**Flashbacking on previous task report:**

In the last task report submitted on 04.01.2015 I reported that, I downloaded data set from http://www.esrl.noaa.gov/psd/gcos_wgsp/Timeseries/SOI/ for analyzing time series and forecasting. This data is about normalized surface pressure difference between Tahiti and Darwin islands for every month of the year starting from 1866 to 2012.

Then I inserted the data set into database in a table named "soi1866_2012" which contains columns named "year", "january", "february", "march", "april", "may", "june", "july", "august", "september", "october", "november", " december".

After that I made forecast for next 20 steps. Forecasts were made for every month of the year separately. For example, for April, I made forecast on the point that what amount of surface oscillation can occur in April 2013, April 2014 ... ... April 2032. I used one data for each April from the year 1886 to 2012 as training data set. In this way I made forecast for other months also.

Here forecasts are for next 20 steps.



I said you that I was planning to read documentation http://msdn.microsoft.com/en-us/library/bb895174.aspx on 'Cross Validation' to meaesure the accuracy of mining models.

## New report

I read whole documentation on "Cross validation" and partly on "Lift Chart" , "Classification Matrix". But none of these data model accuracy measuring tools is not valid for measuring accuracy of time series models.



If you perform cross-validation by using the stored procedures, you have the additional option of choosing the source of testing data. If you perform cross-validation by using the Data Mining Designer, you must use the testing data set that is associated with the model or structure, if any. Generally, if you want to specify advanced settings, you should use the cross-validation stored procedures.

Cross-validation cannot be used with time series or sequence clustering models. Specifically, no model that contains a KEY TIME column or a KEY SEQUENCE column can be included in cross-validation.

▲ Related Content



Validation

and Validation Tasks and
s

You can add multiple models to a lift chart, as long as the models all have the same predictable attribute. Models that do not share the attribute will be unavailable for selection in the **Input** tab.

You cannot display time series models in a lift chart or profit chart. A common practice for measuring the accuracy of time series predictions is to reserve a portion of historical data and compare that data to the predictions. For more information, see Microsoft Time Series Algorithm.
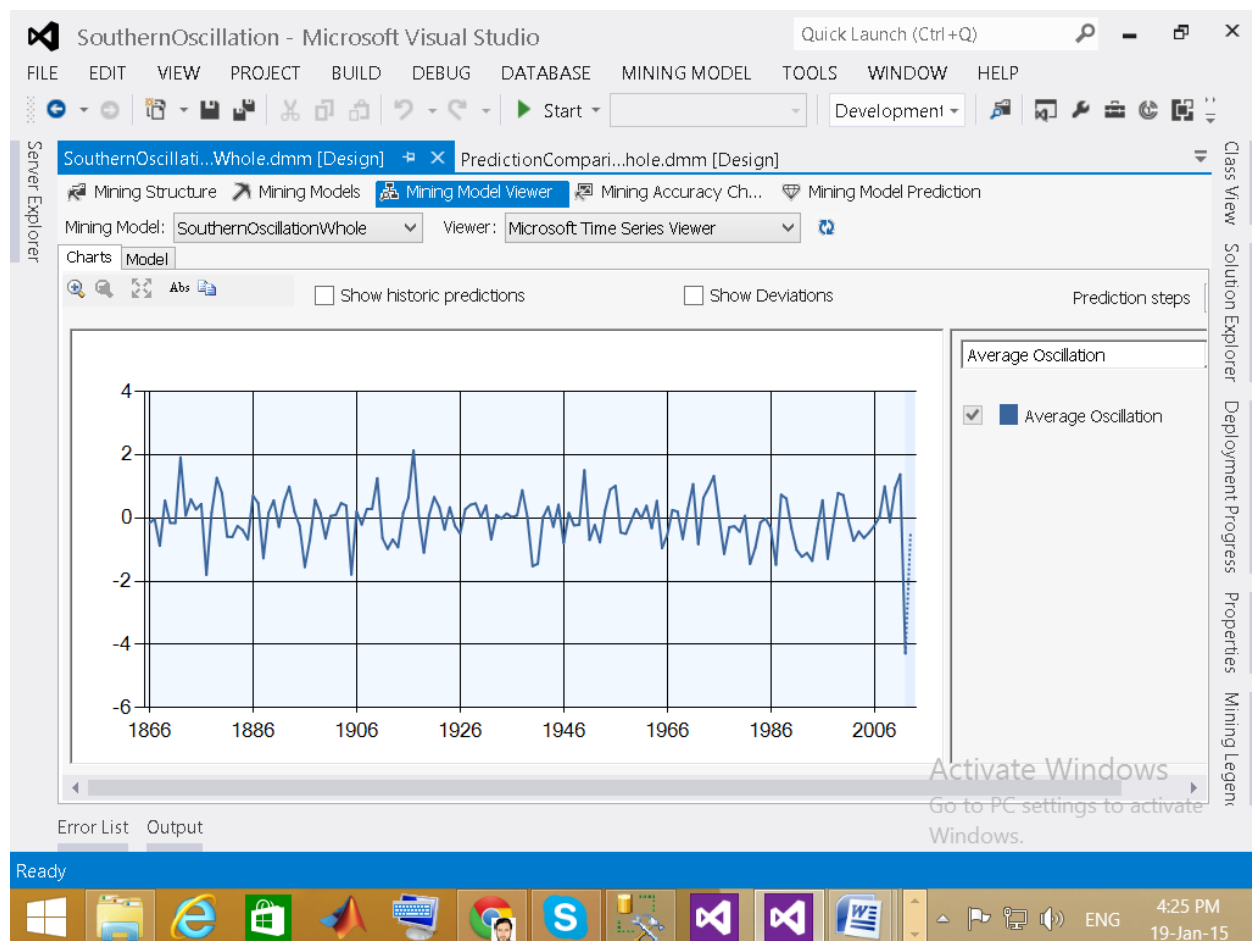
**Related Content**
Back to Top

