

STAT 443 - Forecasting

Cameron Roopnarine

Last updated: March 13, 2021

Contents

Contents	1
1 Week 1	3
1.1 What is a time series?	3
1.2 Basic Principles of Forecasting	6
1.3 Definitions of Stationary	7
1.4 White Noise and Stationary Examples	9
1.5 Weak versus Strong Stationary	11
1.6 † Theoretical L2 Framework for Time Series	13
1.7 Signal and Noise Models	14
1.8 Time Series Differencing	16
2 Week 2	20
2.1 Autocorrelation and Empirical Autocorrelation	20
2.2 Modes of Convergence of Random Variables	23
2.3 † M-dependent CLT	27
2.4 † Two Plus Delta Moment Calculation	31
2.5 † Linear Process CLT	32
2.6 Asymptotic Properties of Empirical ACF	34
2.7 Interpreting the Autocorrelation Function (Non-stationary)	36
3 Week 3	40
3.1 Moving Average Processes	40
3.2 Autoregressive Processes	43
3.3 ARMA Processes	49
3.4 ARMA Process Examples and ACF	51
4 Week 4	55
4.1 Stationary Process Forecasting	55
4.2 Best Linear Prediction	56
4.3 Partial ACF	58
4.4 ARMA Forecasting	60
4.5 ARMA Forecasting Example 1: Cardiovascular Mortality	62
4.6 ARMA Forecasting Example 2: Johnson and Johnson	68
5 Week 5	69
5.1 ARMA Parameter Estimation: AR Case	69
5.2 ARMA Parameter Estimation: MLE	71
5.3 Model Selection Diagnostic Tests	72
5.4 Model Selection Information Criteria	74
5.5 ARIMA Models	75
5.6 ARIMA Modelling Example	77

6	Week 6	78
6.1	SARIMA Models	78
6.2	SARIMA Cardiovascular Mortality Example	79
6.3	Time Series Cross-Validation	80
6.4	Cross-Validation Example	81
6.5	Simulated and Bootstrapped Prediction Intervals	81
6.6	Bootstrap Prediction Intervals Example	82
7	Week 7	83
7.1	Exponential Smoothing Models Introduction	83
7.2	Exponential Smoothing as a State Space Model	84
7.3	Multiplicative Exponential Smoothing Models	85
7.4	Exponential Smoothing Model Selection	87
7.5	J and J Exponential Smoothing Forecast	87

Chapter 1

Week 1

1.1 What is a time series?

In classical statistics, we normally consider $X_1, \dots, X_n \in \mathbf{R}^p$, a **simple random sample**.

In particular,

- (1) X_1, \dots, X_n are i.i.d. (independent and identically distributed)
- (2) $X_i \sim F_\theta$ which is a common distribution characterized by θ .

Examples:

1. $X_i \stackrel{\text{iid}}{\sim} \mathcal{N}(\mu, \sigma^2)$, and we wish to estimate and perform inference on μ and σ^2 .
2. $X_i = \begin{bmatrix} Y_i \\ Z_i \end{bmatrix}$ where Y_i is a dependent variable, and Z_i is an independent variable. Perhaps we happen to observe Y_i and Z_i in pairs, and we posit a model:

$$Y_i = \beta^\top Z_i + \varepsilon_i, \quad \varepsilon_i \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma_\varepsilon^2)$$

REMARK 1.1.1

The relationship between Y_i and Z_i doesn't depend on i , it only depends upon the common parameter β , and it assumes that ε_i has fixed variance for each i .

3. In such settings, one is typically interested in:
 - (a) Prediction: based on the data, how can we predict the behaviour of these variables in the future?
 - (b) Inference: how do we use the data to try to estimate and better understand the underlying mechanism which generates the data? For example, a linear model or simple Gaussian model.

DEFINITION 1.1.2: Time series

We say X_1, \dots, X_T is an (observed) **time series** of length T if X_t denotes an observation obtained at time t . In particular, the observations are ordered in time.

DEFINITION 1.1.3: Real-valued time series

If $X_t \in \mathbf{R}$, we say X_1, \dots, X_T is a **real-valued (scalar) time series**.

DEFINITION 1.1.4: Multivariate time series

If $X_t \in \mathbb{R}^p$, we say X_1, \dots, X_T is a **multivariate (vector-valued) time series**.



Figure 1.1: Quarterly Johnson and Johnson Earnings

Figure 1.1

```
plot(jj, type = "o", ylab = "Quarterly Earnings per Share")
```

Observe that in Figure 1.1:

- The earnings are steadily increasing over time.
- There is heterogeneity in the variance over time.

With time series data, we are typically concerned with the same goals as in classical statistics (prediction and inference). However, in contrast with time series, the data often exhibit:

(1) **Heterogeneity**

- Time trends $\rightarrow \mathbb{E}[X_t] \neq \mathbb{E}[X_{t+h}]$.
- Heteroskedasticity $\rightarrow \mathbb{V}(X_t) \neq \mathbb{V}(X_{t+h})$.

In classical statistics, it's assumed that all the observations have the same distribution which is clearly not the case in time series.

(2) **Serial Dependence (Serial Correlation)**

- Observations that are temporally close appear to depend on each other.

In classical statistics, each successive observation is assumed to be independent which is clearly not the case in time series.

Figure 1.2

```
plot(gtemp, type = "o", ylab = "Global Temperature Deviations")
```

Observe that in Figure 1.2:

- The global temperature is steadily increasing over time.
- Heterogeneity exists within the mean over time.



Figure 1.2: x_t is the deviation of global mean yearly temperature from the mean computed from 1951 to 1980

- Heterogeneity exists within the variance over time, although it is not very apparent.
- Serial dependence occurs.

Let's formally define a time series.

DEFINITION 1.1.5: Time series, Observed stretch

We say $\{X_t\}_{t \in \mathbf{Z}}$ is a **time series** if $\{X_t : t \in \mathbf{Z}\}$ is a stochastic process indexed by \mathbf{Z} . In other words, there is a common probability space $(\Omega, \mathcal{F}, \mathbb{P})$ such that $X_t : \Omega \rightarrow \mathbf{R}$ is a random variable for all t . In relation to the original definition, we say X_1, \dots, X_T is an **observed stretch (realization, simple path)** of length T from $\{X_t\}_{t \in \mathbf{Z}}$.

Formally speaking, we think of a time series as being a little snippet of one long sample path the stochastic process for which would characterize all the serial dependence, time trends, and heteroskedasticity that exist within a time series as can be seen in 1.3.



Figure 1.3: Time Series

1.2 Basic Principles of Forecasting

Consider a time series of length T , namely X_1, \dots, X_T . Based on X_1, \dots, X_T , we would like to produce a “best guess” for X_{T+h} :

$$\hat{X}_{T+h} = \hat{X}_{T+h|T} = f_h(X_T, \dots, X_1)$$

DEFINITION 1.2.1: Forecast, Horizon

For $h \geq 1$, our “best guess”

$$\hat{X}_{T+h} = f_h(X_T, \dots, X_1)$$

is called a **forecast** of X_{T+h} at **horizon** h .

Goals of Forecasting

Goal 1

- Choose f_h “optimally.” Normally, we or the practitioner have some measure, say $L(\cdot, \cdot)$, in mind for determining how “close” \hat{X}_{T+h} is to the true value, X_{T+h} . We then wish to choose f_h so that $L(X_{T+h}, f_h(X_T, \dots, X_1))$ is minimized, where $L(\cdot, \cdot)$ is a loss function.

EXAMPLE 1.2.2

The most common measure of $L(\cdot, \cdot)$ is the **mean-squared error** (MSE), defined by

$$L(X, Y) = \mathbb{E}[(X - Y)^2]$$

Goal 2

- Quantify the uncertainty in the forecast. This entails providing some description of how close we expect \hat{X}_{T+h} to be to X_{T+h} .

EXAMPLE 1.2.3: Why is it important to quantify uncertainty?

Suppose every minute, we flip a coin and denote

- (Heads): $H \rightarrow 1$
- (Tails): $T \rightarrow -1$
- X_t = outcome in minute t , where $t = 1, \dots, T$.

This produces a time series of length T , which is a random sequence of (1)’s and (−1)’s. Note $\mathbb{E}[X_t] = 0$ for all t . If we wish to forecast for $h \geq 1$, consider $\hat{X}_{T+h} = f(X_T, \dots, X_1)$, thus

$$\begin{aligned} L(X_{T+h}, \hat{X}_{T+h}) &= \mathbb{E}[(X_{T+h} - \hat{X}_{T+h})^2] \\ &= \mathbb{E}[X_{T+h}^2] + \mathbb{E}[\hat{X}_{T+h}^2] - 2\mathbb{E}[X_{T+h}\hat{X}_{T+h}] \\ &= \mathbb{E}[X_{T+h}^2] + \mathbb{E}[\hat{X}_{T+h}^2] - 2\mathbb{E}[X_{T+h}]\mathbb{E}[\hat{X}_{T+h}] \\ &= \mathbb{E}[X_{T+h}^2] + \mathbb{E}[\hat{X}_{T+h}^2] \end{aligned}$$

Note that we can write $\mathbb{E}[X_{T+h}\hat{X}_{T+h}] = \mathbb{E}[X_{T+h}]\mathbb{E}[\hat{X}_{T+h}]$ since \hat{X}_{T+h} is a function of the data X_T, \dots, X_1 , and hence independent of X_{T+h} .

Furthermore, note that $\mathbb{E}[X_{T+h}^2] = \mathbb{V}(X_t)$ since $\mathbb{E}[X_{T+h}] = 0$.

We can minimize this by taking $\hat{X}_{T+h} = 0$. There’s nothing “wrong” with this forecast, but ideally we would also be able to say that the sequence appears to be random, and that we don’t expect this forecast to be close to the actual value.

Furthermore, for this basic reason, one can always argue that any forecast that’s not accompanied by some type of quantification of how close we expect the forecast to be, is at very least hard to

interpret; at worst, meaningless because it doesn't describe the accuracy for which we expect the forecast to perform.

How can we quantify the uncertainty in forecasting?

Ideal: The predictive distribution, that is,

$$X_{T+h} \mid X_T, \dots, X_1$$

Excellent: Predictive intervals/sets, that is, for some $\alpha \in (0, 1)$ find an interval I_α such that

$$\mathbb{P}(X_{T+h} \in I_\alpha \mid X_T, \dots, X_1) = \alpha$$

A common example is with $\alpha = 0.95$. Often times, such intervals take the form

$$I_\alpha = (\hat{X}_{T+h} - \hat{\sigma}_h, \hat{X}_{T+h} + \hat{\sigma}_h)$$

Concluding Remarks

1. Estimating predictive distribution leads one towards *estimating* the joint distribution of

$$X_{T+h}, X_T, \dots, X_1$$

For example, the ARMA and ARIMA models.

2. It is important that we acknowledge that some things cannot be predicted!

“It's tough to make predictions, especially about the future.”—Yogi Berra

1.3 Definitions of Stationary

Given a time series X_1, \dots, X_T , we are frequently interested in estimating the joint distribution of

$$X_{T+h}, X_T, \dots, X_1$$

which is useful for forecasting and inference.

The joint distribution is a feature of the process $\{X_t\}_{t \in \mathbb{Z}}$

$$X_1, \dots, X_T \xrightarrow{\text{infer}} \{X_t\}_{t \in \mathbb{Z}}$$

- X_1, \dots, X_T : Observed data.
- $\{X_t\}_{t \in \mathbb{Z}}$: Stochastic process.

The worst case: $X_t \sim F_t$, where F_t is a *changing* function of t . If so, it is hard to pool the data X_1, \dots, X_T to estimate F_t . If **serial dependence** occurs; that is, if the distribution of (X_t, X_{t+h}) depends strongly on t , then we have a similar problem in estimating e.g., $\text{Cov}(X_t, X_{t+h})$.

DEFINITION 1.3.1: Strictly stationary

We say that a time series $\{X_t\}_{t \in \mathbb{Z}}$ is **strictly stationary (strongly stationary)** if for each $k \geq 1$, $i_1, \dots, i_k, h \in \mathbb{Z}$,

$$(X_{i_1}, \dots, X_{i_k}) \equiv (X_{i_1+h}, \dots, X_{i_k+h})$$

If we look at the k -dimensional joint distribution $(X_{i_1}, \dots, X_{i_k})$ of the series at points i_1, \dots, i_k , then **strict stationary means this is shift-invariant**. That is, shifting the window on which you view the data, does not change its distribution. This implies that if $F_t = \text{CDF of } X_t$, then $F_t = F_{t+h} = F$; that is, all variables have a common distribution function.

**DEFINITION 1.3.2: Mean function**

For a time series $\{X_t\}_{t \in \mathbf{Z}}$, with $\mathbb{E}[X_t^2] < \infty$ for all $t \in \mathbf{Z}$, we denote the **mean function** of the time series as

$$\mu_t = \mathbb{E}[X_t]$$

DEFINITION 1.3.3: Autocovariance function

The **autocovariance** function of the time series $\{X_t\}_{t \in \mathbf{Z}}$ is defined as

$$\gamma(t, s) = \mathbb{E}[(X_t - \mu_t)(X_s - \mu_s)] = \text{Cov}(X_t, X_s)$$

DEFINITION 1.3.4: Weakly stationary, Lag

We say that a time series $\{X_t\}_{t \in \mathbf{Z}}$ is **weakly stationary** if $\mathbb{E}[X_t] = \mu$ which does not depend on t , and if

$$\gamma(t, s) = f(|t - s|)$$

that is, $\gamma(t, s)$ is a function of $|t - s|$. In this case, we usually write

$$\gamma(h) = \text{Cov}(X_t, X_{t+h})$$

where we call the input h the **lag** parameter.

Additional Terminology

- The property when $\mathbb{E}[X_t] = \mu$ which does not depend on t is often called **first order stationary**.
- The property when $\gamma(t, s) = f(|t - s|)$ only depends on the lag $|t - s|$ is called **second order stationary**.
- For a second order stationary process,

$$\begin{aligned} \gamma(h) &= \text{Cov}(X_t, X_{t+h}) \\ &= \text{Cov}(X_{t-h}, X_{t-h+h}) & t \rightarrow (t-h) \\ &= \text{Cov}(X_t, X_{t-h}) \\ &= \gamma(-h) \end{aligned}$$

Since $\gamma(h) = \gamma(-h)$, we normally only record $\gamma(h)$ for $h \geq 1$.

1.4 White Noise and Stationary Examples

DEFINITION 1.4.1: Strong white noise

We say $\{X_t\}_{t \in \mathbb{Z}}$ is a **strong white noise** if $\mathbb{E}[X_t] = 0$ and the $\{X_t\}_{t \in \mathbb{Z}}$ are i.i.d.

DEFINITION 1.4.2: Weak white noise

We say $\{X_t\}_{t \in \mathbb{Z}}$ is a **weak white noise** if $\mathbb{E}[X_t] = 0$ and

$$\gamma(t, s) = \text{Cov}(X_t, X_s) = \begin{cases} \sigma^2 & |t - s| = 0 \\ 0 & |t - s| > 0 \end{cases}$$

DEFINITION 1.4.3: Gaussian white noise

We say $\{X_t\}_{t \in \mathbb{Z}}$ is a **Gaussian white noise** if $X_t \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma^2)$.



Figure 1.4: Gaussian White Noise of Length 500

```
# Figure 1.4
plot.ts(rnorm(500), main = "Gaussian White Noise", ylab = "w")
```

Figure 1.4 is a Gaussian *white* noise series. **White** comes from spectral analysis, in which a white noise series shares the same spectral properties as white light: all periodicities occur with equal strength.

EXAMPLE 1.4.4

Suppose $\{W_t\}_{t \in \mathbb{Z}}$ is a strong white noise, then $\mathbb{E}[W_t] = 0$; that is, the mean of W_t doesn't depend on t .

$$\gamma(t, s) = \text{Cov}(W_t, W_s) = \mathbb{E}[W_t W_s] = \begin{cases} \sigma_W^2 & |t - s| = 0 \\ 0 & |t - s| > 0 \end{cases}$$

$\gamma(t, s)$ only depends on $|t - s|$. Therefore, $\{W_t\}_{t \in \mathbb{Z}}$ is **weakly stationary**. Furthermore, we claim that

$\{W_t\}_{t \in \mathbf{Z}}$ is **strictly stationary**. Let $k \geq 1$, $i_1, \dots, i_k, h \in \mathbf{Z}$ with $i_1 < \dots < i_k$, then

$$\begin{aligned} \mathbb{P}(W_{i_1} \leq t_1, \dots, W_{i_k} \leq t_k) &= \prod_{j=1}^k \mathbb{P}(W_{i_j} \leq t_j) && \text{independence} \\ &= \prod_{j=1}^k \mathbb{P}(W_{i_j+h} \leq t_j) \\ &= \mathbb{P}(W_{i_1+h} \leq t_1, \dots, W_{i_k+h} \leq t_k) \end{aligned}$$

EXAMPLE 1.4.5

Suppose $\{W_t\}_{t \in \mathbf{Z}}$ is a strong white noise. Define $X_t = W_t + \theta W_{t-1}$ for $\theta \in \mathbf{R}$. Since $\{W_t\}_{t \in \mathbf{Z}}$ is a strong white noise, we have $\mathbb{E}[W_t] = 0$ for all t , hence we have $\mathbb{E}[X_t] = \mathbb{E}[W_t + \theta W_{t-1}] = \mathbb{E}[W_t] + \theta \mathbb{E}[W_{t-1}] = 0$ which is first order stationary.

$$\gamma(t, s) = \text{Cov}(X_t, X_s) = \begin{cases} (1 + \theta^2)\sigma_W^2 & |t - s| = 0 \\ \theta\sigma_W^2 & |t - s| = 1 \\ 0 & |t - s| > 1 \end{cases}$$

We obtain these calculations as follows:

- $|t - s| = 0$.

$$\mathbb{E}[(W_t + \theta W_{t-1})^2] = \mathbb{E}[W_t^2] + \theta^2 \mathbb{E}[W_{t-1}^2] + 2\mathbb{E}[\theta W_t W_{t-1}] = (1 + \theta^2)\sigma_W^2$$

since W_t is independent of W_{t-1} . The calculation is easy to verify.

- $t = s + 1$ (or $s = t + 1$).

$$\mathbb{E}[(W_{s+1} + \theta W_s)(W_s + \theta W_{s-1})] = \theta \mathbb{E}[W_s^2] = \theta \sigma_W^2$$

since W_{s+1} is independent of W_s and W_{s-1} . The calculation is easy to verify.

- $|t - s| > 1$. $W_t + \theta W_{t-1}$ is independent of $W_s + \theta W_{s-1}$.

We claim that $\{X_t\}_{t \in \mathbf{Z}}$ is also strictly stationary. Let $k \geq 1$, $i_1, \dots, i_k, h \in \mathbf{Z}$ with $i_1 < \dots < i_k$, then

$$\begin{aligned} \mathbb{P}(X_{i_1} \leq t_1, \dots, X_{i_k} \leq t_k) &= \mathbb{P}(W_{i_1} + \theta W_{i_1-1} \leq t_1, \dots, W_{i_k} + \theta W_{i_k-1} \leq t_k) \\ &= \mathbb{P}\left(\begin{bmatrix} W_{i_1-1} \\ W_{i_1} \\ \vdots \\ W_{i_k} \end{bmatrix} \in \mathcal{B}\right) \\ &= \mathbb{P}\left(\begin{bmatrix} W_{i_1-1+h} \\ \vdots \\ W_{i_k+h} \end{bmatrix} \in \mathcal{B}\right) \\ &= \mathbb{P}(X_{i_1+h} \leq t_1, \dots, X_{i_k+h} \leq t_k) \end{aligned}$$

where \mathcal{B} is some subset of $\mathbf{R}^{i_k - i_1 + 1}$, and hence is shift-invariant.

DEFINITION 1.4.6: Bernoulli shift

Suppose $\{\varepsilon_t\}_{t \in \mathbf{Z}}$ is a strong white noise. If $X_t = g(\varepsilon_t, \varepsilon_{t-1}, \dots)$ for some function $g : \mathbf{R}^\infty \rightarrow \mathbf{R}$, we say that $\{X_t\}_{t \in \mathbf{Z}}$ is a **Bernoulli shift**.

REMARK 1.4.7

We can also make a more general definition for a Bernoulli shift. Suppose $\{\varepsilon_t\}_{t \in \mathbb{Z}}$ is a strong white noise. If $X_t = g(\dots, \varepsilon_{t-1}, \varepsilon_t, \varepsilon_{t+1}, \dots)$ for some function $g : \mathbb{R}^{\mathbb{Z}} \rightarrow \mathbb{R}$, we say that $\{X_t\}_{t \in \mathbb{Z}}$ is a **Bernoulli shift**.

