

Time series Problems

Friday, January 19, 2024 11:35 AM

Very important notebook by daniel: https://dev.mdbourke.com/tensorflow-deep-learning/02_neural_network_classification_in_tensorflow/Improving_a_model

<https://www.uber.com/en-DE/blog/forecasting-introduction/>

This particular links tell about the how Uber tries to predict forecasting and what they have used for market demand

What can we predict into the future; because prediction can always go wrong
<https://otexts.com/fpp3/> this books tell about forecasting principles and fundamentals go through it.

<https://arxiv.org/abs/1905.10437> (Nbeats model)

- 1. Horizon: Think of the horizon as the future you're looking at. If you're trying to predict the weather for the next seven days, your horizon is seven days. It's like gazing into the future to see how things will unfold.
- 2. Window: Now, imagine you have a window that you can slide over your historical data. This window helps you see what's been happening in the past. For instance, if you want to predict tomorrow's weather, you might look at the weather patterns in the past month or year. The window is like a tool that helps you focus on a specific period in the past to make your predictions for the future.

From <https://whatnews.com/86341495-5d6c-4850-8d75-d8b8b6e61a16>

Time series models :

Moving average <https://imachivalearningmastery.com/moving-average-smoothing-for-time-series-forecasting-python>.

ARIMA (Autoregression Integrated Moving Average) <https://imachivalearningmastery.com/arima-for-time-series-forecasting-with-python>.

sktime (Scikit-Learn for time series) <https://github.com/alien-turing-institute/sktime>

TensorFlow Decision Forests (random forest, gradient boosting trees)

<https://www.tensorflow.org/decision-forests>

Facebook Kats (purpose-built forecasting and time series analysis library by Facebook analysis library by Facebook) <https://github.com/facebookresearch/Kats>

LinkedIn Greyscale (flexible, intuitive and fast forecasts) <https://github.com/linkedin/greyscale>

To know about auto correlation and how it benefits naive model or forecasting model (simple model): <https://towardsdatascience.com/how-not-to-use-machine-learning-for-time-series-forecasting-avoiding-the-pitfalls-1d95d7a6d544>

- `if keras.callbacks.EarlyStopping()` - stop the model from training if it doesn't improve validation loss for 200 epochs and restore the best performing weights using `restore_best_weights=True` (this'll prevent the model from training for loooooooooooooooong Period of time without improvement)
- `if keras.callbacks.ReduceLROnPlateau()` - if the model's validation loss doesn't improve for 100 epochs, reduce the learning rate by 10x to try and help it make incremental improvements (the smaller the learning rate, the smaller updates a model tries to make)

Some important points form the book Forecasting >Principles and Practice

```
We create this data as a pandas object using the following function:
def create_data():
    n = 1000000
    data = pd.Series(
        data=np.random.randn(n),
        index=pd.date_range("2000-01-01", periods=n, freq="D"),
    )
    return data
```

There are exactly the same data as were shown earlier, but now the data from each season are overlaid, a seasonal plot allows the underlying seasonal pattern to be seen more clearly, and is especially useful in identifying trends in which the pattern changes.

Lag plots are a type of statistical visualization used in time series analysis to examine the relationship between a variable and its lagged values (past observations). These plots help assess autocorrelation, which is the correlation of a time series with its own past values. The basic idea is to scatter plot each data point against its lagged values to identify patterns or dependencies.

- 1. Axis Values:
 - The x-axis represents the current observations (or time points) of the time series.
 - The y-axis represents the lagged observations, usually with a lag of 1 (previous value), 2 (value two time points ago), and so on.
- 2. Scatter Plot:
 - Each point on the lag plot corresponds to a pair of observations, where one is the current value, and the other is its lagged value.
 - The scatter plot visually displays how closely related the current observation is to its lagged values.
- 3. Interpretation:
 - If there is a strong correlation between the current observation and its lagged values, you will observe a recognizable pattern or trend in the lag plot.
 - Common patterns include diagonal lines, curves, or clusters of points, indicating a relationship between the variable and its past values.
- 4. Autocorrelation Assessment:
 - Lag plots are particularly useful for identifying autocorrelation in time series data. Autocorrelation is the correlation of a signal with a delayed copy of itself. If there is a strong autocorrelation, the lag plot will reveal a structured pattern.
- 5. Randomness Check:
 - In contrast, a lag plot with no apparent pattern suggests randomness or lack of autocorrelation in the time series data.

correlation measures the extent of a linear relationship

Autocorrelation

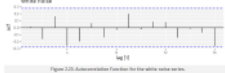
When there is a trend, the autocorrelation for small lags tend to be large and positive because observations nearby to those are likely to follow the same trend. As the lag increases, the correlation tends to decrease, often approaching zero as the lag increases.

When there is no trend, the autocorrelation for small lags tend to be small, indicating that the current value is not strongly related to its lagged values. As the lag increases, the correlation tends to decrease, often approaching zero as the lag increases.



correlogram.

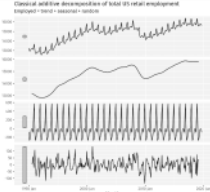
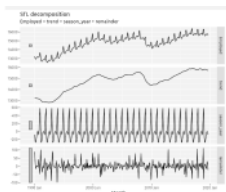
A time series data which does not have an autocorrelation is known as white noise. For white noise we consider that each autocorrelation is equal to zero. But we know that the relation cannot be exactly zero, so we consider a boundary of 95% that is if the autocorrelation is between that boundary then it is a white noise. If one or 5 % of the autocorrelation is out of that boundary then it is not considered as a white noise example the figure below :



Additive and multiplicative decomposition

Time series decomposition is the process of breaking down a time series into its constituent parts. The most common decomposition is the additive decomposition, which assumes that the time series is composed of a trend, a seasonal component, and a random component. The multiplicative decomposition is another common method, which assumes that the time series is composed of a trend, a seasonal component, and a random component, but the components are multiplied together rather than added.

