
SPI-BIRDS: STANDARD DATA QUALITY CHECK

VERSION 1.1

FEBRUARY 17, 2022

PRODUCED BY:

SPI-BIRDS TEAM

STEFAN J.G. VRIEND, LIAM D. BAILEY, CHRIS TYSON,
ANTICA CULINA, & MARCEL E. VISSER

Table of Contents

1	Introduction	3
2	Check overview	5
2.1	Brood data	5
2.2	Capture data	6
2.3	Individual data	6
2.4	Location data	7
3	Brood data	8
3.1	Check B1	8
3.2	Check B2	8
3.3	Check B3	9
3.4	Check B4	9
3.5	Check B5a-d	9
3.6	Check B6	10
3.7	Check B7	10
3.8	Check B8	10
3.9	Check B9	10
3.10	Check B10	11
3.11	Check B11	11
3.12	Check B12	11
3.13	Check B13	11
3.14	Check B14	11
3.15	Check B15	12
4	Capture data	13
4.1	Check C1a-b	13
4.2	Check C2	14
4.3	Check C3	14
4.4	Check C4	14
4.5	Check C5	14
4.6	Check C6	15
5	Individual data	16
5.1	Check I1	16
5.2	Check I2	16
5.3	Check I3	16
5.4	Check I4	17
5.5	Check I5	17

6	Location data	18
6.1	Check L1	18

Introduction

This document is created as part of the **SPI-Birds** project. **SPI-Birds** aims to create a global network of Studies on Populations of Individuals - Birds (**SPI-Birds**), with the aim to archive data, improve data accessibility and transparency, and to facilitate collaboration. Within this project, we are building robust code pipelines that convert data stored in different formats (i.e., primary data) into data in a standard format (that consists of Brood data, Capture data, Individual data, Location data). This standard format ([version 1.1](#)) is aimed at facilitating greater collaboration by allowing data from multiple populations to be easily collated and compared.

In this document, we outline the standard data quality checks which are carried out as a part of the **SPI-Birds** workflow. The aim of the data quality check procedure is to increase the integrity of the data by highlighting all the values within the dataset that are unlikely or impossible. As quality checks are performed on data that have been already converted into the standard format, and because all the datasets are checked in the same way, homogeneity in data quality across different datasets is further increased.

Quality checks identify suspicious data records and flag them in two detailed quality check reports:

1. The first and main document contains records flagged as *potential errors*. Potential errors are values that are considered impossible

(e.g., negative values for clutch size). Whenever a record is flagged as a potential error, a line is added to this report with information on the type of check that was violated and the row number of the corresponding record. This row number refers to the column Row in the corresponding data table in the standard format and does not refer to the row number in the primary data.

2. The second document contains records flagged as *warnings*. Also in the case of warnings, flagged records are added to the report with information on the type of check and its row number in the data in the standard format. In addition, this report contains *verified records*. Some of the flagged records may be confirmed by a data owner to be true observations. We do not want these verified values to be flagged each time a new quality check is re-run. To avoid this, we have implemented an ‘approve-listing’ procedure that will prevent validated records (i.e., flagged but subsequently verified by the data owner) from appearing in future quality check reports.

In addition to these two reports, the quality check procedure adds two new columns (Warning and Error) to each of the data tables in the standard format to allow data users to easily identify and filter out potentially spurious records.

The standard data quality checks are run on three occasions. First, when the SPI-Birds team has created a new tailored pipeline and derived the data in the standard format, the quality checks are run on these standard data. This initial quality report is shared with the data owner. Second, when pipelines have been re-run following an update of the primary data by the data owner (e.g., when a new year of breeding information has been collected), the quality checks are re-run and the report shared with the data owner again. Data owners can decide to cross-check the records highlighted by the quality checks with their field notes, and if possible, correct some of the flagged entries, which will then be added to the list of verified records in the second quality check report. Third, when a user’s data request is approved, the quality checks are run on the requested data, and the corresponding quality check report sent to the user alongside the requested data.

Check overview

The standard data quality checks are continuously updated to improve existing checks and include new checks. The following checks are currently included in **version 1.1** of the standard data quality checks.

2.1 Brood data

A more detailed description of *Brood data* checks can be found in section 3.

CheckID	Description
B1	Compare ClutchSize_observed and BroodSize_observed
B2	Compare BroodSize_observed and NumberFledged_observed
B3	Compare LayDate_observed and HatchDate_observed
B4	Compare HatchDate_observed and FledgeDate_observed
B5	Compare breeding variables against reference values
B6	Compare BroodSize_observed and the number of chicks recorded in <i>Capture data</i>
B7	Check that BroodID is unique within populations
B8	Check that the order of ClutchType_observed is correct
B9	Compare the species of the parents of a brood
B10	Compare the species of the parents and the species of the brood
B11	Compare the species of the brood and the species of the chicks
B12	Check the sex of the mother
B13	Check the sex of the father
B14	Check that parents appear in <i>Capture data</i>
B15	Check that nest locations appear in <i>Location data</i>

2.2 Capture data

A more detailed description of *Capture data* checks can be found in section [4](#).

