
FINAL_ROUTE	VARCHAR2(22 BYTE)	Final Route to Diagnosis using an adapted version of the standard NCRAS methodology
QUINTILE_2019	VARCHAR2(26 BYTE)	Income deprivation quintile defined using the standard NCRAS methodology
CHRL_TOT_27_03	NUMBER	Charlson score defined using the standard NCRAS methodology
TUMOUR_MORPHOLOGY	VARCHAR2(255 BYTE)	Tumour morphology as recorded in the COSD system

## Appendix 3 - Cancer groups used for matching

Table A3: Rapid Registration ICD-10 tumour inclusion list

ICD	CANCER_GROUP	ICD	CANCER_GROUP
C00	Head & Neck	C54	Gynae
C01	Head & Neck	C55	Gynae
C02	Head & Neck	C56	Gynae
C03	Head & Neck	C57	Gynae
C04	Head & Neck	C58	Gynae
C05	Head & Neck	C59	Other
C06	Head & Neck	C60	Urology
C07	Head & Neck	C61	Prostate
C08	Head & Neck	C62	Urology
C09	Head & Neck	C63	Urology
C10	Head & Neck	C64	Urology
C11	Head & Neck	C65	Urology
C12	Head & Neck	C66	Urology
C13	Head & Neck	C67	Urology
C14	Head & Neck	C68	Urology
C15	O-G	C69	Brain & CNS
C16	O-G	C70	Brain & CNS
C17	Upper GI	C71	Brain & CNS
C18	Colorectal	C72	Brain & CNS

ICD	CANCER_GROUP	ICD	CANCER_GROUP
C19	Colorectal	C73	Endocrine
C20	Colorectal	C74	Endocrine
C21	Colorectal	C75	Endocrine
C22	Upper GI	C76	Unknown Primary
C23	Upper GI	C77	Unknown Primary
C24	Upper GI	C78	Unknown Primary
C25	Upper GI	C79	Unknown Primary
C26	Upper GI	C80	Unknown Primary
C27	Other	C81	Haematological
C28	Other	C82	Haematological
C29	Other	C83	Haematological
C30	Head & Neck	C84	Haematological
C31	Head & Neck	C85	Haematological
C32	Head & Neck	C86	Haematological
C33	Lung	C87	Haematological
C34	Lung	C88	Haematological
C35	Other	C89	Haematological
C36	Other	C90	Haematological
C37	Other	C91	Haematological
C38	Lung	C92	Haematological
C39	Lung	C93	Haematological
C40	Bone & ST	C94	Haematological
C41	Bone & ST	C95	Haematological
C42	Other	C96	Haematological
C43	Melanoma	C97	Unknown Primary
C44	NMSC	D05	Breast
C45	Lung	D06	Gynae
C46	Bone & ST	D09	Urology
C47	Brain & CNS	D32	Brain & CNS
C48	Gynae	D33	Brain & CNS
C49	Bone & ST	D35	Brain & CNS
C50	Breast	D41	Urology
C51	Gynae	D42	Brain & CNS
C52	Gynae	D43	Brain & CNS
C53	Gynae	D44	Brain & CNS

## Appendix 4 - Alternative defining events

Several options were considered as to the defining events for the Rapid Registrations. Both standalone datasets, subsets of standalone datasets, and combined datasets were explored and their FNE and FPE figures quantified. A subset of these alternatives are presented below as a demonstration of the process but the majority of this exploratory work is out of scope for this document.

Candidates for diagnosis events from the three main datasets that are rapidly available and have nominally full coverage of cancer patients are shown below (SACT and RTDS were also examined but data is not presented). Of the three, the CWT data has the best FPE but the FNE is substantially higher than the COSD dataset. HES produced the worst results in both measures. A filtering process was applied to the standalone COSD data to remove apparently new diagnoses that were actually recurrences of prior tumours. This improved the FPE at a cost of increasing the FNE. We continue to test whether this process can be further refined to improve the combined FPE and FNE figures, and monitor changes in the underlying datasets that might also give new opportunities to do so.