THEOREM 1.4.8

If $\{X_t\}_{t \in \mathbb{Z}}$ is a Bernoulli shift, then $\{X_t\}_{t \in \mathbb{Z}}$ is strictly stationary.

REMARK 1.4.9

Norbert Wiener conjectured that **every** stationary sequence is a Bernoulli shift, which is not true. The truth is, almost every one is.

EXERCISE 1.4.10

Suppose $\{W_t\}_{t \in \mathbb{Z}}$ is a strong white noise. The **two-sided random walk** is defined as

$$X_t = \sum_{i=0}^t W_i + \sum_{i=t}^{-1} W_i$$

Show that $\{X_t\}_{t \in \mathbb{Z}}$ is first order stationary, but $\{X_t\}_{t \in \mathbb{Z}}$ is not second order stationary.

Solution. $\{X_t\}_{t \in \mathbb{Z}}$ is first order stationary since

$$\begin{aligned} \mathbb{E}[X_t] &= \mathbb{E}\left[\sum_{i=0}^t W_i + \sum_{i=t}^{-1} W_i\right] \\ &= \mathbb{E}[W_0 + W_1 + \dots + W_{t-1} + W_t + W_t + W_{t-1} + \dots + W_0 + W_{-1}] \\ &= \mathbb{E}[W_{-1}] + \mathbb{E}[2W_0] + \mathbb{E}[2W_1] + \dots + \mathbb{E}[2W_{t-1}] \\ &= 0 + 2(0) + \dots + 2(0) \\ &= 0 \end{aligned}$$

since $\{W_t\}_{t \in \mathbb{Z}}$ is a strong white noise; that is, $\mathbb{E}[W_t] = 0$ for all t .

$\{X_t\}_{t \in \mathbb{Z}}$ is not second order stationary since if $t > 0$ the second sum is simply $\sum_{i=t}^{-1} W_i = 0$, and we have

$$\begin{aligned} \mathbb{E}[(X_t - \mu_t)(X_t - \mu_t)] &= \mathbb{E}[X_t^2] \\ &= \mathbb{E}\left[\left(\sum_{i=0}^t W_i\right)^2\right] \\ &= \mathbb{E}[W_0^2] + \dots + \mathbb{E}[W_t^2] && \text{since } W_i \perp\!\!\!\perp W_j \text{ for } i \neq j \\ &= t\sigma_W^2 \end{aligned}$$

which depends on t .

1.5 Weak versus Strong Stationary

Sadly, $\{X_t\}_{t \in \mathbb{Z}}$ is strictly stationary does not imply $\{X_t\}_{t \in \mathbb{Z}}$ is weakly stationary.

EXAMPLE 1.5.1

Suppose $X_t \stackrel{\text{iid}}{\sim}$ Cauchy Random Variables; that is,

$$\mathbb{P}(X_t \leq s) = \int_{-\infty}^s \frac{1}{\pi(1+x^2)} dx$$

Then, $\mathbb{E}[X_t]$ does not exist, and hence $\{X_t\}_{t \in \mathbb{Z}}$ cannot be weakly stationary. However, $\{X_t\}_{t \in \mathbb{Z}}$ is strictly stationary in this case since $\{X_t\}_{t \in \mathbb{Z}}$ is a strong white noise.

THEOREM 1.5.2

If $\{X_t\}_{t \in \mathbb{Z}}$ is strongly stationary and $\mathbb{E}[X_0^2] < \infty$, then $\{X_t\}_{t \in \mathbb{Z}}$ is weakly stationary.

Proof of: Theorem 1.5.2

Note that if $\{X_t\}_{t \in \mathbb{Z}}$ is strictly stationary, then

$$(X_t) \equiv (X_0)$$

so that $\mathbb{E}[X_t] = \mathbb{E}[X_0] = \mu$ which does not depend on t , and also

$$\mathbb{V}(X_t) = \mathbb{V}(X_0)$$

By the Cauchy-Schwarz inequality,

$$\begin{aligned} \gamma(t, s) &= \text{Cov}(X_t, X_s) \\ &= \mathbb{E}[(X_s - \mu)(X_t - \mu)] \\ &\leq \left\{ \mathbb{E}[(X_s - \mu)^2] \right\}^{1/2} \left\{ \mathbb{E}[(X_t - \mu)^2] \right\}^{1/2} \\ &= \sqrt{\mathbb{V}(X_s)} \sqrt{\mathbb{V}(X_t)} \\ &= \mathbb{V}(X_t) < \infty \end{aligned}$$

If $t < s$, then

$$\text{Cov}(X_t, X_s) = \text{Cov}(X_0, X_{s-t}) = f(|s-t|)$$

since it is shift-invariant, and hence if we shift everything over by t ,

$$(X_t, X_s) \equiv (X_{t-t}, X_{s-t}) \equiv (X_0, X_{s-t})$$

DEFINITION 1.5.3: Gaussian process

$\{X_t\}_{t \in \mathbb{Z}}$ is said to be a **Gaussian process (Gaussian time series)** if for each $k \in \mathbb{Z}_{\geq 1}$, $i_1 < i_2 < \dots < i_k$ we have

$$(X_{i_1}, \dots, X_{i_k}) \sim \text{MVN}(\boldsymbol{\mu}_k(i_1, \dots, i_k), \Sigma_{k \times k}(i_1, \dots, i_k))$$

$$\boldsymbol{\mu}_k = \begin{bmatrix} \mathbb{E}[X_{i_1}] \\ \vdots \\ \mathbb{E}[X_{i_k}] \end{bmatrix} \quad \Sigma_{k \times k} = \text{Cov}(X_{i_j}, X_{i_r})_{1 \leq j, r \leq k}$$

THEOREM 1.5.4

If $\{X_t\}_{t \in \mathbb{Z}}$ is weakly stationary and is a Gaussian process, then $\{X_t\}_{t \in \mathbb{Z}}$ is strictly stationary.

Proof of: Theorem 1.5.4

If $\{X_t\}_{t \in \mathbb{Z}}$ is weakly stationary, then $\mathbb{E}[X_t] = \mu$ for all t .

$$(X_{i_1}, \dots, X_{i_k}) \rightarrow \begin{bmatrix} \mathbb{E}[X_{i_1}] \\ \vdots \\ \mathbb{E}[X_{i_k}] \end{bmatrix} = \begin{bmatrix} \mu \\ \vdots \\ \mu \end{bmatrix} = \boldsymbol{\mu} = \begin{bmatrix} \mathbb{E}[X_{i_1+h}] \\ \vdots \\ \mathbb{E}[X_{i_k+h}] \end{bmatrix}$$

Also,

$$\begin{aligned} \mathbb{V}(X_{i_1}, \dots, X_{i_k}) &= \text{Cov}(X_{i_j}, X_{i_r})_{1 \leq j, r \leq k} \\ &= \text{Cov}(X_0, X_{i_r - i_j})_{1 \leq j, r \leq k} \\ &= \text{Cov}(X_0, X_{i_r + h - (i_j + h)})_{1 \leq j, r \leq k} \\ &= \text{Cov}(X_{i_j + h}, X_{i_r + h})_{1 \leq j, r \leq k} \\ &= \mathbb{V}(X_{i_1+h}, \dots, X_{i_k+h}) \end{aligned}$$

Using the Gaussian assumption

$$(X_{i_1}, \dots, X_{i_k}) \equiv \text{MVN}(\boldsymbol{\mu}, \boldsymbol{\Sigma}_{k \times k}) \equiv (X_{i_1+h}, \dots, X_{i_k+h})$$

Hence $\{X_t\}_{t \in \mathbb{Z}}$ is strictly stationary in this case.

EXERCISE 1.5.5

Prove that if $\{X_t\}_{t \in \mathbb{Z}}$ is not weakly stationary; that is, either $\mathbb{E}[X_t]$ depends on t or $\gamma(t, s)$ does not depend on the lag, and has a finite mean and variance, then $\{X_t\}_{t \in \mathbb{Z}}$ is not strictly stationary.

1.6 † Theoretical L2 Framework for Time Series

- $X_t = \lim_{h \rightarrow \infty} X_{h,t}$. In what sense does this limit exist?
- How “close” are two random variables X and Y ?
- Is there a random variable that achieves

$$\inf_{y \in S} d(Y, S)$$

DEFINITION 1.6.1: L^2 space

Consider a probability space $(\Omega, \mathcal{F}, \mathbb{P})$. The space L^2 is the set of random variables $X : \Omega \rightarrow \mathbb{R}$ measurable so that $\mathbb{E}[X^2] < \infty$.

DEFINITION 1.6.2: L^2 -time series

We say that $\{X_t\}_{t \in \mathbb{Z}}$ is an L^2 -time series if $X_t \in L^2$ for all $t \in \mathbb{Z}$.

L^2 is a Hilbert space when equipped with inner product, $X, Y \in L^2$.

$$\langle X, Y \rangle = \mathbb{E}[XY]$$

$\langle \cdot, \cdot \rangle$ is an inner product since it is

(1) Linear: $\langle aX + bY, Z \rangle = a\langle X, Z \rangle + b\langle Y, Z \rangle$.

(2) “Almost” Positive Definite: $\langle X, X \rangle = \mathbb{E}[X^2] = 0 \iff X = 0$ almost surely; that is, $\mathbb{P}(X = 0) = 1$.

(3) Symmetric: $\langle X, Y \rangle = \langle Y, X \rangle$.

L^2 is complete with this inner product; that is, whenever $X_n \in L^2$ so that $\mathbb{E}[(X_n - X_m)^2] \rightarrow 0$ as $n, m \rightarrow \infty$, then there exists $X \in L^2$ so that $X_n \rightarrow X$; that is, $\mathbb{E}[(X_n - X)^2] \rightarrow 0$. This follows from the “famous” Riesz-Fischer Theorem.

Useful Tools for Time Series

(1) **Existence of Limits**

$$X_{t,n} = \sum_{j=0}^n \psi_j \varepsilon_{t-j}$$

$\{\varepsilon_t\}_{t \in \mathbb{Z}}$ is a strong white noise. Since for $n > m$,

$$\mathbb{E}[(X_{t,n} - X_{t,m})^2] = \mathbb{E}\left[\left(\sum_{j=m+1}^n \psi_j \varepsilon_{t-j}\right)^2\right] = \sum_{j=m+1}^n \psi_j^2 \sigma_\varepsilon^2 \rightarrow 0 \text{ as } n, m \rightarrow \infty$$

only if $\sum_{j=0}^{\infty} \psi_j^2 < \infty$, then there **must** exist a random variable X_t (by the completeness of L^2), so that

$$X_t = \lim_{n \rightarrow \infty} X_{t,n} = \sum_{j=0}^{\infty} \psi_j \varepsilon_{t-j}$$

(2) **Projection Theorem and Forecasting.** Forecasting can be often cast as finding a random variable Y among a collection of possible forecasts \mathcal{M} (e.g., $\mathcal{M} = \text{Span}(X_T, \dots, X_1)$) so that

$$Y = \arg \inf_{Z \in \mathcal{M}} \mathbb{E}[(X_{T+h} - Z)^2]$$

When \mathcal{M} is a closed linear subspace of L^2 , the Projection Theorem guarantees that such a Y exists, and it must satisfy

$$\langle X_{T+h} - Y, Z \rangle = 0 \quad \forall Z \in \mathcal{M}$$

must be in the orthogonal complement.

1.7 Signal and Noise Models

“Ideally,” a time series that we are considering was generated from a stationary process. If so, we can pool data to estimate the processes underlying structure (e.g., its marginal distribution, and serial dependence structure).

Most time series are evidently not stationary.

Looking back at Figure 1.1:

- Mean appears to increase, so it is not first order stationary;
- Variability also appears to increase, so it is not second order stationary;
- Therefore, it is not strictly stationary.

Signal and Noise Model: $X_t = s_t + \varepsilon_t$

- s_t is the **deterministic** “signal” or “trend” of the series.
- ε_t is the “noise” added to the signal satisfying $\mathbb{E}[\varepsilon_t] = 0$, hence $\mathbb{E}[X_t] = \mathbb{E}[s_t + \varepsilon_t] = \mathbb{E}[s_t]$. There exists a (strong) white noise $\{W_t\}_{t \in \mathbb{Z}}$ so that

$$\varepsilon_t = g(W_t, W_{t-1}, \dots) \quad [\text{Stationary Noise}]$$

$$\varepsilon_t = g_t(W_t, W_{t-1}, \dots) \quad [\text{Non-stationary Noise}]$$

The terms $\{W_t\}_{t \in \mathbb{Z}}$ are often called the “innovations” or “shocks” driving the random behaviour of X_t .
 g is used to try to capture noise that can potentially have serial dependence.

EXAMPLE 1.7.1

An example of a function g so that $\varepsilon_t = g_t(W_t, W_{t-1}, \dots)$ might be a **random walk**; that is, $\varepsilon_t = \sum_{j=0}^t W_j$. Another example could be the **changing variance models**; that is, $\varepsilon_t = \sigma(t)W_t$.

Our goal is to estimate s_t , and then infer the structure of ε_t .

In Figure 1.2, the model appears to be non-stationary (trending upwards over time), so we might try the signal and noise model. We might posit a linear trend, or even higher order functions.

For the temperature data, we may posit that

$$s_t = \beta_0 + \beta_1 t \quad [\text{Linear Trend}]$$

The trend may be estimated by ordinary least squares (OLS). We choose β_0 and β_1 to minimize

$$\sum_{t=1}^T [X_t - (\beta_0 + \beta_1 t)]^2$$

This can be done in R using the `lm()` command, and can easily be computed with calculus. Figure 1.5 is a small example of the global temperature data superimposed with the `lm()` estimate.

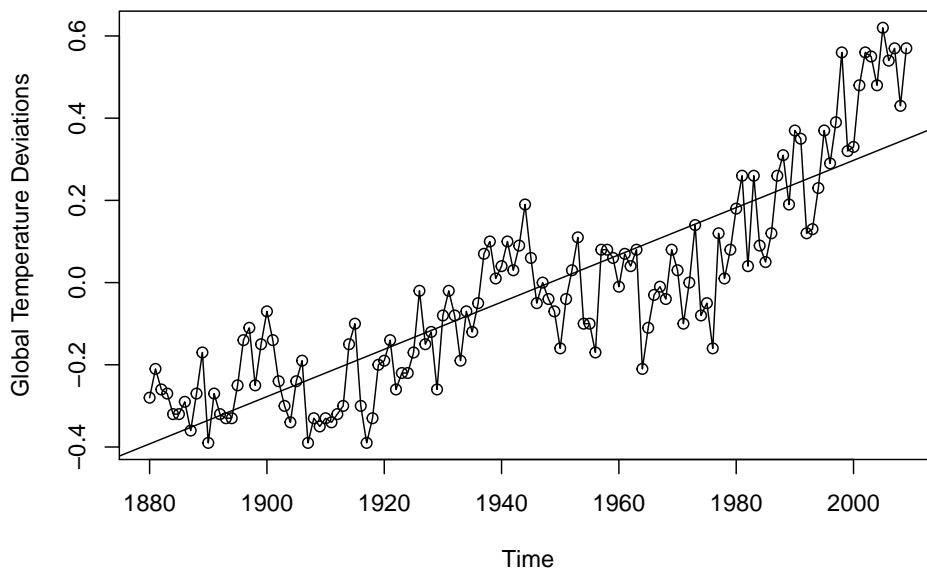


Figure 1.5: OLS estimate of linear trend

```
# Figure 1.5
fit <- lm(gtemp ~ time(gtemp), na.action = NULL)
plot.ts(gtemp, type = "o", ylab = "Global Temperature Deviations")
abline(fit)
```

Let's introduce some terminology about trends.

DEFINITION 1.7.2: Detrended time series

Detrending a time series constitutes computing the residuals based on an estimate for the signal/trend. A **detrended time series** is a time series of such residuals.

1. Estimate $s_t \rightarrow \hat{s}_t$.
2. Detrend series: $X_t - \hat{s}_t = Y_t$ where Y_t is the “detrended” series.



Figure 1.6: Residuals of OLS fit.

```
# Figure 1.6
plot(resid(fit), type = "o", main = "detrended")
```

In Figure 1.6: If trend is now zero, there appears to be a substantial serial dependence remaining in the time series.

1.8 Time Series Differencing

Signal and Noise Model: $X_t = s_t + \varepsilon_t$. Hopefully, upon estimating s_t with \hat{s}_t , we find $X_t - \hat{s}_t = \hat{\varepsilon}_t$ (detrended series) which looks reasonably stationary. If the residuals were reasonably stationary, we might proceed in estimating their underlying structure of $\{\hat{\varepsilon}_t\}_{t=1, \dots, T}$ as if it were stationary. In particular, we might try to estimate their marginal distributions and/or their serial dependence structure. If we thought those estimates were reasonably good, we would have a good idea of how the time series X_t behaves.

Random Walk with Drift Model. Let ε_t be a strong white noise.

$$\begin{aligned}
 X_t &= \delta + X_{t-1} + \varepsilon_t \\
 &= \delta + \delta + X_{t-2} + \varepsilon_{t-1} + \varepsilon_t \\
 &= \delta + \delta + \delta + X_{t-3} + \varepsilon_{t-2} + \varepsilon_{t-1} + \varepsilon_t \\
 &\vdots \\
 &= t\delta + X_0 + \sum_{j=1}^t \varepsilon_j
 \end{aligned}
 \qquad t \text{ times}$$

where we note that $t\delta + X_0 = s_t$ is a linear signal, and $\sum_{j=1}^t \varepsilon_j$ is a random walk noise.

Notice that under the Random Walk Model.

$$X_t - X_{t-1} = \nabla X_t = \delta + \varepsilon_t$$

So, if X_t follows a random walk model, the series $Y_t = \nabla X_t$ should behave like a white noise shifted by δ .



Figure 1.7: First differenced series. Average of first differenced series is $\hat{\delta} \approx 0.0066$

```
# Figure 1.7
plot(diff(gtemp), type = "o", main = "first difference")
```

In Figure 1.7: To see what this looks like in this temperature example, here is a plot of $\nabla X_t = X_t - X_{t-1}$ for Figure 1.2. As you can see if you look at this compared to the detrended series using linear trend, I would say this series looks much more like a white noise (there does not appear to be any discernible patterns in this first difference). If you calculate the mean of this first difference series, that would be an estimator for the drift term in the random walk model which here is ≈ 0.0066 .

DEFINITION 1.8.1: Differenced time series

Differencing a time series constitutes computing the difference between successive terms. A **differenced time series** is a time series of such differences. The first differenced series is denoted

$$\nabla X_t = X_t - X_{t-1}$$

and is the series of length $T - 1$, namely

$$X_2 - X_1, X_3 - X_2, \dots, X_T - X_{T-1}$$

Higher order differences are calculated recursively, so

$$\nabla^d X_t = \nabla^{d-1} \nabla X_t$$

where ∇^d is the d^{th} order difference, and we define $\nabla^0 X_t = X_t$.

Detrending and Differencing are both ways of reducing a (potentially non-stationary) time series to an approximately stationary series.

Differencing vs. Detrending

Pros:

- Differencing does not require the parameter estimation (don't need to estimate s_t).
- Higher order differencing can reduce even very “trendy” series to look more like noise.

Cons:

- Differencing can “wash away” features of the series, and introduce more complicated structures.
- The trend is often of interest, and good estimates of the trend lead to improved long-range forecasts.

EXAMPLE 1.8.2: Potentially Complicating Series with Differencing

$X_t = W_t$ where W_t is a strong white noise.

$$\nabla X_t = W_t - W_{t-1} = Y_t$$

$$\gamma_X(h) = \text{Cov}(X_t, X_{t+h}) = \begin{cases} \sigma_W^2 & h = 0 \\ 0 & h \geq 1 \end{cases}$$

More complicated:

$$\gamma_Y(h) = \text{Cov}(Y_t, Y_{t+h}) = \begin{cases} 2\sigma_W^2 & h = 0 \\ -\sigma_W^2 & h = 1 \\ 0 & h \geq 2 \end{cases}$$



Figure 1.8: First Difference and White Noise

```
# Figure 1.8
par(mfrow = c(2, 1))
```

```
plot(diff(gtemp), main = "first difference Temp data")
plot(rnorm(gtemp),
     type = "l",
     main = "white noise",
     ylab = "w")
```

In Figure 1.8: If these two series behave in the same way, then it stands to reason that

$$g(\varepsilon_t, \varepsilon_{t-1}, \dots) = \varepsilon_t \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma_{\text{temp}}^2)$$

Chapter 2

Week 2

2.1 Autocorrelation and Empirical Autocorrelation

Usually through either detrending or differencing, we arrive at a series $\{X_t\}_{t \in \mathbb{Z}}$ that we may consider as stationary.