TIME SERIES DECOMPOSITION

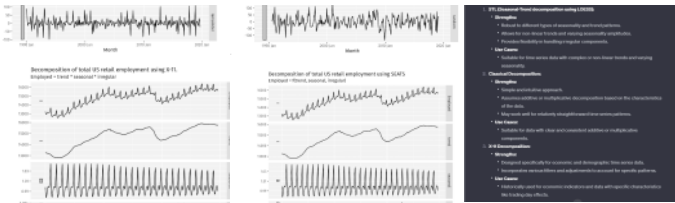


1. Additive decomposition

- The "trend" and "seasonal" components are added together.
- The "trend" component is the long-term movement of the data.
- The "seasonal" component is the periodic movement of the data.
- The "random" component is the noise in the data.

2. Multiplicative decomposition

- The "trend" and "seasonal" components are multiplied together.
- The "trend" component is the long-term movement of the data.
- The "seasonal" component is the periodic movement of the data.
- The "random" component is the noise in the data.



Residuals or difference between the predicted values and actual values:

1. Residuals are calculated as the difference between the observed values and the predicted values. They are used to assess the quality of the model fit. 2. Residuals are plotted against the predicted values to check for any patterns. 3. Residuals are used to calculate the standard error of the estimate. 4. Residuals are used to calculate the coefficient of determination. 5. Residuals are used to calculate the F-statistic. 6. Residuals are used to calculate the t-statistic. 7. Residuals are used to calculate the p-value. 8. Residuals are used to calculate the confidence interval. 9. Residuals are used to calculate the standard deviation. 10. Residuals are used to calculate the standard error of the mean. 11. Residuals are used to calculate the standard error of the regression. 12. Residuals are used to calculate the standard error of the estimate.

If a transformation has been used in the model, then it is often useful to look at residuals on the transformed scale. We call these "innovation residuals".

Useful properties of residuals is that: They have constant variance and they have normal distribution. This does not mean that all the residuals should display this property, but it is useful to have. And makes life easier.

Linear regression.

Assumptions

When we use a linear regression model, we are implicitly making some assumptions about the relationship in Equation (1.1).

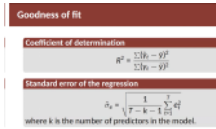
First, we assume that the model is a reasonable approximation to reality, that is, the relationship between the linear variable and the predictor variable is linear.

Second, we make the following assumptions about the errors $(\epsilon_1, \dots, \epsilon_n)$:

- 1. They have a mean of zero.
- 2. They are uncorrelated.
- 3. They have a constant variance.
- 4. They are normally distributed.

It is also useful to assume the errors follow a normal distribution with a constant variance σ^2 in order to make prediction intervals.

Another important assumption is that the linear regression model is the best predictor of the response variable. If we were predicting a variable that was not linear, we could not use the linear model. In this case, we would need to use a non-linear model. In this case, we would need to use a non-linear model. In this case, we would need to use a non-linear model.



Why there is necessary that we consider intercept for liner regression?

In layman terms, adding an intercept to a linear regression model is like including a baseline value that represents the starting point or average level of the dependent variable when all the independent variables are set to zero.

Imagine you are trying to predict someone's weight based on their height. The linear regression equation might look like this:

$$\hat{y} = \text{Intercept} + (\text{Coefficient} \times \text{Height})$$

Here:

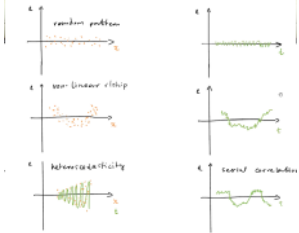
- The Intercept is a constant term, and it represents the expected weight when the person's height is zero. In reality, height can't be zero, but the intercept helps capture the baseline weight that is not influenced by height.

- The Coefficient is the slope of the line, indicating how much the weight changes for each unit increase in height.

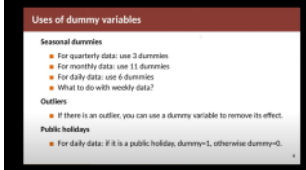
Without the intercept, the line would have to pass through the origin (0,0), and the model would assume that when height is zero, weight is also zero. This might not make sense in many real-world scenarios. The intercept allows the line to start at a more meaningful point, reflecting the baseline value when the independent variables are all zero or not applicable.

So, the intercept provides a starting point for the relationship between the variables and allows the model to better capture the real-world situation where some variables may not be relevant or have a meaningful value of zero.

Here are some plots of the error vs the predicted values and time series data.



Using predictors as dummy variable in time series data



How to select the best predictors to use in regression models.

A common approach is to use a stepwise selection procedure. This involves adding predictors to the model one by one, and at each step, the model with the best fit is selected.

Another common approach is to use a stepwise selection procedure. This involves adding predictors to the model one by one, and at each step, the model with the best fit is selected.

Adjusted R squared: The normal R2 tells us how well a model fits to historical data, but not how well the model forecast in future.

And there is no degree of freedom, adding a variable will increase the value of R2, for these reasons R2 should be used to determine whether the model will give good predictions.

An alternative which is designed to overcome these problems is the adjusted R2 (also called "R-bar-squared"): $R^2_{adj} = 1 - (1 - R^2) \frac{n}{n - k - 1}$.

From <https://datacamp.com/topics/evaluating-predictive-models>

Akaike's Information Criterion

Akaike's information criterion (AIC) is a measure of the relative quality of statistical models for a given dataset.

$$AIC = -2 \ln \left(\frac{L(\hat{\theta})}{L(\hat{\theta}_0)} \right) + 2k$$

where $L(\hat{\theta})$ is the likelihood of the model with parameters $\hat{\theta}$ and $L(\hat{\theta}_0)$ is the likelihood of the model with parameters $\hat{\theta}_0$.

The model with the lowest AIC is the best model. The AIC is a measure of the relative quality of statistical models for a given dataset.

The model with the lowest AIC is the best model. The AIC is a measure of the relative quality of statistical models for a given dataset.

The model with the lowest AIC is the best model. The AIC is a measure of the relative quality of statistical models for a given dataset.

The model with the lowest AIC is the best model. The AIC is a measure of the relative quality of statistical models for a given dataset.

The model with the lowest AIC is the best model. The AIC is a measure of the relative quality of statistical models for a given dataset.

The model with the lowest AIC is the best model. The AIC is a measure of the relative quality of statistical models for a given dataset.

