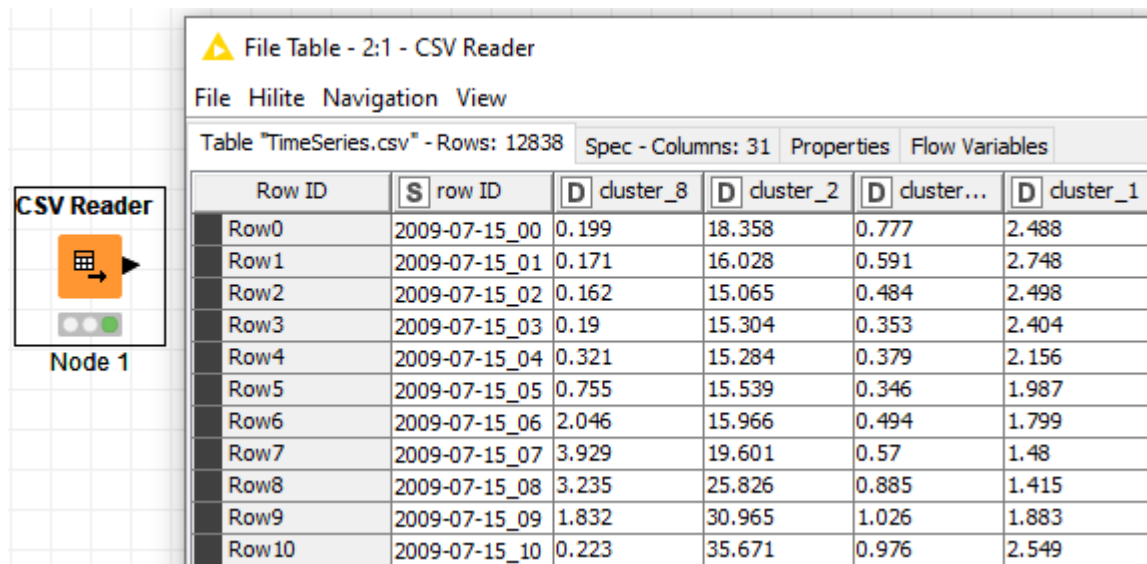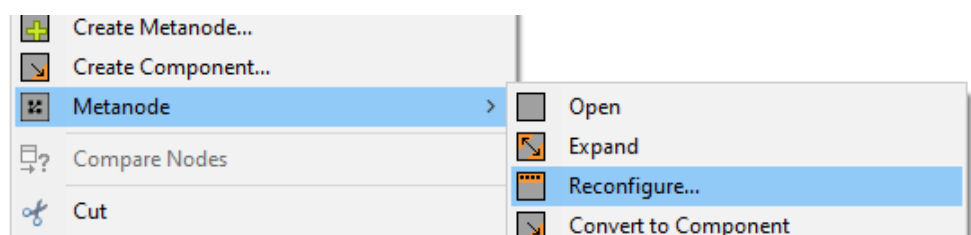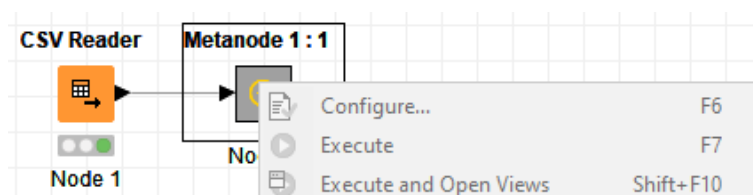# Regression on Time Series Data

In this demonstration, we will apply linear regression on time series data.

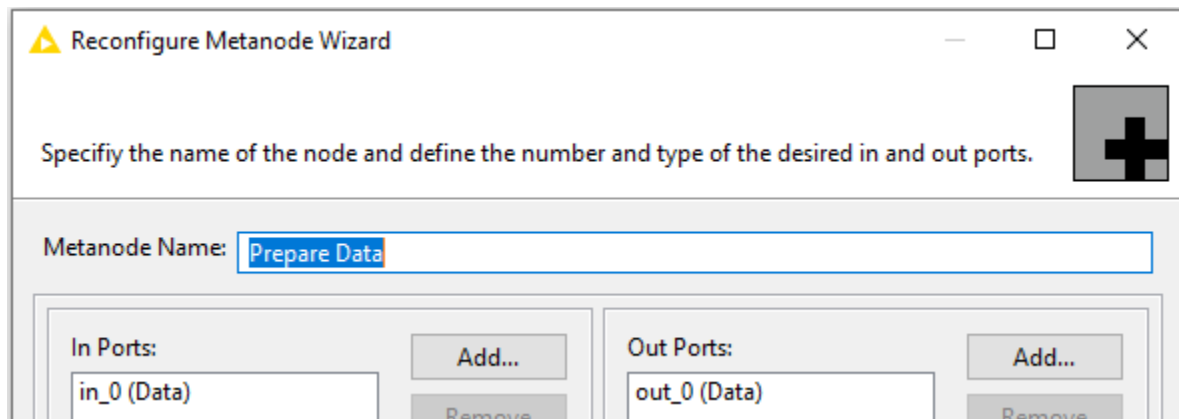1. Firstly, let us read in the data. Use a *CSV Reader* node to read in the datasets/TimeSeries.csv file.
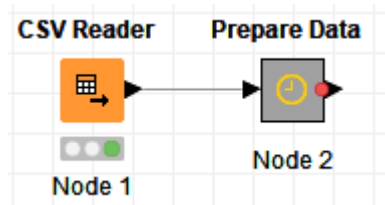


