

Removing Duplicate Values in Excel

Problem Statement:

Student data downloaded from Red Raider Registry can contain duplicate values. These duplicate Unique Identifiers can result from searching courses over time or result from multiple course attempts by a student. This tutorial will provide steps on how to identify if the data set contains duplicate values and how to remove those duplicate values from the data set.

This example shows how to review and identify duplicate unique identifiers (students) in downloaded course data over the last ten (10) years.

Step 1- Download Dataset from Red Raider Registry

A. Select the fields required for the dataset

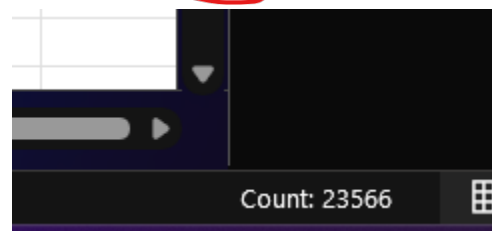
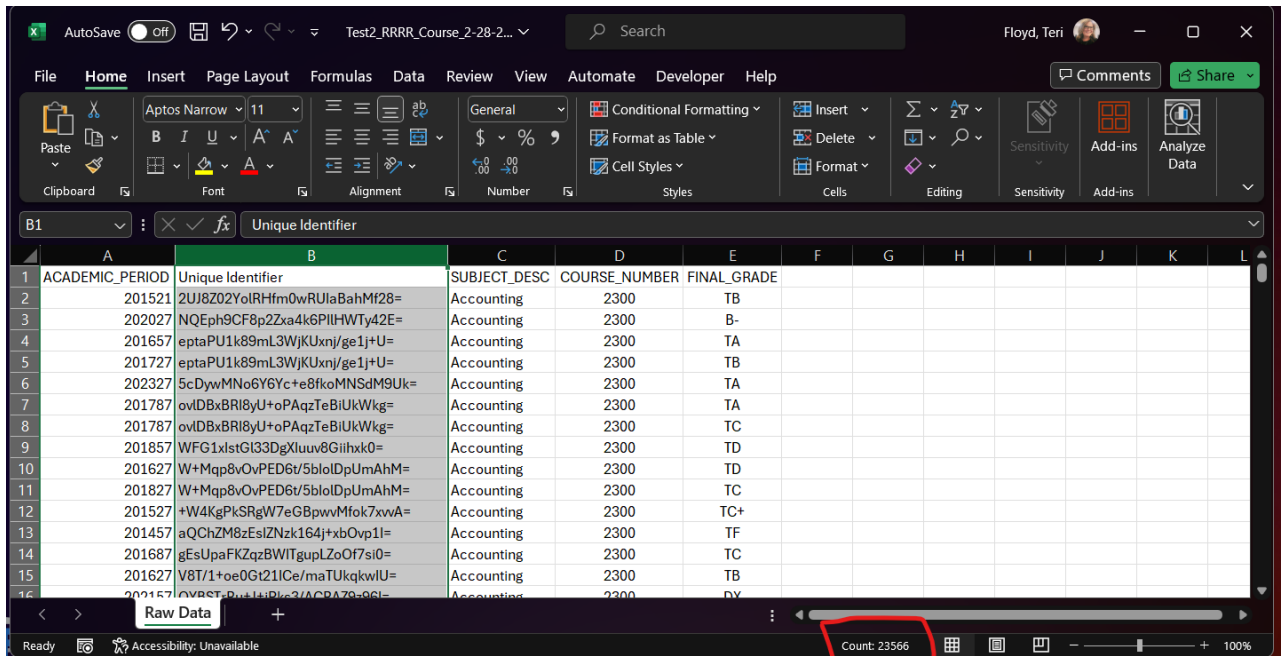
The screenshot displays the search interface for the Red Raider Registry. On the left, there are five filter categories: 'Course' (highlighted in blue), 'Demographics', 'Program', 'Admission Test Attributes', and 'Student Academic Attributes'. The 'Course' filter is selected, leading to a search form on the right. The form is titled 'Course' and includes the following fields: 'Subject*' (a dropdown menu with 'Accounting (ACCT)' selected), 'Course Number (Optional)' (a text input field with '2300'), 'From' (a dropdown menu with 'Spring 2014 TTU' selected), and 'To' (a dropdown menu with 'Spring 2024 TTU' selected). Below these fields, there is a section titled 'Select Fields for your dataset' with several buttons: 'College', 'Department', 'Subject', 'Course Number', 'Final Grade', and 'Mid-term Grade'. The 'Subject', 'Course Number', and 'Final Grade' buttons are highlighted in red. At the bottom of this section, there are two buttons: 'Select All' and 'Deselect All'.