Given such a series, we wish to estimate a function g , so that

$$X_t = g(W_t, W_{t-1}, \dots)$$

$\{W_t\}_{t \in \mathbb{Z}}$ is a “innovation” sequence (strong white noise) which could admit serial dependence, etc.

In a first pass, it’s reasonable to assume that g is a linear function.

DEFINITION 2.1.1: Linear process

A time series $\{X_t\}_{t \in \mathbb{Z}}$ is said to be a **linear process** if there exists a strong white noise $\{W_t\}_{t \in \mathbb{Z}}$ and coefficient $\{\psi_\ell\}_{\ell \in \mathbb{Z}}$ where $\psi_\ell \in \mathbb{R}$, so that

$$\sum_{\ell=-\infty}^{\infty} |\psi_\ell| < \infty$$

and

$$X_t = \sum_{\ell=-\infty}^{\infty} \psi_\ell W_{t-\ell}$$

Note that the sum defining X_t is well-defined as a limit in L^2 . Also, we must require that $\mathbb{V}(W_{t-\ell}) < \infty$.

DEFINITION 2.1.2: Causal linear process

We say $\{X_t\}_{t \in \mathbb{Z}}$ is a **causal linear process** if

$$X_t = \sum_{\ell=0}^{\infty} \psi_\ell W_{t-\ell}$$

Note that X_t only depends on W ’s in the “past.”

EXAMPLE 2.1.3

$X_t = W_t$ is a linear process, so all ψ ’s are 0, except for $\psi_0 = 1$ which is a strong white noise sequence.

REMARK 2.1.4

Linear processes are **strictly stationary** since they can be written as Bernoulli-shifts.

EXAMPLE 2.1.5

$X_t = W_t + \theta W_{t-1}$ where $\{W_t\}_{t \in \mathbb{Z}}$ is a strong white noise with finite variance. X_t is a linear process.

$$\gamma_X = \begin{cases} (1 + \theta^2)\sigma_W^2 & h = 0 \text{ always non-zero} \\ \theta\sigma_W^2 & h = 1 \\ 0 & h \geq 2 \end{cases}$$

$\gamma_X(h)$ non-zero for $h \geq 1$ only where “lagged” terms in the linear process are non-zero. Suggests a way of sleuthing out what

$$g(W_t, W_{t-1}, \dots) = \sum_{\ell=0}^{\infty} \psi_{\ell} W_{t-\ell}$$

must look like.

DEFINITION 2.1.6: Autocorrelation function

Suppose $\{X_t\}_{t \in \mathbb{Z}}$ is weakly stationary. The **autocorrelation function** (ACF) of $\{X_t\}_{t \in \mathbb{Z}}$ is

$$\rho_X(h) = \frac{\gamma(h)}{\gamma(0)} \quad (h \geq 0)$$

Note since $\gamma(0) = \mathbb{V}(X_t) = \mathbb{V}(X_0)$ (since the process is stationary),

$$|\gamma(h)| = |\text{Cov}(X_t, X_{t+h})| \stackrel{\text{CS}}{\leq} \sqrt{\mathbb{V}(X_t)\mathbb{V}(X_{t+h})} = \mathbb{V}(X_0)$$

Same # by stationarity

Hence, $|\rho(h)| \leq 1 \implies -1 \leq \rho(h) \leq 1$.

Estimating $\gamma(h)$ and $\rho(h)$

$$\gamma(h) = \text{Cov}(X_t, X_{t+h}) = \mathbb{E}[(X_t - \mu)(X_{t+h} - \mu)]$$

where $\mu = \mathbb{E}[X_t]$. Hence, a sensible estimator is

$$\hat{\mu} = \frac{1}{T} \sum_{t=1}^T X_t = \bar{X}$$

which is the **sample mean (time series average)**.

$$\hat{\gamma}(h) = \frac{1}{T} \sum_{t=1}^{T-h} (X_t - \bar{X})(X_{t+h} - \bar{X}) \approx \frac{1}{T-h} \sum_{t=1}^{T-h} (X_t - \bar{X})(X_{t+h} - \bar{X})$$

where $(X_t - \bar{X})(X_{t+h} - \bar{X})$ is the averaging over centred terms h -time steps apart.

$$\hat{\rho}(h) = \frac{\hat{\gamma}(h)}{\hat{\gamma}(0)}$$

EXAMPLE 2.1.7

$X_t = W_t$ where $\{W_t\}_{t \in \mathbb{Z}}$ is a strong white noise with $\mathbb{V}(W_t) = \sigma_W^2 < \infty$.

$$\gamma_X(h) = \begin{cases} \sigma_W^2 & h = 0 \\ 0 & h \geq 1 \end{cases}$$

Therefore,

$$\rho_X(h) = \begin{cases} 1 & h = 0 \\ 0 & h \geq 1 \end{cases}$$

Note that it's always the case that

$$\rho(0) = \frac{\gamma(0)}{\gamma(0)} = 1$$



Figure 2.1: ACF of white noise, sample length 130

```
# Figure 2.1
acf(rnorm(500))
```

In Figure 2.1: Let's then have a look at what the empirical autocorrelation function looks like when we apply it to a strong white noise sample. In this case, we are considering a strong Gaussian white noise with variance 1. This is what the sample ACF looks like. What we're plotting here is on the x -axis we have the lags h , and on the y -axis we have the magnitudes of the autocorrelation $\hat{\rho}(h)$. What we're seeing here is $\hat{\rho}(0) = 1$ (by definition). However, for lags other than zero, for the other autocorrelations plotted, we can see that they are relatively small compared to $\hat{\rho}(0) = 1$, which is the point of the blue lines (explained in the next lecture). The basic interpretation of blue lines is that if an autocorrelation would go inside the blue lines then you could imagine that it would be consistent with the series being a strong white noise, which is what we observe here. There are small violations that can occur by sheer chance.

2.2 Modes of Convergence of Random Variables

$\hat{\gamma}(h)$ is an estimator of $\gamma(h)$ when the data is stationary, and we want to discuss the asymptotic properties of this estimator.

Review/Introduce

- (1) Stochastic Boundedness (convergence of random variables): $\mathcal{O}(p)$ and $o(p)$
- (2) Convergence in Probability
- (3) Convergence in Distribution

DEFINITION 2.2.1: Bounded in probability

Suppose $\{X_n\}_{n \geq 1}$ is a sequence of random variables. We say that X_n is **bounded in probability** by Y_n if for all $\varepsilon > 0$, there exists real numbers M, N , so that for all $n \geq N$,

$$\mathbb{P}\left(\left|\frac{X_n}{Y_n}\right| > M\right) \leq \varepsilon$$

Notation: $X_n = \mathcal{O}_p(Y_n)$, and in English, we say “ X_n is on the order of Y_n .”

DEFINITION 2.2.2: Converges in probability

We say X_n **converges in probability** to X if for all $\varepsilon > 0$,

$$\lim_{n \rightarrow \infty} \mathbb{P}(|X_n - X| > \varepsilon) = 0$$

If a_n is a sequence of scalars, we abbreviate $\frac{X_n}{a_n}$ converges in probability to zero as

$$X_n = o_p(a_n) \iff \mathbb{P}\left(\left|\frac{X_n}{a_n}\right| > \varepsilon\right) \xrightarrow{n \rightarrow \infty} 0 \quad (\forall \varepsilon > 0)$$

Hence, X_n converges in probability to zero is denoted $X_n = o_p(1)$. Likewise, we also write $X_n \xrightarrow{p} X$ to denote X_n converges in probability to X .

DEFINITION 2.2.3: Converges in distribution

We say that the sequence of scalar random variables X_n with respective CDF's $F_n(x)$ **converges in distribution** to X with CDF $F(x)$ if for all continuity points of F ,

$$\lim_{n \rightarrow \infty} |F_n(y) - F(y)| = 0$$

REMARK 2.2.4

When $F(x)$ is the CDF of a continuous random variable (e.g., a normal CDF), then

$$\lim_{n \rightarrow \infty} |F_n(y) - F(y)| = 0 \quad (\forall y \in \mathbf{R})$$

THEOREM 2.2.5: Markov's Inequality

If $\mathbb{E}[Y^2] < \infty$, then

$$\mathbb{P}(|Y| \geq m) \leq \frac{\mathbb{E}[Y^2]}{m^2}$$

Proof of: Theorem 2.2.5

$$\begin{aligned} \mathbb{E}[Y^2] &= \mathbb{E}\left[Y^2 \mathbb{I}\{|Y| \geq m\} + Y^2 \mathbb{I}\{|Y| < m\}\right] \\ &= \mathbb{E}\left[Y^2 \mathbb{I}\{|Y| \geq m\}\right] + \mathbb{E}\left[Y^2 \mathbb{I}\{|Y| < m\}\right] \\ &\geq \mathbb{E}\left[Y^2 \mathbb{I}\{|Y| \geq m\}\right] \\ &\geq m^2 \mathbb{E}\left[\mathbb{I}\{|Y| \geq m\}\right] && \text{since } Y^2 \geq m^2 \\ &= m^2 \mathbb{P}(|Y| \geq m) \end{aligned}$$

REMARK 2.2.6: Generalization of Markov's Inequality

If $\mathbb{E}[Y^k] < \infty$, then

$$\mathbb{P}(|Y| \geq m) \leq \frac{\mathbb{E}[|Y|^k]}{m^k}$$

EXAMPLE 2.2.7

Suppose X_n is a strong white noise in L^2 ($\mathbb{E}[X_0^2] < \infty$), and let

$$\bar{X}_T = \frac{1}{T} \sum_{t=1}^T X_t$$

Then,

$$(1) |\bar{X}_T| = o_p(1).$$

$$\begin{aligned} \mathbb{V}(\bar{X}_T) &= \mathbb{E}[\bar{X}_T^2] \\ &= \frac{1}{T^2} \mathbb{E}\left[\left(\sum_{t=1}^T X_t\right)^2\right] \\ &= \frac{1}{T^2} \sum_{t=1}^T \sum_{s=1}^T \mathbb{E}[X_t X_s] \\ &= \frac{1}{T^2} \sum_{t=1}^T \mathbb{E}[X_t^2] \\ &= \frac{1}{T^2} \sum_{t=1}^T \mathbb{E}[X_0^2] \\ &= \frac{\sigma^2}{T} && \text{since } \sigma^2 = \mathbb{E}[X_0^2] \end{aligned}$$

Therefore, for $\varepsilon > 0$, by Markov's Inequality we have

$$\mathbb{P}(|\bar{X}_T| > \varepsilon) \leq \frac{\mathbb{E}[|\bar{X}_T|^2]}{\varepsilon^2} = \frac{\sigma^2/T}{\varepsilon^2} \xrightarrow{T \rightarrow \infty} 0$$

Hence, $|\bar{X}_T| \xrightarrow{p} 0$
 (2) $\bar{X}_T = \mathcal{O}_p(1/\sqrt{T})$, as before

$$\mathbb{V}\left(\frac{\bar{X}_T}{1/\sqrt{T}}\right) = \mathbb{V}(\sqrt{T}\bar{X}_T) = T\mathbb{V}(\bar{X}_T) = \sigma^2$$

So by Markov's Inequality, for $M > 0$

$$\mathbb{P}(|\sqrt{T}\bar{X}_T| > M) \leq \frac{\mathbb{V}(\sqrt{T}\bar{X}_T)}{M^2} = \frac{\sigma^2}{M^2} \xrightarrow{M \rightarrow \infty} 0$$

Hence $\sqrt{T}\bar{X}_T = \mathcal{O}_p(1) \implies \bar{X}_T = \mathcal{O}_p(1/\sqrt{T})$.

REMARK 2.2.8

Alternatively, we can show this using the CLT. By the CLT,

$$\sqrt{T}\bar{X}_T \xrightarrow{D} \mathcal{N}(0, \sigma^2)$$

Therefore, if $F_T \sim$ CDF of $\sqrt{T}\bar{X}_T$ and $\Phi \sim$ CDF of $\mathcal{N}(0, 1)$ random variable we have

$$\left|F_T(x) - \Phi\left(\frac{x}{\sigma}\right)\right| \xrightarrow{T \rightarrow \infty} 0 \quad (\forall x \in \mathbf{R})$$

For $\varepsilon > 0$, choose M such that

$$\Phi\left(-\frac{M}{\sigma}\right) = 1 - \Phi\left(\frac{M}{\sigma}\right) \leq \frac{\varepsilon}{4}$$

For this M , choose T_0 such that if $T \geq T_0$, then

$$\left|F_T(-M) - \Phi\left(-\frac{M}{\sigma}\right)\right| \leq \frac{\varepsilon}{4}$$

and

$$\left|F_T(M) - \Phi\left(\frac{M}{\sigma}\right)\right| \leq \frac{\varepsilon}{4}$$

Then,

$$\begin{aligned} \mathbb{P}(|\sqrt{T}\bar{X}_T| \geq M) &= F_T(-M) + (1 - F_T(M)) \\ &= \Phi\left(-\frac{M}{\sigma}\right) + \left[1 - \Phi\left(\frac{M}{\sigma}\right)\right] + F_T(-M) - \Phi\left(-\frac{M}{\sigma}\right) + \Phi\left(\frac{M}{\sigma}\right) - F_T(M) \\ &\leq \frac{\varepsilon}{4} + \frac{\varepsilon}{4} + \frac{\varepsilon}{4} + \frac{\varepsilon}{4} \\ &= \varepsilon \end{aligned}$$

REMARK 2.2.9

In general,

$$\frac{X_n}{a_n} \xrightarrow{D} \text{non-degenerate random variable} \implies X_n = \mathcal{O}_p(a_n)$$

REMARK 2.2.10: Algebra of \mathcal{O}_p and $o(p)$ notation

1. If $X_n = \mathcal{O}_p(a_n)$ and $Y_n = \mathcal{O}_p(b_n)$, then

$$X_n + Y_n = \mathcal{O}_p(\max(a_n, b_n))$$

2. If $X_n = o_p(1)$ and $Y_n = o_p(1)$, then

$$X_n + Y_n = o_p(1)$$

3. If $X_n = o_p(1)$ and $Y_n = o_p(1)$, then

$$X_n Y_n = o_p(1)$$

EXAMPLE 2.2.11

Suppose W_t is a strong white noise in L^2 with $\mathbb{E}[W_t^4] < \infty$. Let $X_t = W_t + \theta W_{t-1}$ for $\theta \in \mathbf{R}$. Show that

$$\hat{\gamma}(1) \xrightarrow{p} \theta \sigma_W^2$$

Solution.

$$\begin{aligned} \bar{X}_T &= \frac{1}{T} \sum_{t=1}^T X_t \\ &= \frac{1}{T} \sum_{t=1}^T (W_t + \theta W_{t-1}) \\ &= \frac{1}{T} \sum_{t=1}^T W_t + \frac{\theta}{T} \sum_{t=1}^T W_{t-1} \\ &= o_p(1) \end{aligned} \quad \text{by WLLN}$$

$$\begin{aligned} \hat{\gamma}(1) &= \frac{1}{T} \sum_{t=1}^{T-1} (X_t - \bar{X}_T)(X_{t+1} - \bar{X}_T) \\ &= \frac{1}{T} \sum_{t=1}^{T-1} [X_t X_{t+1} - X_t \bar{X}_T - \bar{X}_T X_{t+1} + (\bar{X}_T)^2] \\ &= \frac{1}{T} \sum_{t=1}^{T-1} X_t X_{t+1} - \frac{\bar{X}_T}{T} \sum_{t=1}^{T-1} X_t - \frac{\bar{X}_T}{T} \sum_{t=1}^{T-1} X_{t+1} + \frac{T-1}{T} (\bar{X}_T)^2 \\ &= \frac{1}{T} \sum_{t=1}^{T-1} X_t X_{t+1} + R_1 + R_2 + R_3 \end{aligned}$$

Notice that $R_i = o_p(1)$ for $i = 1, 2, 3$ since, for example, $\bar{X}_T = o_p(1)$ and $\sum_{t=1}^T X_t = o_p(1)$ so their product is $o_p(1)$; so we only need to focus on the first term.

$$\begin{aligned} \frac{1}{T} \sum_{t=1}^{T-1} X_t X_{t+1} &= \frac{1}{T} \sum_{t=1}^{T-1} (W_t + \theta W_{t-1})(W_{t+1} + \theta W_t) \\ &= \frac{1}{T} \sum_{t=1}^{T-1} \theta W_t^2 + G_1 + G_2 + G_3 \end{aligned}$$

Now,

$$\frac{1}{T} \sum_{t=1}^{T-1} \theta W_t^2 \xrightarrow{\text{a.s.}} \theta \mathbb{E}[W_t^2] = \theta \sigma_W^2$$

by strong law of large numbers. We now wish to calculate the variance of

$$\begin{aligned}
 G_1 &= \frac{1}{T} \sum_{t=1}^{T-1} W_t W_{t+1}. \\
 \mathbb{E}[G_1] &= \frac{1}{T} \sum_{t=1}^{T-1} \mathbb{E}[W_t W_{t+1}] = 0 \\
 \mathbb{V}(G_1) &= \mathbb{E}[G_1^2] \\
 &= \frac{1}{T^2} \sum_{t=1}^{T-1} \sum_{s=1}^{T-1} \underbrace{\mathbb{E}[W_t W_{t+1} W_s W_{s+1}]}_{\neq 0 \Leftrightarrow s=t} \\
 &= \frac{1}{T^2} \sum_{t=1}^{T-1} \mathbb{E}[W_t^2 W_{t+1}^2] \\
 &= \frac{T-1}{T^2} \sigma_W^4 \xrightarrow{T \rightarrow \infty} 0
 \end{aligned}$$

By Markov's Inequality: $G_1 = o_p(1)$, and similarly, for G_2 and G_3 .

2.3 † M-dependent CLT

Suppose X_t is a mean zero strictly stationary time series with $\mathbb{E}[X_t^2] < \infty$. We are frequently faced with the problems:

- (1) What is the approximate distribution of

$$\frac{1}{\sqrt{T}} \sum_{t=1}^T X_t = \sqrt{T} \bar{X}_T \stackrel{D}{\approx} \mathcal{N}(0, \sigma_X^2)$$

- (2) If X_t is a strong white noise, what the approximate distribution of

$$\hat{\gamma}(h) = \frac{1}{T} \sum_{t=1}^{T-h} \underbrace{X_t X_{t+h}}_{\text{not iid}} + o_p(1)$$

$X_t X_{t+h} = Y_t$ is strictly stationary.

- Only way to understand how $\{X_t\}_{t \in \mathbb{Z}}$ behaves, we have to observe replicates of the process.
- If process is suitably “weakly dependent,” then we can observe replicates of the process by viewing in on overlapping windows.

DEFINITION 2.3.1: m -dependent

We say a time series $\{X_t\}_{t \in \mathbb{Z}}$ is **m -dependent** for a positive integer m , if for all

$$t_1 < t_2 < \dots < t_{d_1} < s_1 < s_2 < \dots < s_{d_2} \in \mathbb{Z}$$

so that $t_{d_1+m} \leq s_1$, then

$$(X_{t_1}, \dots, X_{t_{d_1}})$$

is **independent of**

$$(X_{s_1}, \dots, X_{s_{d_2}})$$

EXAMPLE 2.3.2

$X_t = W_t + \theta W_{t-1}$ for $\theta \in \mathbf{R}$ where W_t is a strong white noise is 2-dependent.

THEOREM 2.3.3: Generalization of the standard CLT to m -dependent

Suppose X_t is a strictly stationary and m -dependent time series for $m \in \mathbf{Z}_{>0}$ with $\mathbb{E}[X_t] = 0$ and $\mathbb{E}[X_t^2] < \infty$, then if

$$S_T = \frac{1}{\sqrt{T}} \sum_{t=1}^T X_t = \sqrt{T} \bar{X}_T \xrightarrow[T \rightarrow \infty]{D} \mathcal{N}(0, \sigma_m^2)$$

where

$$\sigma_m^2 = \sum_{h=-m}^m \gamma(h) = \gamma(0) + 2 \sum_{h=1}^m \gamma(h)$$

Note that σ_m^2 is just the variance of S_T and can be easily calculated.