2. To make workflows look cleaner and to consolidate related nodes together for easier reusability, metanodes are sometimes used. Create a new 1:1 metanode by clicking on this symbol in the top menu bar and reconfigure it:
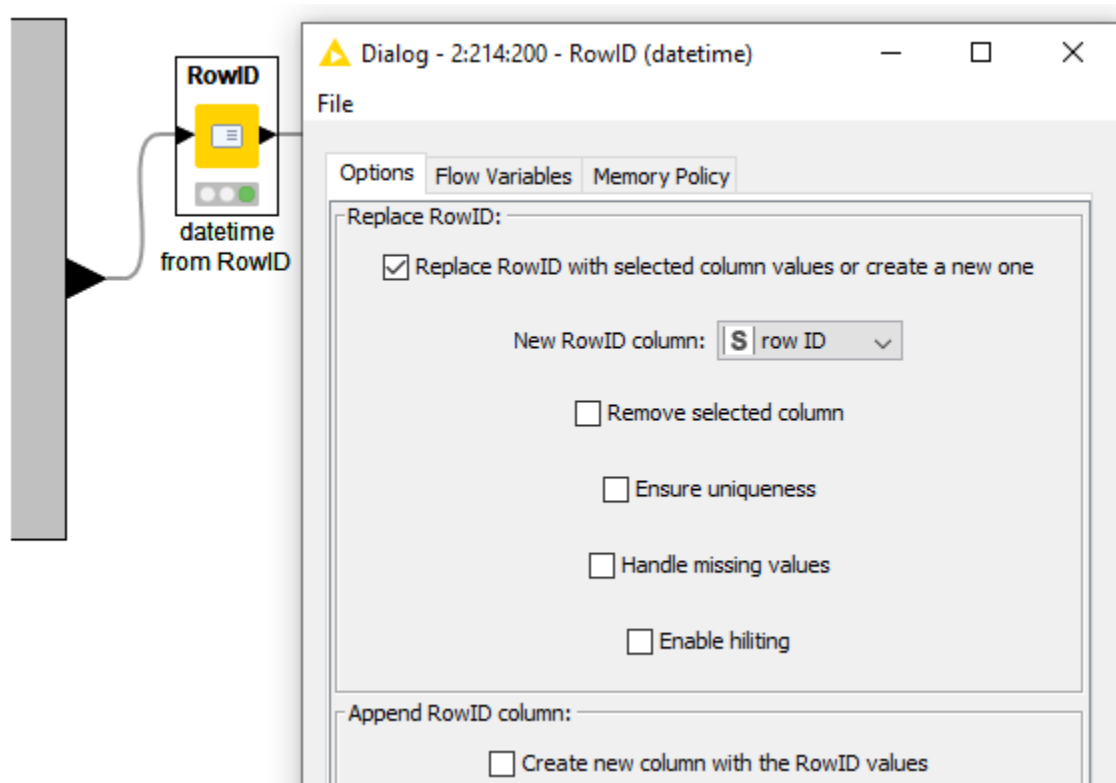
3. Rename the metanode to "Prepare Data":



Renamed node:



4. Double-click into the node and create a RowID node to make a row id using the datetime "row ID" column:

5. You can see that the datetime column is copied as a row id:



6. We will next change the datetime column (not the newly created row id) into an actual datetime variable.





Scroll to the end of the dataset to see the datetime variable:

7.  Sort the data into an ascending order using the *Sorter* node:



8.  We shall now arbitrarily choose one of the cluster i.e. cluster_26, as our data of interest:



9.  Just to make life easier, we shall rename "cluster_26" to simply "cluster":

10. Let's go back to the main workflow (by clicking the correct tab) and introduce a Partitioning node to split the data. We are going to learn from 90% of the past data to predict the latest 10%:
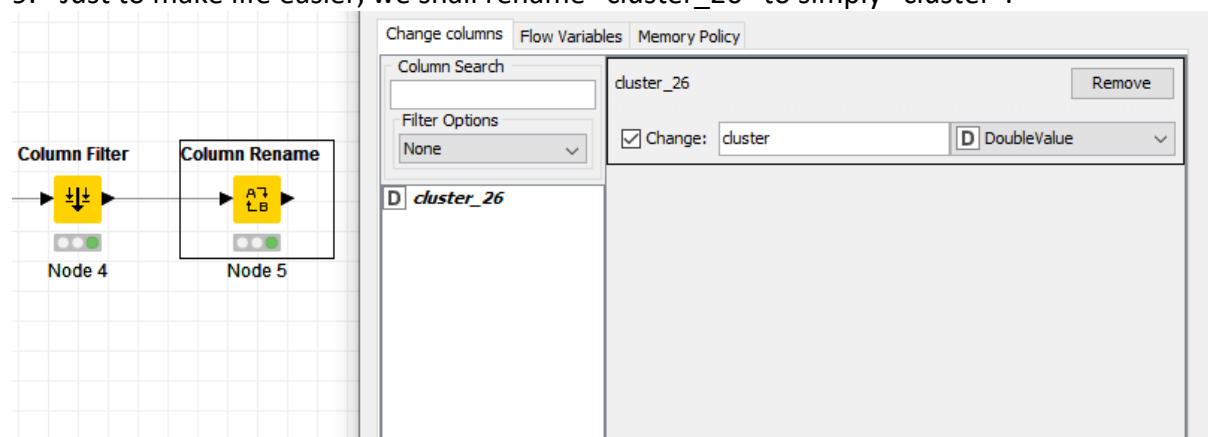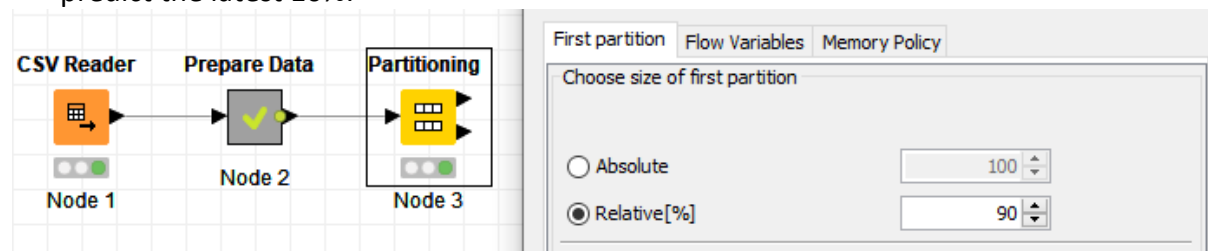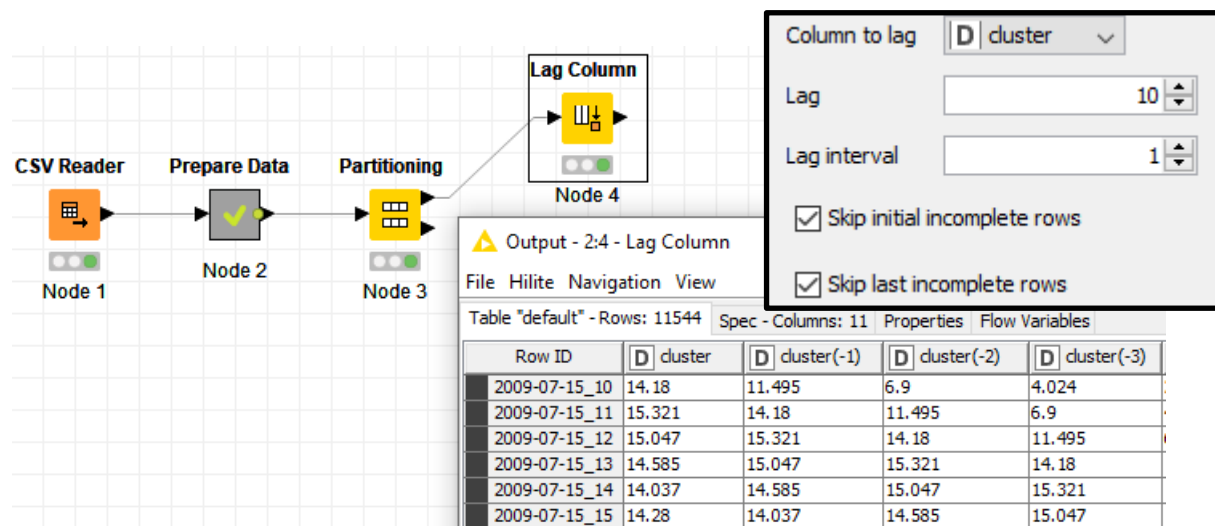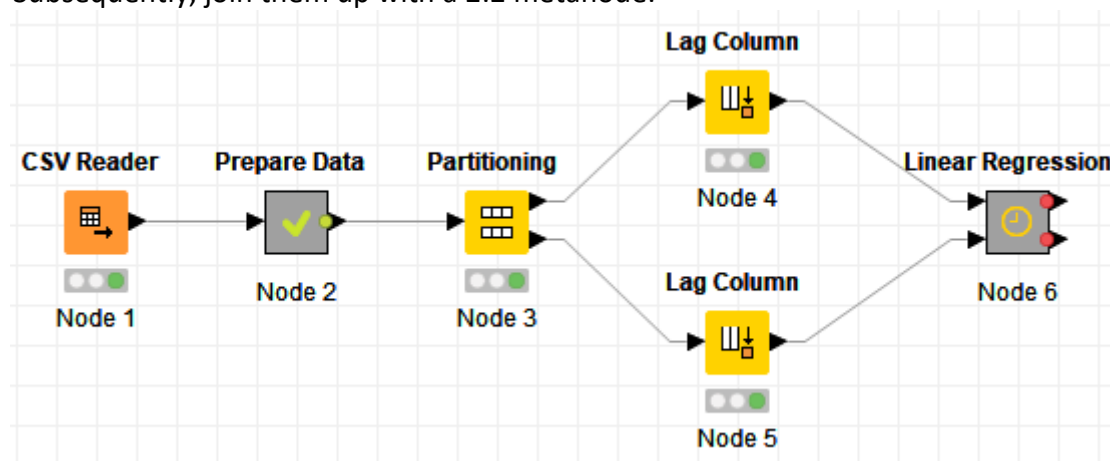


