

# Learn How to Calculate Mahalanobis Distance Using SPSS

Authored by  
**Mohammed looti**

November 8, 2025

## RECOMMENDED CITATION

Mohammed looti (2025). *Learn How to Calculate Mahalanobis Distance Using SPSS*. PSYCHOLOGICAL STATISTICS. Retrieved from <https://statistics.arabpsychology.com/?p=12760>

The [Mahalanobis distance](#) is recognized as an exceptionally powerful metric within the realm of **statistical analysis**. Unlike the simple measurement provided by standard Euclidean distance, this measure fundamentally quantifies the separation between a specific observation (a point) and the center of a data cluster (the mean of a distribution), crucially adjusting for the inherent [correlation](#) and **variance structure** of the dataset. By utilizing the [covariance matrix](#) to appropriately weigh the influence of each variable, the Mahalanobis measure becomes an indispensable tool when operating within a [multivariate space](#).

This sophisticated distance measure holds immense value, particularly in the critical task of identifying [outliers](#)--those observations that significantly deviate from the central tendencies of the data cloud. In complex modeling environments, such as those employing regression or other **multivariate techniques**, an observation exhibiting an unusually large Mahalanobis distance can exert disproportionate or undue influence on the resulting model's parameters and conclusions. Therefore, calculating this specific metric is a foundational step in rigorous **data cleaning** and comprehensive model diagnostics before proceeding to inferential testing.

This comprehensive guide is designed to detail the precise methodology required for calculating the [Mahalanobis distance](#) using [SPSS](#) (Statistical Package for the Social Sciences). We will systematically navigate the steps necessary to generate this metric for every single observation in your dataset and subsequently determine if any of these distances are statistically significant. This process provides empirical confirmation of the presence and identity of **multivariate outliers** within your sample.

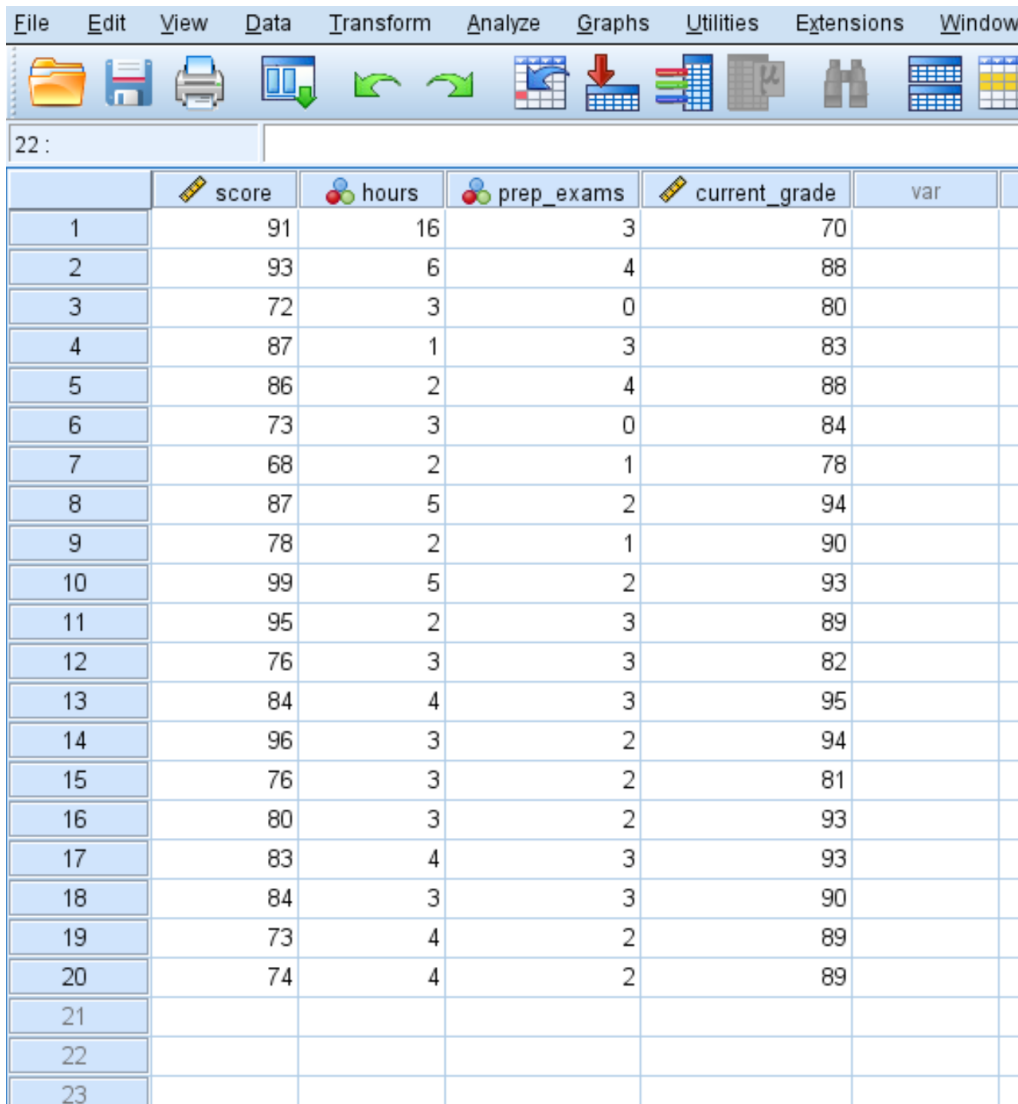
## Preparing the Dataset for Multivariate Analysis