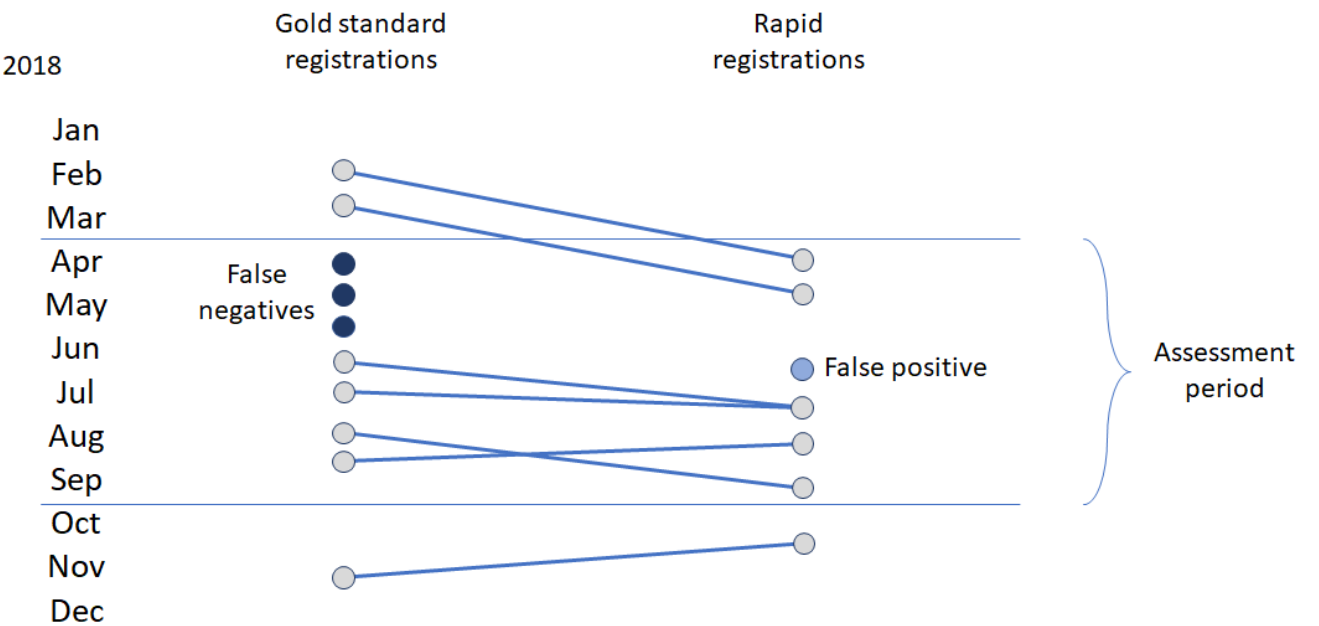
Table A4: Rapid Cancer Registrations: alternative defining events

Event	FPE	FNE
Event 52 - standalone CWT	7.6%	28.3%
Event 53 - standalone HES	13.2%	38.9%
Event 54 - standalone COSD	8.1%	15.8%
Event 101 (up to cas2106) - filtered COSD	5.2%	17.7%
Event 101 (cas2107) - filtered combined COSD/CWT	5.6%	16.4%
Event 101 (cas2108) - filtered combined COSD/CWT	5.1%	16.5%
Event 101 (cas2109) - filtered combined COSD/CWT	5.1%	16.6%
Event 101 (cas2110) - filtered combined COSD/CWT/HES	5.1%	14.7%

## Appendix 5 - Counts and error tabulations

Figure A1 shows an example for a very small dataset of how counts and error proportions are derived. This dataset has 10 Gold Standard Registrations and 7 Rapid Registrations overall (both indicated by the dots in the figure, with time running vertically over the course of 2018 and Gold Standard vs Rapid Registrations divided horizontally). Successful linkages between Gold Standard and Rapid Registrations are indicated by blue lines. False negatives and false positives are indicated. Only tumours in the 6-month assessment period are included in the tabulations below, although these can link to tumours outside the period as shown, and many-to-one linkages are also allowed. The false negative rate is therefore 3 in 7 and the false positive rate 1 in 6 below.

Figure A1: Illustration of counts and errors tabulation



Tables A5 and A6 below tabulate counts of Gold Standard and Rapid Registrations together with the numbers of false positive and false negative errors. When considering comparisons between figures the nature of the linkage and relationships displayed in the diagram above should be kept in mind.