11. We have to properly structure the data to make current data (i.e. cluster) dependent cluster data that comes before it in time (i.e. cluster(-1), cluster(-2) etc…). The *Lag Column* node is used for this purpose. Note how in the final output, the first 10 rows are skipped as they do not have complete lag data, as we have checked the corresponding checkbox.



12. Use the same Lag Column node for both the training 90% data and testing 10% data. Subsequently, join them up with a 2:2 metanode:

13. In this metanode, put in a Linear Regression Learner to learn from the training 90% data:
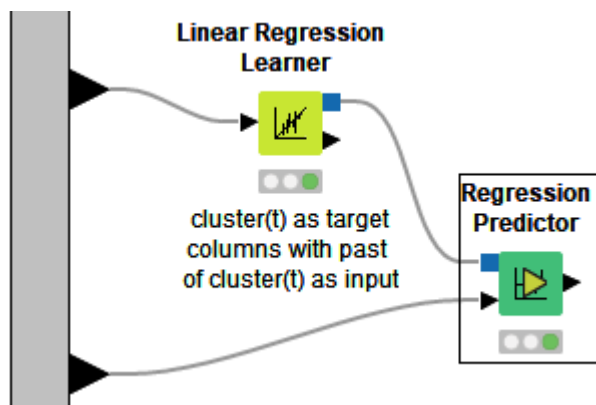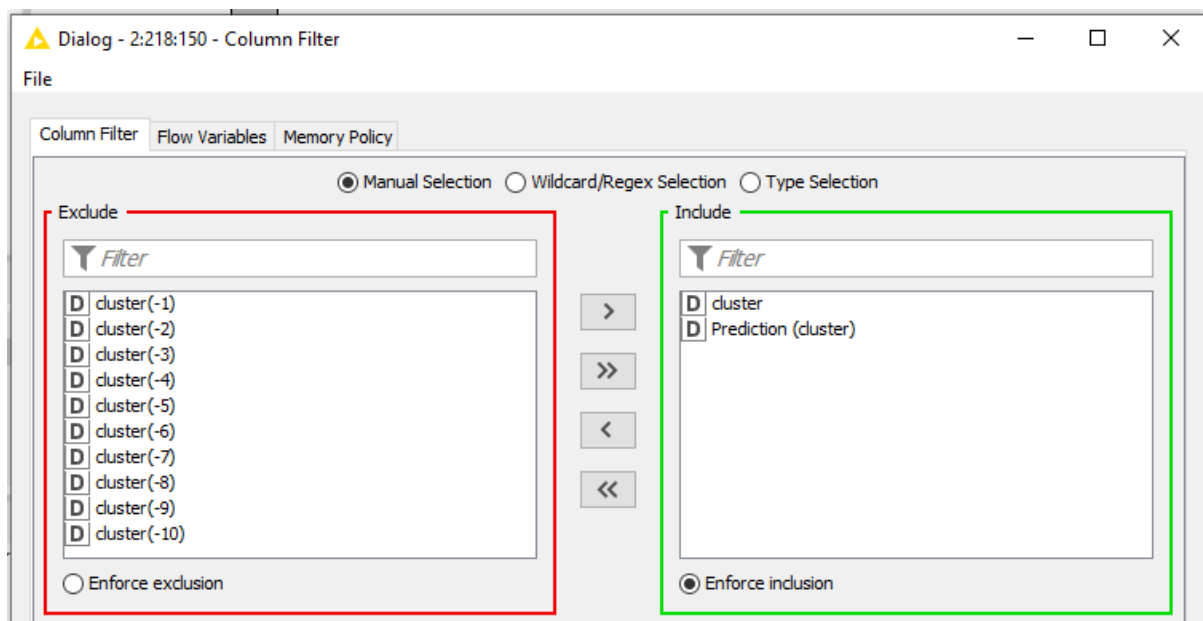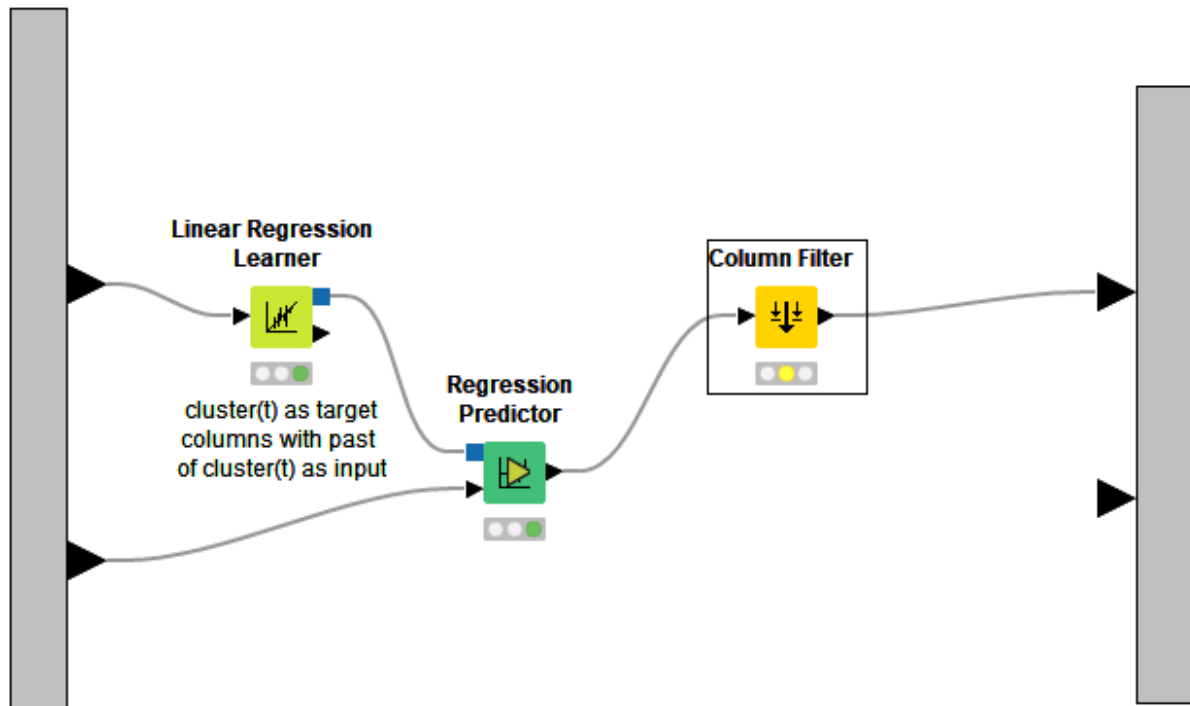