To effectively illustrate the precise calculation of the Mahalanobis distance within a practical context, we will utilize a simulated sample dataset focused on student performance metrics. This dataset incorporates four core variables: the student's exam **score** (designated as our dependent variable), alongside three powerful predictor variables: the **number of hours they spent studying**, the **number of preparatory exams they undertook**, and their **current academic grade in the course**. This structure allows us to examine the combined influence of predictors.

Our sample comprises 20 distinct observations, and the primary analytical objective is to determine if the unique combination of predictor variables (hours, prep exams, and grade) for any individual student is statistically unusual when compared against the overall group mean and covariance structure. Identifying such highly unusual or extreme observations is paramount before implementing standard inferential procedures, such as **multiple regression analysis**, which are highly sensitive to these extreme points.

The structure of our example dataset, exactly as it is loaded into the [SPSS](#) Data View, is presented in the image below. We will use this established structure as the basis for the subsequent steps,

which involve calculating the Mahalanobis distance for each case. This systematic approach enables us to accurately and efficiently identify potential **multivariate outliers** within the sample, ensuring the integrity of downstream statistical models.



	score	hours	prep_exams	current_grade	var
1	91	16	3	70	
2	93	6	4	88	
3	72	3	0	80	
4	87	1	3	83	
5	86	2	4	88	
6	73	3	0	84	
7	68	2	1	78	
8	87	5	2	94	
9	78	2	1	90	
10	99	5	2	93	
11	95	2	3	89	
12	76	3	3	82	
13	84	4	3	95	
14	96	3	2	94	
15	76	3	2	81	
16	80	3	2	93	
17	83	4	3	93	
18	84	3	3	90	
19	73	4	2	89	
20	74	4	2	89	
21					
22					
23					