B. Click on Create Dataset

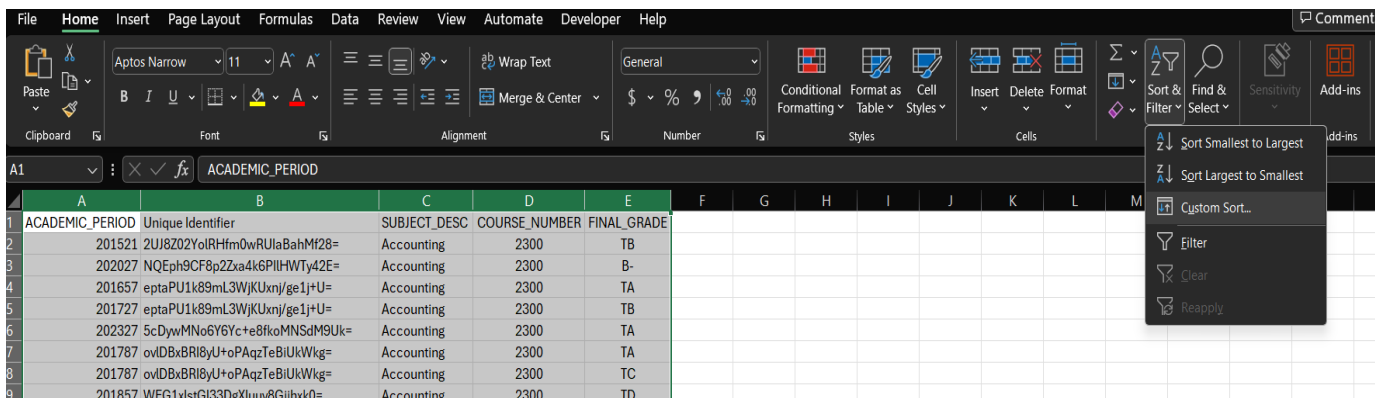
This screenshot shows a grey rectangular box containing text and a button. At the top, it says 'Explanation of Grades and their Meanings' in red. Below that, it says 'Once your dataset is created, a CSV file will be downloaded to your device.' in black. Further down, it says 'Combining and importing files for analysis' in red. At the bottom, there is a black button with the text 'Create Dataset' in white.

Step 2- Review Dataset and Sort

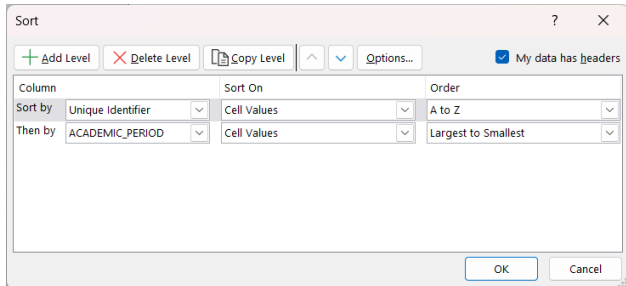
- A. Review the row count. Click on column B to see how many rows were counted. In this example, there is a count of 23,566 rows. Note that this includes the column names in row one, so there are 23,565 rows of student data.



- B. To identify students that have retaken a course it is important to sort the data. Select all columns in the worksheet and then select Home > Editing > Sort & Filter > Custom Sort.



- C. Sort by unique identifier (A to Z) and academic period in descending order (largest to smallest). The reason for sorting first by alphabetical order then by academic period in descending order is to identify the last TTU term that a student attempted the class. That way once duplicate values are removed, the data shows the most recent student information.