DEFINITION 2.3.4: Triangular array

We say $\{X_{i,j}, 1 \leq j \leq n_i, 1 \leq i < \infty\}$ forms a **triangular array** of mean zero L^2 random variables, if $\mathbb{E}[X_{i,j}] = 0$, $\mathbb{E}[X_{i,j}^2] < \infty$, and for each i -fixed we have $X_{i,1}, \dots, X_{i,n_i}$ are independent with $n_i < n_{i+1}$.

Visually, row-wise random variables are independent:

$$\begin{array}{cccc} X_{1,1} & \cdots & X_{1,n_1} & \\ X_{2,1} & \cdots & \cdots & X_{2,n_2} \\ \vdots & \ddots & \ddots & \ddots \end{array}$$

THEOREM 2.3.5: Linderberg-Feller CLT for Triangular Arrays

Let $\{X_{i,j}, 1 \leq j \leq n_i, 1 \leq i < \infty\}$ be a triangular array of mean zero L^2 random variables. Define

$$\sigma_i^2 = \sum_{j=1}^{n_i} \mathbb{V}(X_{i,j})$$

and

$$S_i = \frac{1}{\sigma_i} \sum_{j=1}^{n_i} X_{i,j}$$

If for $\varepsilon > 0$,

$$\frac{1}{\sigma_i^2} \sum_{j=1}^{n_i} \mathbb{E} \left[X_{i,j}^2 \mathbb{I}\{|X_{i,j}| > \varepsilon \sigma_i\} \right] \xrightarrow{i \rightarrow \infty} 0$$

then

$$S_i \xrightarrow{D} \mathcal{N}(0, 1)$$

Proof of: Theorem 2.3.3

Bernstein Blocking Argument: we take a given time series of length T .

Let a_T = big block size and m = little block size. We assume $a_T \rightarrow \infty$ as $T \rightarrow \infty$, but $\frac{a_T}{T} \rightarrow 0$. Then,

$$N = \text{number of blocks} = \left\lfloor \frac{T}{M + a_T} \right\rfloor$$

$$B_j = \{i : (j-1)(a_T + m) + 1 \leq i \leq ja_T + (j-1)m\}$$

$$b_j = \{i : ja_T + (j-1)m + 1 \leq i \leq j(a_T + m)\}$$

Since a_T is increasing up to infinity, for T sufficiently large, $a_T > m$, and so by m -dependence,

$$\sum_{t \in B_j} X_t$$

is independent of

$$\sum_{t \in B_k} X_t \quad (j \neq k)$$

similarly for $B_j, B_k \rightarrow b_j, b_k$.

$$\frac{1}{\sqrt{T}} \sum_{t=1}^T X_t = \frac{1}{\sqrt{T}} \sum_{j=1}^N \sum_{t \in B_j} X_t + \underbrace{\frac{1}{\sqrt{T}} \sum_{j=1}^N \sum_{t \in b_j} X_t}_{\text{iid}} + \text{Remainder} = G_1 + G_2 + G_3$$

We want to show the big blocks dominate.

$$\mathbb{V}(G_2) = \frac{1}{T} \sum_{j=1}^N \mathbb{E} \left[\left(\sum_{t \in b_j} X_t \right)^2 \right] = \frac{N}{T} \mathbb{E} \left[\left(\sum_{t=1}^m X_t \right)^2 \right]$$

due to strict stationarity.

Also,

$$\mathbb{E} \left[\left(\sum_{t=1}^m X_t \right)^2 \right] = \sum_{t=1}^m \sum_{s=1}^m \mathbb{E}[X_t X_s] = \sum_{t=1}^m \sum_{s=1}^m \gamma(|t-s|)$$

Let $h = t - s$, set of possible values for h is $m - |h|$, so

$$= \sum_{h=1-m}^{m-1} (m - |h|) \gamma(h) < \infty$$

noting that $\gamma(h) = \gamma(-h)$, therefore for C as a constant, we have

$$\mathbb{V}(G_2) = \frac{N}{T} C = \frac{\left\lfloor \frac{T}{a_T + m} \right\rfloor}{T} (C) \xrightarrow{a_T \rightarrow \infty} 0$$

and hence $G_2 = o_p(1)$.

Let's deal with the big block terms. Notice

$$G_1 = \frac{1}{\sqrt{T}} \sum_{j=1}^N \sum_{t \in B_j} X_t = \sum_{j=1}^N \frac{\sum_{t \in B_j} X_t}{\sqrt{T}} = \sum_{j=1}^N Y_j$$

where Y_j is a triangular array. So, $\mathbb{V}(G_1) = \sum_{j=1}^N \mathbb{V}(Y_j)$.

$$\begin{aligned} \mathbb{V}(Y_j) &= \mathbb{V}(Y_1) \\ &= \frac{1}{T} \mathbb{E} \left[\left(\sum_{t=1}^{a_T} X_t \right)^2 \right] \\ &= \frac{1}{T} \sum_{t=1}^{a_T} \sum_{s=1}^{a_T} \mathbb{E}[X_t X_s] \\ &= \frac{1}{T} \sum_{h=1-a_T}^{a_T-1} (a_T - |h|) \gamma(h) \end{aligned}$$

Note that since the process is m -dependent, $\gamma(h) = 0$ if $|h| \geq m$. Continuing,

$$\frac{1}{T} \sum_{h=1-a_T}^{a_T-1} (a_T - |h|) \gamma(h) = \sum_{h=-m}^m (a_T - |h|) \gamma(h)$$

Therefore,

$$\mathbb{V}(G_1) = \frac{N}{T} \sum_{h=-m}^m (a_T - |h|) \gamma(h) \xrightarrow{T \rightarrow \infty} \sum_{h=-m}^m \gamma(h)$$

$\approx 1/a_T$

Therefore, the variance of G_1 is bounded. We showed $\sigma_N^2 = \mathbb{V}(G_1) \approx \text{constant}$. So, we must show

$$\sum_{j=1}^N \mathbb{E} \left[\underbrace{Y_j^2}_{\text{iid}} \mathbb{I}\{|Y_j| > \varepsilon \sigma_N\} \right] = N \mathbb{E} \left[Y_1^2 \mathbb{I}\{|Y_1| > \varepsilon \sigma_N\} \right] \xrightarrow{T \rightarrow \infty} 0$$

Aside: For $\delta > 0$,

$$\begin{aligned} \mathbb{E}[|Y|^{2+\delta}] &\geq \mathbb{E}[|Y|^{2+\delta} \mathbb{I}\{|Y| > \varepsilon\}] \\ &\geq \varepsilon^\delta \mathbb{E}[|Y|^2 \mathbb{I}\{|Y| > \varepsilon\}] \\ \implies \mathbb{E}[|Y|^2 \mathbb{I}\{|Y| > \varepsilon\}] &\leq \frac{\mathbb{E}[|Y|^{2+\delta}]}{\varepsilon^\delta} \end{aligned}$$

It may be shown that for $C > 0$

$$\mathbb{E}[|Y_j|^{2+\delta}] \leq C \left(\frac{a_T}{T} \right)^{\frac{2+\delta}{2}}$$

So

$$\begin{aligned} N \mathbb{E}[Y_1^2 \mathbb{I}\{|Y_1| > \varepsilon \sigma_N\}] &\leq \frac{N}{(\varepsilon \sigma_N)^\delta} C \left(\frac{a_T}{T} \right)^{\frac{2+\delta}{2}} \\ &= \frac{C}{(\varepsilon \sigma_N)^\delta} \frac{N a_T}{T} \left(\frac{a_T}{T} \right)^{\delta/2} \xrightarrow{T \rightarrow \infty} 0 \end{aligned}$$

Therefore, by Theorem 2.3.3

$$\frac{G_1}{\sigma_N} \xrightarrow{T \rightarrow \infty} \mathcal{N}(0, 1)$$

and since

$$\sigma_N^2 \rightarrow \sum_{j=-m}^m \gamma(j)$$

we have

$$G_1 \xrightarrow{D} \mathcal{N}\left(0, \sum_{h=-m}^m \gamma(h)\right)$$

Since $G_2 = o_p(1)$ we have

$$\frac{1}{\sqrt{T}} \sum_{t=1}^T X_t \xrightarrow{D} \mathcal{N}\left(0, \sum_{h=-m}^n \gamma(h)\right)$$

2.4 † Two Plus Delta Moment Calculation

We want to show

$$\mathbb{E}[|Y_1|^{2+\delta}] \leq C \left(\frac{a_T}{T} \right)^{\frac{2+\delta}{2}}$$

where

$$Y_1 = \frac{1}{\sqrt{T}} \sum_{t=1}^{a_T} X_t$$

$a_T = \text{big block size} \rightarrow \infty$ as $T \rightarrow \infty$

$$\frac{a_T}{T} \rightarrow 0$$

X_t are m -dependent random variables.

$$\mathbb{E}[|X_i|^{2+\delta}] < \infty \quad (\delta > 0) \iff \mathbb{E} \left[\left| \sum_{t=1}^{a_T} X_t \right|^{2+\delta} \right] \leq C a_T^{\frac{2+\delta}{2}}$$

THEOREM 2.4.1: Rosenthal's Inequality

If X_1, \dots, X_n are independent random variables with $\mathbb{E}[|X_i|^{2+\delta}] < \infty$ for $\delta > 0$, then

$$\mathbb{E} \left[\left| \sum_{i=1}^n X_i \right|^{2+\delta} \right] < c_p n^{\delta/2} \sum_{i=1}^n \mathbb{E}[|X_i|^{2+\delta}]$$

In particular, if X_1, \dots, X_n are i.i.d., then

$$\mathbb{E} \left[\left| \sum_{i=1}^n X_i \right|^{2+\delta} \right] \leq c_p n^{\frac{2+\delta}{2}} \mathbb{E}[|X_1|^{2+\delta}]$$

Proof of: Theorem 2.4.1

See Petrov, Limit theorems of Probability Theory, p.g. 59.

PROPOSITION 2.4.2

For arbitrary random variables X_1, \dots, X_n ,

$$\mathbb{E} \left[\left| \sum_{i=1}^n X_i \right|^{2+\delta} \right] \leq n^{(2+\delta)-1} \sum_{i=1}^n \mathbb{E}[|X_i|^{2+\delta}]$$

Proof of: Proposition 2.4.2

Since $\varphi(x) = |x|^{2+\delta}$ is convex where $a_1, \dots, a_n \in \mathbf{R}$, by Jensen's Inequality,

$$\left| \frac{1}{n} \sum_{i=1}^n a_i \right|^{2+\delta} \leq \frac{1}{n} \sum_{i=1}^n |a_i|^{2+\delta}$$

Rearranging yields

$$\left| \sum_{i=1}^n a_i \right|^{2+\delta} \leq n^{(2+\delta)-1} \sum_{i=1}^n |a_i|^{2+\delta}$$

Replace $a_i \sim X_i$, take expectation.

$$\sum_{t=1}^{a_T} X_t = \sum_{j=0}^m \sum_{\substack{t \equiv j \pmod{m+1} \\ 1 \leq t \leq a_T}} X_t$$

Variables in the second sum are separated by at least m -time steps, and hence i.i.d. Therefore,

$$\begin{aligned} \mathbb{E} \left[\left| \sum_{t=1}^{a_T} X_t \right|^{2+\delta} \right] &\leq (m+1)^{(2+\delta)-1} \mathbb{E} \left[\left| \sum_{\substack{t \equiv j \pmod{m+1} \\ 1 \leq t \leq a_T}} X_t \right|^{2+\delta} \right] && \text{by Proposition 2.4.2} \\ &\leq (m+1)^{(2+\delta)-1} \left(\frac{a_T}{m+1} \right)^{\frac{2+\delta}{2}} \mathbb{E}[|X_1|^{2+\delta}] && \text{by Theorem 2.4.1} \\ &= C a_T^{\frac{2+\delta}{2}} \end{aligned}$$

where C is the same constant as in Section 2.3.

2.5 † Linear Process CLT

EXAMPLE 2.5.1

$X_t = \sum_{\ell=0}^m \psi_\ell W_{t-\ell}$ where $\{W_t\}_{t \in \mathbb{Z}}$ is a strong white noise in L^2 .

A general linear process $X_t = \sum_{\ell=0}^m \psi_\ell W_{t-\ell}$ is not m -dependent.

THEOREM 2.5.2: Basic Approximation Theorem (BAT)

Suppose X_n is a sequence of random variables so that there exists an array

$$\{Y_{m,n} : m, n \in \mathbb{Z}_{\geq 1}\}$$

so that:

- (1) For each fixed m , $Y_{m,n} \xrightarrow{D} Y_m$ as $n \rightarrow \infty$.
- (2) $Y_m \xrightarrow{D} Y$ as $m \rightarrow \infty$ for some random variable Y .
- (3) For all $\varepsilon > 0$,

$$\lim_{m \rightarrow \infty} \left[\limsup_{n \rightarrow \infty} \mathbb{P}(|X_n - Y_{m,n}| > \varepsilon) \right] = 0$$

Then $X_n \xrightarrow{D} Y$ as $n \rightarrow \infty$.

REMARK 2.5.3

$Y_{m,n}$ is often an “ m -dependent” approximation to X_n

Proof of: Theorem 2.5.2

Shumway and Stoffer using characteristic functions.

THEOREM 2.5.4: Linear Process CLT

Suppose $X_t = \sum_{\ell=0}^{\infty} \psi_{\ell} W_{t-\ell}$ is a causal linear process with $\sum_{\ell=0}^{\infty} |\psi_{\ell}| < \infty$ with $\{W_t\}_{t \in \mathbb{Z}}$ is a strong white noise in L^2 . If

$$S_t = \frac{1}{\sqrt{T}} \sum_{t=1}^T X_t$$

then

$$S_T \xrightarrow[T \rightarrow \infty]{D} \mathcal{N}\left(0, \sum_{\ell=-\infty}^{\infty} \gamma(\ell)\right)$$

Proof of: Theorem 2.5.4

X_t is strictly (and weakly) stationary.

$$\begin{aligned} \gamma(h) &= \mathbb{E}[X_t X_{t+h}] \\ &= \mathbb{E}\left[\left(\sum_{\ell=0}^{\infty} \psi_{\ell} W_{t-\ell}\right) \left(\sum_{j=0}^{\infty} \psi_j W_{t+h-j}\right)\right] \\ &= \sum_{\ell=0}^{\infty} \sum_{j=0}^{\infty} \psi_{\ell} \psi_j \mathbb{E}[W_{t-\ell} W_{t+h-j}] && \text{Fubini's Theorem} \\ &= \sum_{\ell=0}^{\infty} \psi_{\ell} \psi_{\ell+h} \sigma_W^2 \end{aligned}$$

Then,

$$\sum_{h=-\infty}^{\infty} \gamma(h) = \sum_{h=-\infty}^{\infty} \left| \sum_{\ell=0}^{\infty} \psi_{\ell} \psi_{\ell+h} \sigma_W^2 \right| \leq \sum_{\ell=0}^{\infty} |\psi_{\ell}| \sum_{h=-\infty}^{\infty} |\psi_h| \sigma_W^2 < \infty$$

by the Triangle Inequality. So $\sum_{h=-\infty}^{\infty} \gamma(h)$ is well-defined. Note that $\mathbb{E}[S_T] = 0$ since $\mathbb{E}[X_t] = 0$. Also,

$$\mathbb{V}(S_T) = \frac{1}{T} \sum_{t=1}^T \sum_{s=1}^T \mathbb{E}[X_t X_s] = \frac{1}{T} \sum_{h=1-T}^{T-1} (T - |h|) \gamma(h) = \sum_{h=1-T}^{T-1} \left(1 - \frac{|h|}{T}\right) \gamma(h)$$

Note that $\left(1 - \frac{|h|}{T}\right) \leq |\gamma(h)|$ since $\{\gamma(h)\}$ is summable by Dominated Convergence Theorem (DCT).

Define

$$\begin{aligned} X_{t,m} &= \sum_{\ell=0}^m \psi_{\ell} W_{t-\ell} \\ S_{T,m} &= \frac{1}{\sqrt{T}} \sum_{t=1}^T X_{t,m} \end{aligned}$$

is an m -dependent approximation to S_T .

(1) By the m -dependent CLT,

$$S_{T,m} \xrightarrow{D} \mathcal{N}\left(0, \sum_{h=-m}^m \gamma_m(h)\right) := S'_m$$

and $\gamma_m(h) = \mathbb{E}[X_{t,m} X_{t+h,m}]$.

(2) By DCT,

$$\sum_{h=-m}^m \gamma_m(h) \xrightarrow{m \rightarrow \infty} \sum_{h=-\infty}^{\infty} \gamma(h)$$

and hence

$$S'_m \xrightarrow{D} \mathcal{N}\left(0, \sum_{h=-\infty}^{\infty} \gamma(h)\right)$$

(3)

$$\begin{aligned}
\mathbb{E}[(S_{T,m} - S_T)^2] &= \frac{1}{T} \mathbb{E} \left[\left(\sum_{t=1}^T (X_t - X_{t,m}) \right)^2 \right] \\
&\leq \sum_{h=1-T}^{T-1} \left(1 - \frac{|h|}{T} \right) \sum_{\ell=m+1}^{\infty} |\psi_{\ell}| |\psi_{\ell+h}| \sigma_W^2 \\
&\leq \sum_{\ell=m+1}^{\infty} |\psi_{\ell}| \left(\sum_{h=-\infty}^{\infty} |\psi_h| \right) \sigma_W^2 \xrightarrow{m \rightarrow \infty} 0
\end{aligned}$$

So condition (3) of the BAT is satisfied using Markov's Inequality. Therefore,

$$S_T = \frac{1}{\sqrt{T}} \sum_{t=1}^T X_t \xrightarrow{D} \mathcal{N} \left(0, \sum_{h=-\infty}^{\infty} \gamma(h) \right)$$

2.6 Asymptotic Properties of Empirical ACF

If X_1, \dots, X_T is an observed time series in which we think was generated by a stationary process, then $\gamma(h) = \text{Cov}(X_t, X_{t+h})$ does not depend on t . Recall that

$$\hat{\gamma}(h) = \frac{1}{T} \sum_{t=1}^{T-h} (X_t - \bar{X})(X_{t+h} - \bar{X})$$

$$\rho(h) = \text{Corr}(X_t, X_{t+h}) = \frac{\gamma(h)}{\gamma(0)}$$

$$\hat{\rho}(h) = \frac{\hat{\gamma}(h)}{\hat{\gamma}(0)}$$

Questions

- (1) Are $\hat{\gamma}$ and $\hat{\rho}$ **consistent**?
- (2) What is the approximate distribution of $\hat{\gamma}(h)$ and $\hat{\rho}(h)$.

Consistency

By adding and subtracting μ in the definition of $\hat{\gamma}(h)$, we may assume without loss of generality that $\mathbb{E}[X_t] = 0$.

Suppose $\{X_t\}_{t \in \mathbb{Z}}$ is strictly stationary, and

$$X_t = g(W_t, W_{t-1}, \dots)$$

We first need to establish the consistency of

$$\bar{X} = \frac{1}{T} \sum_{t=1}^T X_t$$

where X_t 's are not i.i.d. so Law of Large numbers does not work. Instead, we would use the Ergodic Theorem, but we will not cover it here. Therefore,

$$\bar{X} \xrightarrow{P} 0$$

Furthermore,

$$\begin{aligned}\hat{\gamma}(h) &= \frac{1}{T} \sum_{t=1}^{T-h} (X_t - \bar{X})(X_{t+h} - \bar{X}) \\ &= \frac{1}{T} \sum_{t=1}^{T-h} X_t X_{t+h} - \bar{X} \frac{1}{T} \sum_{t=1}^{T-h} X_t - \bar{X} \frac{1}{T} \sum_{t=1}^{T-h} X_{t+h} + \frac{T-h}{T} (\bar{X})^2\end{aligned}$$

where we note that the last three terms converge in probability to 0 by the Ergodic Theorem.

Also, note that $\mathbb{E}[X_t X_{t+h}] = \gamma(h)$ and $X_t X_{t+h} = g_h(W_{t+h}, W_{t+h-1}, \dots)$.

Again, by the Ergodic Theorem,

$$\frac{1}{T} \sum_{t=1}^{T-h} X_t X_{t+h} \xrightarrow{P} \gamma(h)$$

Therefore, $\hat{\gamma}(h) \xrightarrow{P} \gamma(h)$ and $\hat{\rho}(h) = \frac{\hat{\gamma}(h)}{\hat{\gamma}(0)} \xrightarrow{P} \rho(h)$ under strict stationarity and $\mathbb{E}[X_t^2] < \infty$.