The model with the lowest AIC is the best model. The AIC is a measure of the relative quality of statistical models for a given dataset.

The model with the lowest AIC is the best model. The AIC is a measure of the relative quality of statistical models for a given dataset.

The model with the lowest AIC is the best model. The AIC is a measure of the relative quality of statistical models for a given dataset.

The model with the lowest AIC is the best model. The AIC is a measure of the relative quality of statistical models for a given dataset.

The model with the lowest AIC is the best model. The AIC is a measure of the relative quality of statistical models for a given dataset.

The model with the lowest AIC is the best model. The AIC is a measure of the relative quality of statistical models for a given dataset.

The model with the lowest AIC is the best model. The AIC is a measure of the relative quality of statistical models for a given dataset.

The model with the lowest AIC is the best model. The AIC is a measure of the relative quality of statistical models for a given dataset.

While R^2 is widely used, neither has several limits than the other measures, its tendency to select in-sample predictors over variables even if they actually help forecasting.

Many studies have also used the Akaike Information Criterion (AIC) to select the best model. However, the AIC will select the model given enough data. However, in reality, there is rarely a true underlying model, and even if there was a true underlying model, selecting that model will not necessarily give the best forecasts (because the parameter estimates may not be accurate).

Consequently, we recommend that one of the MIC , AIC , or CF statistics be used, each of which has forecasting as their objective. If the value of T is large enough, they will all lead to the same model. In most of the examples in this book, we use the MIC value to select the forecasting model.

Imagine that you have 40 predictors and you can fit 2^{40} models which is tough so a method of forward stepwise regression can be used, where you start with only the intercept and then keep on adding each one of the predictors and see whether the accuracy is improving this is done till the model shows no improvement.

* **Laيمان's Explanation:** Ex-ante forecasts are predictions or estimates made before an event or period has occurred. It's like trying to guess the outcome of a game before it starts. For example, if you predict the score of a soccer match before the game begins, that's an ex-ante forecast.

- **• Layman's Explanation:** Ex-post forecasts, on the other hand, are predictions made after the event or period has already taken place. It's like making a guess about the score of a soccer match after it has ended. If you predict the score once the game is over, that's an ex-post forecast.

The simplest way of modelling a non linear relationship is to transform either the forecast variable or the predictor variable. The most commonly used variable is log functions.

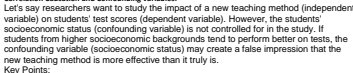
- ****Slope β_1 :**** This is a coefficient in the linear regression equation that represents the change in the dependent variable (y) for a one-unit change in the independent variable.

- **Elasticity Concept:** Elasticity is a measure of the responsiveness of one variable to changes in another variable. In this case, it's the responsiveness of $\frac{Y}{Y}$ to changes in the independent variable.

This interpretation is useful when dealing with percentage changes, especially in situations where the relationship between variables is better understood in terms of proportional or percentage adjustments rather than absolute changes.

When log transformation is not suitable or does not adequately address the characteristics of the data, piecewise transformations may be considered. A piecewise transformation involves dividing the data into different segments and applying distinct transformations to each segment.

- It's important to note that the choice of piecewise transformations should be guided by an understanding of the data and the underlying relationships between variables. Segmenting the data and selecting appropriate transformations require domain knowledge and careful exploration of the dataset.



- In summary, confounding variables are external factors that, if not properly accounted for, can lead to incorrect conclusions about the relationship between the variables being studied.

Having correlated predictors is not really a problem for forecasting, as we can still compute forecasts without needing to separate out the effects of the predictors. However, it becomes a problem with scenario forecasting as the scenarios should take account of the relationships between predictors. It is also a problem if some historical analysis of the contributions of various predictors is required.

Multicollinearity
Fortunately, if your purpose is primarily to predict or forecast Y , strong multicollinearity may not be a problem because a careful multiple regression program can still produce the best (least-squares) forecasts of Y based on all of the X variables.

From <<https://www.sciencedirect.com/topics/mathematics/multicollinearity-problem>>

Forecasts will be unreliable if the values of the future predictors are outside the range of the historical values of the predictors. For example, suppose you have fitted a regression model with predictors X_1 and X_2 which are highly correlated with each other, and suppose that the values of X_1 and X_2 in the training data ranged between 0 and 100. Then forecasts based on $X_1 > 100$ or $X_2 > 100$ or $X_1 < 0$ will be unreliable. It is always a little dangerous when future values of the predictors lie much outside the historical range, but it is especially problematic when multicollinearity is present. For example, if you are using gender as a predictor, and you are not interested in the specific contributions of each predictor, and if the values of your predictor variables are within their historical ranges, there is nothing to worry about – multicollinearity is not a problem except when there is perfect correlation.

from <https://l0tests.com/1003/causality.html>

Forecasts produced using exponential smoothing methods are weighted averages of past observations, with the weights decaying exponentially as the observations get older. In other words, the more recent the observation the higher the associated weight. This framework generates reliable forecasts quickly and for a wide range of time series, which is a great advantage and of major importance to applications in industry.

From <<https://otexts.com/fon3/exnsmooth.html>

When to be used :

- No clear trend or seasonal patterns

It falls between naive an average method, one giving only weight to the last one and another giving average weights to all the observation.

Exponential smoothing falls between these two models,

$$\hat{y}_{t+h|t} = \alpha y_t + \alpha(1-\alpha)y_{t-1} + \alpha(1-\alpha)^2 y_{t-2} + \dots, \quad (9.1)$$

where $0 \leq \alpha \leq 1$ is the smoothing parameter. The one-step-ahead forecast for time $T+1$ is a weighted average of all of the observations in the series y_1, \dots, y_T . The rate at which the weights decrease is controlled by the parameter α .

For any alpha between 0 and 1 the weights attached to the observation decrease exponentially as we go back in time and hence the name exponential smoothing.

Component form
where there is forecasting equation and smoothing equation

Component form

Now how to estimate alpha or the smoothing parameters?
Choosing an appropriate value for α involves a trade-off between bias and variance.

recent changes and stability based on past observations. The following methods for determining the value of α :

- The choice of α is often based on domain knowledge and the characteristics of the time series data.
- Smaller values of α (e.g., 0.1) result in smoother forecasts that react more slowly to changes.
- Larger values of α (e.g., 0.5 to 0.9) lead to more responsive forecasts that adapt quickly to recent changes.