Table A5: Counts and errors tabulation by cancer group

Cancer group	Gold Standard (GS) Registrations	Rapid Registrations	Difference	Percentage Rapid/GS	FPE	FNE
Brain & CNS	5458	4225	1233	77.4%	567	1643
Breast	28885	24760	4125	85.7%	403	2895
Colorectal	18919	17319	1600	91.5%	939	2055
Endocrine	1890	1490	400	78.8%	130	471
Gynae	9753	8665	1088	88.8%	533	1400
Haematological	13793	11984	1809	86.9%	803	2492
Head & Neck	5268	4868	400	92.4%	399	698
Lung	21586	19476	2110	90.2%	720	2528
Melanoma	8233	7749	484	94.1%	876	1047
O-G	6617	6343	274	95.9%	398	581
Prostate	26920	24917	2003	92.6%	469	2307
Bone & Soft Tissue	1132	1228	-96	108.5%	512	354
Unknown Primary	3423	2532	891	74.0%	1283	2108
Upper GI	9194	8198	996	89.2%	837	1749
Urology	16914	13326	3588	78.8%	658	3808

Table A6: Counts and errors tabulation by cancer site

Cancer site	Gold Standard (GS) Registrations	Rapid Registrations	Difference	Percentage Rapid/GS	FPE	FNE
C00	109	147	-38	134.9%	61	22
C01	643	453	190	70.5%	10	71
C02	604	613	-9	101.5%	16	83
C03	233	106	127	45.5%	5	67
C04	252	241	11	95.6%	11	30
C05	214	186	28	86.9%	8	32
C06	270	283	-13	104.8%	17	50
C07	236	276	-40	116.9%	88	51
C08	81	89	-8	109.9%	16	13
C09	913	754	159	82.6%	14	69
C10	150	239	-89	159.3%	9	30
C11	110	106	4	96.4%	6	15
C12	155	99	56	63.9%	1	10
C13	142	130	12	91.5%	11	22
C14	24	61	-37	254.2%	13	12
C15	3997	4204	-207	105.2%	118	275
C16	2620	2139	481	81.6%	236	306
C17	805	690	115	85.7%	143	245
C18	12407	11402	1005	91.9%	612	1519

Cancer site	Gold Standard (GS) Registrations	Rapid Registrations	Difference	Percentage Rapid/GS	FPE	FNE
C19	993	857	136	86.3%	25	118
C20	4874	4437	437	91.0%	102	368
C21	645	623	22	96.6%	76	50
C22	2616	2383	233	91.1%	258	543
C23	473	452	21	95.6%	26	66
C24	642	517	125	80.5%	29	95
C25	4509	3983	526	88.3%	137	683
C26	149	173	-24	116.1%	131	117
C30	162	148	14	91.4%	20	28
C31	92	64	28	69.6%	5	25
C32	878	873	5	99.4%	51	68
C33	13	11	2	84.6%	1	3
C34	20132	18144	1988	90.1%	503	2324
C37	166	87	79	52.4%	11	57
C38	72	349	-277	484.7%	42	23
C39	NA	13	NA	NA%	4	NA
C40	118	109	9	92.4%	12	23
C41	116	150	-34	129.3%	80	44
C43	8233	7749	484	94.1%	764	1047
C45	1203	872	331	72.5%	10	121
C46	68	41	27	60.3%	4	26
C47	26	18	8	69.2%	9	19
C48	283	392	-109	138.5%	111	84
C49	830	928	-98	111.8%	370	261
C50	25069	22100	2969	88.2%	214	2409
C51	641	500	141	78.0%	25	141
C52	92	94	-2	102.2%	10	19
C53	1318	1215	103	92.2%	36	164
C54	4093	3613	480	88.3%	79	252
C55	72	314	-242	436.1%	17	24
C56	2976	2219	757	74.6%	121	663
C57	268	295	-27	110.1%	25	51
C58	10	23	-13	230.0%	17	2
C60	302	296	6	98.0%	40	46
C61	26920	24917	2003	92.6%	243	2307
C62	1052	1027	25	97.6%	71	88