CheckID	Description
C1	Compare capture variables against reference values
C2	Check that chick age values are within the range of expectations
C3	Check that adults caught on a nest during the breeding season are recorded as the parents of that nest
C4	Check the age of subsequent captures
C5	Check that individuals in <i>Capture data</i> appear in <i>Individual data</i>
C6	Check that capture locations appear in <i>Location data</i>

2.3 Individual data

A more detailed description of *Individual data* checks can be found in section [5](#).

CheckID	Description
I1	Check that IndvID is unique within populations
I2	Check that BroodIDLaid and BroodIDFledged match the brood in <i>Capture data</i>
I3	Check for uncertainty in Sex_calculated
I4	Check for uncertainty in Species
I5	Check that all individuals in <i>Individual data</i> appear in <i>Capture data</i>

2.4 Location data

There are currently no checks on Location data. A more detailed description of *Location data* checks can be found in section 6.

CheckID	Description
L1	Check coordinates of locations

Brood data

3.1 Check B1

Compare clutch size (ClutchSize_observed) and brood size (BroodSize_observed) within each brood. We expect that clutch size is larger than or equal to brood size. Broods that do not meet this assumption are flagged. Broods that were not subject to experimental manipulation are flagged as a *potential error*. Broods that were subject to experimental manipulation are flagged as a *warning*, as clutch size might be smaller than brood size as a result of the experimental procedure.

3.2 Check B2

Compare brood size (BroodSize_observed) and number of fledglings (NumberFledged_observed) within each brood. We expect that brood size is larger than or equal to number of fledglings. Broods that do not meet this assumption are flagged. Broods that were not subject to experimental manipulation are flagged as a *potential error*. Broods that were subject to experimental manipulation are flagged as a *warning*, as brood size might be smaller than number of fledglings as a result of the experimental procedure.

3.3 Check B3

Compare lay date (LayDate_observed) and hatch date (HatchDate_observed) within each brood. We expect that lay date is earlier than hatch date. Expected length of incubation period is not considered (e.g., hatch date one day after lay date would not be flagged). Broods that do not meet this assumption are flagged as a *potential error*. This check does not flag records as a *warning*.

3.4 Check B4

Compare hatch date (HatchDate_observed) and fledge date (FledgeDate_observed) within each brood. We expect that hatch date is earlier than fledge date. Expected length of fledging period is not considered (e.g., fledge date one day after hatch date would not be flagged). Broods that do not meet this assumption are flagged as a *potential error*. This check does not flag records as a *warning*.

3.5 Check B5a-d

Compare breeding parameters (i.e., ClutchSize_observed, BroodSize_observed, NumberFledged_observed, LayDate_observed) against expected upper and lower limits of each variable. These reference values are generated by population- and species-specific data if the number of observations is sufficiently large ($n \geq 100$). Population-species combinations that have too few records are evaluated for negative values or other values that are considered impossible. This check does not flag records as a *warning*.

For ClutchSize_observed, BroodSize_observed and NumberFledged_observed, records are considered impossible if they are negative or larger than 2 times the 99th percentile and are flagged as a *potential error*. For population-species combinations that have too few observations ($n < 100$), records are considered impossible if they are negative and are flagged as a *potential error*.

For LayDate_observed, records are considered impossible if they are earlier than January 1st or later than December 31st for the current breeding season and are flagged as a *potential error*. For population-species combinations that have too few observations ($n < 100$), records are considered impossible if they are earlier than January 1st or later than December 31st for the current breeding season and are flagged as a *potential error*.

3.6 Check B6

Compare brood size (BroodSize_observed) of a specific brood with the number of chicks recorded in Individual data. We expect that these numbers should be equal. Records where BroodSize_observed is larger than the number of chicks recorded in Individual data are flagged as a *warning* because chicks might have died before ringing and measuring. In experimentally manipulated broods BroodSize_observed might be smaller than the number of chicks in Individual data. If so, the record is flagged as a *warning*. In non-manipulated broods BroodSize_observed should never be smaller than the number of chicks in Individual data. If so, the record is flagged as a *potential error*.

3.7 Check B7

Check that brood identity (BroodID) is unique within a population. Records with non-unique brood identifiers are flagged as *potential error*. Brood identities are not required to be unique across populations. This check does not flag records as a *warning*.

3.8 Check B8

Check that the order of clutch types (ClutchType_observed) per breeding female per season is correct. Replacement and second clutches can occur in any order, but never before first clutches. First clutches that are not the earliest brood in the season are flagged as a *potential error*. This check does not flag records as a *warning*.

3.9 Check B9

Check that the parents of a brood are of the same species. We expect that the parents of the majority of broods are marked as the same species. Common, biologically possible multi-species broods (i.e., PARMAJ – CYACAE and FICHYP – FICALB) are flagged as a *warning*. Other combinations of species are flagged as a *potential error*.

3.10 Check B10

Check whether the parents of a brood and the brood itself are of the same species. We expect that the parents and their broods are marked as the same species. Broods with a combination of species for which brood fostering is known to exist (i.e., PARMAJ – CYACAE and FICHYP – FICALB) are flagged as a *warning*. Other combinations of species are flagged as a *potential error*.