- Grid Search: You can perform a grid search over a range of potential ϕ values and evaluate their performance on a validation set or through cross-validation.
- Iterate through different values of ϕ (e.g., 0.1, 0.2, ..., 0.9) and select the one that provides the best forecasting accuracy.

Optimization Algorithms:

- Some optimization algorithms can be used to automatically find the optimal value of ϕ that minimizes a forecasting error metric (e.g., Mean Squared Error) on a training dataset.
- Techniques like gradient descent or other optimization algorithms can be employed to search for the best ϕ value.

4. Exponential Smoothing Libraries:

- Many time series forecasting libraries, such as statsmodels in Python, provide functions for automatic parameter selection. These functions may use optimization algorithms or heuristics to determine the optimal smoothing parameters.

From <https://chat.openai.com/c/86d4149f-5d46-489d-8d78-c9d8b94eaf16>

Holt's linear trend:

Holt extended the smooth exponential with a trend

$$\begin{aligned} \text{Forecast equation: } & \hat{y}_{t+h|t} = a_t + b_t h \\ \text{Level equation: } & a_t = \alpha y_t + (1-\alpha) \hat{a}_{t-1} \\ \text{Trend equation: } & b_t = \beta (a_t - a_{t-1}) + (1-\beta) b_{t-1} \end{aligned}$$

where \hat{a}_t denotes an estimate of the level at time t , a_t denotes the smoothed component of the level (output of the filter at time t), α is the smoothing parameter for the level, $0 \leq \alpha \leq 1$, and β is the smoothing parameter for the trend, $0 \leq \beta \leq 1$. The filter takes as input the observed data that will be represented by y_t .

Disadvantage of holt's linear trend >

The forecasts generated by Holt's linear method display a constant trend (increasing or decreasing) indefinitely into the future. Empirical evidence indicates that these methods tend to over-forecast, especially for longer forecast horizons.

From <https://chat.openai.com/c/3a3c3eb7-ba9c-4000-8078-c9d8b94eaf16>

So for this we introduce a damped parameter known as dampers

When $\phi = 1$ (damped parameter) then it's equal to Holt's equation

In practice, ϕ rarely less than 0.8 as the damping has a very strong effect for smaller values. Values of ϕ closer to 1 will mean that a damped model is not able to be distinguished from a non-damped model. For these reasons, we usually restrict ϕ to a minimum of 0.8 and a maximum of 0.98.

From <https://chat.openai.com/c/3a3c3eb7-ba9c-4000-8078-c9d8b94eaf16>

Holt's winters seasonal method consist of all three smoothing equation one for level and another one for trend and another for seasonal.

- There are two variations to this method that differ in nature of the seasonal component
- Additive method - when seasonal variation are constant throughout the seasons and
 - Multiplicative when the seasonal variations proportional to the level of the series.
- With the additive method, the seasonal component is expressed in absolute terms in the scale of the observed series, and in the level equation the series is seasonally adjusted by subtracting the seasonal component. Within each year, the seasonal component will add up to approximately zero. With the multiplicative method, the seasonal component is expressed in relative terms (percentages), and the series is seasonally adjusted by dividing through by the seasonal component. Within each year, the seasonal component will sum up to approximately 1.

From <https://chat.openai.com/c/3a3c3eb7-ba9c-4000-8078-c9d8b94eaf16>

Short-hand	Method
α, β, γ	Simple exponential smoothing
α, β	Holt's linear method
$\alpha, \beta, \gamma, \delta$	Holt's linear trend method
α, β, γ	Winters' seasonal method
$\alpha, \beta, \gamma, \delta$	Winters' trend method
α, β, γ	Winters' seasonal method

Table 16.1 Describes the recursive calculations and prediction formulas. In each case, \hat{a}_t denotes the series level at time t , \hat{b}_t denotes the slope at time t , \hat{c}_t denotes the seasonal component of the series at time t , and \hat{d}_t denotes the slope of the trend at time t . $\alpha, \beta, \gamma, \delta$ are smoothing parameters. $\hat{a}_t = \alpha y_t + (1-\alpha) \hat{a}_{t-1}$, $\hat{b}_t = \beta (\hat{a}_t - \hat{a}_{t-1}) + (1-\beta) \hat{b}_{t-1}$, $\hat{c}_t = \gamma y_t + (1-\gamma) \hat{c}_{t-1}$, $\hat{d}_t = \delta (\hat{b}_t - \hat{b}_{t-1}) + (1-\delta) \hat{d}_{t-1}$.

Method	Level	Trend	Seasonal
SES	$\hat{a}_t = \alpha y_t + (1-\alpha) \hat{a}_{t-1}$	$\hat{b}_t = 0$	$\hat{c}_t = 0$
LES	$\hat{a}_t = \alpha y_t + (1-\alpha) \hat{a}_{t-1}$	$\hat{b}_t = \beta (\hat{a}_t - \hat{a}_{t-1}) + (1-\beta) \hat{b}_{t-1}$	$\hat{c}_t = 0$
LTES	$\hat{a}_t = \alpha y_t + (1-\alpha) \hat{a}_{t-1}$	$\hat{b}_t = \beta (\hat{a}_t - \hat{a}_{t-1}) + (1-\beta) \hat{b}_{t-1}$	$\hat{c}_t = \gamma y_t + (1-\gamma) \hat{c}_{t-1}$
WSES	$\hat{a}_t = \alpha y_t + (1-\alpha) \hat{a}_{t-1}$	$\hat{b}_t = \beta (\hat{a}_t - \hat{a}_{t-1}) + (1-\beta) \hat{b}_{t-1}$	$\hat{c}_t = \gamma y_t + (1-\gamma) \hat{c}_{t-1}$
WTES	$\hat{a}_t = \alpha y_t + (1-\alpha) \hat{a}_{t-1}$	$\hat{b}_t = \beta (\hat{a}_t - \hat{a}_{t-1}) + (1-\beta) \hat{b}_{t-1}$	$\hat{c}_t = \gamma y_t + (1-\gamma) \hat{c}_{t-1}$
WTES	$\hat{a}_t = \alpha y_t + (1-\alpha) \hat{a}_{t-1}$	$\hat{b}_t = \beta (\hat{a}_t - \hat{a}_{t-1}) + (1-\beta) \hat{b}_{t-1}$	$\hat{c}_t = \gamma y_t + (1-\gamma) \hat{c}_{t-1}$