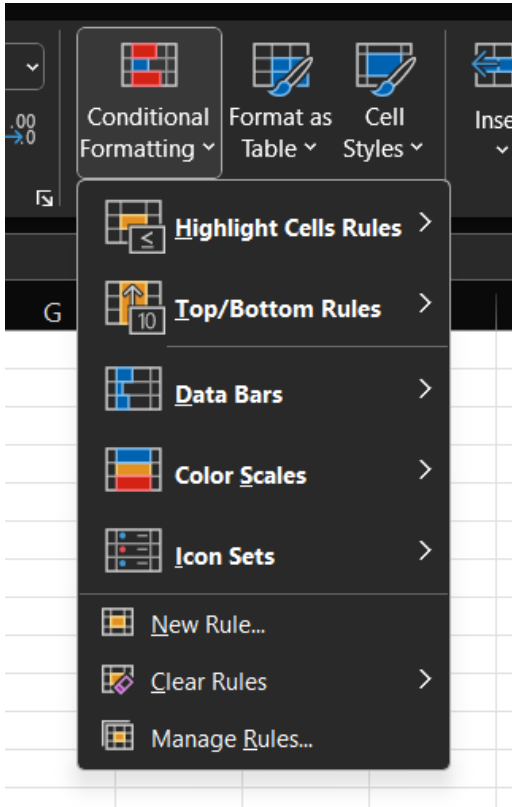
- D. The data is now displayed in the proper order to identify duplicate values (repeat students) and when they took the course (most recent shows first).

Step 3- Identify Duplicates in the Data Set

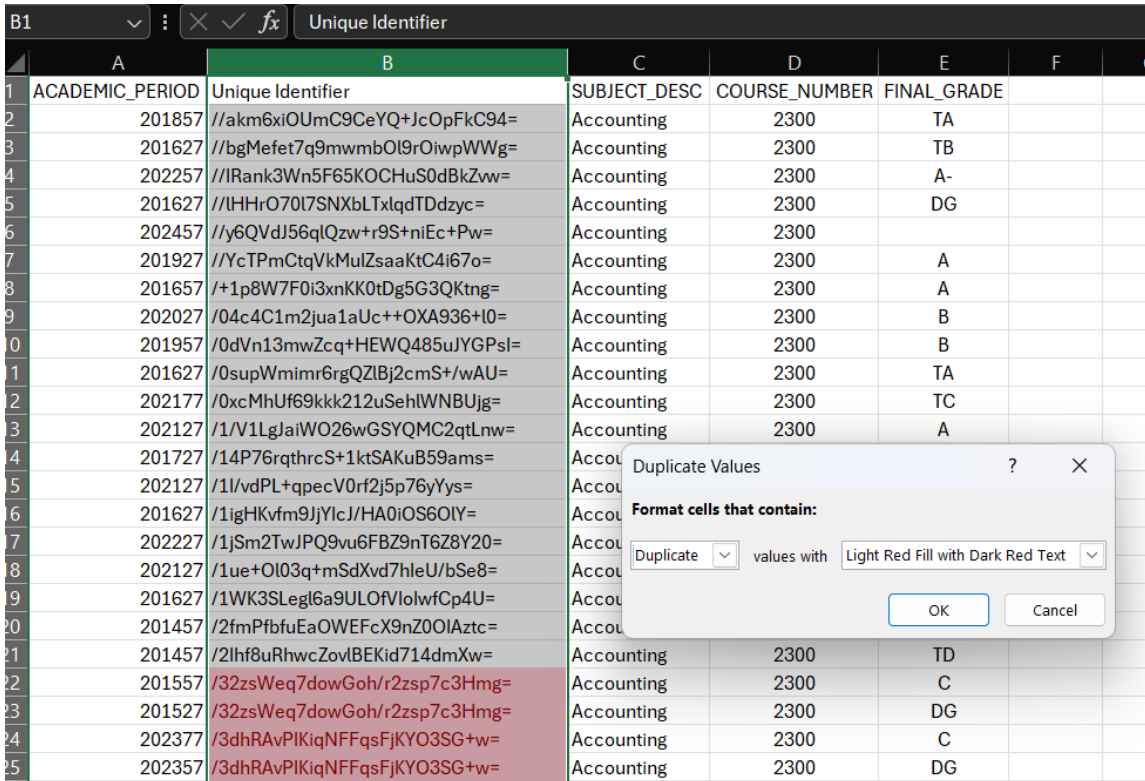
- A. Select the cells you want to check for duplicates. In this example, we are selecting the Unique Identification column to check for duplicate values.

