

# Vetting data

---

The printable version is no longer supported and may have rendering errors. Please update your browser bookmarks and please use the default browser print function instead.

## Vetting Data

Vetting is the process of checking to assure the quality of the data we import into IFs. Even when data is imported through an automated application there could be errors or missing data in the imported series. Vetting the new data you are bringing into IFs requires you to compare the new series to the existing historical series in IFs.

There are a number of reasons for vetting data. One reason is that a country (or a number of countries) could be missed when importing due to not properly concurring countries with the proper country concordance table. Make sure the IFs country concordance list (i.e., the relevant column in the **Country** Translation table in IFs.mdb) is correct and up to date and you are using the correct concordance list as listed in the DataDict. Additionally, make sure there are no missing years in the new data. Sometimes the year column is missed because the data is formatted as text and not as numbers, or the source simply did not provide an update for that year. If this is the case, be sure to blend in the missing years in the IFs vetting tool. One more general issue is around the units. Make sure you perform any unit conversion required before you import the data. Comparing the new data with the current IFs series will give you an idea on any unit conversion that might be needed.

## Vetting Data Checklist:

- Data in Access file and source Excel file match
- Use IFs vetting tool to compare new data to historical data
- Large discrepancies between new and old data are documented
- Check to see if there are big spikes in the new, imported data
- Blend countries/years with missing data in IFs vetting tool
- Zeroes in Access file are actually zeroes and not null values
- No missing countries or years
- DataDict contains no errors, Original Source has website name and Name in Source contains variable source
- Initials added to both Access and Excel file names, as well as in the DataDict notes
- Send Word file with vetting notes, along with original files to puller and project leads

## Process:

1. After the data series are imported into IFs, the IFsDataImport.mdb file is passed to another RA (assigned as a vetter) who brings the data from the IFsDataImport.mdb file into the model via the "Vet Imported Data" feature in the Extended Features menu

2. The vetter uses the vetting tool to determine if there are any remarkable inconsistencies between the old and new data. This is a bit subjective, but necessarily so as the threshold for concern will be different depending on which series from which source is being vetted
3. If significant errors are found that need to be corrected, and these errors had occurred during the import process, these are documented in the RA's vetting notes and the files are passed back to the original data puller to correct and re-send.
4. If no errors are found, the vetter blends columns for new years and preserves historical data points as appropriate, and saves the changes in the IFsDataImport file.
5. The vetter passes the IFsDataImport file (renamed to reflect series/batch update name and completion date) back to the project lead and data team supervisor
6. Data team supervisor stores data import file for consolidation process

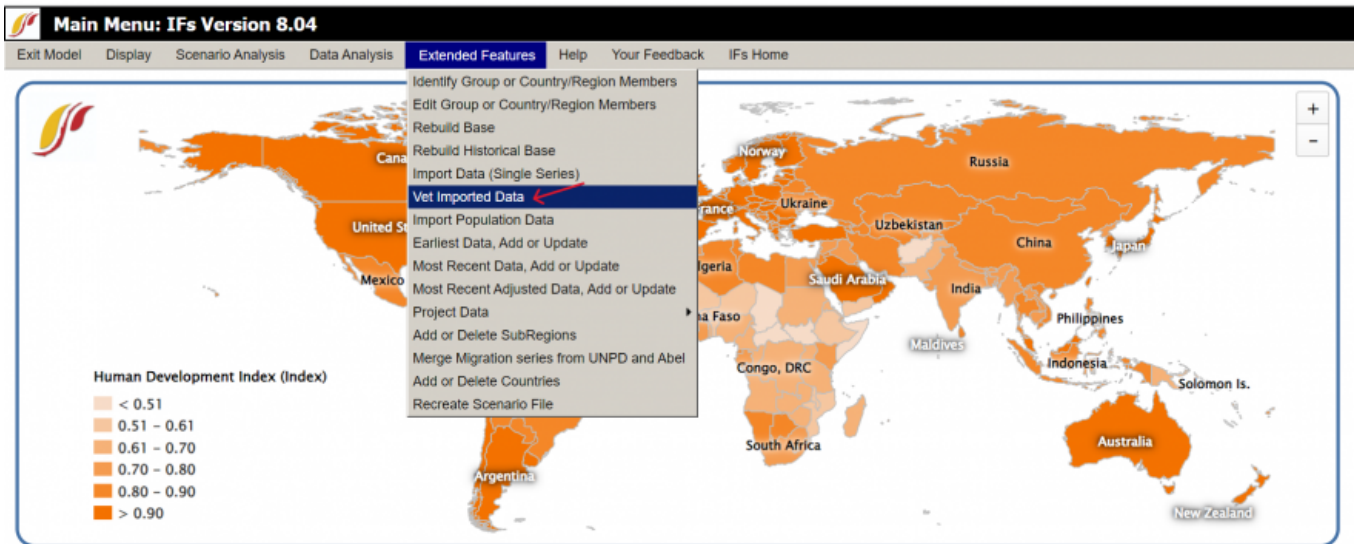
### **You can ensure data quality by following certain procedures:**