ETS(A,N,N): simple exponential smoothing with additive errors

From <https://chat.openai.com/c/3a3c3eb7-ba9c-4000-8078-c9d8b94eaf16>

- estimates.
 - Error Correction Form:
 - By rearranging the smoothing equation, we get the error correction form.
 - It shows that the current level (\hat{a}_t) is adjusted based on the error in the previous forecast ($\hat{a}_{t-1} - y_{t-1}$).
 - The error ($\hat{a}_{t-1} - y_{t-1}$) is the difference between the actual observation (y_{t-1}) and the previous forecast (\hat{a}_{t-1}).
 - The adjustment is proportional to ϕ and the error ($\hat{a}_{t-1} - y_{t-1}$). If the error is large, the adjustment is larger, and if the error is small, the adjustment is smaller.
 - The adjustment corrects the estimated level based on how well the previous forecast aligned with the actual data.
- In simpler terms, the model continuously learns and adjusts its estimate of the underlying level based on the errors in its past predictions. If the model consistently overestimates or underestimates, the adjustments become larger or smaller, respectively. The smoothing parameter (ϕ) controls the sensitivity to new observations, and a smaller ϕ leads to a smoother estimate.

The innovation state space model is a type of statistical model used for time series forecasting, and it belongs to the broader class of state space models. It shares some similarities with exponential smoothing models but also introduces additional components to capture various aspects of time series behavior. Here's a brief overview of both concepts:

- Exponential Smoothing Models:**
 - Exponential smoothing models, such as Simple Exponential Smoothing (SES), Double Exponential Smoothing (Holt's method), and Triple Exponential Smoothing (Holt-Winters method), are based on the idea of exponentially decreasing weights for past observations.
 - These models typically involve components like level, trend, and seasonality. Simple Exponential Smoothing considers only the level, Holt's method adds a trend, and Holt-Winters method incorporates both trend and seasonality.
 - The smoothing parameters control the weight given to recent observations, and the models make predictions based on weighted averages.
- Innovation State Space Model:**
 - The innovation state space model extends the state space modeling framework to capture innovations or shocks in the time series.
 - It involves two key components: the state equation and the observation equation.
 - "State Equation:" Describes how the underlying state of the system evolves over time, including potential trend or other latent components.
 - "Observation Equation:" Relates the observed data to the underlying state, accounting for any innovations or shocks that might occur.
 - The model explicitly models the innovations, allowing it to capture sudden changes or unexpected events in the time series.

- Differences:**
 - While exponential smoothing models focus on exponentially weighted averages and may include trend and seasonality components, the innovation state space model provides a more general framework that explicitly models innovations.
 - The state space model is more flexible and can handle various types of underlying structures and dependencies.
 - The innovation state space model is often used when there are irregular or abrupt changes in the time series that might not be well-captured by traditional exponential smoothing models.

In summary, the innovation state space model introduces a broader framework that allows for a more explicit representation of innovations or shocks in the time series, providing greater flexibility in capturing complex patterns and sudden changes. The choice between these models depends on the characteristics of the specific time series being analyzed.

- Innovation state space model as explained to 5 year old
- Imagine you have a magical toy that moves around and does different things. This toy has two special parts:
- Magic Blueprint (State Equation):** This is like the secret plan or recipe for how the toy changes over time. It says if the toy is going up, down, or staying the same. So, if the toy is a car, the blueprint might say if it's speeding up, slowing down, or cruising at the same speed.
 - Watching the Toy (Observation Equation):** This part is like having a friend who watches the toy and tells you what's happening. Sometimes, the friend might see the toy doing unexpected things, like suddenly zooming faster or slowing down out of the blue. Now, let's imagine the toy sometimes does really surprising stuff—maybe it suddenly jumps or spins around. The innovation state space model is like having a special way to understand and explain these surprises. It helps us figure out why the toy did something unexpected and understand the plan behind it.
- In simpler terms, it's like having a cool way to understand and predict the surprises or changes in how the toy moves around. It's like having a magical guide that helps us make sense of the toy's actions, even when it does things we didn't expect! So, the innovation state space model is a helpful tool to understand and predict unexpected or surprising events in a way that even a 5-year-old can grasp.

From <https://chat.openai.com/c/86d4149f-5d46-489d-8d78-c9d8b94eaf16>

- Innovation Residuals:**
 - Definition: Innovation residuals, also known as innovations or shocks, represent the unanticipated components of the time series. They capture the difference between the observed values and the expected or predicted values from the model.
 - Significance: Innovation residuals are essential for capturing sudden changes or unexpected events in the time series. They provide insight into the unexplained variations that the model hasn't accounted for explicitly.
- Regular Residuals:**
 - Definition: Regular residuals, often referred to simply as residuals, are the differences between the observed values and the predicted values of the time series based on the chosen model.
 - Significance: Regular residuals are useful for assessing the overall goodness-of-fit of the model. They indicate how well the model captures the observed data, and examining the distribution of residuals helps identify patterns or systematic errors.

From <https://chat.openai.com/c/86d4149f-5d46-489d-8d78-c9d8b94eaf16>

Estimation and model selection

- Smoothing parameters, α , β , and γ , and the initial states ($\hat{a}_0, \hat{b}_0, \hat{c}_0$), are estimated by maximizing the "likelihood" of the data arising from the smoothed model.
- For models with additive error variances to minimizing SSE.
- For models with multiplicative errors, see equivalent to minimizing SE.

Innovation means shock or sudden changes in the time series data

This emphasizes that ETS models, through their statistical framework, explicitly account for errors, making it possible to generate prediction intervals. The concept of prediction intervals is crucial for understanding the uncertainty associated with forecasts, and ETS models provide a more structured way to handle this uncertainty compared to traditional exponential smoothing methods.