Distribution of $\hat{\gamma}(h)$

Consider simple (but most important case) when $\{X_t\}_{t \in \mathbb{Z}}$ is a strong white noise with $\mathbb{E}[X_t^4] < \infty$. The finite 4th moment assumption is not really assumed here, but this will be explained why it's classically assumed.

$$\hat{\gamma}(h) \xrightarrow{P} 0$$

Similarly,

$$\hat{\gamma}(h) = \frac{1}{T} \sum_{t=1}^{T-h} X_t X_{t+h} + \underbrace{R}_{\tilde{\gamma}(h)}$$

Note that $\mathbb{E}[\tilde{\gamma}(h)] = 0$ for $h \geq 1$. Also,

$$\mathbb{V}(\tilde{\gamma}(h)) = \mathbb{E}[\tilde{\gamma}^2(h)] = \frac{1}{T^2} \sum_{t=1}^{T-h} \sum_{s=1}^{T-h} \mathbb{E}[X_t X_{t+h} X_s X_{s+h}]$$

is non-zero only when $t = s$, so

$$\mathbb{V}(\tilde{\gamma}(h)) = \frac{1}{T^2} \sum_{t=1}^{T-h} \mathbb{E}[X_t^2 X_{t+h}^2] = \frac{T-h}{T^2} \sigma_X^4$$

where $\mathbb{E}[X_t^2] = \sigma_X^2$. Therefore,

$$\mathbb{V}(\sqrt{T} \tilde{\gamma}(h)) \xrightarrow{T \rightarrow \infty} \sigma_X^4$$

THEOREM 2.6.1

If $\{X_t\}_{t \in \mathbb{Z}}$ is a strong white noise with $\mathbb{E}[X_t^4] < \infty$, then

$$\sqrt{T} \tilde{\gamma}(h) = \frac{1}{\sqrt{T}} \sum_{t=1}^{T-h} X_t X_{t+h} \xrightarrow{D} \mathcal{N}(0, \sigma_X^4)$$

Proof of: Theorem 2.6.1

Using Martingale CLT which is derived from m -dependent CLT.

COROLLARY 2.6.2

It follows that if

$$\sqrt{T}\hat{\gamma} \xrightarrow{D} \mathcal{N}(0, \sigma_X^4)$$

and $\hat{\gamma}(0) \xrightarrow{P} \sigma_X^2$ (SLLN), then by Slutsky's Theorem,

$$\sqrt{T} \frac{\hat{\gamma}(h)}{\hat{\gamma}(0)} = \sqrt{T} \hat{\rho}(h) \xrightarrow{D} \mathcal{N}(0, 1)$$

If $\{X_t\}_{t \in \mathbb{Z}}$ is a strong white noise,

$$\left(-\frac{z_{\alpha/2}}{\sqrt{T}}, \frac{z_{\alpha/2}}{\sqrt{T}} \right)$$

is a $(1 - \alpha)$ prediction interval for $\hat{\rho}(h)$ for all h with T large where $\Phi(z_{\alpha/2}) = 1 - \alpha$. Hence,

$$\left(-\frac{1.96}{\sqrt{T}}, \frac{1.96}{\sqrt{T}} \right)$$

is an approximate 95% prediction interval for $\hat{\rho}(h)$ assuming the data is generated by a strong white noise process.

Now, we know that the blue boundaries are $\pm \frac{1.96}{\sqrt{T}}$ in Figure 2.1. Also, we might be able to say that exists mild serial correlation at lag 1 of the ACF for Figure 2.2 since there are lines that go outside the blue boundaries.



Figure 2.2: ACF of first differenced temperature data

```
# Figure 2.2
plot(acf(diff(gtemp)))
```

2.7 Interpreting the Autocorrelation Function (Non-stationary)

We have an excellent understanding of how $\hat{\rho}(h)$ behaves when X_1, \dots, X_T is a strong white noise.

- Consistency:

$$\hat{\rho}(h) \xrightarrow{P} 0 \quad (h \geq 1)$$

- Distribution:

$$\hat{\rho}(h) \stackrel{D}{\approx} \mathcal{N}\left(0, \frac{1}{T}\right) \quad (T \text{ is large})$$

What happens when we calculate the empirical ACF for a non-stationary time series?

EXAMPLE 2.7.1

$X_t = t + W_t$ where W_t is a strong white noise. Note that X_t has a linear trend, and hence not stationary. First,

$$\bar{X} = \frac{1}{T} \sum_{t=1}^T [t + W_t] = \frac{1}{T} \frac{[T(T+1)]}{2} + \bar{W} = \frac{T+1}{2} + \bar{W}$$

Also,

$$\begin{aligned} \hat{\gamma}(h) &= \frac{1}{T} \sum_{t=1}^{T-h} \left(t + W_t - \frac{T+1}{2} - \bar{W} \right) \left(t + h + W_{t+h} - \frac{T+1}{2} - \bar{W} \right) \\ &= \frac{1}{T} \sum_{t=1}^{T-h} \left(t - \frac{T+1}{2} \right) \left(t + h - \frac{T+1}{2} \right) + R \\ &= \frac{1}{T} \sum_{t=1}^{T-h} \left(t - \frac{T+1}{2} \right)^2 + \frac{1}{T} \sum_{t=1}^{T-h} h \left(t - \frac{T+1}{2} \right) \\ &= \frac{1}{T} \sum_{t=1}^{T/2} t^2 + \frac{h}{T} \left[\frac{(T-h)(T-h+1)}{2} - \frac{(T+1)(T-h)}{2} \right] \\ &\approx \mathcal{O}(T^2) + \mathcal{O}(T) \end{aligned}$$

where R is the remainder with the white noise terms. Note that the dominant term; that is, the $\mathcal{O}(T^2)$ doesn't depend on h .

It follows that in this case that

$$\frac{\hat{\gamma}(h)}{T^2} \xrightarrow{T \rightarrow \infty} C \quad (\forall h)$$

Hence

$$\hat{\rho}(h) = \frac{\hat{\gamma}(h)}{\hat{\gamma}(0)} \frac{T^2}{T^2} = \frac{\hat{\gamma}(h)}{\hat{\gamma}(0)} \frac{T^2}{T^2} \xrightarrow{P} 1 \quad (\forall h)$$

Moral: If X_t has a trend that is not properly removed, $\hat{\rho}(h)$ is likely to be large.

Figure 2.3

`acf(gtemp)`

Figure 2.4

`plot(as.ts(cumsum(rnorm(100))), main = "autoregression, phi=1")`

Figure 2.5

`acf(as.ts(cumsum(rnorm(100))))`

- Looking back at Figure 1.2, we see that this time series has an upwards trend. Therefore, based on what we just did, we expect that the ACF should be very large (close to 1) at each lag for this time series. Clearly, Figure 2.3 is indicative of a strong trend or non-stationarity.
- In Figure 2.4, we are plotting

$$X_t = X_{t-1} + W_t$$

with $X_0 = 0$ and $X_t = \sum_{j=1}^t W_j$ which is non-stationary. Some people say it has a “stochastic trend.”



Figure 2.3: ACF of raw temperature data, sample length 130

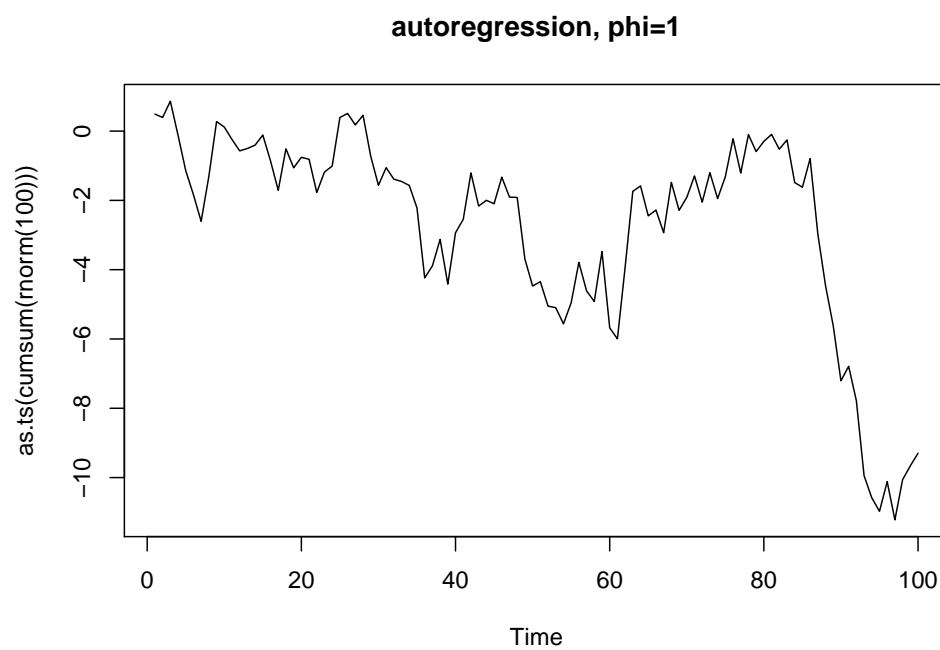
Figure 2.4: Realization of an AR(1) with $\phi = 1$ starting from $x_0 = 0$



Figure 2.5: ACF of an AR(1) with $\phi = 1$ starting from $x_0 = 0$

- In Figure 2.5 there exists a similar pattern which is indicative of non-stationarity.

Chapter 3

Week 3

3.1 Moving Average Processes

Suppose X_t is stationary. Identify serial dependence using ACF $\hat{\rho}(h)$. If the lines go out of the dotted blue boundaries, namely $\pm \frac{1.96}{\sqrt{T}}$, within the ACF plot of $\hat{\rho}(h)$, then we suspect serial dependence.

Posit

$$X_t = g(W_t, W_{t-1}, \dots) = \sum_{\ell=0}^{\infty} \psi_{\ell} W_{t-\ell} \quad [\text{Linear Process}]$$

Not feasible to estimate infinitely many parameters $\{\psi\}_{\ell=0}^{\infty}$. Assume coefficients arise from a *parsimonious* linear model for X_t .

DEFINITION 3.1.1: Moving average process

Suppose $\{W_t\}_{t \in \mathbb{Z}}$ is a strong white noise with $\mathbb{V}(W_t) = \sigma_W^2 < \infty$. We say X_t is a **moving average process** of order q or $\text{MA}(q)$, if there exists $\theta_1, \dots, \theta_q \in \mathbf{R}$ with $\theta_q \neq 0$ such that

$$X_t = W_t + \theta_1 W_{t-1} + \dots + \theta_q W_{t-q} = \sum_{\ell=0}^q \theta_{\ell} W_{t-\ell}$$

where $\theta_0 = 1$. In other words, we've truncated the linear process representation at the level q .

DEFINITION 3.1.2: Backshift operator

The **backshift operator**, B , is defined by

$$B^j X_t = X_{t-j}$$

B is assumed further to be linear in the sense that for $a, b \in \mathbf{R}$

$$(aB^j + bB^k)X_t = aB^j X_t + bB^k X_t = aX_{t-j} + bX_{t-k}$$

EXAMPLE 3.1.3

- First difference of X_t :

$$\nabla X_t = (1 - B)X_t = X_t - BX_t = X_t - X_{t-1}$$

- Second difference of X_t :

$$\nabla^2 X_t = (1 - B)^2 X_t = (1 - 2B + B^2)X_t = X_t - 2X_{t-1} + X_{t-2}$$

DEFINITION 3.1.4: Moving average operator

The **moving average operator** is defined by

$$\theta(B) = 1 + \theta_1 B + \theta_2 B^2 + \cdots + \theta_q B^q$$

DEFINITION 3.1.5: Moving average polynomial

The **moving average polynomial** is defined as

$$\theta(x) = 1 + \theta_1 x + \cdots + \theta_q x^q$$

If $X_t \sim \text{MA}(q)$, then

$$X_t = W_t + \theta_1 W_{t-1} + \cdots + \theta_q W_{t-q} = \theta(B)W_t$$

which is a succinct expression defining $\text{MA}(q)$.

Properties of $\text{MA}(q)$ Processes

- (1) $\text{MA}(0)$ process is a strong white noise.
- (2) If $X_t \sim \text{MA}(q)$, then

$$\mathbb{E}[X_t] = \mathbb{E}\left[\sum_{\ell=0}^q \theta_\ell W_{t-\ell}\right] = 0 \quad (\theta_0 = 1)$$

$$\mathbb{V}(X_t) = \mathbb{E}\left[\left(\sum_{\ell=0}^q \theta_\ell W_{t-\ell}\right)^2\right] = \sum_{\ell=0}^q \theta_\ell^2 \sigma_W^2$$

$$\begin{aligned} \gamma(h) &= \text{Cov}(X_t, X_{t+h}) \\ &= \mathbb{E}\left[\left(\sum_{\ell=0}^q \theta_\ell W_{t-\ell}\right)\left(\sum_{k=0}^q \theta_k W_{t+h-k}\right)\right] \quad t - \ell = t + h - k \implies k = \ell + h \\ &= \begin{cases} \sigma_W^2 \sum_{j=0}^{q-h} \theta_j \theta_{j+h} & 1 \leq h \leq q \\ 0 & h > q \end{cases} \end{aligned}$$

Recall that $\gamma(h) = \gamma(-h)$, so we will only display the values for $h \geq 0$. Note that $\gamma(q)$ cannot be zero because $\theta \neq 0$. The cutting off of $\gamma(h)$ after q lags is the signature of the $\text{MA}(q)$ model. Therefore,

$$\rho(h) = \frac{\gamma(h)}{\gamma(0)} = \begin{cases} \frac{\sum_{j=0}^{q-h} \theta_j \theta_{j+h}}{\sum_{j=0}^q \theta_j^2} & 1 \leq h \leq q \\ 0 & h > q \end{cases}$$

REMARK 3.1.6

By choosing $\theta_1, \dots, \theta_q$ appropriately, we can get any ACF we want $\rho(h)$ where $1 \leq h \leq q$.

- (3) If $X_t \sim \text{MA}(q)$, then X_t is q -dependent.

In Figure 3.1, let's look an example now of what a moving average process would actually look like if we were to realize a moving average process. On the top of Figure 3.1, I've plotted a moving average process of order 0, which is just a strong white noise. Then, as we progress down to panel 2 and panel 3 I've calculated moving averages of orders 1 and 2 based on this strong white noise sequence. In the second panel, $X_t = W_t + W_{t-1}$, so this is a moving average process of order 1, in which $\theta_1 = 1$. In the third panel, we have a moving average process of order 2, in which $X_t = W_t + W_{t-1} + W_{t-2}$, which is a moving average process of order 2 where $\theta_1 = \theta_2 = 1$. One thing to observe when going from a moving average process of order 0 to 2 is that the time series is getting "smoother."



Figure 3.1: Realizations of MA processes with coefficients equal to 1

```
# Figure 3.1
par(mfrow = c(3, 1))
```

```

ma0.sim <- arima.sim(list(order = c(0, 0, 0), ma = c()), n = 134)
plot(ma0.sim, ylab = "x", main = "white noise")

ma1.sim <- arima.sim(list(order = c(0, 0, 1), ma = c(1)), n = 134)
plot(ma1.sim, ylab = "v", main = (expression(MA(1) ~ ~ ~ theta[1] == 1)))

ma2.sim <-
  arima.sim(list(order = c(0, 0, 2), ma = c(1, 1)), n = 134)
plot(ma2.sim, ylab = "y", main = (expression(paste(
  MA(2), ~ ~ ~ theta[1], " = ", theta[2], " = ", 1
))))

```

In Figure 3.2, the difference is apparent since going from MA(0) to MA(1) shows that MA(1) has significant serial correlation at lag 1. Similarly, for MA(2) there is significant serial correlation at lag 2.

```

# Figure 3.2
acf(ma0.sim)
acf(ma1.sim)
acf(ma2.sim)

```

3.2 Autoregressive Processes

DEFINITION 3.2.1: Autoregressive process

Suppose $\{W_t\}_{t \in \mathbb{Z}}$ is a strong white noise with $\mathbb{V}(W_t) = \sigma_W^2 < \infty$. We say X_t is an **autoregressive process** of order 1, or AR(1), if there exists a constant ϕ such that

$$X_t = \phi X_{t-1} + W_t \quad (t \in \mathbb{Z})$$

Using the backshift operator, this may also be expressed as

$$(1 - \phi B)X_t = W_t$$

Interpretation

Prediction: Form a linear model (regression) predicting X_t as

$$X_t = \phi X_{t-1} + W_t$$

where X_t is the dependent variable and X_{t-1} is the covariant/independent variable.

Markov Property:

$$X_t \mid (X_{t-1}, X_{t-2}, \dots) = X_t \mid X_{t-1}$$

Question: Does there exist a stationary process X_t satisfying the following?

$$X_t = \phi X_{t-1} + W_t$$

Let's see.

$$\begin{aligned}
 X_t &= \phi X_{t-1} + W_t \\
 &= \phi(\phi X_{t-2} + W_{t-1}) + W_t \\
 &= \phi^2 X_{t-2} + \phi W_{t-1} + W_t \\
 &\vdots \\
 &= \phi^k X_{t-k} + \sum_{j=0}^{k-1} \phi^j W_{t-j}
 \end{aligned}$$

k times
 if $|\phi| > 1$, the sum diverges



Figure 3.2: ACF plots of corresponding moving average series.

Suppose $|\phi| < 1$, then

$$\xrightarrow[k \rightarrow \infty]{L^2\text{-sense}} 0 + \sum_{j=0}^{\infty} \phi^j W_{t-j}$$

which is a causal linear process. Moreover, if $X_t = \sum_{j=0}^{\infty} \phi^j W_{t-j}$, then X_t is strictly stationary, and

$$\begin{aligned} X_t &= \sum_{j=0}^{\infty} \phi^j W_{t-j} \\ &= \sum_{j=1}^{\infty} \phi^j W_{t-j} + W_t \\ &= \phi \sum_{j=1}^{\infty} \phi^{j-1} W_{t-j} + W_t & j \rightarrow j-1 \\ &= \phi \sum_{j=0}^{\infty} \phi^j W_{t-1-j} + W_t \\ &= \phi X_{t-1} + W_t \end{aligned}$$

Therefore, X_t satisfies AR(1) equation.

THEOREM 3.2.2

If $|\phi| < 1$, then there exists a strictly stationary and causal linear process X_t such that

$$X_t = \phi X_{t-1} + W_t$$

What if $|\phi| > 1$? If $X_t = \phi X_{t-1} + W_t$ for $t \in \mathbb{Z}$, then that implies

$$\begin{aligned} X_t &= \phi^{-1} X_{t+1} - \phi^{-1} W_{t+1} \\ &= \phi^{-1} (\phi^{-1} X_{t+1} - \phi^{-1} W_{t+1}) - \phi^{-1} W_{t+1} \\ &\vdots \\ &= \phi^{-k} X_{t+k} - \sum_{j=1}^{k-1} \phi^{-j} W_{t+j} \end{aligned} \quad k \text{ times}$$

Therefore,

$$X_t = \frac{X_{t+k}}{\phi^k} - \sum_{j=1}^{k-1} \frac{W_{t+j}}{\phi^j} \xrightarrow[k \rightarrow \infty]{L^2\text{-sense}} - \sum_{j=1}^{\infty} \frac{W_{t+j}}{\phi^j}$$

since $\sum_{j=1}^{\infty} \frac{1}{\phi^j} < \infty$. This sequence is strictly stationary since it is a Bernoulli shift. However, what we have derived is not desirable as this model is future dependent, normally we try to avoid this.