ACADEMIC_PERIOD	Unique Identifier	SUBJECT_DESC	COURSE_NUMBER	FINAL_GRADE
201521	2UJ8Z02YolRHfm0wRUlaBahMf28=	Accounting	2300	TB
202027	NQEph9CF8p2Zxa4k6PIlHWTy42E=	Accounting	2300	B-
201657	eptaPU1k89mL3WjKUxnj/ge1j+U=	Accounting	2300	TA
201727	eptaPU1k89mL3WjKUxnj/ge1j+U=	Accounting	2300	TB
202327	5cDywMNo6Y6Yc+e8fkoMNSdM9Uk=	Accounting	2300	TA
201787	ovlDBxBRI8yU+oPAqzTeBiUkWkg=	Accounting	2300	TA
201787	ovlDBxBRI8yU+oPAqzTeBiUkWkg=	Accounting	2300	TC
201857	WFG1xlstGI33DgXluuv8Giihvk0=	Accounting	2300	TD
201627	W+Mqp8vOvPED6t/5blolDpUmAhM=	Accounting	2300	TD
201827	W+Mqp8vOvPED6t/5blolDpUmAhM=	Accounting	2300	TC
201527	+W4KgPkSRgW7eGBpwwMfok7xvA=	Accounting	2300	TC+
201457	aQChZM8zEslZNzk164j+xbOvp1l=	Accounting	2300	TF
201687	gEsUpaFKZqzBWITgupLZoOf7si0=	Accounting	2300	TC
201627	V8T/1+oeOGt21lCe/maTUKqkwIU=	Accounting	2300	TB
202157	QYBSTRu+J+jPks3/ACRAZ9z96l=	Accounting	2300	DX

B. Select Home > Conditional Formatting > Highlight Cells Rules > Duplicate Values



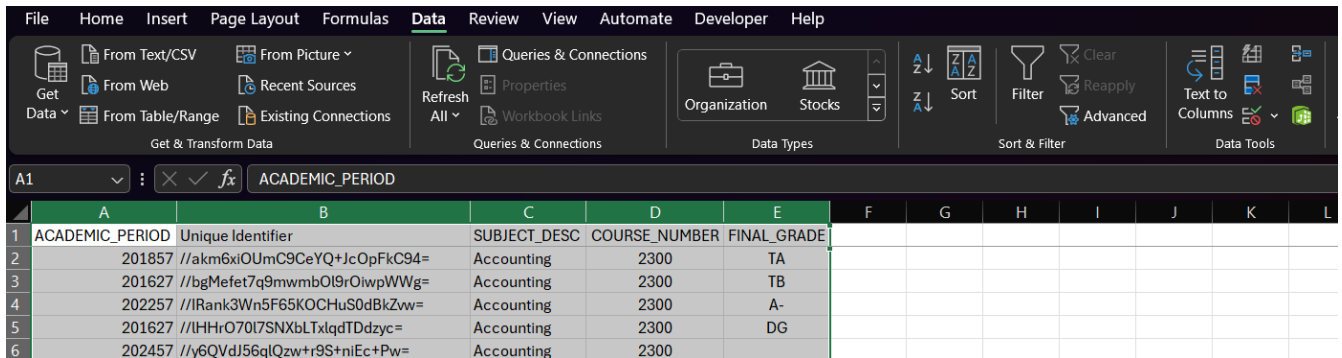
C. In the box next to values with, pick the formatting you want to apply to the duplicate values, and then select OK.