Some source for python time series data :
<https://towardsdatascience.com/time-series-libraries-for-python-4e41e9d0328b>

ARIMA MODELS

Stationary time Series : Time series whose statistical properties do not depend on time, that is they do not have trends or seasonality. Some times it may be confusing as cyclic time series without trends and seasonality can be called stationary data

In general, a stationary time series will have no predictable patterns in the long-term.

Differencing: This is a technique to make non-stationary time series into stationary one because it simplifies the modeling process and makes it easier to identify patterns and trends.

- Definition:**
 - Differencing involves computing the difference between consecutive observations in a time series.
- First Order Difference:**

- The most common form of differencing is the first-order difference, denoted as $y_t - y_{t-1}$
 - This operation calculates the change between each data point and its immediate predecessor.
3. Purpose
- The goal of differencing is to remove the trend or seasonality present in the original time series, making it stationary.
 - Stationarity implies that the statistical properties of the time series, such as mean and variance, do not change over time.
- (Plot such ACF helps identify non-stationary time series.)

For a stationary time series, the ACF will drop to zero relatively quickly, while the ACF of non-stationary data decreases slowly. Also, for non-stationary data, the value of $f(1)$ is often large and positive.

Random Walk

Random Walking: Random walking is a concept used to describe a process where an entity, such as a variable or a position, moves in a sequence of random steps. Each step is independent of the previous ones, and the direction or size of the step is determined by a random element.

In financial terms, a random walk is often used to describe the movement of stock prices. The idea is that, in an efficient market, future price changes are not predictable, and each price movement is like a random step. The unpredictable nature of these steps is akin to a person taking random steps without any clear pattern.

Here's a simple analogy:
Imagine standing at a point and taking steps forward or backward, where the decision to move forward or backward is based on a flip of a coin. If it's heads, you take a step forward; if it's tails, you take a step backward. The sequence of steps you take forms a random walk.

Key Points:

- Independence: Each step is not influenced by the previous steps; it's independent.
- Unpredictability: The overall path becomes unpredictable over time, even though each step is determined randomly.
- Used in Finance: Random walks are often used as a theoretical model for the unpredictable nature of financial markets.
- In financial theory, the idea of a random walk supports the efficient market hypothesis, suggesting that it's challenging to consistently outperform the market by predicting future price movements, as they are unpredictable and follow a random pattern.

From: <https://dataquest.com/blog/random-walk-data-science/>

- A random walk is often considered a non-stationary time series because its mean and variance change over time, and it doesn't exhibit consistent patterns.
 - Conversely, a stationary time series doesn't necessarily have to follow a random walk. It simply means that statistical properties remain constant, making it more amenable to modeling and forecasting.
- In summary, while a random walk is a type of non-stationary time series, not all stationary time series behave like a random walk. Stationary series can have various patterns, trends, or seasonality without the unpredictable nature associated with a strict random walk.

From: <https://dataquest.com/blog/random-walk-data-science/>

Financial data often exhibits a random walk or non-stationary behavior. To make it stationary, differencing is a common technique. Here's a breakdown:

4. First Differencing
- If the financial data follows a random walk, taking the first difference (subtracting each value from its previous one) can help remove the trend.
 - First differencing is often sufficient for making the series stationary if the underlying trend is linear.
5. Second Differencing
- In some cases, especially if the first difference still doesn't yield stationarity, you might need to take a second difference.
 - Second differencing involves differencing the already differenced series, providing a more pronounced effect on removing trends.
- Remember that the goal is to achieve stationarity, where statistical properties remain constant over time. This makes it easier to apply time series models and draw meaningful insights. However, it's important to be cautious with excessive differencing. If you find yourself needing a high order of differencing, it might indicate that the data is highly unpredictable, and modeling challenges could persist. Always assess the results and choose the differencing order that best suits your specific dataset.

From: <https://dataquest.com/blog/random-walk-data-science/>

There are Seasonal differences which can be done when there is strong seasonality in the data

But too much differentiation is also not good

Now how we check whether we need a differencing or not to make it stationary and for that we use Unit root test

In this test, the null hypothesis is that the data are stationary, and we look for evidence that the null hypothesis is false. Consequently, small p-values (e.g., less than 0.05) suggest that differencing is required. From: <https://datascience.com/blog/unit-root-test/>

The KPSS test p-value is reported as a number between 0.01 and 0.1. If the actual p-value is less than 0.01, it is reported as 0.01; and if the actual p-value is greater than 0.1, it is reported as 0.1. In this case, the p-value is shown as 0.01 (and therefore it may be smaller than that), indicating that the null hypothesis is rejected. That is, the data are not stationary. We can difference the data, and apply the test again. From: <https://datascience.com/blog/unit-root-test/>

A similar feature for determining whether seasonal differencing is required is `unitroot_seasonal`, which uses the measure of seasonal strength introduced in Section 4.3 to determine the appropriate number of seasonal differences required. No seasonal differences are suggested if $F_S < 0.64$, otherwise one seasonal difference is suggested.

From: <https://dataquest.com/blog/unit-root-test/>

Backshift Operator or notation:

With back shift operator we are differencing we can write like this:

The backward shift operator B is a useful notational device when working with time series lags: $By_t = y_{t-1}$.

In general, a d th-order difference can be written as $(1-B)^d y_t$.

For example, a seasonal difference followed by a first difference can be written as $(1-B)(1-B_n)^d y_t$

From: <https://datascience.com/backshift-operator/>

Auto regression Models

Autoregression, often denoted as AR, is a statistical modeling technique used in time series analysis. In simple terms, autoregression refers to modeling the relationship between a variable and its own past values. In other words, the value of the variable at the current time step is modeled as a linear combination of its previous values.

For example, an autoregressive model of order 1, denoted as AR(1), expresses the current value (y_t) as a function of the previous value (y_{t-1}):

$$y_t = \phi_1 y_{t-1} + \epsilon_t$$

Here:

- y_t is the value at time t .
 - ϕ_1 is a coefficient representing the impact of the previous value on the current value.
 - y_{t-1} is the value at the previous time step.
 - ϵ_t is the error term, representing the random noise or unobserved factors affecting y_t .
- The autoregressive process continues for higher orders (AR(2), AR(3), and so on), involving more past values to predict the current value. Autoregressive models are widely used in time series forecasting to capture the temporal dependencies and patterns present in the data.