What if $|\phi| = 1$? In this case we claim that there is no stationary process such that $X_t = \phi X_{t-1} + W_t$. Let's prove this. Suppose $|\phi| = 1$. If $X_t = X_{t-1} + W_t$, then

$$X_t = \sum_{j=1}^t W_j + X_0 \quad (\text{by iterating}) \implies X_t - X_0 = \sum_{j=1}^t W_j \quad [\text{Random Walk}]$$

Now, $|\text{Cov}(X_t, X_0)|^2 \leq \mathbb{V}(X_t)\mathbb{V}(X_0) = (\mathbb{V}(X_0))^2$, so we get

$$|\text{Cov}(X_t, X_0)| \leq \sqrt{\mathbb{V}(X_t)\mathbb{V}(X_0)} = \sqrt{(\mathbb{V}(X_0))^2} = \mathbb{V}(X_0)$$

Therefore, $-2\text{Cov}(X_t, X_0) \leq 2|\text{Cov}(X_t, X_0)| \leq 2\mathbb{V}(X_0)$. Finally,

$$\mathbb{V}(X_t - X_0) = \mathbb{V}(X_t) + \mathbb{V}(X_0) - 2\text{Cov}(X_t, X_0) \leq 4\mathbb{V}(X_0)$$

where in the last inequality we used the fact that X_t is stationary. Furthermore,

$$\mathbb{V}\left(\sum_{j=1}^t W_j\right) = t\sigma_W^2 \xrightarrow{t \rightarrow \infty} \infty$$

Properties of Causal AR(1) for $|\phi| < 1$

(1) The span of dependence of X_t is “infinite”

$$X_t = \sum_{\ell=0}^{\infty} \phi^\ell W_{t-\ell}$$

(2) ACF.

$$\mathbb{V}(X_t) = \mathbb{E}\left[\left(\sum_{\ell=0}^{\infty} \phi^\ell W_{t-\ell}\right)^2\right] = \sum_{\ell=0}^{\infty} \phi^{2\ell} \sigma_W^2 = \frac{\sigma_W^2}{1-\phi^2}$$

$$\begin{aligned} \gamma(h) &= \text{Cov}(X_t, X_{t+h}) \\ &= \mathbb{E}\left[\left(\sum_{\ell=0}^{\infty} \phi^\ell W_{t-\ell}\right)\left(\sum_{k=0}^{\infty} \phi^k W_{t+h-k}\right)\right] \\ &= \sum_{\ell=0}^{\infty} \phi^\ell \phi^{\ell+h} \sigma_W^2 \\ &= \phi^h \sum_{\ell=0}^{\infty} \phi^{2\ell} \sigma_W^2 \\ &= \phi^h \left(\frac{\sigma_W^2}{1-\phi^2}\right) \end{aligned}$$

where in the first sum we let $t-\ell = t+h-k$ and in the second sum we let $k = \ell+h$ for $\ell = 0, 1, 2, \dots$. Hence,

$$\rho(h) = \frac{\gamma(h)}{\gamma(0)} = \phi^h \quad (h \geq 0)$$

Note: this decays geometrically in the lag parameter.

Figure 3.3

```
ar0.sim <- arima.sim(list(order = c(1, 0, 0), ar = c(0.5)), n = 134)
plot(ar0.sim, ylab = "x", main = (expression(AR(1) ~ ~ ~ phi[1] == 0.5)))

ar1.sim <- arima.sim(list(order = c(1, 0, 0), ar = c(0.9)), n = 134)
plot(ar1.sim, ylab = "y", main = (expression(AR(1) ~ ~ ~ phi[1] == 0.9)))

ar2.sim <-
  arima.sim(list(order = c(1, 0, 0), ar = c(-0.9)), n = 134)
plot(ar2.sim, ylab = "z", main = (expression(AR(1) ~ ~ ~ phi[1] == -0.9)))
```

Figure 3.4

```
acf(ar0.sim)
acf(ar1.sim)
acf(ar2.sim)
```



Figure 3.3: Realizations of AR(1) processes



Figure 3.4: Corresponding ACF plots

DEFINITION 3.2.3: Autoregressive process, Autoregressive polynomial

We say X_t follows an **autoregressive process** of order p , or $\text{AR}(p)$, if there exists coefficients $\phi_1, \dots, \phi_p \in \mathbb{R}$ with $\phi_p \neq 0$ such that

$$X_t = \phi_1 X_{t-1} + \dots + \phi_p X_{t-p} + W_t$$

We also define the **autoregressive polynomial** to be

$$\phi(x) = 1 - \phi_1 x - \dots - \phi_p x^p$$

$X_t \sim \text{AR}(p)$ if $\phi(B)X_t = W_t$.

3.3 ARMA Processes

We've seen the moving average polynomial:

$$\theta(x) = 1 + \theta_1 x + \dots + \theta_q x^q \quad (\theta_q \neq 0)$$

and the autoregressive polynomial:

$$\phi(x) = 1 - \phi_1 x - \dots - \phi_p x^p \quad (\phi_p \neq 0)$$

If $W_t \sim$ strong white noise

$$X_t = \theta(B)W_t \quad (X_t \sim \text{MA}(q))$$

$$\phi(B)X_t = W_t \quad (X_t \sim \text{AR}(p))$$

Why not combine the two?

DEFINITION 3.3.1: Autoregressive moving average

Given a strong white noise sequence W_t , we say that X_t is an **autoregressive moving average process** of orders p and q , or $\text{ARMA}(p, q)$, if X_t is stationary and

$$\phi(B)X_t = \theta(B)W_t$$

$$\phi(z) = 1 - \phi_1 z - \dots - \phi_p z^p \quad (\phi_p \neq 0)$$

$$\theta(z) = 1 + \theta_1 z + \dots + \theta_q z^q \quad (\theta_q \neq 0)$$

This implies that the model is

$$X_t = \phi_1 X_{t-1} + \dots + \phi_p X_{t-p} + W_t + \theta_1 W_{t-1} + \dots + \theta_q W_{t-q}$$

Using ARMA models to model autocorrelation: ARMA combines the following two points.

- $\text{MA}(q)$: ACF may be specified at lags $1, \dots, q$
- $\text{AR}(p)$: ACF has geometric decay/oscillations

REMARK 3.3.2: Parameter redundancy

Consider $X_t = W_t$ where $X_t \sim \text{MA}(0)$, then

$$0.5X_{t-1} = 0.5W_{t-1}$$

Therefore,

$$X_t - 0.5X_{t-1} = W_t - 0.5W_{t-1} \implies X_t \sim \text{ARMA}(1, 1)$$

$$\phi(z) = 1 - 0.5z \implies \text{zero of } \phi \text{ is } z_0 = 2$$

$$\theta(z) = 1 - 0.5z \implies \text{zero of } \theta \text{ is } z_0 = 2$$

Parameter redundancy manifests as shared zeros in ϕ and θ . We always assume the models are “reduced” by factoring and diving away common zeros in ϕ .

DEFINITION 3.3.3: Causal ARMA

We say an $\text{ARMA}(p, q)$ is **causal** if there exists $\{X_t\}_{t \in \mathbb{Z}}$ satisfying $\phi(B)X_t = \theta(B)W_t$ and

$$X_t = \sum_{\ell=0}^{\infty} \psi_{\ell} W_{t-\ell} = \psi(B)W_t \quad [\text{Causal Linear Process Solution}]$$

where $\psi(B) = \sum_{\ell=0}^{\infty} \psi_{\ell} B^{\ell}$ and $\sum_{\ell=0}^{\infty} |\psi_{\ell}| < \infty$ with $\psi_0 = 1$.

DEFINITION 3.3.4: Invertible ARMA

An $\text{ARMA}(p, q)$ is **invertible** if there exists $\{X_t\}_{t \in \mathbb{Z}}$ satisfying $\phi(B)X_t = \theta(B)W_t$ and

$$W_t = \sum_{\ell=0}^{\infty} \pi_{\ell} X_{t-\ell} = \pi(B)X_t$$

where $\pi(B) = \sum_{\ell=0}^{\infty} \pi_{\ell} B^{\ell}$ and $\sum_{\ell=0}^{\infty} |\pi_{\ell}| < \infty$ with $\pi_0 = 1$.

REMARK 3.3.5

Causality + Invertibility \implies Information in $\{X_t\}_{t \leq T}$ is the same as Information in $\{W_t\}_{t \leq T}$ where $\{X_t\}_{t \leq T}$ is an observed time series.

THEOREM 3.3.6: Causality

By the fundamental theorem of algebra, the autoregressive polynomial $\phi(z)$ has p roots, say $z_1, \dots, z_p \in \mathbf{C}$. If $\rho = \min_{1 \leq j \leq p} |z_j| > 1$, then there exists a stationary and causal X_t to the ARMA equations: $\phi(B)X_t = \theta(B)W_t$.

$$X_t = \sum_{\ell=0}^{\infty} \psi_{\ell} W_{t-\ell}$$

The coefficients $\{\psi_{\ell}\}_{\ell=0}^{\infty}$ satisfy

$$\sum_{\ell=0}^{\infty} |\psi_{\ell}| < \infty$$

in fact,

$$|\psi_{\ell}| \leq \frac{1}{\rho^{\ell}}$$

which is the geometric decay. Also,

$$\psi(z) = \sum_{\ell=0}^{\infty} \psi_{\ell} z^{\ell} = \frac{\theta(z)}{\phi(z)} \quad (|z| \leq 1)$$

In essence,

$$X_t = \frac{\theta(B)}{\phi(B)} W_t = \sum_{j=0}^{\infty} \psi_j B^j W_t$$

Key: $\frac{1}{\phi(z)} = \sum_{j=0}^{\infty} \phi_j z^j$ where $|z| \leq 1$ so $\frac{1}{\phi}$ has a convergent power series representation for $|z| \leq 1$.

THEOREM 3.3.7: Invertibility

If z_1, \dots, z_q are the zeros of $\theta(z)$ and $\min_{1 \leq i \leq q} |z_i| > 1$, then X_t is invertible,

$$W_t = \sum_{\ell=0}^{\infty} \pi_{\ell} X_{t-\ell}$$

Coefficients $\{\pi_{\ell}\}_{\ell=0}^{\infty}$ satisfy

$$\pi(z) = \sum_{\ell=0}^{\infty} \pi_{\ell} z^{\ell} = \frac{\phi(z)}{\theta(z)} \quad (|z| \leq 1)$$

Moral: When we look for coefficients $\phi_1, \dots, \phi_p, \theta_1, \dots, \theta_q$, we want to do so in such a way that

$$\phi(z), \theta(z) \neq 0 \quad (|z| \leq 1)$$

3.4 ARMA Process Examples and ACF**EXAMPLE 3.4.1**

Consider the ARMA(2, 2) model

$$X_t = \frac{1}{4}X_{t-1} + \frac{1}{8}X_{t-2} + W_t - \frac{5}{6}W_{t-1} + \frac{1}{6}W_{t-2}$$

Questions:

- Is there a stationary and causal solution to X_t ?

- Is it invertible?
- Is there parameter redundancy?

AR polynomial:

$$\phi(z) = 1 - \frac{1}{4}z - \frac{1}{8}z^2$$

MA polynomial:

$$\theta(z) = 1 - \frac{5}{6}z + \frac{1}{6}z^2$$

Roots for ϕ :

$$\frac{2 \pm \sqrt{4 + 4(8)}}{-2} = -1 \pm 3 = -4, 2$$

Roots for θ : 2, 3

$$\Rightarrow \phi(z) = -\frac{1}{8}(z+4)(z-2), \quad \theta(z) = \frac{1}{6}(z-2)(z-3)$$

Note that $\phi(z)$ and $\theta(z)$ share common $(z-2)$ which indicates that the parameters are redundant. Therefore, X_t satisfies an ARMA(1, 1) with

$$\phi(z) = -\frac{1}{8}(z+4), \quad \theta(z) = \frac{1}{6}(z-3)$$

Since the roots of ϕ and θ are outside the unit circle in \mathbb{C} , X_t is stationary, causal, and invertible.

EXAMPLE 3.4.2

Suppose

$$X_t = -\frac{1}{4}X_{t-1} + W_t - \frac{1}{3}W_{t-1}$$

where $X_t \sim \text{ARMA}(1, 1)$.

$$\phi(z) = 1 + \frac{1}{4}z \Rightarrow \text{Root is } -4.$$

So X_t is stationary and causal, and can be represented as a linear process:

$$X_t = \sum_{\ell=0}^{\infty} \psi_{\ell} W_{t-\ell}$$

We need to calculate the coefficients ψ_{ℓ} .

We know

$$\begin{aligned} \psi(z) &= \sum_{\ell=0}^{\infty} \psi_{\ell} z^{\ell} = \frac{\theta(z)}{\phi(z)} \quad (|z| \leq 1) \\ \Rightarrow \psi(z)\phi(z) &= \theta(z) \end{aligned}$$

Note that both $\psi(z)\phi(z)$ and $\theta(z)$ are power series, therefore we can calculate ψ_{ℓ} by matching coefficients.

- $\phi(z) = 1 + \frac{1}{4}z$
- $\theta(z) = 1 - \frac{1}{3}z$
- $\psi(z)\phi(z) = \theta(z)$

Let's compute it.

$$\begin{aligned}
 z^0 : \quad \psi_0 &= 1 \\
 z^1 : \quad \frac{\psi_0}{4} + \psi_1 &= -\frac{1}{3} & \implies \psi_1 &= -\frac{7}{12} \\
 z^2 : \quad \frac{\psi_1}{4} + \psi_2 &= 0 & \implies \psi_2 &= \frac{7}{12} \left(\frac{1}{4} \right) \\
 &\vdots \\
 z^\ell : \quad \frac{\psi_{\ell-1}}{4} + \psi_\ell &= 0 & \implies \psi_\ell &= (-1)^\ell \frac{7}{12} \left(\frac{1}{4} \right)^{\ell-1} \quad (\ell \geq 1)
 \end{aligned}$$

Simplifying,

$$\psi_j = \begin{cases} 1 & j = 0 \\ \frac{7}{3} \left(-\frac{1}{4} \right)^j & j \geq 1 \end{cases}$$

We can automate ψ_j in R with `ARMAtoMA()`.

```
library(astsa)
ARMAtoMA(ar=-1/4, ma=-1/3, 10)
```

If X_t is a stationary and causal solution to the $\text{ARMA}(p, q)$ model.

$$X_t = \sum_{j=0}^{\infty} \psi_j W_{t-j}$$

$$\gamma_X(h) = \mathbb{E}[X_t X_{t+h}] = \mathbb{E} \left[\left(\sum_{j=0}^{\infty} \psi_j W_{t-j} \right) \left(\sum_{k=0}^{\infty} \psi_k W_{t+h-k} \right) \right]$$

Note that

$$t - j = t + h - k, \implies k = h + j, \quad j = 0, 1, 2, \dots \quad \mathbb{E}[X_{t-j}^2] = \sigma_W^2$$

Therefore,

$$\gamma_X(h) = \sigma_W^2 \sum_{j=0}^{\infty} \psi_j \psi_{j+h}$$

We can automate $\gamma_X(h)$ in R with `ARMAacf()`.

For $h \geq 1$, we have

$$\begin{aligned}
 \gamma_X(h) &= \sum_{j=0}^{\infty} \psi_j \psi_{j+h} \\
 &= \psi_0 \psi_h + \sum_{j=1}^{\infty} \psi_j \psi_{j+1} \\
 &= \frac{7}{3} \left(-\frac{1}{4} \right)^h + \sum_{j=1}^{\infty} \left[\frac{7}{3} \left(-\frac{1}{4} \right)^j \frac{7}{3} \left(-\frac{1}{4} \right)^{j+1} \right] \\
 &= \frac{91}{135} (-1)^h 4^{1-h}
 \end{aligned}$$

Then,

$$\begin{aligned}
 \gamma_X(0) &= \sum_{j=0}^{\infty} \psi_j^2 \\
 &= (1)^2 + \sum_{j=1}^{\infty} \psi_j^2 \\
 &= 1 + \sum_{j=1}^{\infty} \frac{7}{3} \left(-\frac{1}{4}\right)^j \\
 &= \frac{184}{135}
 \end{aligned}$$

Therefore, the ACF for $h \geq 1$ is given by

$$\rho_X(h) = \begin{cases} 1 & h = 0 \\ \frac{\gamma_X(h)}{\gamma_X(0)} = \frac{\frac{91}{135}(-1)^h 4^{1-h}}{\frac{184}{135}} = \frac{91}{23}(-1)^h 2^{-2h-1} & h \geq 1 \end{cases}$$

Let's verify this in R.

```
round(ARMAacf(ar = -1 / 4, ma = -1 / 3, 5), 6)
h <- seq(1, 10, by = 1)
round((91 / 23) * (-1) ^ h * 2 ^ (-2 * h - 1), 6)
```

Output:

```

      0          1          2          3          4          5
1.000000 -0.494565  0.123641 -0.030910  0.007728 -0.001932
      -0.494565  0.123641 -0.030910  0.007728 -0.001932
```

As we can see, this is correct.

Chapter 4

Week 4

4.1 Stationary Process Forecasting

Suppose we observe a time series X_1, \dots, X_T that we believe has been generated by an underlying stationary process. We would like to produce an h -step ahead forecast

$$\hat{X}_{T+h} = \hat{X}_{T+h|T} = f(X_T, \dots, X_1)$$

forecasting X_{T+h} . Ideally, \hat{X}_{T+h} would minimize the prediction error

$$L(X_{T+h}, \hat{X}_{T+h}) = \min_f L(X_{T+h}, f(X_T, \dots, X_1))$$

where L is a loss function.

Frequently, the loss function is taken to be the *mean-squared error* (MSE)

$$L(X_{T+h}, \hat{X}_{T+h}) = \mathbb{E}[(X_{T+h} - \hat{X}_{T+h})^2]$$

When using MSE, it is natural to consider

$$L^2 = \{\text{Random variables } X : \mathbb{E}[X^2] < \infty\}$$

L^2 is a Hilbert space when equipped with the inner product

$$\langle X, Y \rangle = \mathbb{E}[XY]$$

Hilbert spaces are generalizations of Euclidean space (\mathbf{R}^d) in which the geometry and notation of projection are preserved.

$$\text{Proj}(X \rightarrow Y) = \langle X, Y \rangle Y$$

DEFINITION 4.1.1: Closed Linear Subspace

We say $\mathcal{M} \subseteq L^2$ is a **closed linear subspace**, if

- (i) Linearity: $X, Y \in \mathcal{M}$, $\alpha, \beta \in \mathbf{R}$ then $\alpha X + \beta Y \in \mathcal{M}$
- (ii) Closed: If $X_n \rightarrow X$ (in the sense that $\mathbb{E}[(X_n - X)^2] \rightarrow 0$), and $X_n \in \mathcal{M}$, then $X \in \mathcal{M}$.

THEOREM 4.1.2: Projection Theorem

If \mathcal{M} is a closed linear subspace in L^2 and $x \in L^2$, then there exists a unique $\hat{X} \in \mathcal{M}$ such that

$$\mathbb{E}[(X - \hat{X})^2] = \inf_{Y \in \mathcal{M}} \mathbb{E}[(X - Y)^2]$$

Moreover, \hat{X} satisfies the prediction equations/normal equations:

$$(X - \hat{X}) \in \mathcal{M}^\perp \implies \mathbb{E}[(X - \hat{X})Y] = 0 \quad (\forall Y \in \mathcal{M})$$

In MSE forecasting, we want to choose \hat{X}_{T+h} satisfying

$$\mathbb{E}[(X_{T+h} - \hat{X}_{T+h})^2] = \inf_{Y \in \mathcal{M}} \mathbb{E}[(X_{T+h} - Y)^2]$$

where \mathcal{M} is a closed linear subspace based on the available data.

(1) $\mathcal{M}_1 = \{z : z = f(X_T, \dots, X_1), f \text{ is any Borel Measurable function}\}$. In this case

$$\hat{X}_{T+h} = \mathbb{E}[X_{T+h} \mid X_T, \dots, X_1]$$

Unfortunately \mathcal{M}_1 is enormous and complicated!

(2) $\mathcal{M}_2 = \overline{\text{Span}}(1, X_T, \dots, X_1) = \{Y : Y = \alpha_0 + \sum_{j=1}^T \alpha_j X_j, \alpha_0, \dots, \alpha_T \in \mathbf{R}\}$ which is the linear functions of X_1, \dots, X_T . \hat{X}_{T+h} is called the **best linear predictor** (BLP).

4.2 Best Linear Prediction

Suppose X_t is a (weakly) stationary time series. Best linear prediction entails finding \hat{X}_{T+h} so that

$$\mathbb{E}[(X_{T+h} - \hat{X}_{T+h})^2] = \inf_{Y \in \mathcal{M}_2} \mathbb{E}[(X_{T+h} - Y)^2]$$

\hat{X}_{T+h} is the best prediction among all linear functions of X_T, \dots, X_1 .

DEFINITION 4.2.1: Projection

If \hat{X} satisfies

$$\mathbb{E}[(X - \hat{X})^2] = \inf_{Y \in \mathcal{M}} \mathbb{E}[(X - Y)^2]$$

we say that \hat{X} is the **projection** of X onto \mathcal{M} , and we write $\hat{X} = \text{Proj}(X \mid \mathcal{M})$.