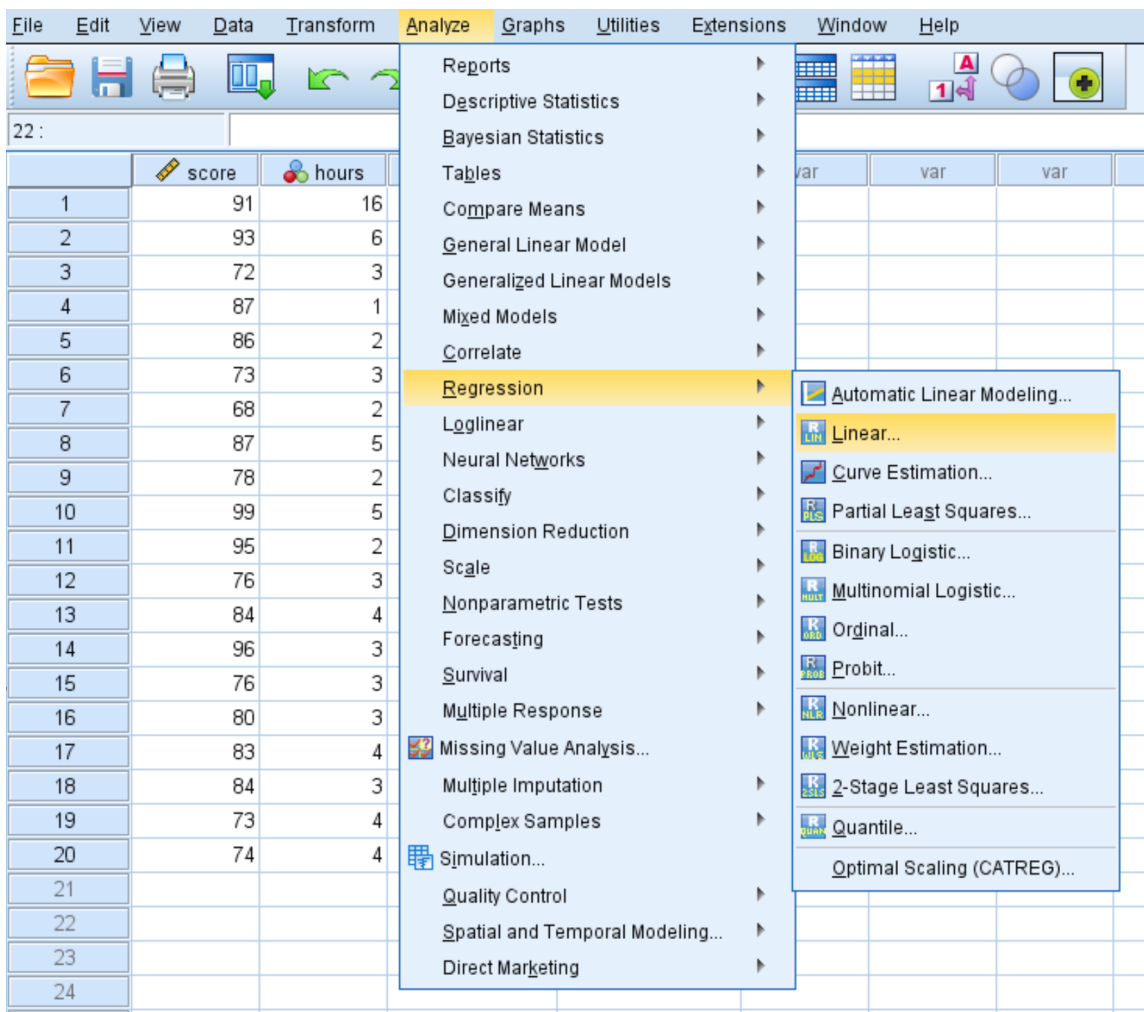
## Step 1: Initiating the Linear Regression Procedure

In the [SPSS](#) environment, the most efficient and standard methodology for calculating the Mahalanobis distance involves leveraging the comprehensive diagnostic capabilities embedded within the **linear regression** dialog box. While our immediate goal is not necessarily to fit a predictive regression model, the software uses the powerful underlying **matrix algebra** inherent to regression analysis to generate this essential distance metric quickly and accurately.

To commence the process, navigate directly to the main menu bar interface. Click on the **Analyze** tab, hover your cursor over the **Regression** option, and then select **Linear** from the cascading

submenu. Executing this sequence opens the primary [Linear Regression](#) dialog box. This specific interface is where we will define our variables and, critically, request the necessary diagnostic statistics that include the Mahalanobis distance.

Following these precise steps in the menu navigation is absolutely critical, as it grants access to the advanced diagnostic options required to successfully compute and save the distance metric. The visual guide below confirms the correct menu path.