From: <https://dataquest.com/blog/random-walk-data-science/>

Moving Average Models

A moving average model, often denoted as MA, is a statistical modeling technique used in time series analysis to capture patterns or dependencies in the data. Unlike autoregressive models that consider the relationship with past values of the variable itself, moving average models focus on the relationship between the current value and the past error terms.

The basic idea of a moving average model of order q , denoted as MA(q), is to express the current value (y_t) as a weighted sum of the past q error terms (ϵ_t):

$$y_t = \epsilon_t + \theta_1 \epsilon_{t-1} + \theta_2 \epsilon_{t-2} + \dots + \theta_q \epsilon_{t-q}$$

Here:

- y_t is the value at time t .
- ϵ_t is the mean of the time series.
- ϵ_t is the current error term, representing the random noise or unobserved factors affecting y_t .
- $\theta_1, \theta_2, \dots, \theta_q$ are coefficients representing the weights assigned to the past error terms.

The moving average process helps capture short-term fluctuations or random patterns in the time series. In combination with autoregressive models, moving average models contribute to creating more sophisticated models like ARMA (Autoregressive Moving Average) and ARIMA (Autoregressive Integrated Moving Average) for time series forecasting.

From: <https://dataquest.com/blog/random-walk-data-science/>

Invertibility

Imagine you have a sequence of numbers that represent some kind of data, let's say daily temperatures. A Moving Average (MA) model helps you understand how the current temperature depends on past random fluctuations.

Now, the term "invertible" is like saying you want to make this information more manageable and easy to understand. It's like taking a complex puzzle and turning it into a simpler one.

In the context of an MA model, being invertible means that you can look at the past random fluctuations in a way that makes sense. It's like saying, "Okay, these past influences don't have an endless impact on today's temperature." In other words, you can summarize the past in a neat way, making it easier to interpret and work with.

For example, if you find out that your MA model is not invertible, it's like saying, "Oops, the past influences are too chaotic, and I can't simplify them." It's a bit like having a messy puzzle that's hard to put together.

So, making an MA model invertible is like tidying up the information, making it more organized and easier to use for understanding and predicting future temperatures. It's about simplifying the complexity in a way that makes sense for analysis.

From: <https://dataquest.com/blog/random-walk-data-science/>

ARIMA

a) Non-seasonal ARIMA models

In time series analysis, ARIMA (Autoregressive Integrated Moving Average) models are used to model and forecast time series data. ARIMA models are a combination of three components: Autoregressive (AR), Integrated (I), and Moving Average (MA).

The general form of an ARIMA model is $ARIMA(p, d, q)$, where:

- p : Autoregressive order (number of lagged values of the time series used in the model).
- d : Integrated order (number of times the time series is differenced to make it stationary).
- q : Moving Average order (number of lagged values of the error term used in the model).

Let's consider an example of an ARIMA(1, 1, 1) model. The model can be written as:

$$(1 - \phi_1 B)(1 - B)(1 - \theta_1 B) y_t = \epsilon_t$$

where ϵ_t is the error term, B is the backward shift operator, ϕ_1 is the autoregressive coefficient, θ_1 is the moving average coefficient, and y_t is the time series value at time t .

The model can be simplified to:

$$(1 - \phi_1 B)(1 - B)(1 - \theta_1 B) y_t = \epsilon_t$$

where ϵ_t is the error term, B is the backward shift operator, ϕ_1 is the autoregressive coefficient, θ_1 is the moving average coefficient, and y_t is the time series value at time t .

The model can be simplified to:

$$(1 - \phi_1 B)(1 - B)(1 - \theta_1 B) y_t = \epsilon_t$$

where ϵ_t is the error term, B is the backward shift operator, ϕ_1 is the autoregressive coefficient, θ_1 is the moving average coefficient, and y_t is the time series value at time t .

The model can be simplified to:

$$(1 - \phi_1 B)(1 - B)(1 - \theta_1 B) y_t = \epsilon_t$$

where ϵ_t is the error term, B is the backward shift operator, ϕ_1 is the autoregressive coefficient, θ_1 is the moving average coefficient, and y_t is the time series value at time t .

The model can be simplified to:

$$(1 - \phi_1 B)(1 - B)(1 - \theta_1 B) y_t = \epsilon_t$$

where ϵ_t is the error term, B is the backward shift operator, ϕ_1 is the autoregressive coefficient, θ_1 is the moving average coefficient, and y_t is the time series value at time t .

The model can be simplified to:

$$(1 - \phi_1 B)(1 - B)(1 - \theta_1 B) y_t = \epsilon_t$$

where ϵ_t is the error term, B is the backward shift operator, ϕ_1 is the autoregressive coefficient, θ_1 is the moving average coefficient, and y_t is the time series value at time t .

The model can be simplified to:

$$(1 - \phi_1 B)(1 - B)(1 - \theta_1 B) y_t = \epsilon_t$$

where ϵ_t is the error term, B is the backward shift operator, ϕ_1 is the autoregressive coefficient, θ_1 is the moving average coefficient, and y_t is the time series value at time t .

The model can be simplified to:

$$(1 - \phi_1 B)(1 - B)(1 - \theta_1 B) y_t = \epsilon_t$$

where ϵ_t is the error term, B is the backward shift operator, ϕ_1 is the autoregressive coefficient, θ_1 is the moving average coefficient, and y_t is the time series value at time t .

The model can be simplified to:

$$(1 - \phi_1 B)(1 - B)(1 - \theta_1 B) y_t = \epsilon_t$$

where ϵ_t is the error term, B is the backward shift operator, ϕ_1 is the autoregressive coefficient, θ_1 is the moving average coefficient, and y_t is the time series value at time t .

The model can be simplified to:

$$(1 - \phi_1 B)(1 - B)(1 - \theta_1 B) y_t = \epsilon_t$$

where ϵ_t is the error term, B is the backward shift operator, ϕ_1 is the autoregressive coefficient, θ_1 is the moving average coefficient, and y_t is the time series value at time t .