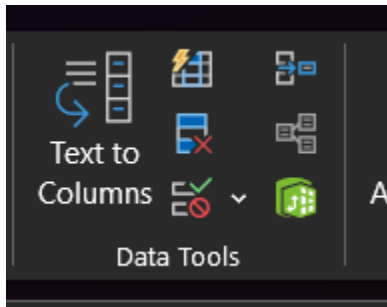
Step 4 – Remove the Duplicate Values

Please note that when you use Remove Duplicates, the duplicate data is permanently deleted. It is good practice to copy the original data to another worksheet before deleting duplicates, so you don't accidentally lose information.

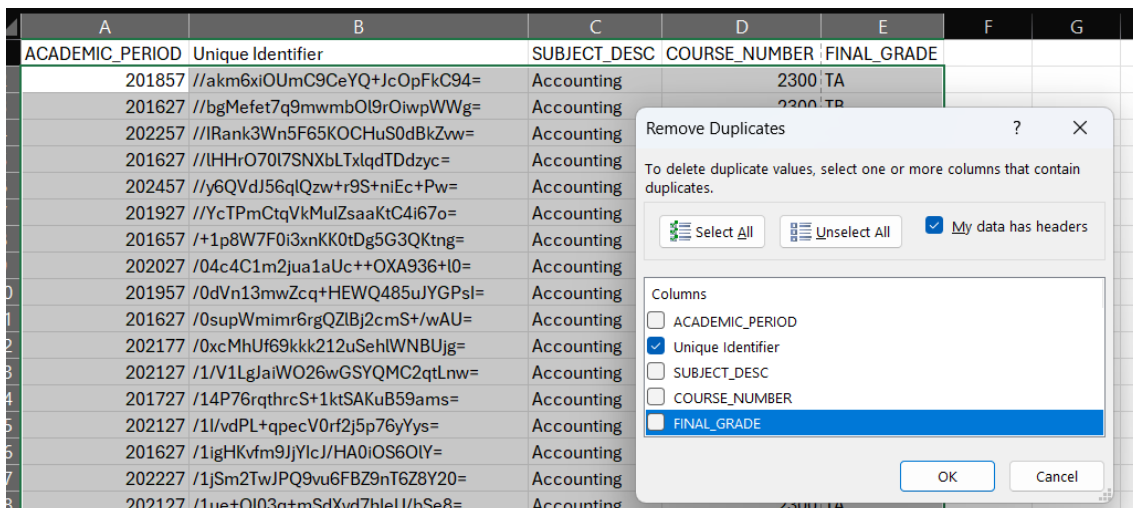
- A. Select the range of cells that has duplicate values you want to remove.



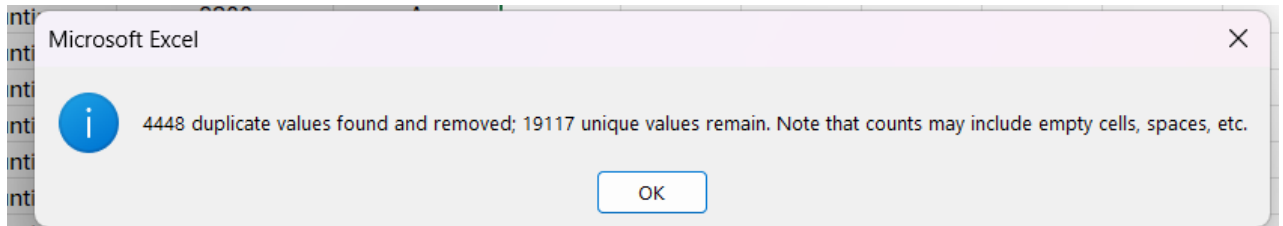
- B. Select Data > Data Tools > Remove Duplicates



- C. Click on the Remove Duplicates icon and select the columns that contain duplicates. In this test, we wanted to remove duplicate Unique Identifiers so that was the only column checked. Click OK.



- D. Excel will remove all duplicates that occur after the first occurrence and a notification will pop up indicating how many duplicates were found / removed and how many unique values remain. Click OK.



- E. Now the data has been deduplicated, it is another good practice to double check row counts to ensure that there was no additional data lost.

23,565 rows of raw data (contains duplicates)
- 4,448 duplicate values

19,117 unique values