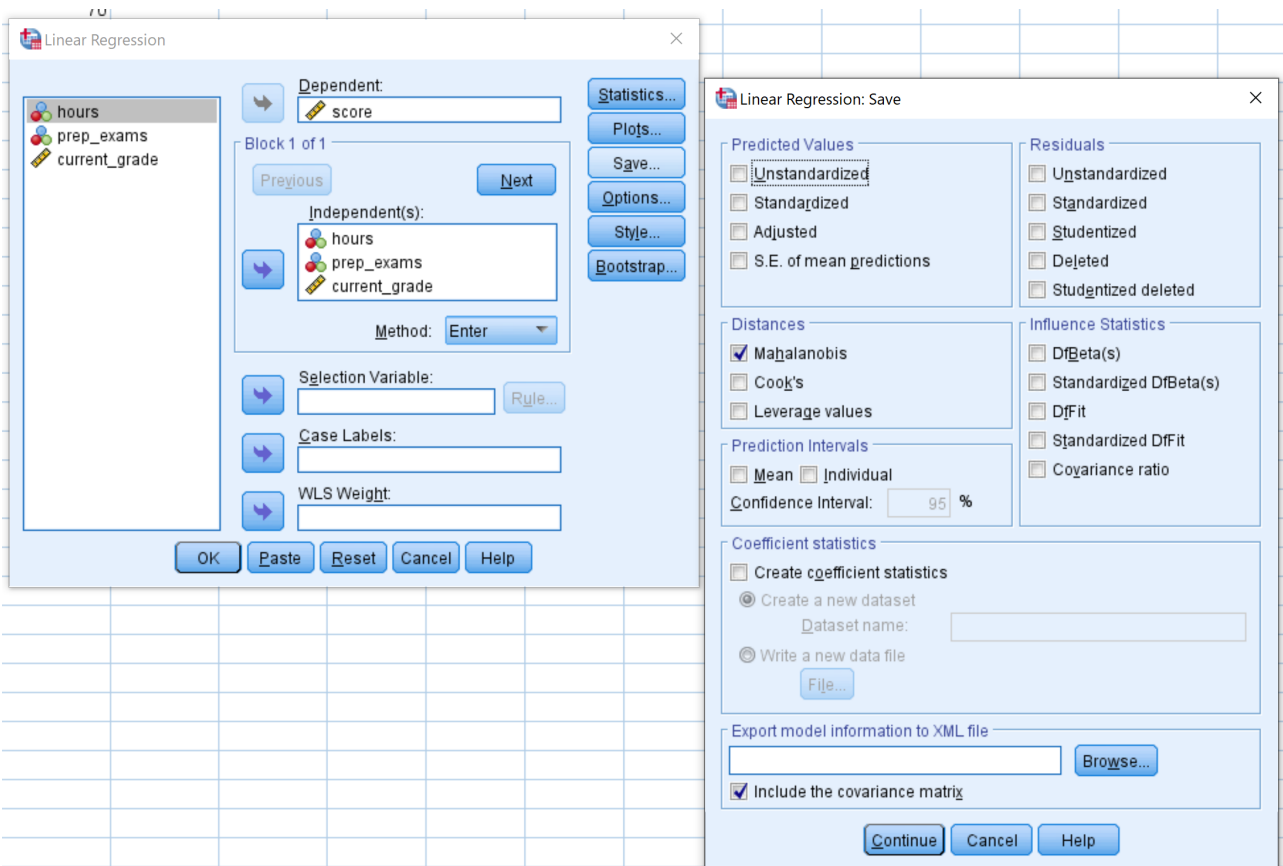


## Step 2: Generating the Mahalanobis Distance via Save Options

Once inside the Linear Regression dialog box, variables must be assigned to their correct roles. Begin by dragging the continuous outcome variable, *score*, into the designated box labeled **Dependent**. Subsequently, transfer all three of the identified predictor variables (hours studied, prep exams, and current grade) into the box labeled **Independent(s)**. Correct variable specification is vital as the calculation depends on the covariance structure of these independent variables.

After the variables have been precisely specified, the calculation of the [Mahalanobis distance](#) is initiated through the procedural save options. Locate and click the **Save** button situated on the right side of the dialog box. This action immediately launches the "Linear Regression: Save" secondary window, which provides a comprehensive list of various diagnostic and residual statistics that can be saved directly as new variables within your active dataset.

Within this Save window, locate the specialized "Distances" section. It is imperative that you place a checkmark next to the option labeled **Mahalanobis**. This crucial selection instructs the statistical software to compute the Mahalanobis distance for every single case, basing the calculation strictly on the observed covariance matrix of the previously specified predictor variables. After confirming this selection, click **Continue** to exit the Save window, and then click **OK** in the main Linear Regression dialog box to execute the analysis and generate the new variable.