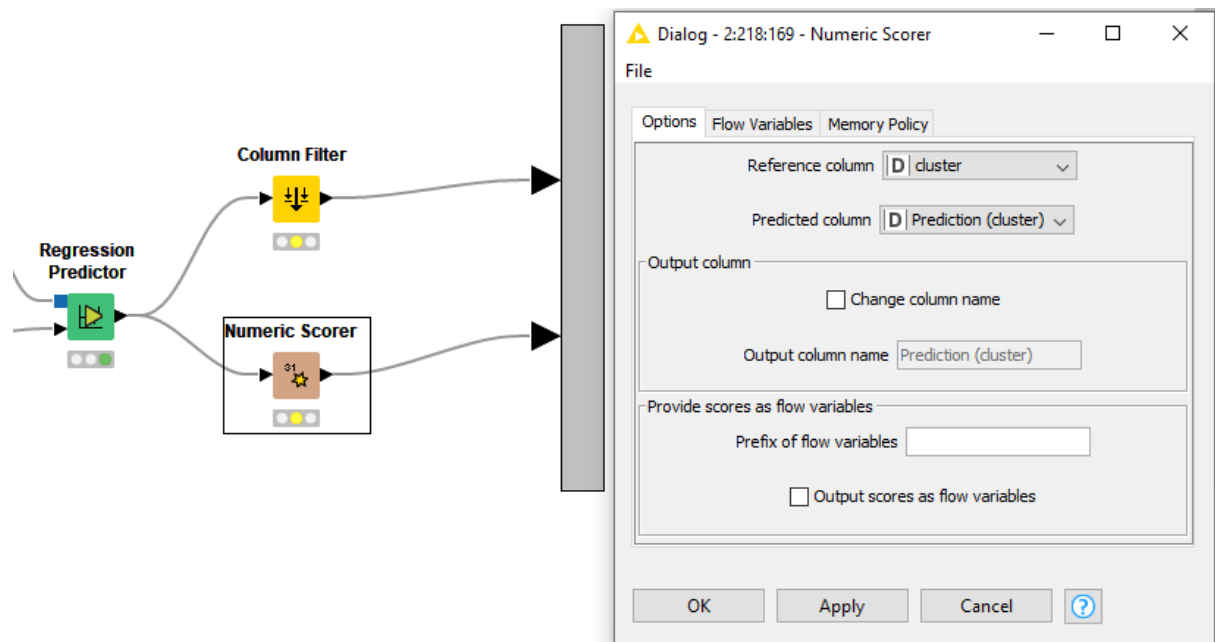14. Next, apply the learner's model to the test 10% data:

pages

15. Now, since we are only interested in current data and its prediction (all the cluster(-i) data are redundant as the model is