In particular, the BLP is

$$\hat{X}_{T+h} = \text{Proj}(X_{T+h} \mid \mathcal{M}_2)$$

Consider the case when $h = 1$. From the Projection Theorem, the BLP is of the form

$$\hat{X}_{T+1} = \phi_{T,0} + \sum_{j=1}^T \phi_{T,j} X_j \approx \phi_{T,0} + \sum_{j=0}^T \phi_{T,j} (X_j - \mu)$$

where $\mu = \mathbb{E}[X_t]$. \hat{X}_{T+1} must satisfy the **prediction equations**,

$$\mathbb{E}[(X_{T+1} - \hat{X}_{T+1})Y] = 0 \quad (\forall Y \in \mathcal{M}_2)$$

In particular,

$$\mathbb{E}[(X_{T+1} - \hat{X}_{T+1})1] = 0 \quad (Y = 1)$$

$$\mathbb{E}[(X_{T+1} - \hat{X}_{T+1})X_j] = 0 \quad (1 \leq j \leq T, Y = X_j)$$

We have $T + 1$ equations. Since $\mathbb{E}[X_j - \mu] = 0$,

$$0 = \mathbb{E}[X_{T+1} - \hat{X}_{T+1}] = \mu - \phi_{T,0} + 0 \implies \phi_{T,0} = \mu$$

Before proceeding, note that this implies

$$\mathbb{E}[(X_{T+1} - \hat{X}_{T+1})X_j] = \mathbb{E}[(X_{T+1} - \mu - (\hat{X}_{T+1} - \mu))(X_j - \mu)]$$

So we may assume without loss of generality that $\mu = 0$, therefore $\mathbb{E}[X_i X_j] = \gamma(j - i)$. Therefore,

$$0 = \mathbb{E}[(X_{T+1} - \hat{X}_{T+1})X_k] = \gamma(T + 1 - k) - \sum_{j=1}^T \phi_{T,j} \gamma(j - k) \quad (1 \leq k \leq T)$$

Therefore, we have linear system of equations for $\phi_{T,1}, \dots, \phi_{T,T}$:

$$\sum_{j=1}^T \phi_{T,j} \gamma(j - k) = \gamma(T + 1 - k)$$

Let

$$\gamma_T = \begin{pmatrix} \gamma(T) \\ \vdots \\ \gamma(1) \end{pmatrix} \in \mathbf{R}^T$$

$$\Gamma_T = [\gamma(j - k), 1 \leq j, k \leq T] \in \mathbf{R}^{T \times T}$$

$$\phi_T = (\phi_{T,1}, \dots, \phi_{T,T})^\top \in \mathbf{R}^T$$

this linear system may be expressed as

$$\Gamma_T \phi_T = \gamma_T \implies \phi_T = \Gamma_T^{-1} \gamma_T$$

given that Γ_T is invertible.

The BLP is of the form

$$\hat{X}_{T+1} = \phi_T^\top \mathbf{X}_T = (\Gamma_T^{-1} \gamma_T)^\top \mathbf{X}_T$$

where $\mathbf{X}_T = (X_T, \dots, X_1)^\top \in \mathbf{R}^T$.

When is Γ_T non-singular?

THEOREM 4.2.2

If $\gamma(0) > 0$, and $\gamma(h) \rightarrow 0$ as $h \rightarrow \infty$, then Γ_T is non-singular.

Takeaway: Most stationary processes (those whose serial dependence decays over time) have non-singular Γ_T .

Note that

$$\hat{X}_{T+1}^2 = \gamma_T^\top \Gamma_T^{-1} \mathbf{X}_T \mathbf{X}_T^\top \Gamma_T^{-1} \gamma_T$$

Note that $\mathbb{E}[\mathbf{X}_T \mathbf{X}_T^\top] = \Gamma_T$. Therefore, $\mathbb{E}[\hat{X}_{T+1}^2] = \gamma_T^\top \Gamma_T^{-1} \gamma_T$. Also, since

$$\mathbb{E}[X_{T+1} \mathbf{X}_T] = \gamma_T \implies \mathbb{E}[X_{T+1} \hat{X}_{T+1}] = \gamma_T^\top \Gamma_T^{-1} \gamma_T$$

It follows that the mean-squared prediction error is

$$\begin{aligned} P_{T+1}^T &= \mathbb{E}[(X_{T+1} - \hat{X}_{T+1})^2] \\ &= \mathbb{E}[X_{T+1}^2 - 2X_{T+1} \hat{X}_{T+1} + \hat{X}_{T+1}^2] \\ &= \gamma(0) - 2\gamma_T^\top \Gamma_T^{-1} \gamma_T + \gamma_T^\top \Gamma_T^{-1} \gamma_T \\ &= \gamma(0) - \gamma_T^\top \Gamma_T^{-1} \gamma_T \end{aligned}$$

The mean-squared prediction error has a simple, computable form depending on $\gamma(h)$ for $1 \leq h \leq T$.

4.3 Partial ACF

If $X_t \sim \text{ARMA}(p, q)$, then we might be able to identify p, q by looking at the ACF.

$$X_t \sim \text{AR}(p) \implies \text{ACF has a geometric decay}$$

$$X_t \sim \text{MA}(q) \implies \text{ACF is non-zero at the first } q \text{ lags, then zero beyond}$$

ACF of an $\text{ARMA}(p, q)$ model can be calculated by calculating the linear process coefficients $\{\psi\}_{\ell=0}^{\infty}$. Automated in R using `ARMAacf()`.

In Figure 4.1, it looks like geometric decay. However, it is hard to tell the difference between the $\text{ARMA}(1, 1)$ process and the $\text{AR}(p)$ process via the ACF. Therefore, we want to define the *partial autocorrelation function*.

```
# Figure 4.1 (Omitted the PACF)
ACF = ARMAacf(ar = c(.8), ma = 1, 24)[-1]
PACF = ARMAacf(ar = c(.8),
               ma = 1,
               24,
               pacf = TRUE)
par(mfrow = c(1, 2))
plot(ACF,
     type = "h",
     xlab = "lag",
     ylim = c(-.8, 1))
abline(h = 0)
plot(PACF,
     type = "h",
     xlab = "lag",
     ylim = c(-.8, 1))
abline(h = 0)
```



Figure 4.1: $\text{ARMA}(1, 1)$: $X_t = 0.9X_{t-1} + W_t + 0.5W_{t-1}$

DEFINITION 4.3.1: Partial autocorrelation function

The **partial autocorrelation function** of a stationary process $\{X_t\}_{t \in \mathbb{Z}}$ is

$$\phi_{h,h} = \text{Corr}(X_{t+h} - \text{Proj}(X_{t+h} \mid X_{t+h-1}, \dots, X_{t+1}), X_t - \text{Proj}(X_t \mid X_{t+h-1}, \dots, X_{t+1}))$$

Interpretation: Autocorrelation between X_t and X_{t+h} after removing the linear dependence on the intervening variables $X_{t+h-1}, \dots, X_{t+1}$.

REMARK 4.3.2

If $X_t \sim \text{AR}(p)$, then $\phi_{h,h} = 0$ for $h \geq p + 1$.

Proof of: Remark 4.3.2

If $X_t \sim \text{AR}(p)$, then $X_{t+h} = \sum_{j=1}^p \phi_j X_{t+h-j} + W_{t+h}$.

$$\text{Proj}(X_{t+h} \mid X_{t+h-1}, \dots, X_{t+1}) = \sum_{k=1}^{h-1} \beta_k X_{t+h-k}$$

and minimizes

$$\begin{aligned} \mathbb{E} \left[\left(X_{t+h} - \sum_{k=1}^{h-1} \beta_k X_{t+h-k} \right)^2 \right] &= \mathbb{E} \left[\left(W_{t+h} + \sum_{j=1}^p \phi_j X_{t+h-j} - \sum_{k=1}^{h-1} \beta_k X_{t+h-k} \right)^2 \right] \\ &= \sigma_W^2 + \mathbb{E} \left[\left(\sum_{j=1}^p \phi_j X_{t+h-j} - \sum_{k=1}^{h-1} \beta_k X_{t+h-k} \right)^2 \right] \end{aligned}$$

where the second term is minimized by setting $\beta_j = \phi_j$ for $1 \leq j \leq p$ and $\beta_j = 0$ for $j \geq p$. Note that W_{t+h} is independent of other terms. Hence,

$$X_{t+h} - \text{Proj}(X_{t+h} \mid X_{t+h-1}, \dots, X_{t+1}) = W_{t+h} \quad (h \geq p + 1)$$

Therefore,

$$\phi_{h,h} = \text{Corr}(W_{t+h}, X_t - \text{Proj}(X_t \mid X_{t+h-1}, \dots, X_{t+1}))$$

which is independent by causality. Therefore, $\phi_{h,h} = 0$.

REMARK 4.3.3

It can be shown that if $X_t \sim \text{MA}(q)$ (invertible), then

$$\phi_{h,h} \neq 0$$

$$|\phi_{h,h}| = \mathcal{O}(r^h) \quad (0 < r < 1)$$

which is geometric decay.

	ACF	PACF
MA(q)	Cuts off after lag q	Geometric decay
AR(p)	Geometric decay	Cuts off after lag p

Estimating the PACF

Using the BLP theory,

$$\hat{\phi}_{h,h} = (\hat{\Gamma}_h^{-1} \hat{\gamma}_h)(h)$$

where

$$\begin{aligned} \hat{\Gamma}_h &= [\hat{\gamma}(j-k), 1 \leq j, k \leq h] \in \mathbf{R}^{h \times h} \\ \hat{\gamma}_h &= [\hat{\gamma}(1), \dots, \hat{\gamma}(h)] \in \mathbf{R}^h \end{aligned}$$

4.4 ARMA Forecasting

Suppose X_t follows a stationary and invertible ARMA(p, q) model so that $\phi(B)X_t = \theta(B)W_t$. Having observed X_T, \dots, X_1 , we wish to predict X_{T+h} .

$$\hat{X}_{T+h} = \text{Proj}(X_{T+h} \mid \mathcal{M}_2) \approx \mathbb{E}[X_{T+h} \mid X_T, \dots, X_1]$$

by causality and invertibility $X_t \sim$ linear function of W_t .

Furthermore,

$$\hat{X}_{T+h} \approx \tilde{X}_{T+h} = \mathbb{E}[X_{T+h} \mid X_T, \dots, X_1, X_0, \dots]$$

which is geometric decay of the dependence on past values.

Since X_t is casual and invertible,

$$X_t = \sum_{\ell=0}^{\infty} \psi_{\ell} W_{t-\ell}$$

$$W_t = \sum_{\ell=0}^{\infty} \pi_{\ell} X_{t-\ell}$$

where $\psi_0 = \pi_0 = 1$. Note that ψ 's and π 's are computable by solving homogeneous linear difference equations.

These representations imply,

$$\text{Information in } (X_T, X_{T-1}, \dots) = \text{Information in } (W_T, W_{T-1}, \dots)$$

So

$$\tilde{X}_{T+h} = \mathbb{E}[X_{T+h} \mid X_T, X_{T-1}, \dots] = \mathbb{E}[X_{T+h} \mid W_T, W_{T-1}, \dots]$$

$$\begin{aligned} \tilde{X}_{T+h} &= \mathbb{E} \left[\sum_{\ell=0}^{\infty} \psi_{\ell} W_{T+h-\ell} \mid W_T, W_{T-1}, \dots \right] \\ &= \mathbb{E} \left[\sum_{\ell=0}^{h-1} \psi_{\ell} W_{T+h-\ell} \mid W_T, \dots \right] + \mathbb{E} \left[\sum_{\ell=h}^{\infty} \psi_{\ell} W_{T+h-\ell} \mid W_T, \dots \right] \\ &= \sum_{\ell=h}^{\infty} \psi_{\ell} W_{T+h-\ell} \quad \text{since } \psi_{\ell} W_{T+h-\ell} = 0 \text{ for } 0 \leq \ell \leq h-1 \end{aligned}$$

Also, using invertibility,

$$0 = \mathbb{E}[W_{T+h} \mid X_T, X_{T-1}, \dots] = \mathbb{E} \left[\sum_{\ell=0}^{\infty} \pi_{\ell} X_{T+h-\ell} \mid X_T, \dots \right]$$

by independence, and furthermore, with $\pi_0 = 1$ we have

$$0 = \tilde{X}_{T+h} + \sum_{\ell=1}^{h-1} \pi_{\ell} \tilde{X}_{T+h-\ell} + \sum_{\ell=h}^{\infty} \pi_{\ell} X_{T+h-\ell}$$

Therefore,

$$\tilde{X}_{T+h} = - \sum_{\ell=1}^{h-1} \pi_{\ell} \tilde{X}_{T+h-\ell} - \sum_{\ell=h}^{\infty} \pi_{\ell} X_{T+h-\ell}$$

Truncated ARMA Prediction

$$\hat{X}_{T+h} = - \sum_{j=1}^{h-1} \pi_j \hat{X}_{T+h-j} - \sum_{j=h}^{T+h-1} \pi_j X_{T+h-j}$$

Residuals:

$$\hat{W}_t = \phi(B)\hat{X}_t - \theta_1 \hat{W}_{t-1} - \dots - \theta_q \hat{W}_{t-q}$$

Mean initialization:

- $\hat{W}_t = 0$ for $t \leq 0$ and $t \geq T$.
- $\hat{X}_t = 0$ for $t \leq 0$ and $t \geq T+1$.
- $\hat{X}_t = X_t$ for $1 \leq t \leq T$.

Estimator for σ_W^2 :

$$\hat{\sigma}_W^2 = \frac{1}{T} \sum_{t=1}^T \hat{W}_t^2$$

Mean Squared Prediction Error: Since $\hat{X}_{T+h} \approx \sum_{j=h}^{\infty} \psi_j W_{T+h-j}$,

$$P_{T+h}^T = \mathbb{E}[(X_{T+h} - \hat{X}_{T+h})^2] = \mathbb{E}\left[\left(\sum_{j=0}^{h-1} \psi_j W_{T-j}\right)^2\right] = \sigma_W^2 \sum_{j=0}^{h-1} \psi_j^2$$

Estimated Mean Squared Prediction Error:

$$\hat{P}_{T+h}^T = \hat{\sigma}_W^2 \sum_{j=0}^{h-1} \psi_j^2$$

Construction of Prediction Intervals: Since $\hat{X}_{T+h} \approx \mathbb{E}[X_{T+h} | X_T, \dots]$,

$$\mathbb{E}[\hat{X}_{T+h} - X_{T+h}] = 0 \quad (\text{Tower Property})$$

$$\mathbb{E}[(\hat{X}_{T+h} - X_{T+h})^2] = P_{T+h}^T$$

Hence

$$\frac{\hat{X}_{T+h} - X_{T+h}}{\sqrt{\hat{P}_{T+h}^T}}$$

is an approximately mean zero and unit variance random variable.

Suppose c_α is the α -critical value of this random variable, then

$$\hat{X}_{T+h} \pm c_{\alpha/2} \sqrt{\hat{P}_{T+h}^T}$$

is an approximate $(1 - \alpha)$ prediction interval for X_{T+h} .

Choices for c_α :

- (1) z_α (standard normal critical value).

Motivation: If W_t is Gaussian, then $X_t = \sum_{\ell=0}^{\infty} \psi_\ell W_{t-\ell}$ is Gaussian.

- (2) Empirical critical value of residuals (standardized)

$$\frac{\hat{W}_t}{\sigma_W} \quad (1 \leq t \leq T)$$

- (3) t -distribution, Pareto, or skewed distribution fit to standardized residuals.

Long Range Behaviour of ARMA Forecasts

Suppose $Y_t = s_t + X_t$ where $X_t \sim \text{ARMA}(p, q)$.

$$\hat{Y}_{T+h} = \hat{s}_{T+h} + \hat{X}_{T+h} = \hat{s}_{T+h} + \underbrace{\sum_{j=h}^{\infty} \psi_j W_{T+h-j}}_{\rightarrow 0 \text{ (geometrically)}}$$

\hat{Y}_{T+h} is converging fast to \hat{s}_{T+h} . Therefore, when we are doing ARMA forecasting in a trend + noise framework, we better get the trend correct for long range forecasts. Long range forecasts are only going to depend on the trend, and very little on the noise because we know that ARMA processes have a geometric decay to their dependent structure.

$$P_{T+h}^T = \sigma_W^2 \sum_{\ell=0}^{h-1} \psi_{\ell}^2 \xrightarrow{h \rightarrow \infty} \sigma_W^2 \sum_{\ell=0}^{\infty} \psi_{\ell}^2 = \gamma_X(0) = \sigma_W^2$$

In the long run, the MSE is the variance of X_t .

4.5 ARMA Forecasting Example 1: Cardiovascular Mortality

[R Code] Cardiovascular Mortality

Slide 1

Let's give ARMA forecasting a try on real data.

Slide 2



Figure 4.2: Weekly cardiovascular mortality, LA County.

Slide 3

Let X_t = cardiovascular mortality series. Our model is

$$X_t = s_t + Y_t$$

where $Y_t \sim \text{ARMA}(p, q)$.

$$X_t = \underbrace{\beta_0 + \beta_1 t + \beta_2 t^2 + \beta_3 t^3}_{\text{polynomial}} + \underbrace{\beta_4 \sin\left(\frac{2\pi}{52}t\right) + \beta_5 \cos\left(\frac{2\pi}{52}t\right) + \beta_6 \sin\left(\frac{2\pi}{26}t\right) + \beta_7 \cos\left(\frac{2\pi}{26}t\right)}_{\text{seasonal}}$$

where the first four terms are the polynomial trends, the next two terms are the yearly cycle, and the last two are the half-yearly cycle.

Decided on the trend using AIC, which will be discussed next week.

Slide 4

s_t estimated using ordinary least squares.



Slide 5



Series residuals(reg2)



- $\hat{Y}_t = X_t - \hat{s}_t$ “seems reasonably stationary.”
- Mild serial correlation in \hat{Y}_t — Might be well modelled by MA(2) or ARMA(1, 1).

Slide 6

Normal Q-Q Plot



- \hat{Y}_t seems reasonably normal, suggests using

$$\pm Z_{\alpha/2} \sqrt{P_{T+h}^T}$$

to construct prediction bounds.

Slide 7

Considering the PACF: On the first two lags these are large which is indicative of an autoregressive 2 structure, that is, AR(2) structure.

Slide 8

Model \hat{Y}_t as ARMA(2, 1).

$$Y_t = 0.0885Y_{t-1} + 0.3195Y_{t-2} + W_t + 0.1328W_{t-1}$$

parameters estimated by MLE.

Slide 9



Slide 10



$\hat{Y}_{T+h|T}$, $h = 1, \dots, 10$.

$$\hat{Y}_{T+h|T} \pm 1.96\sqrt{\hat{P}_{T+h}^T}$$

where 1.96 is the 97.5% critical value of $\mathcal{N}(0, 1)$.

Slide 11



Figure 4.3: 30 weeks of data with predicted trend

Slide 12



Figure 4.4: Forecasts with 95% prediction intervals

Fluctuations attribute to serial dependence. Red lines show that forecasts quickly converge to trend.

4.6 ARMA Forecasting Example 2: Johnson and Johnson

[R Code] Johnson and Johnson

X_t Johnson and Johnson Earnings.

$$X_t = e^{s_t + Y_t}$$

where Y_t is stationary. In this case,

$$\log(X_t) = s_t + Y_t$$

where $Y_t \sim \text{ARMA}(p, q)$.

Chapter 5

Week 5

5.1 ARMA Parameter Estimation: AR Case

Suppose we observe a time series $X_1, \dots, X_T \sim \text{ARMA}(p, q)$

$$\phi(B)X_t = \theta(B)X_t$$

$$\phi(z) = 1 - \phi_1 z - \dots - \phi_p z^p$$

$$\theta(z) = 1 + \theta_1 z + \dots + \theta_q z^q$$

Our goal is to estimate

- ϕ_1, \dots, ϕ_p (AR parameters)
- $\theta_1, \dots, \theta_q$ (MA parameters)
- σ_W^2 (white noise variance)

AR(1) case: $X_t = \phi X_{t-1} + W_t$ with $\mathbb{E}[W_t^2] = \sigma_W^2$. The idea is to use OLS.

$$\hat{\phi} = \arg \min_{|\phi| < 1} \sum_{t=2}^T (X_t - \phi X_{t-1})^2$$

This leads to (upon some calculations):

$$\hat{\phi} = \frac{\frac{1}{T} \sum_{t=2}^T X_t X_{t-1}}{\frac{1}{T} \sum_{t=2}^T X_t^2} \approx \frac{\hat{\gamma}(1)}{\hat{\gamma}(0)} = \hat{\rho}(1) \xrightarrow{T \rightarrow \infty} \phi$$

$$\hat{\sigma}_W^2 = \frac{1}{T-1} \sum_{t=2}^T (X_t - \hat{\phi} X_{t-1})^2$$

where $X_t - \hat{\phi} X_{t-1}$ is estimated W_t and $\hat{\sigma}_W^2$ is the sample variance of residuals.