Cancer site	Gold Standard (GS) Registrations	Rapid Registrations	Difference	Percentage Rapid/GS	FPE	FNE
C63	30	17	13	56.7%	6	24
C64	4804	4151	653	86.4%	219	848
C65	411	307	104	74.7%	18	92
C66	356	245	111	68.8%	9	127
C67	4463	4939	-476	110.7%	134	808
C68	95	50	45	52.6%	4	43
C69	368	334	34	90.8%	36	55
C70	20	40	-20	200.0%	6	8
C71	2250	2076	174	92.3%	174	268
C72	78	78	0	100.0%	33	24
C73	1722	1372	350	79.7%	74	380
C74	114	77	37	67.5%	23	58
C75	54	41	13	75.9%	28	33
C76	94	317	-223	337.2%	234	79
C77	276	223	53	80.8%	146	107
C78	599	184	415	30.7%	144	405
C79	229	258	-29	112.7%	178	151
C80	2225	1550	675	69.7%	542	1366
C81	895	853	42	95.3%	10	77
C82	1203	1033	170	85.9%	10	143
C83	3150	2654	496	84.3%	34	370
C84	390	227	163	58.2%	12	123
C85	1341	907	434	67.6%	48	362
C86	NA	98	NA	NA%	3	NA
C88	203	367	-164	180.8%	13	43
C90	2518	2077	441	82.5%	42	499
C91	2205	1798	407	81.5%	72	492
C92	1748	1477	271	84.5%	201	332
C93	23	172	-149	747.8%	17	1
C94	29	140	-111	482.8%	118	10
C95	50	53	-3	106.0%	5	21
C96	38	128	-90	336.8%	105	19
D05	3816	2660	1156	69.7%	41	486
D09	4894	413	4481	8.4%	33	1497
D32	1340	729	611	54.4%	34	603
D33	416	480	-64	115.4%	64	197



Cancer site	Gold Standard (GS) Registrations	Rapid Registrations	Difference	Percentage Rapid/GS	FPE	FNE
D35	450	261	189	58.0%	37	231
D41	507	1881	-1374	371.0%	24	235
D42	139	7	132	5.0%	2	55
D43	261	176	85	67.4%	34	108
D44	110	26	84	23.6%	14	75

## Appendix 6 - False negative errors and basis of diagnosis

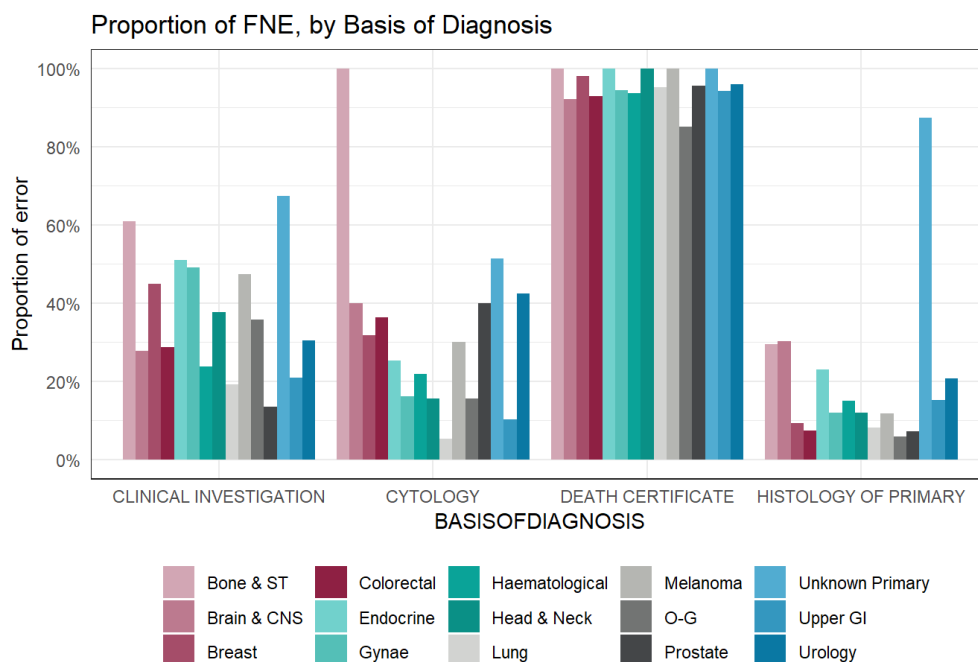
This appendix explores the reason for the overall age-dependence of the false negative error rate.