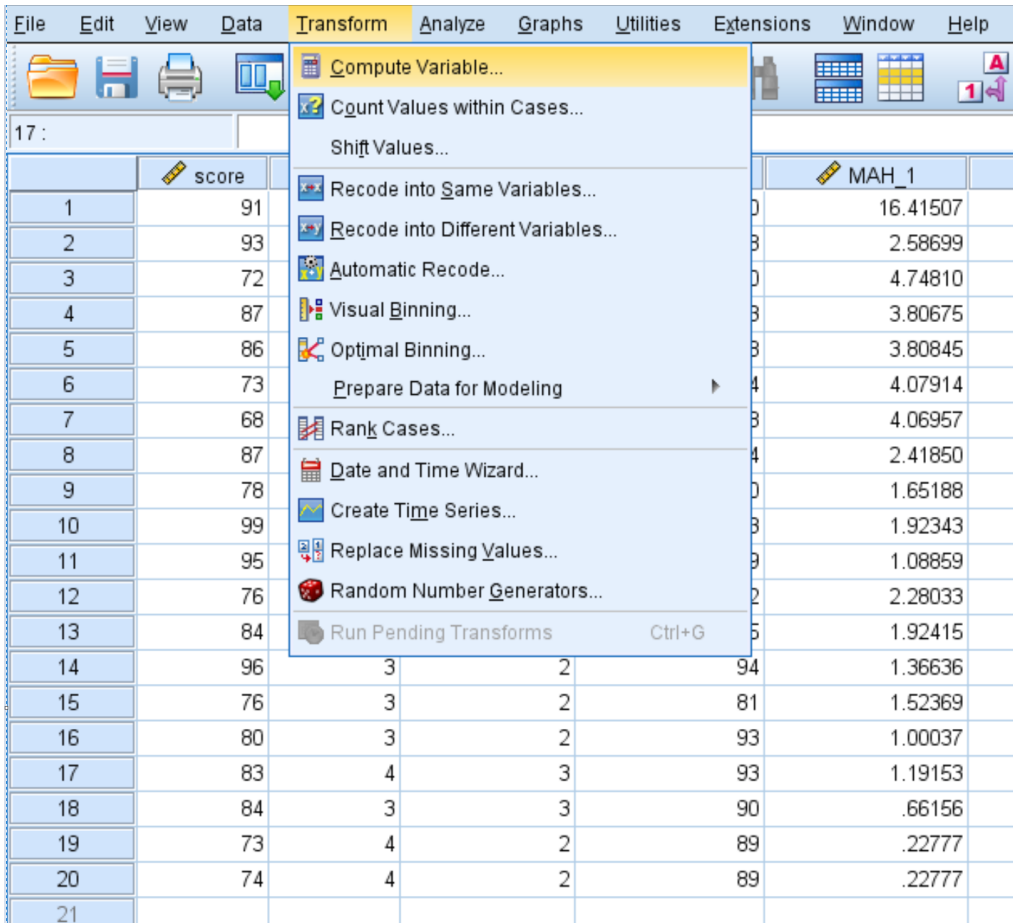
Following the execution of the procedure, a new column, which is typically and automatically named **MAH\_1**, will instantly appear in your Data View tab. This newly created column contains the raw, calculated Mahalanobis distance value for each observation in the dataset. A preliminary examination of these raw values serves as the essential first step in the identification process, as cases displaying notably larger distances signal potential [outliers](#) that require subsequent, more rigorous statistical scrutiny.

	score	hours	prep_exams	current_grade	MAH_1
1	91	16	3	70	16.41507
2	93	6	4	88	2.58699
3	72	3	0	80	4.74810
4	87	1	3	83	3.80675
5	86	2	4	88	3.80845
6	73	3	0	84	4.07914
7	68	2	1	78	4.06957
8	87	5	2	94	2.41850
9	78	2	1	90	1.65188
10	99	5	2	93	1.92343
11	95	2	3	89	1.08859
12	76	3	3	82	2.28033
13	84	4	3	95	1.92415
14	96	3	2	94	1.36636
15	76	3	2	81	1.52369
16	80	3	2	93	1.00037
17	83	4	3	93	1.19153
18	84	3	3	90	.66156
19	73	4	2	89	.22777
20	74	4	2	89	.22777

### Step 3: Assessing Outlier Significance using P-Values

While the raw Mahalanobis distance values provide an excellent measure of relative separation, they do not inherently offer the statistical context needed to declare an observation a significant outlier. To move from a descriptive measure to an inferential one, we must convert the calculated distance into a corresponding statistical [p-value](#). This conversion is possible because the Mahalanobis distance adheres to a known probability distribution, specifically the [Chi-Square distribution](#). This allows us to formally test the probability of observing a distance of that magnitude, or greater, purely by random chance.