trained), use a *Column Filter* node to remove them.
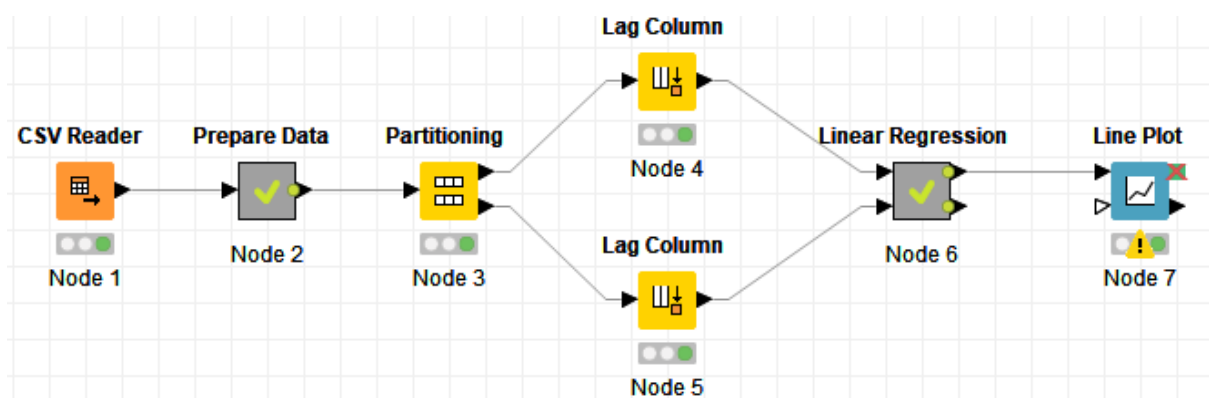
16. We can put in a *Numeric Scorer* node to assess the accuracy of the model on the test data:
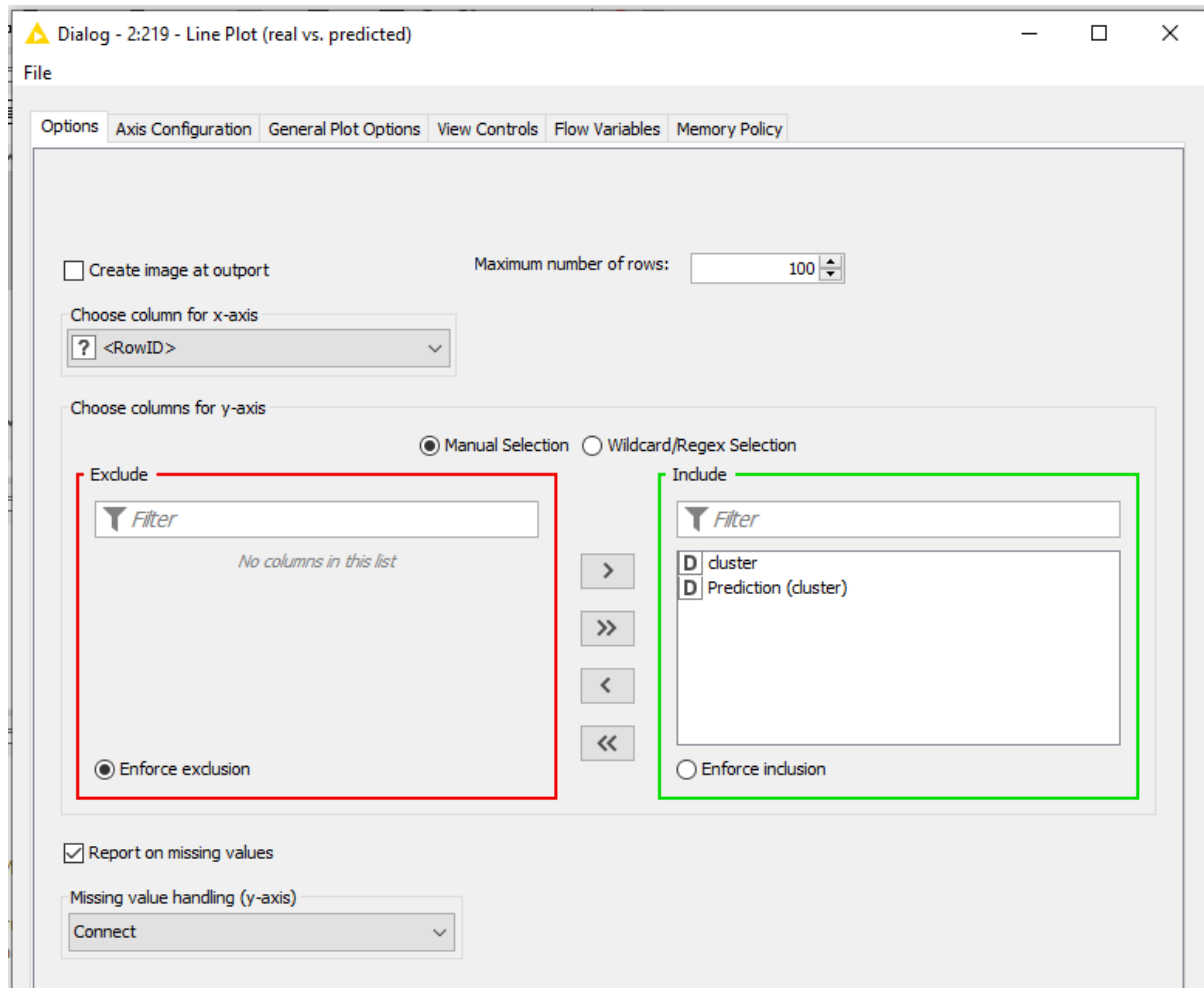