The most common methods of confirming a diagnosis (histology and cytology) account for the lowest proportion of false negatives (Figure A2). Where diagnosis comes from specific tumour markers, the Rapid Registrations are much more likely to "miss" the significant event or events. Patients diagnosed clinically (from imaging, consultation by a doctor but without a pathological sample being taken) are also more likely to be "missed" in the Rapid Registrations dataset.

Those patients for whom a diagnosis method cannot be determined (unknown) or died before they could be offered cancer treatment (death certificate), are most likely to be "missed" in the Rapid Registrations dataset. As Figure A3 indicates though, these account for a small proportion of those falsely omitted from the Rapid Registrations.

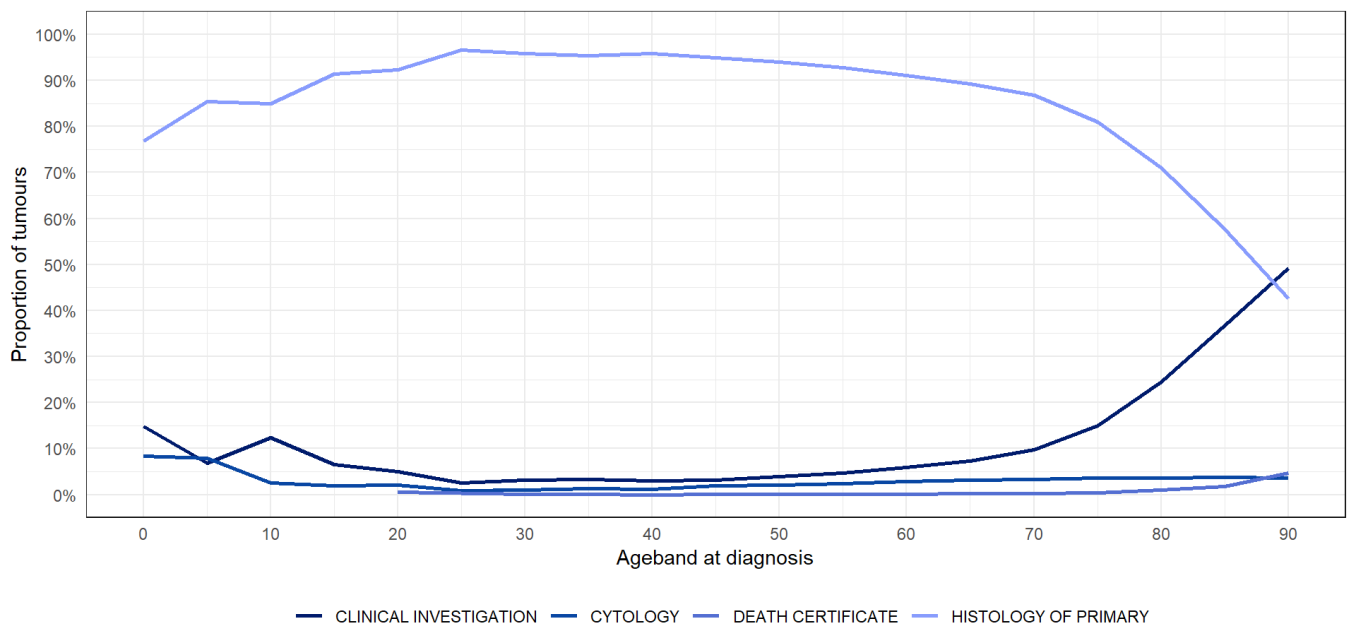
The marked reduction in the proportion of patients having their diagnosis confirmed from a pathological specimen (histology or cytology) explains the increase often observed at older ages in Figure A3, from the age of around 70, reflecting fewer patients having an invasive procedure performed on them as age increases. This is likely to be the reason behind the increasing false negative proportions by age observed overall and in most tumour groups (Figures 5 and 6).

Figure A2: The proportion of false negative Rapid Registrations by tumour group and basis of diagnosis, England, 2018



Source: NHS Digital, National Cancer Registration and Analysis Service

Figure A3: The proportion of false negative Rapid Registrations by method of diagnosis, England, 2018 (all tumour types combined)



Source: NHS Digital, National Cancer Registration and Analysis Service