The model can be simplified to:

$$(1 - \phi_1 B)(1 - B)(1 - \theta_1 B) y_t = \epsilon_t$$

where ϵ_t is the error term, B is the backward shift operator, ϕ_1 is the autoregressive coefficient, θ_1 is the moving average coefficient, and y_t is the time series value at time t .

The model can be simplified to:

$$(1 - \phi_1 B)(1 - B)(1 - \theta_1 B) y_t = \epsilon_t$$

where ϵ_t is the error term, B is the backward shift operator, ϕ_1 is the autoregressive coefficient, θ_1 is the moving average coefficient, and y_t is the time series value at time t .

The model can be simplified to:

$$(1 - \phi_1 B)(1 - B)(1 - \theta_1 B) y_t = \epsilon_t$$

where ϵ_t is the error term, B is the backward shift operator, ϕ_1 is the autoregressive coefficient, θ_1 is the moving average coefficient, and y_t is the time series value at time t .

Difference between Moving Average and Moving Average Smoothing

Both moving average (MA) and moving average smoothing models are related concepts used in time series analysis, but they serve different purposes.

1. Moving Average (MA) Model:

- Purpose: The MA model is a statistical model that describes the relationship between the current value of a time series and past error terms (residuals).
- Equation: The general form of an MA model of order q is expressed as a weighted sum of past error terms.
- Example Equation: $y_t = \epsilon_t + \theta_1 \epsilon_{t-1} + \theta_2 \epsilon_{t-2} + \dots + \theta_q \epsilon_{t-q}$
- Usage: It is used to capture short-term fluctuations or random patterns in the data.

2. Moving Average Smoothing:

- Purpose: Moving average smoothing, on the other hand, is a technique used to reduce noise and highlight trends in a time series by computing the average of adjacent values.
 - Calculation: It involves taking the average of a specific number of consecutive observations (a window or period) to smooth out fluctuations.
 - Example Calculation: For a simple moving average of order k , each value in the smoothed series is the average of the current observation and the $k-1$ previous observations.
 - Usage: It is used for trend analysis and noise reduction, making it easier to identify underlying patterns.
- In summary, the MA model focuses on modeling the relationship between the current value and past error terms to capture short-term dependencies, while moving average smoothing is a data preprocessing technique that involves averaging neighboring values to create a smoother representation of the time series.

From: <https://dataquest.com/blog/random-walk-data-science/>

seasonal AR processes are stationary and invertible

This criterion is then an important effect on the long-term forecasts obtained from these models.

- If $\phi = 0$ and $d = 0$, the long-term forecasts will go to zero.
- If $\phi = 0$ and $d = 1$, the long-term forecasts will go to a non-zero constant.
- If $\phi = 0$ and $d = 2$, the long-term forecasts will follow a straight line.
- If $\phi \neq 0$ and $d = 0$, the long-term forecasts will go to a value of the data.
- If $\phi \neq 0$ and $d = 1$, the long-term forecasts will follow a straight line.
- If $\phi \neq 0$ and $d = 2$, the long-term forecasts will follow a quadratic trend. (This is not recommended, and ϕ and γ will not persist.)

The value of d also has an effect on the prediction intervals. As higher the value of d , the more rapidly the prediction interval increases in size. If $d = 0$, the long-term forecast prediction interval will grow at the constant rate of the forecast error, so the prediction intervals will all be the same width for the same.

This behavior is seen in Figure 1 where $d = 0$ and $d = 1$ in this figure, the prediction intervals are about the same width for the last few forecast horizons, and the first few forecasts are close to the mean of the data.

Partial Autocorrelation (PACF) and ACF graph

Now you have to select a pure AR model or pure MA model
i.e. ARIMA(p,0,0) and ARIMA(0,0,q)

We check the ACF and PACF graphs,

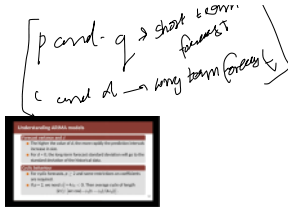
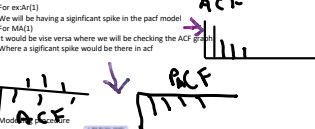
For ex: A(1)

We will be having a significant spike in the pacf model

For MA(1)

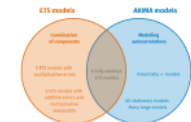
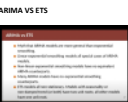
It would be vice versa where we will be checking the ACF graph

Where a significant spike would be there in acf

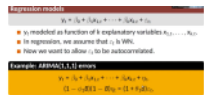


Seasonal ARIMA models

We will also have a part where we have a seasonal models.
It can be represented by the (P,D,Q)m

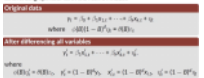


DYNAMIC REGRESSION MODELS



1. The estimated coefficients $\hat{\beta}_0, \dots, \hat{\beta}_p$ are no longer the best estimates, as more information has been gained for the coefficients.
2. The estimated coefficients $\hat{\beta}_0, \dots, \hat{\beta}_p$ are no longer the best estimates, as more information has been gained for the coefficients.
3. The ARIMA model of the first model is no longer a good model as in which the best model is becoming.
4. In most cases, the regression model with the coefficients will be too small, and so some additional variables will appear to be important when they are not. This is known as "spurious regression".

Any regression with an ARIMA error can be rewritten as an ARIMA model with an ARIMA error by differencing all variables with the same differencing operator as in the ARIMA model.



There are two different ways of modelling a dynamic model. A deterministic trend is obtained using the regression model $y_t = \beta_0 + \beta_1 t + \beta_2 t^2 + \dots + \beta_p t^p + \epsilon_t$ where y_t is an ARIMA process, ϵ_t is white noise, and β_0, \dots, β_p are parameters to be estimated. Alternatively, a stochastic trend is obtained using the model $y_t = \beta_0 + \beta_1 t + \beta_2 t^2 + \dots + \beta_p t^p + \epsilon_t$ where y_t is an ARIMA process with $d = 1$ or more, and ϵ_t is white noise. In this case, the difference operator $(1 - B)$ is applied to the data, and the resulting series is modelled as an ARIMA process. This is known as "stochastic regression".

More in the book

Hierarchical and grouped time series :