To perform this crucial transformation, navigate to the **Transform** tab located in the main menu bar, and then select the **Compute Variable** function. This powerful utility enables the creation of a new variable (our desired p-value) based entirely on a specified mathematical expression or function applied to existing variables in the dataset. This step is mandatory for achieving statistical significance testing.



score	MAH_1
91	16.41507
93	2.58699
72	4.74810
87	3.80675
86	3.80845
73	4.07914
68	4.06957
87	2.41850
78	1.65188
99	1.92343
95	1.08859
76	2.28033
84	1.92415
96	1.36636
76	1.52369
80	1.00037
83	1.19153
84	.66156
73	.22777
74	.22777

Within the Compute Variable dialog box, proceed with the following two main steps:

In the designated **Target Variable** box, assign a clear and meaningful name for the new variable, such as "pvalue" or "Mahal\_P".

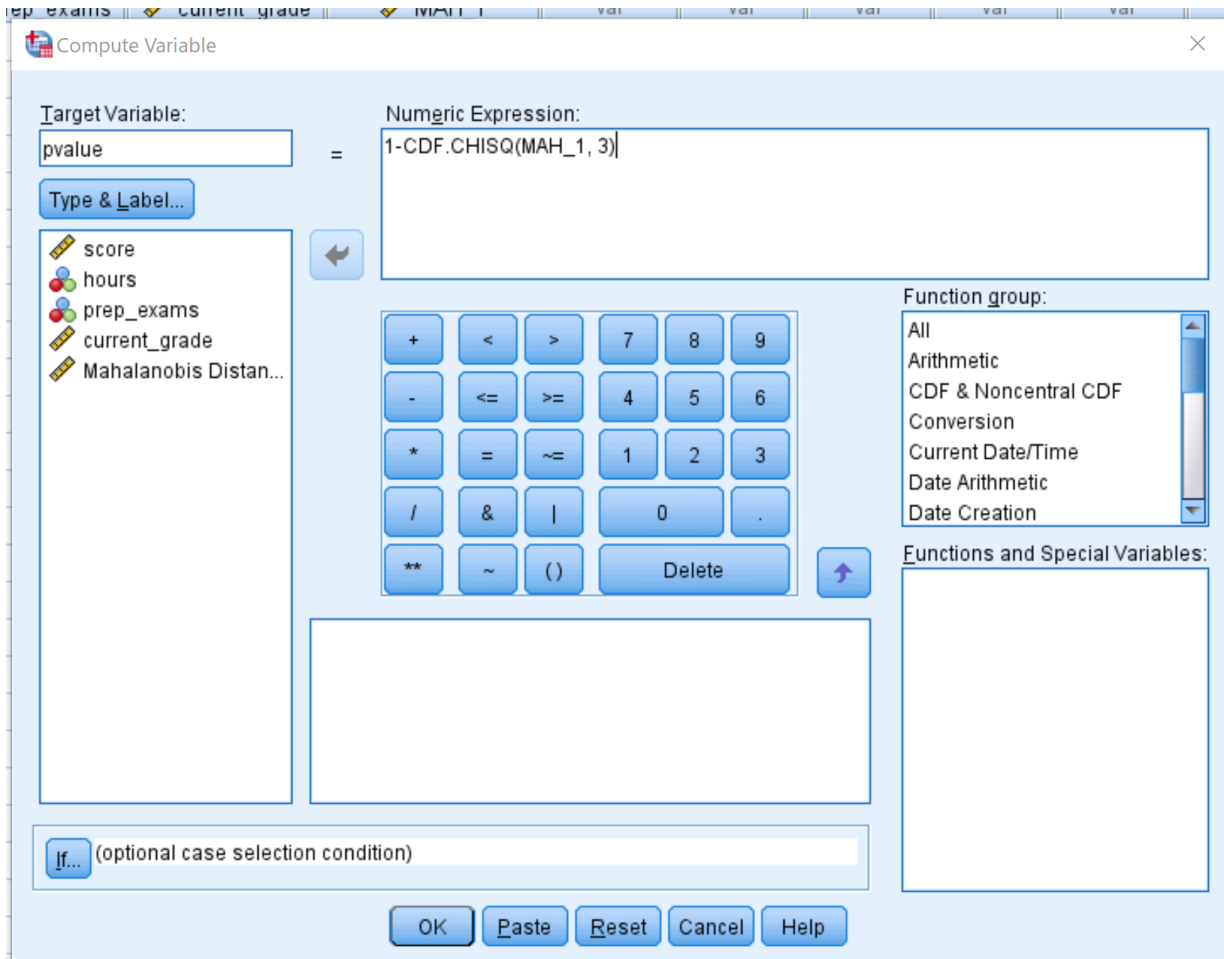
In the **Numeric Expression** box, accurately input the following formula, which utilizes the cumulative distribution function for the Chi-Square distribution:

1 - CDF.CHISQ(MAH\_1, 3)

The structure of the expression, `CDF.CHISQ(MAH_1, 3)`, calculates the cumulative probability of the Chi-Square distribution up to the specific distance value found in the `MAH_1` variable. Since the conventional [p-value](#) fundamentally represents the probability of observing a distance \*greater\* than the calculated Mahalanobis distance, we must subtract the resulting cumulative probability from 1.