The output tatistics of the scorer tells us that the model has quite a high r-square value:

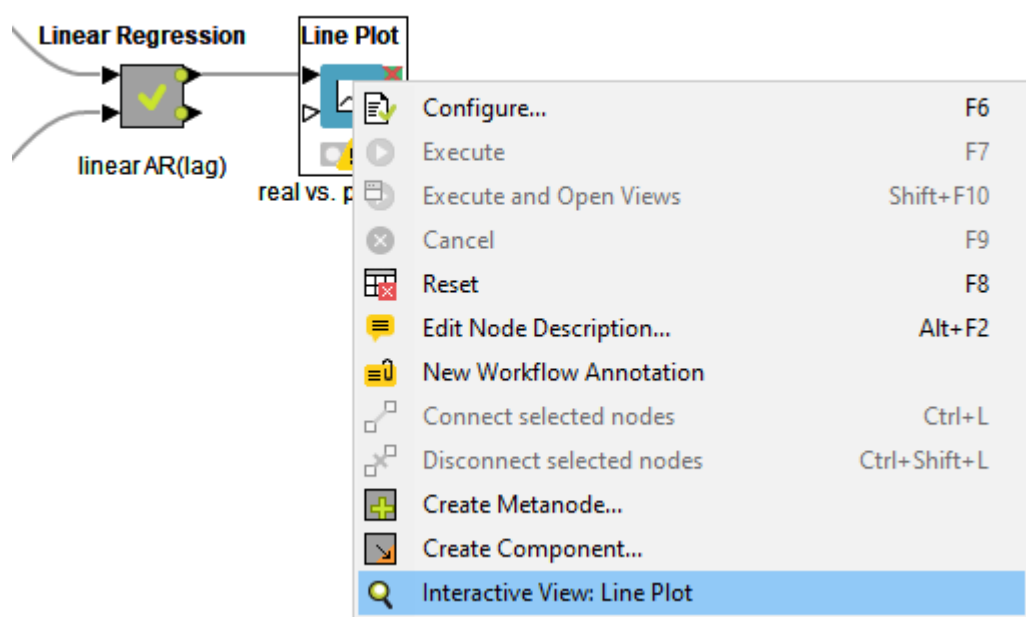| Row ID | D Predicti... |
|---|---|
| R^2 | 0.959 |
| mean absolut... | 0.773 |
| mean square... | 1.366 |
| root mean sq... | 1.169 |
| mean signed ... | -0.115 |
| mean absolut... | 0.105 |