1. For an updated table, open the existing and updated table side by side. Compare data points for all years and all countries to make sure there are no, or very small, differences.
2. For all series, be sure to check large and important countries like USA, China or Germany.
3. Check countries with similar name (like the two Congos and two Koreas), as these can sometimes get mixed up.
4. Check for zeros and make sure they are actually zeros in the source data. We take no data as an empty cell. If we have zeros that need to be data and must be a feasible value (e.g., GDP cannot be 0).
5. Check the variable definition and make sure it makes sense and is in fact the right definition.
6. Make sure percentages are below 100 (there are cases when percentages can be above 100, e.g., gross enrollment rate).
7. Check values: GDP growth rate of more than 10% should raise a flag. For instances like this, check against the source data.
8. Creating line graphs for the countries in a series is a great way to quickly check for transients. This can easily be done in Excel or Tableau.

### **Using the IFs Vetting Tool**

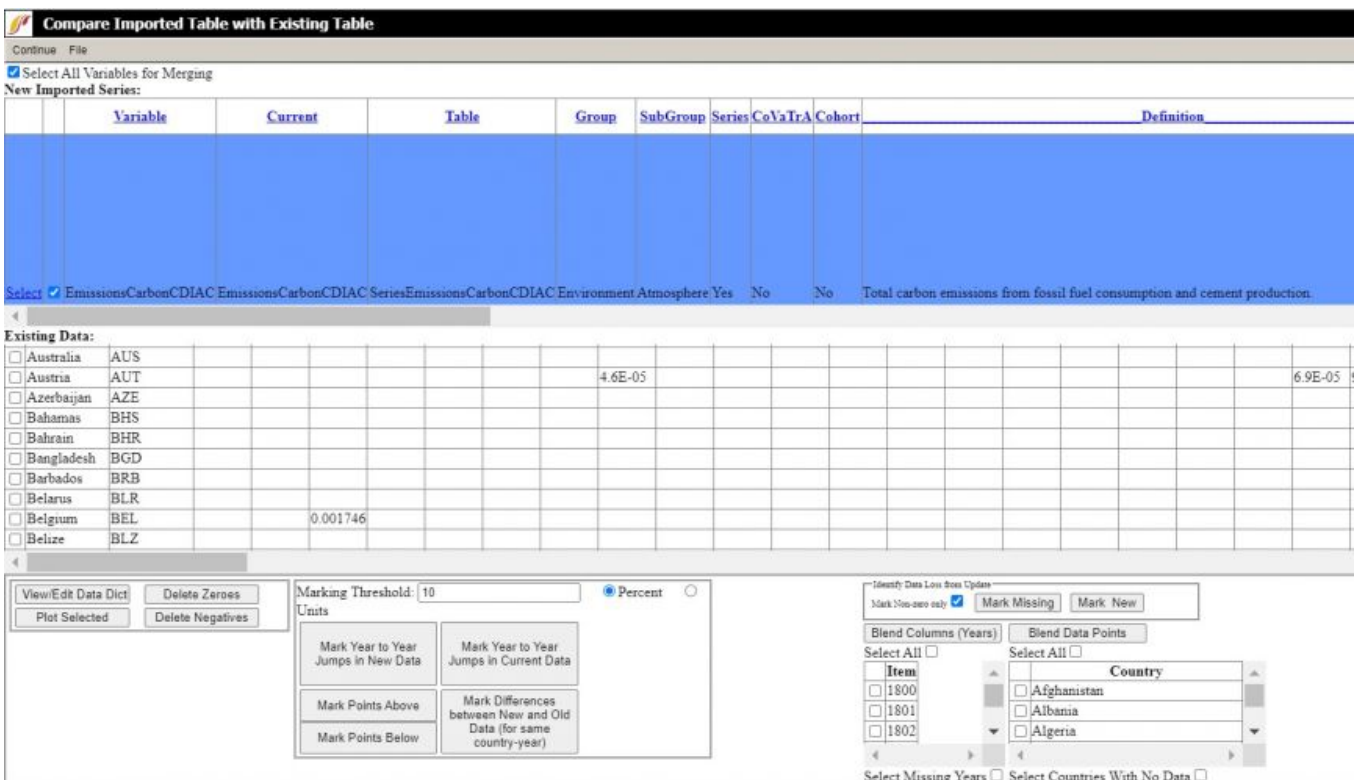
IFs has a feature to do some basic initial vetting and is a tool that should be used in every vetting process. The vetting tool can be found by the following:

**IFs -> MainMenu->Extended Features -> Vet Imported Data.**



IFs Main Page

The Compare Imported Table With Existing Table screen is then opened. Currently, the .NET version of IFs only compares tables that have been imported previously, as opposed to being able to open Access files. Make sure you import data into IFs before vetting. Select a series from "New Imported Series" you would like to vet.



Compare Imported Table with Existing Table

This will then populate the grids below with the new data (the data you are importing) and the old, historical data that is in IFs. The old data is in the first grid and the new data is in the second grid, as shown below.

Existing Data:

Country	FIPS_CODE	1800	1801	1802	1803	1804	1805	1806	1807	1808	1809	1810	1811	1812	1813	1814	1815	1816	1817	1818	1819	
<input type="checkbox"/> Afghanistan	AFG																					
<input type="checkbox"/> Albania	ALB																					
<input type="checkbox"/> Algeria	DZA																					
<input type="checkbox"/> Angola	AGO																					
<input type="checkbox"/> Argentina	ARG																					
<input type="checkbox"/> Armenia	ARM																					
<input type="checkbox"/> Australia	AUS																					
<input type="checkbox"/> Austria	AUT								4.6E-05												6.9E-05	
<input type="checkbox"/> Azerbaijan	AZE																					

Marking Threshold: 10  Percent  Units

Mark New only

Select All 

- 1800
- 1801
- 1802

Select All 
 Country
 

- Afghanistan
- Albania
- Algeria

Select Missing Years
  Select Countries With No Data

New Data Series:

	Country	FIPS_CODE	1949	1950	1951	1952	1953	1954	1955	1956	1957	1958	1959	1960	1961	1962	1963	1964	1965	1966	1967	1968	1969	1970	1971	1972
Edit <input type="checkbox"/>	Afghanistan	AFG	4	23	25	25	29	29	42	50	80	90	105	113	134	188	193	229	275	298	350	334	257	456	517	418
Edit <input type="checkbox"/>	Albania	ALB	277	81	110	102	113	137	181	229	411	328	394	552	622	672	568	550	593	696	731	838	885	1021	1187	153
Edit <input type="checkbox"/>	Algeria	DZA	248	1033	1128	1062	1092	1134	1258	1364	1512	1424	1545	1680	1654	1546	1480	1541	1799	2299	2302	2471	3075	4111	5091	773
Edit <input type="checkbox"/>	Angola	AGO		51	68	85	75	95	113	137	169	162	169	150	124	322	314	334	324	424	271	456	760	977	930	122
Edit <input type="checkbox"/>	Argentina	ARG	4195	8168	9544	9853	9591	10032	10799	12092	12996	12064	13362	13312	13957	14643	13658	15197	16053	17218	17874	18839	21088	22562	24254	245
Edit <input type="checkbox"/>	Armenia	ARM																								
Edit <input type="checkbox"/>	Australia	AUS	10276	14941	16112	16432	16223	18517	19291	19934	20340	21184	22849	24053	24704	25883	27551	29719	32988	32815	35251	36712	38794	40256	41662	429
Edit <input type="checkbox"/>	Austria	AUT	7254	5704	6351	6025	5928	6690	7936	7700	7966	7668	7620	8405	8689	9246	10088	10620	10414	10706	10899	11549	12188	13824	14223	153
Edit <input type="checkbox"/>	Azerbaijan	AZE																								

Vetting Data Page

Here you have a lot of different options to compare the new data with the old. The vetting tool will automatically mark any zeros in the table. You will want to double check to see if, in fact, these are suppose to be zeros or if they are just null values. If they are just null values, you can click the button, **Delete Zeros**.