After that I decided to compare two predictions done by the same model with different amount of traning data to measure accuracy of the model according to this documentation http://msdn.microsoft.com/en-us/library/cc879295.aspx . For doing this I needed to make a cross prediction according to this documentation http://msdn.microsoft.com/en-us/library/cc879284.aspx and after that I compared original prediction which I have done and showed you in the previous task report with cross prediction results.
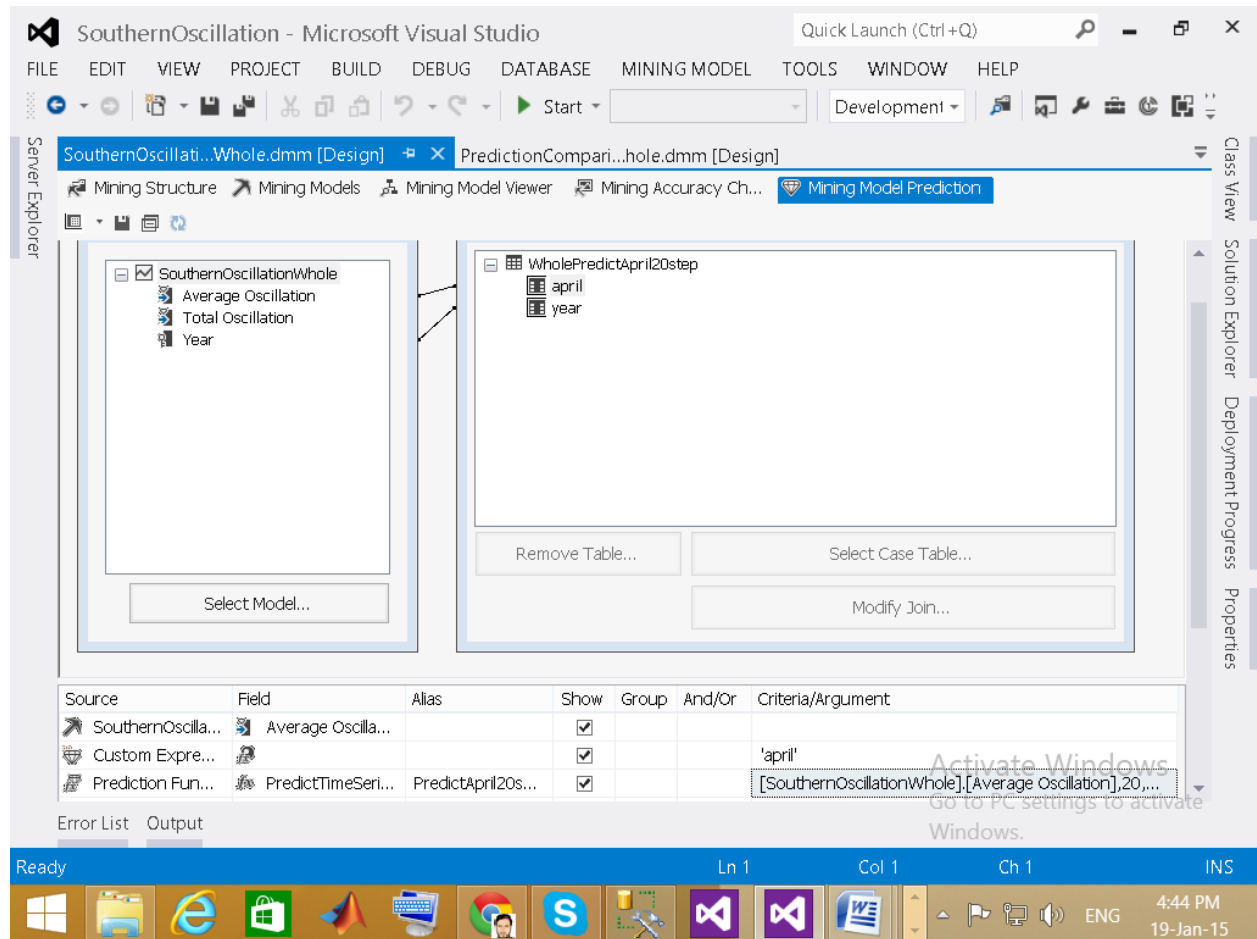
**Step 01: Making a cross prediction**