The integer **3** within this function is a critical and fixed parameter, representing the [degrees of freedom](#) (df). For the Mahalanobis distance calculation, the degrees of freedom must invariably be equal to the total number of predictor variables (independent variables) included in the preceding

Linear Regression analysis. As our example utilized three independent variables (hours, prep exams, and grade), our degrees of freedom is correctly set to 3. Confirm all settings and click **OK** to execute the precise computation.



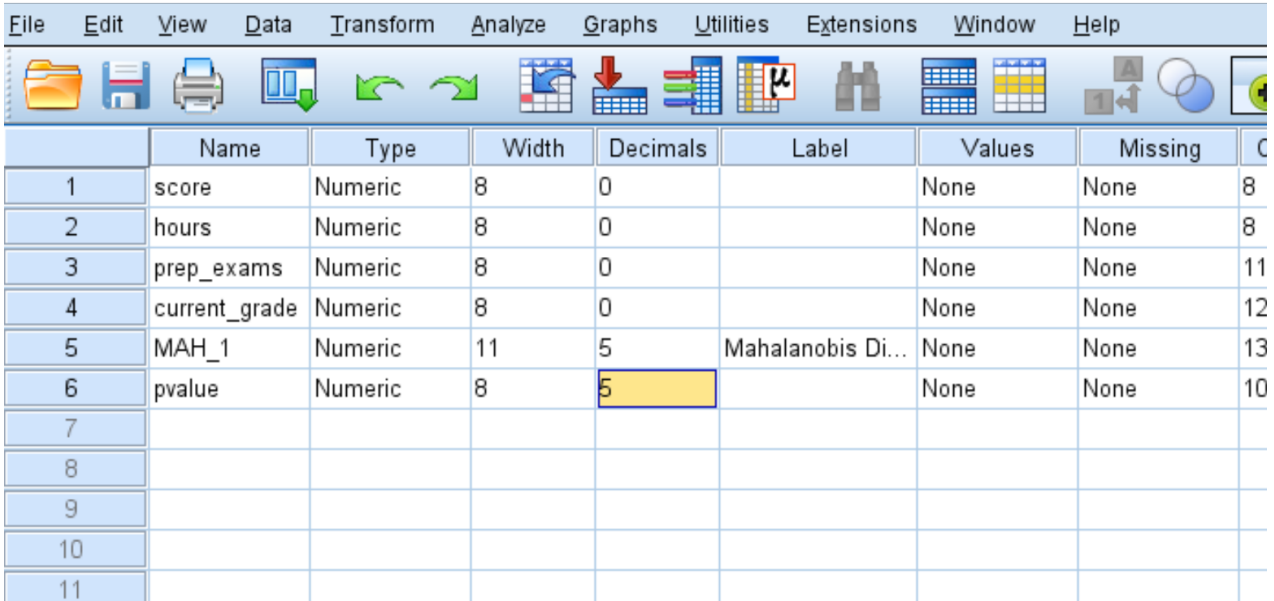
#### Step 4: Interpreting Results and Identifying Significant Outliers

Once the computation is successfully executed, a brand new column containing the calculated p-values will be seamlessly appended to your dataset in the Data View. These resulting p-values directly correspond to the statistical significance level of each respective Mahalanobis distance, effectively quantifying how likely that observation is under the assumption of a multivariate normal distribution.