Next, you need to click the button, **Mark Differences between New and Old Data for same country-year**. This will allow you to easily observe any differences, as the vetting tool will highlight them. You can also change the threshold that it will mark. The default is that it will mark any more than 10% difference. As you can see in the screenshot below, the values for Afghanistan in 2014 and 2015 are marked. You will want to go through all of the countries in each series you are importing to look at these marked differences and write down any that are very large.

You can also click the button **Mark Year to Year Jumps in New Data**, which is helpful to observe any large changes year-to-year.

Existing Data:

Country	FIPS_CODE	1800	1801	1802	1803	1804	1805	1806	1807	1808	1809	1810	1811	1812	1813	1814	1815	1816	1817	1818	1819	
<input type="checkbox"/> Afghanistan	AFG																					
<input type="checkbox"/> Albania	ALB																					
<input type="checkbox"/> Algeria	DZA																					
<input type="checkbox"/> Angola	AGO																					
<input type="checkbox"/> Argentina	ARG																					
<input type="checkbox"/> Armenia	ARM																					
<input type="checkbox"/> Australia	AUS																					
<input type="checkbox"/> Austria	AUT									4.6E-05												6.9E-05
<input type="checkbox"/> Azerbaijan	AZE																					

Marking Threshold: 10  Percent  Units

Mark Year to Year Jumps in New Data

Mark Year to Year Jumps in Current Data

Mark Points Above

Mark Differences between New and Old Data (for same country-year)

Blend Columns (Years)

Select All

Select All

Item  
 1800  
 1801  
 1802

Country  
 Afghanistan  
 Albania  
 Algeria

Select Missing Years

Select Countries With No Data

New Data Series:

Country	FIPS_CODE	1949	1950	1951	1952	1953	1954	1955	1956	1957	1958	1959	1960	1961	1962	1963	1964	1965	1966	1967	1968	1969	1970	1971	1972	
<input type="checkbox"/> Afghanistan	AFG	4	23	25	25	29	29	42	50	80	90	105	113	134	188	193	229	275	298	350	334	257	456	517	418	
<input type="checkbox"/> Albania	ALB	277	81	110	102	113	137	181	229	411	328	394	552	622	672	568	550	593	696	731	838	885	1021	1187	153	
<input type="checkbox"/> Algeria	DZA	248	1033	1128	1062	1092	1134	1258	1364	1512	1424	1545	1680	1654	1546	1480	1541	1799	2299	2302	2471	3075	4111	5091	773	
<input type="checkbox"/> Angola	AGO		51	68	85	75	95	113	137	169	162	169	150	124	322	314	334	324	424	271	456	760	977	930	122	
<input type="checkbox"/> Argentina	ARG	4195	8168	9544	9853	9591	10032	10799	12092	12996	12064	13362	13312	13957	14643	13658	15197	16053	17218	17874	18839	21088	22562	24254	245	
<input type="checkbox"/> Armenia	ARM																									
<input type="checkbox"/> Australia	AUS	10276	14941	16112	16432	16223	18517	19291	19934	20340	21184	22849	24053	24704	25883	27551	29719	32988	32815	35251	36712	38794	40256	41662	429	
<input type="checkbox"/> Austria	AUT	7254	5704	6351	6025	5928	6690	7936	7700	7966	7668	7620	8405	8689	9246	10088	10620	10414	10706	10899	11549	12188	13824	14223	153	
<input type="checkbox"/> Azerbaijan	AZE																									

Vetting Methods

## Blending

This vetting tool is also used to merge country and year columns between updated and existing Access tables. This is a necessity when blending to preserve the existing historical IFs values when there is no corresponding value in the new data series. For example, often times there will be an update from a source but it is missing multiple countries we already have data for. In this case, you can use the blending tool to blend in any historical country values that are not in the new data. The same can be done for missing years in the new data.

Once you have thoroughly vetted a data series, you will want to blend the new and historical series together. This ensures that we do not lose any values from the historical series. To do this for the Country values, you will want to check the **Select All** box in the bottom right hand corner of the screen. This will automatically check all of the countries, as seen below. Once you have check this box, press the **Blend Data Points** button, and click **Yes**. This will update the new Access file with missing Country observations.

You will also want to do the same process for the **Blend Columns (Years)** button, if there are any missing years. This ensures that we retain the historical values.