The process of using data from one series to predict trends in another series is called cross-prediction. That means a model will be trained from a series but will predict for another series. For example, here I trained a time series model named "SouthernOscillationWhole" by the data in the series named "SouthernWholeOscillation" which includes (whole data set) surface pressure occurred in all months (Jan - Dec) from the year 1866 to 2012.

This is model view of the model "SouthernOscillationWhole" trained by the data in series "SouthernWholeOscillation". **The prediction part here shown by dotted line is not important to consider**. This graph is only to observe the data fluctuation throughout the year from 1866 to 2012.

After training from these data I used this trained model to predict for the series named "SeriesApril" which includes one data for each April from the year 1886 to 2012. This is called cross prediction for month of April.

Here in the following picture the result of cross prediction for from the month of April of 2013 to April of 2032. In this way I made cross prediction for other months also.



Then I saved the result of cross prediction for April in a table named "WholePredictApril20step" in the database named "SouthernOscillation" which I attached with the mail for your kind consideration.

**Step 02: Comparing the cross prediction result and original predicrtion result:**

Now summary is following:

In the database "SouthernOscillation" contains two prediction result table for the month of April named "PrimaryPredictApril20step" and "WholePredictApril20step"

The table "PrimaryPredictApril20step" contains prediction result for the month of April done the model named "SouthernOscillationModel" trained by the data in the series named "SeriesApril" which includes one data for each April from the year 1886 to 2012 to predict surface pressure in April. I have done it and showed you in the previous task report. The table "PrimaryPredictApril20step" contains two columns named "PredictApril.$TIME", "PredictApril.April".

The table "WholePredictApril20step" contains cross prediction result for the month of April done by the model "SouthernOscillationWhole" which was trained by the data in the series named "SouthernWholeOscillation" which includes (whole data set) surface pressure occurred in all months (Jan - Dec) from the year 1866 to 2012. The table "WholePredictApril20step" contains two columns named "PredictAprilWhole.$TIME", "PredictAprilWhole.Average Oscillation".

**It is important to mention that, for the model "SouthernOscillationWhole" I used the same algorithm (Time series) and algorithm parameters which I used for the model named "SouthernOscillationModel". So that the prediction results tables "WholePredictApril20step" and "PrimaryPredictApril20step" of two models respectively can be considered the prediction of same model but by different training data. But for the easiness of task I used the model using different names.**
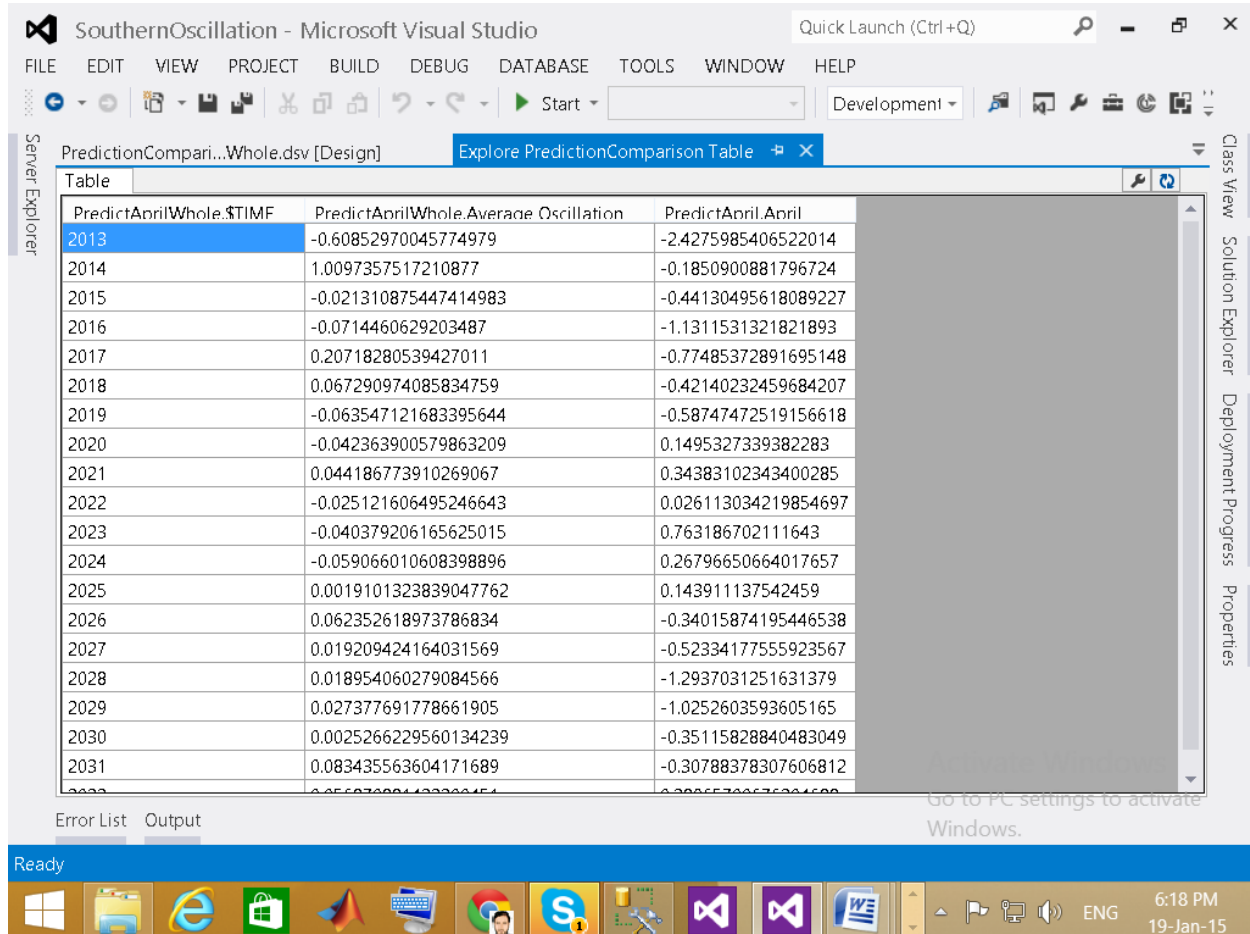
Comparison:

I made another data source view named "PredictionComparisonPrimaryVsWhole" in the project SouthernOscillation in MS data analysis tools using the following sql

```
select
[dbo].[WholePredictApril20step].[PredictAprilWhole.$TIME],[dbo].[WholePredictApril20step]
.[PredictAprilWhole.Average
Oscillation],[dbo].[PrimaryPredictApril20step].[PredictApril.April] from
[dbo].[WholePredictApril20step],[dbo].[PrimaryPredictApril20step] where
[dbo].[WholePredictApril20step].[PredictAprilWhole.$TIME]=[dbo].[PrimaryPredictApril20ste
p].[PredictApril.$TIME]
```

That means new data source view contains three columns named "`PredictAprilWhole.$TIME`", "`PredictAprilWhole.Average Oscillation`", "`PrimaryPredictApril20step`" where the columns "`PredictAprilWhole.Average Oscillation`" and "`PrimaryPredictApril20step`" contain the cross prediction result (done by the model "SouthernOscillationWhole" ) and primary prediction result (done by the model "SouthernOscillationModel") respectively.

The following picture is the explored view of the data source view named "PredictionComparisonPrimaryVsWhole"

Then I created another time series model named "PredictionComparisonPrimaryVsWhole" using the the column "PredictAprilWhole.$TIME" as key time and the columns "PredictAprilWhole.Average Oscillation" and "PrimaryPredictApril20step" as input and predict type. It is impotant to mention the purpose of creating of this model named "PredictionComparisonPrimaryVsWhole" is not to predict but to observe the graph of model view (which is actualy a visualized graph of our cross prediction and primary prediction) so that our conception about the difference between cross prediction and primary prediction.

This is the model view of the model named "PredictionComparisonPrimaryVsWhole". **The dotted line indicates prediction which is not impotant to cosider**. The plain line with blue colour shows primary prediction for the month of April from the year 2013 to 2032 and the plain line with red colour shows the cross prediction for the month of April from the year 2013 to 2032. It is the comparison. From this graph we can understand how the prediction done by the same model differes for the different training data.

**Conclusion:**

We see in the previous picture that, in 2014, there is enough difference between two prediction while the difference is smaller in follwing years. Why so much large difference in 2014? If we observe the original data we see that it is normal that sometimes surface pressure differes largely between the neighbouring times. For example, in the April 1869 surface pressure is 2.12 while at the same time in 1870 it is only 0.47 . Again in June 1884 surface pressure is 0.98 while it is in july in the same year is -0.28. **When we use the whole data set as training data set for cross prediction, the model faces more fluctuation than when use a part of data as training data set for primary prediction. That is the cause, why cross prediction shows more deeper fluctuation than than the primary prediction.**