17. Let's visualize the actual time series data with the prediction using a line plot:
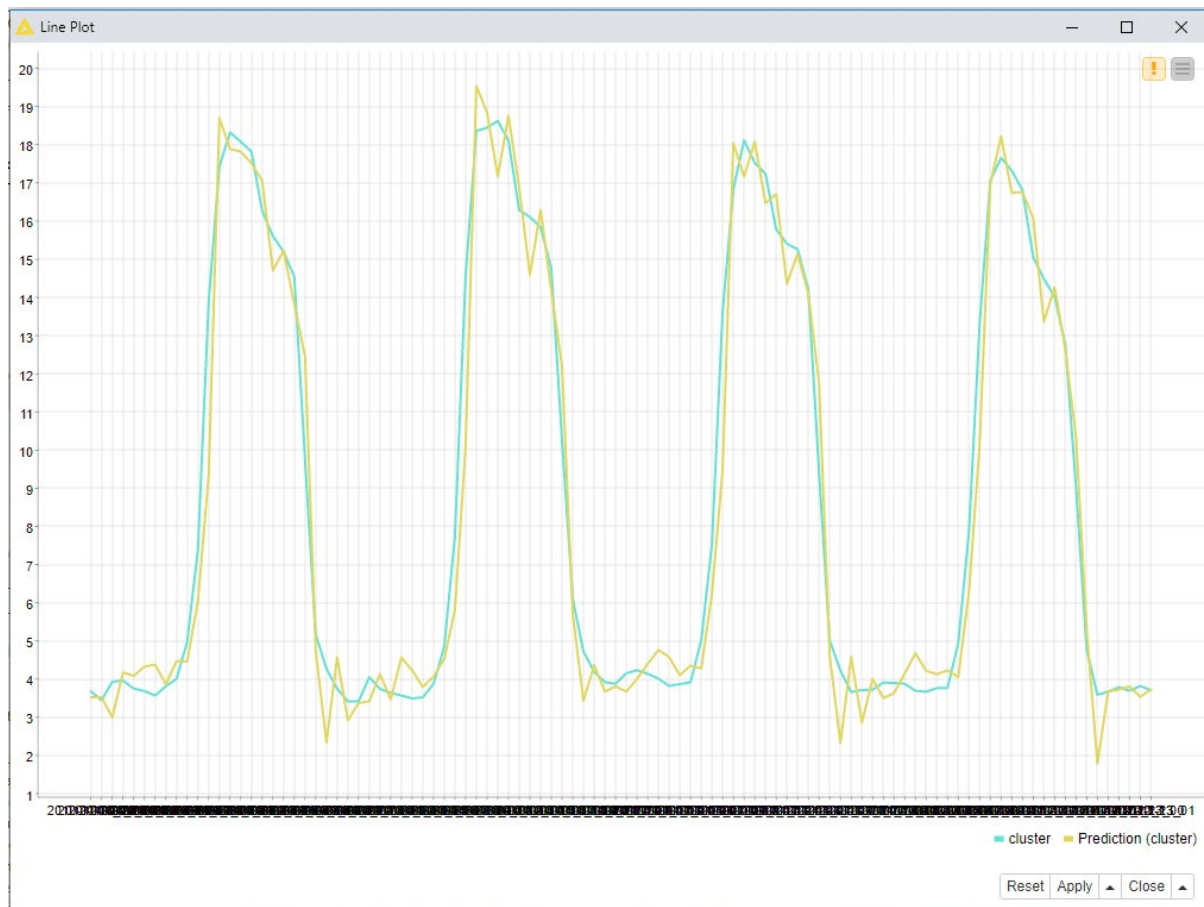
18. We can configure the *Line Plot* node as such:



19. Execcute the *Line Plot* node and open up the "Interactive View: Line Plot" once it is done.

20. The line plot shows that the prediction is quite accurate as compared to the actual cluster data. The prediction seems to have some noise, which one can try to decrease by using more past data for learning.



## Conclusion

We have seen how to use KNIME to use linear regression to predict time series data.