3.11 Check B11

Check that the chicks in a brood and the brood itself are of the same species. We expect that the chicks and the broods are marked as the same species. Broods with a combination of species for which brood fostering is known to exist (i.e., PARMAJ – CYACAE and FICHYP – FICALB) are flagged as a *warning*. Other combinations of species are flagged as a *potential error*.

3.12 Check B12

Check that the individual listed as the mother of a brood is female. Broods where the mother is male are flagged as a *potential error*. This check does not flag records as a *warning*.

3.13 Check B13

Check that the individual listed as the father of the brood is male. Broods where the father is female are flagged as a *potential error*. This check does not flag records as a *warning*.

3.14 Check B14

Check that all individuals recorded as parents of a brood appear at least once in *Capture data*. Broods with either of the parents missing from *Capture data* are flagged as a *potential error*. This check does not flag records as a *warning*.

3.15 Check B15

Check that all nest locations appear in *Location data*. Missing locations should not occur are flagged as a *potential error*. This check does not flag records as a *warning*.

Capture data

4.1 Check C1a-b

Compare chick and adult capture measurements (i.e., Mass, Tarsus) against expected upper and lower limits of each variable. For Mass in adults and Tarsus in adults and chicks, reference values are generated by population- and species-specific data if the number of observations is sufficiently large ($n \geq 100$). Population-species combinations that have too few records are evaluated for negative values or other values that are considered impossible. For Mass in chicks, reference values are calculated for each age (in days). A logistic growth model is fitted to determine reference values for each day. If the model fails, reference values are generated per age if the number of observations for chicks of that age is sufficiently large ($n \geq 100$). Population-species combinations that have too few records are evaluated for negative values or other values that are considered impossible.

Records are considered impossible if they are negative or larger than 2 times the 99th percentile and are flagged as a *potential error*. For population-species combinations that have too few observations ($n < 100$), records are considered impossible if they are negative and are flagged as a *potential error*. This check does not flag records as a *warning*.

4.2 Check C2

Check that chick age values (in number of days since hatching) are within the expected duration of the nestling period. Possible values currently include any number of days between 0 and 30 days since hatching. Values that are outside this range are flagged as a *potential error*. This check does not flag records as a *warning*.

4.3 Check C3

Check that adults captured on a nest during the breeding season are listed as the parents of that nest in *Brood data*. Adults caught on a nest that are not marked as the parents of that nest are flagged as a *warning*. This check does not flag records as a *potential error*.

4.4 Check C4

Check that the observed age of chronologically ordered captures is correct. The age recorded in an individual's subsequent capture is expected to be equal when the capture was in the same year or increase when the capture was in a later year. Records of an individual caught as an adult before records of the same individual caught as a chick are flagged as a *potential error*. Other records where the observed age of a capture is larger than the age of a subsequent capture are flagged as a *warning*.

4.5 Check C5

Check that all individuals in *Capture data* have a record in *Individual data*. *Individual data* is usually a direct product of *Capture data*, so this should never occur. Should it occur nonetheless, it is an indication of problems in the underlying pipelines and missing individuals are flagged as a *potential error*. This check does not flag records as a *warning*. This check is the opposite of check I5 (see 5.5).

4.6 Check C6

Check that all capture locations in *Capture data* have a record in Location data. Missing locations should not occur and are flagged as a *potential error*. This check does not flag records as a *warning*.

Individual data

5.1 Check I1

Check that individual identities (IndvID) are unique. Records with non-unique individual identifiers within populations are flagged as a *potential error*. Individual identities are not required to be unique across populations. This check does not flag records as a *warning*.

5.2 Check I2

Check that brood identities for an individual caught as a chick (BroodIDLaid, BroodIDFledged) match the correct nest in Capture data. Chicks caught on a nest that are not associated with corresponding brood identities are flagged as a *potential error*. This check does not flag records as a *warning*.

5.3 Check I3

Check for uncertainty in the sex of an individual (Sex_calculated). Individuals that have been recorded as both male ('M') and female ('F') in *Capture data*

are marked as conflicted ('C') by the pipeline and flagged as a *potential error* by the check. This check does not flag records as a *warning*.

5.4 Check I4

Check for uncertainty in the species of an individual (Species). Individuals who have been recorded as different species in Capture data are marked as conflicted ('CCCCC') by the pipeline and flagged as a *potential error* by the check. This check does not flag records as a *warning*.

5.5 Check I5

Check that all individuals recorded in *Individual data* also appear in *Capture data*. *Individual data* is usually a direct product of *Capture data*, so this should never occur. Should it occur nonetheless, it is an indication of problems in the underlying pipelines and missing individuals are flagged as a *potential error*. This check does not flag records as a *warning*. This check is the opposite of check C5 (see [4.5](#)).

Location data

6.1 Check L1

Check that coordinates of locations (Longitude, Latitude) are realistic. Records that are 20 km or farther away from the study site centre point are flagged as a *potential error*. Centre points are determined by the maximum kernel density of Longitude and Latitude. This check also produces a map that is printed in the quality check report. This check does not flag records as a *warning*.