It is important to note that by default, [SPSS](#) often displays numerical results limited to only two decimal places. This level of precision is typically insufficient for accurately identifying statistically significant [outliers](#), which frequently possess extremely small probability values requiring greater resolution.

score	hours	prep_exams	current_grade	MAH_1	pvalue
91	16	3	70	16.41507	.00
93	6	4	88	2.58699	.46
72	3	0	80	4.74810	.19
87	1	3	83	3.80675	.28
86	2	4	88	3.80845	.28
73	3	0	84	4.07914	.25
68	2	1	78	4.06957	.25
87	5	2	94	2.41850	.49
78	2	1	90	1.65188	.65
99	5	2	93	1.92343	.59
95	2	3	89	1.08859	.78
76	3	3	82	2.28033	.52
84	4	3	95	1.92415	.59
96	3	2	94	1.36636	.71
76	3	2	81	1.52369	.68
80	3	2	93	1.00037	.80
83	4	3	93	1.19153	.76
84	3	3	90	.66156	.88
73	4	2	89	.22777	.97
74	4	2	89	.22777	.97

To rectify this display issue and improve readability, navigate to the **Variable View** tab at the bottom of the SPSS window. Locate the row corresponding to your newly created p-value variable (e.g., "pvalue") and increase the value specified in the **Decimals** column to five or more digits. This simple adjustment allows for a far more granular and accurate reading of very small probability values, ensuring that potential significant outliers are not overlooked due to rounding.



	Name	Type	Width	Decimals	Label	Values	Missing	C
1	score	Numeric	8	0		None	None	8
2	hours	Numeric	8	0		None	None	8
3	prep_exams	Numeric	8	0		None	None	11
4	current_grade	Numeric	8	0		None	None	12
5	MAH_1	Numeric	11	5	Mahalanobis Di...	None	None	13
6	pvalue	Numeric	8	5		None	None	10
7								
8								
9								
10								
11								

Return to the **Data View** to inspect the refined, high-precision p-values. A widely accepted and statistically conservative threshold utilized for identifying a genuine multivariate outlier via the [Mahalanobis distance](#) is a p-value that is determined to be **less than .001**. Any observation whose probability falls below this stringent criterion is classified as highly unusual and statistically unlikely given the inherent distribution of the other variables in the dataset. In the example illustrated here, we clearly observe that the first case is the sole outlier, as its p-value is markedly smaller than the established .001 criterion.

	score	hours	prep_exams	current_grade	MAH_1	pvalue
1	91	16	3	70	16.41507	.00093
2	93	6	4	88	2.58699	.45977
3	72	3	0	80	4.74810	.19120
4	87	1	3	83	3.80675	.28310
5	86	2	4	88	3.80845	.28291
6	73	3	0	84	4.07914	.25304
7	68	2	1	78	4.06957	.25405
8	87	5	2	94	2.41850	.49020
9	78	2	1	90	1.65188	.64768
10	99	5	2	93	1.92343	.58845
11	95	2	3	89	1.08859	.77983
12	76	3	3	82	2.28033	.51630
13	84	4	3	95	1.92415	.58830
14	96	3	2	94	1.36636	.71344
15	76	3	2	81	1.52369	.67681
16	80	3	2	93	1.00037	.80116
17	83	4	3	93	1.19153	.75504
18	84	3	3	90	.66156	.88221
19	73	4	2	89	.22777	.97299
20	74	4	2	89	.22777	.97299
21						
22						
23						

## Best Practices for Managing Identified Outliers

Once the Mahalanobis distance calculation has definitively identified a statistically significant [outlier](#), the subsequent step requires the analyst to make an informed, methodological decision regarding the appropriate course of action. Handling outliers demands careful, ethical consideration; removing data can potentially compromise the generalizability of results, yet retaining influential outliers risks severely skewing model coefficients and overall statistical inferences.

The decision process typically boils down to two primary strategies:

**Scrutinize the Outlier for Data Entry Errors.** The most crucial initial step is to rigorously verify the accuracy and source of the data point. An extreme observation often arises simply due to a clerical mistake, such as a misplaced decimal, a transposed digit, or an error during the initial recording process. If an error is confirmed, the value must be corrected immediately. If the original source value cannot be definitively verified, the case should either be treated as **missing data** or

potentially removed, depending on the research protocol.

**Evaluate the Impact and Consider Removal.** If the extreme value is confirmed to be an accurate, true data point (a genuine extreme case), the analyst must then assess its specific influence on the planned statistical analysis. If the outlier substantially alters key regression coefficients, inflates standard errors, or significantly shifts the overall model conclusions, its removal may be statistically justifiable. However, this decision must be handled transparently. Should removal be chosen, it is mandatory to clearly document the reasoning and, ideally, report the findings for the analysis both including and excluding the outlier in the final scholarly report or publication to ensure maximum methodological rigor and transparency.