AR(p) case: $X_t = \phi_1 X_{t-1} + \dots + \phi_p X_{t-p} + W_t$. OLS: $\phi = (\phi_1, \dots, \phi_p)^\top \in \mathbf{R}^p$

$$\hat{\phi} = \arg \min_{\hat{\phi}} \sum_{t=p+1}^T (X_t - \phi_1 X_{t-1} - \dots - \phi_p X_{t-p})^2$$

$\hat{\phi}$ admits a stationary and causal solution.

Solve using calculus (take first order partial derivatives set equal to zero), leads to a system of p linear equations of the form

$$\hat{\Gamma}_p \hat{\phi} = \hat{\gamma}_p$$

where

$$\hat{\Gamma}_p = (\hat{\gamma}(j-k), 1 \leq j, k \leq p) \in \mathbf{R}^{p \times p}$$

$$\hat{\gamma}_p = (\hat{\gamma}(1), \dots, \hat{\gamma}(p))^{\top}$$

The resulting OLS estimator takes the form

$$\hat{\phi} = \hat{\Gamma}_p^{-1} \hat{\gamma}_p$$

$$\hat{\sigma}_W^2 = \hat{\gamma}(0) - \hat{\gamma}_p^{\top} \hat{\Gamma}_p^{-1} \hat{\gamma}_p$$

Similar approach: use method of moments (set parameters so that empirical moments match theoretical causal moments induced by the model).

If $X_t \sim \text{AR}(p)$, then for $1 \leq h \leq p$.

$$\begin{aligned} \gamma(h) &= \mathbb{E}[X_t X_{t+h}] \\ &= \mathbb{E}[X_t (\phi_1 X_{t+h-1} + \dots + \phi_p X_{t+h-p} + W_{t+h})] \\ &= \phi_1 \gamma(h-1) + \phi_2 \gamma(h-2) + \dots + \phi_p \gamma(h-p) + 0 \end{aligned}$$

where the 0 occurs since $X_t \perp\!\!\!\perp W_{t+h}$.

This implies the linear system:

$$\gamma_p = \Gamma_p \phi$$

$$\gamma_p = (\gamma(1), \dots, \gamma(p))^{\top} \in \mathbf{R}^p$$

$$\Gamma_p = [\gamma(j-k), 1 \leq j, k \leq p] \in \mathbf{R}^{p \times p}$$

Note that $X_t = \sum_{\ell=0}^{\infty} \psi_{\ell} W_{t-\ell}$ where $\psi_0 = 1$ and $W_t = X_t - \phi_1 X_{t-1} - \dots - \phi_p X_{t-p}$ imply

$$\sigma_W^2 = \mathbb{E}[X_t W_t] = \mathbb{E}[X_t (X_t - \phi_1 X_{t-1} - \dots - \phi_p X_{t-p})] = \gamma(0) - \phi_1 \gamma(1) - \dots - \phi_p \gamma(p)$$

which are **Yule-Walker Equations**.

$$\gamma_p = \Gamma_p \phi$$

Yule-Walker Estimators:

$$\hat{\phi} = \hat{\Gamma}_p^{-1} \hat{\gamma}_p$$

$$\hat{\sigma}_W^2 = \hat{\gamma}(0) - \hat{\gamma}_p^{\top} \hat{\Gamma}_p^{-1} \hat{\gamma}_p$$

EXAMPLE 5.1.1

In the AR(1) case, the YW estimators are

$$\hat{\phi} = \frac{\hat{\gamma}(1)}{\hat{\gamma}(0)} = \hat{\rho}(1)$$

$$\hat{\sigma}_W^2 = \hat{\gamma}(0) - \frac{\hat{\gamma}^2(1)}{\hat{\gamma}(0)}$$

THEOREM 5.1.2

If $X_t \sim \text{AR}(p)$ (causal), then

$$\frac{\hat{\phi}_{OLS, i}}{\hat{\phi}_{YW, i}} \xrightarrow[T \rightarrow \infty]{P} 1$$

OLS and YW estimates are asymptotically equivalent.

THEOREM 5.1.3

$$\sqrt{T}(\hat{\phi}_{YW} - \phi) \xrightarrow[T \rightarrow \infty]{D} MVN(0, \sigma_W^2 \Gamma_p^{-1})$$

$$\hat{\sigma}_W^2 \xrightarrow{P} \sigma_W^2$$

- Optimal variance among all possible (asymptotically) unbiased estimators, hence **efficient**.
- Result can be used to obtain confidence intervals for ϕ .

5.2 ARMA Parameter Estimation: MLE

Ordinary Least Squares and Yule-Walker equation estimators are effective in estimating the $AR(p)$ parameters, but are difficult to apply to fitting $MA(q)$ and general $ARMA(p, q)$ models since the white noises W_t are not observable, and YW equations are not linear in the MA parameters.

Latent Variables (variables associated with W_t) \implies MLE is best.

Suppose $X_t \sim AR(1)$ (causal)

$$X_t = \phi X_{t-1} + W_t$$

where $W_t \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma_W^2)$, then

$$X_t = \sum_{\ell=0}^{\infty} \phi^\ell W_{t-\ell}$$

is Gaussian. L^2 limits of Gaussian random variables are Gaussian. (MGF or characteristic function).

Moreover, X_1, \dots, X_T are jointly Gaussian since

$$a_1 X_1 + \dots + a_T X_T = \sum_{\ell=0}^{\infty} \phi^\ell \underbrace{(a_1 W_{1-\ell} + \dots + a_T W_{t-\ell})}_{\text{Gaussian}}$$

MLE:

$$\mathcal{L}(\phi, \sigma_W^2) = f(X_T, X_{T-1}, \dots, X_1; \phi, \sigma_W^2)$$

where

- $\mathcal{L}(\phi, \sigma_W^2)$ is the likelihood of ϕ and σ_W^2 .
- $f(X_T, X_{T-1}, \dots, X_1; \phi, \sigma_W^2)$ is the joint density of X_T, \dots, X_1 at the observed data. Gaussian Density.

Key idea in Time series: To evaluate the likelihood condition on the path/past!

$$\begin{aligned} f(X_T, \dots, X_1) &= f(X_T | X_{T-1}, \dots, X_1) f(X_{T-1}, \dots, X_1) \\ &\vdots \\ &= f(X_T | X_{T-1}, \dots, X_1) f(X_{T-1} | X_{T-2}, \dots, X_1) \dots f(X_2 | X_1) f(X_1) \\ &= \prod_{i=1}^T f(X_i | X_{i-1}, \dots, X_1) \end{aligned} \quad \text{iterate}$$

According to HW2:

$$X_i | (X_{i-1}, \dots, X_1) \sim \mathcal{N}(\phi X_{i-1}, \sigma_W^2)$$

Note that $X_i | (X_{i-1}, \dots, X_1) = X_i | X_{i-1}, AR(1)$.

Thus,

$$\begin{aligned}\mathcal{L}(\phi, \sigma_W^2) &= \prod_{i=2}^T \frac{1}{\sqrt{2\pi\sigma_W^2}} \exp\left\{-\frac{(X_i - \phi X_{i-1})^2}{2\sigma_W^2}\right\} f(X_1) \\ &= (2\pi\sigma_W^2)^{-\frac{T-1}{2}} \exp\left\{-\frac{\sum_{i=2}^T (X_i - \phi X_{i-1})^2}{2\sigma_W^2}\right\} f(X_1; \phi, \sigma_W^2)\end{aligned}$$

Maximizing $\mathcal{L}(\phi, \sigma_W^2)$ in this case leads to a similar estimator as OLS/YW.

General ARMA(p, q) case: Again, X_T, \dots, X_1 are jointly Gaussian if $W_t \sim \text{Gaussian}$.

$$\begin{aligned}L(\underbrace{\phi_1, \dots, \phi_p, \theta_1, \dots, \theta_q, \sigma_W^2}_{\boldsymbol{\theta} \in \mathbb{R}^{p+q+1}}) &= \prod_{i=1}^T \underbrace{f(X_i | X_{i-1}, \dots, X_1)}_{\text{Gaussian}} \\ X_i | (X_{i-1}, \dots, X_1) &\sim \mathcal{N}(\mathbb{E}[X_i | X_{i-1}, \dots, X_1], \text{MSE}) \\ &\sim \mathcal{N}(\tilde{X}_{i|(i-1)}(\boldsymbol{\theta}), P_{i-1}^i(\boldsymbol{\theta}))\end{aligned}$$

where $P_{i-1}^i(\boldsymbol{\theta})$ is forecast MSE predicting X_i from X_{i-1}, \dots, X_1 .

This likelihood can be maximized using numerical optimization. (Newton-Raphson Algorithm, Conjugate Gradient).

THEOREM 5.2.1: Chapter 8 of Brockwell and Davis, Hannan (1980)

The MLE's of $\phi_1, \dots, \phi_p, \theta_1, \dots, \theta_q, \sigma_W^2$ are \sqrt{T} consistent and asymptotically Normal with asymptotic covariance equal to the inverse of the information matrix. In this sense, they are asymptotically optimal.

REMARK 5.2.2: Takeaway Message

- (1) MLE estimation reduces to OLS, YW equation estimation for AR(p) models.
- (2) For general ARMA(p, q) estimation, MLE is through to be optimal in most situations. (Used as a default/benchmark).

5.3 Model Selection Diagnostic Tests

Using MLE, we can fit an ARMA(p, q) model to an observed series X_1, \dots, X_T .

Question: How do we select the orders p and q of the model?

Usual Methods

- (1) Examine ACF and PACF.
- (2) Model Diagnostics/Goodness-of-Fit tests: Examine the residuals of the ARMA(p, q) model to check for the plausibility of the white noise assumption.
- (3) Model Selection Methods: Information criteria, cross-validation.

Model Diagnostics

If the ARMA(p, q) model fits the data well, then the estimated residuals should behave like white noise.

$$\hat{W}_t = \frac{X_t - \tilde{X}_{t|(t-1)}}{\sqrt{\hat{P}_t^{t-1}}}$$

where

- $\tilde{X}_{t|(t-1)}$ is the truncated predictor of X_t based on X_{t-1}, \dots, X_1 , and
- \hat{P}_t^{t+1} is the estimated MSE.

This can be investigated by considering $\hat{\rho}_W(h)$ which is the empirical ACF of $\hat{W}_1, \dots, \hat{W}_T$.

As a measure of how “white” the residuals are, it is common to evaluate the cumulative significance of $\hat{\rho}_W(h)$ for $1 \leq h \leq H$ by applying a “white noise test.” Suppose W_1, \dots, W_T is a strong white noise, and $\hat{\rho}_W(h)$ is the empirical ACF of this series.

We know that for each fixed h ,

$$\sqrt{T}\hat{\rho}_W(h) \xrightarrow{D} \mathcal{N}(0, 1)$$

Also, for $j \neq h$,

$$\begin{aligned} \text{Cov}(\sqrt{T}\hat{\gamma}_W(h), \sqrt{T}\hat{\gamma}_W(j)) &= T\mathbb{E}\left[\sum_{t=1}^T W_t W_{t+h}\right] \mathbb{E}\left[\sum_{s=1}^T W_s W_{s+j}\right] \\ &= T \sum_{t=1}^T \sum_{s=1}^T \mathbb{E}[W_t W_{t+h} W_s W_{s+j}] \\ &= 0 \end{aligned}$$

Using Martingale, or m -dependent CLT's, it can be shown that

$$\begin{pmatrix} \sqrt{T}\hat{\rho}_W(1) \\ \vdots \\ \sqrt{T}\hat{\rho}_W(H) \end{pmatrix} \xrightarrow{D} \text{MVN}(0, I_{H \times H})$$

Therefore,

$$T \sum_{h=1}^H \hat{\rho}_W^2(h) \xrightarrow{D} \chi^2(H)$$

Box-Ljung-Pierce Test [White Noise Test for ARMA(p, q) Models]

If $X_t \sim \text{ARMA}(p, q)$, and \hat{W}_t are the model residuals with empirical ACF $\hat{\rho}_W(h)$, then if

$$Q(T, H) = T(T+2) \sum_{h=1}^H \frac{\hat{\rho}_W^2(h)}{T-h} \approx T \sum_{h=1}^H \hat{\rho}_W^2(h)$$

$$Q(T, H) \xrightarrow[T \rightarrow \infty]{D} \chi^2(H - (p+q))$$

That is, we lose $p+q$ degrees of freedom for fitting the model.

The BLP test p -value is then computed as

$$P_{\text{BLP}} = \mathbb{P}(\chi^2(H - (p+q)) > Q(T, H))$$

REMARK 5.3.1

If $X_t \sim \text{ARMA}(p, q)$, and \hat{W}_t are calculated based on $\text{ARMA}(p', q')$ model where $p' < p$ or $p' < q$ (model is under specified), then

$$Q(T, H) \xrightarrow[T \rightarrow \infty]{P} \infty$$

Interpretation: If BLP p -values are small, the model is ill-fitting or under specified.

5.4 Model Selection Information Criteria

Suppose we are trying to select the orders p and q of an ARMA(p, q) model to fit X_1, \dots, X_T .

ϕ = AR parameters

θ = MA parameters

σ_W^2 = white noise variance

$$\mathcal{L}(X_1, \dots, X_T; \hat{\phi}, \hat{\theta}, \sigma_W^2)$$

Natural idea: Maximize the likelihood of the data as a function of p and q .

Problem: The likelihood is (monotonically) increasing as a function of p and q . Maximizing would lead to overfitting.

Solution: Maximize the likelihood subject to a penalty term on the number of parameters (complexity) of the model. Let the number of parameters in the ARMA(p, q) model be denoted by $k = p + q + 1$.

$$-2 \log(\mathcal{L}(X_1, \dots, X_T; \hat{\phi}, \hat{\theta}, \sigma_W^2)) + P(T, k)$$

where $P(T, k)$ is an increasing function of k .

Optimal p and q balance model fit with the penalty for complexity.

Common Penalty Term Choices

- $\text{AIC}(p, q) = -2 \log(\mathcal{L}(X_1, \dots, X_T; \hat{\phi}, \hat{\theta}, \sigma_W^2)) + \frac{2k + T}{T}$.
 - Comes from estimating the Kullback–Leibler distance from the fitted model to the “true” model.
- $\text{BIC}(p, q) = -2 \log(\mathcal{L}(X_1, \dots, X_T; \hat{\phi}, \hat{\theta}, \sigma_W^2)) + \frac{k \log(T)}{T}$.
 - Comes from approximating and maximizing the posterior distribution of the model given the data.

Interpretation: Small AIC/BIC mean a better model.

Information criteria are also used in trend fitting. Suppose

$$X_t = s_t + Y_t = f_t(\beta) + Y_t$$

where $\beta \in \mathbf{R}^k$ and $f_t(\beta)$ is the trend we fit.

Estimate β with $\hat{\beta}$ using ordinary least squares.

$$\text{SS}(\text{Res})_T = \sum_{t=1}^T (X_t - f_t(\hat{\beta}))^2$$

Information criteria typically calculated assuming Y_t is a Gaussian white noise.

$$\text{SS}(\text{Res})_T + P(T, k)$$

where for $P(T, k)$ we use AIC or BIC penalty.

REMARK 5.4.1

- (1) In trend fitting, the assumption of Gaussian white noise residuals is often in doubt.
- (2) AIC/BIC are not perfect! They are but one of many tools useful in model selection.
Strengths:
 - (i) Easy to compute.
 - (ii) Facilitates comparing many models quickly.**Weaknesses:**
 - (i) Likelihood must be specified.
 - (ii) There is a degree of “arbitrariness” to the choice of penalty.
- (3) It can be shown that minimizing the AIC is related to minimizing the 1-step forecast MSE, and so when the application is forecasting, AIC is more common.

5.5 ARIMA Models

We have seen that many time series appear stationary after differencing.

DEFINITION 5.5.1: Integrated

We say a time series X_t is **integrated** to order d if $\nabla^d X_t$ is stationary, but $\nabla^j X_t$ for $1 \leq j \leq d$ is not stationary.

Motivation: If Y_t is stationary, and $X_t = \sum_{j=1}^t Y_j$, X_t is integrated to order 1. $Z_t = \sum_{j=1}^t X_j$ is integrated to order 2, and so on.

DEFINITION 5.5.2: ARIMA

We say X_t follows an **Autoregressive Integrated Moving Average Process** (ARIMA) of orders p, d, q if

$$\phi(B)(1-B)^d X_t = \theta(B)W_t$$

and write $X_t \sim \text{ARIMA}(p, d, q)$. Note that $\nabla^d X_t$ follows an $\text{ARMA}(p, q)$ model.

Forecasting $\text{ARIMA}(p, d, q)$ Processes

- (1) $Y_t = \nabla^d X_t$ follows an $\text{ARMA}(p, q)$ model and so can be forecast using truncated ARMA prediction.
- (2) Forecasts $\hat{Y}_{T+h|T}$ can be used to forecast X_{T+h} by reversing the differencing.

EXAMPLE 5.5.3

For $d = 1$, $Y_{T+1} = X_{T+1} - X_T$ so $\hat{X}_{T+1|T} = X_T + \hat{Y}_{T+1|T}$. This can be iterated to produce longer Horizon forecasts.

Predicting MSE is approximately of the form

$$P_{T+h}^T \approx \sigma_W^2 \sum_{j=0}^{h-1} \psi_{j,*}^2$$

where $\psi_{j,*}^2$ is the coefficient of z^j in the power series expansion (centred at zero) of

$$\frac{\theta(z)}{\phi(z)(1-z)^d} \quad (|z| < 1)$$

Idea:

$$X_t \approx \frac{\theta(B)}{\phi(B)(1-B)^d} W_t$$

EXAMPLE 5.5.4

Let $X_t \sim \text{ARIMA}(0, 1, 0)$.

$$X_t - X_{t-1} = (1-B)X_t = W_t \implies X_t = X_{t-1} + W_t \implies X_t = \sum_{j=1}^t W_j$$

if we iterate t -times. If $Y_t = \nabla X_t$, then $\hat{Y}_{T+h|T} = 0$ (forecasting W_t 's). Therefore,

$$\hat{X}_{T+1|T} = X_t + \hat{Y}_{T+1|T} = X_T$$

Similarly, $\hat{X}_{T+h|T} = X_T$. The best predictor of random walk is last known location.
Prediction MSE:

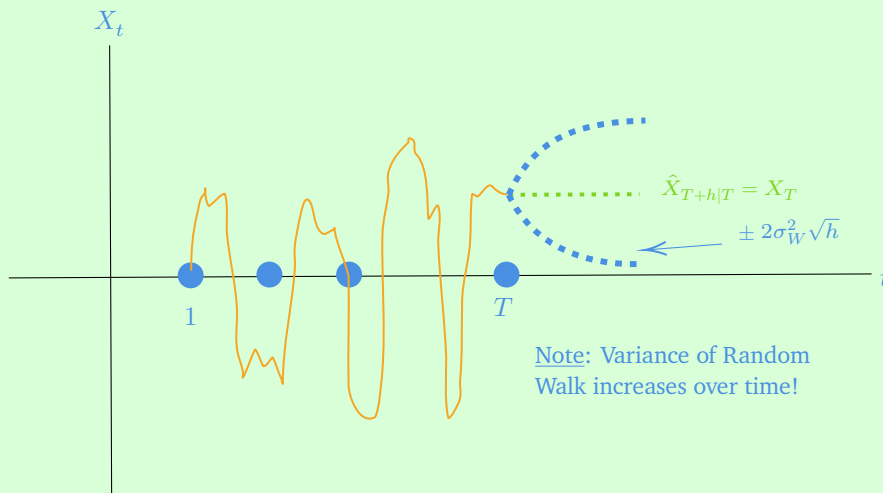
$$\frac{\theta(z)}{\phi(z)(1-z)^d} = \frac{1}{1-z} = \sum_{j=0}^{\infty} z^j \quad (|z| < 1)$$

$$\implies \psi_{j,*} = 1 \quad (\forall j)$$

$$\implies P_{T+h}^T = \sigma_W^2 \sum_{j=0}^{h-1} \psi_{j,*}^2 = h\sigma_W^2$$

Note that

$$\mathbb{E}[(\hat{X}_{T+h|T} - X_{T+h})^2] = \mathbb{E}\left[\left(\sum_{j=T+1}^{T+h} W_j\right)^2\right] = h\sigma_W^2$$



If we forecast into the future, the forecast will be the last observed value. Also, if we plot prediction intervals, they would be of the form $\pm 2\sigma_W^2 \text{MSE}$ where MSE which is on the order of \sqrt{h} . In particular, these are increasing as a function of h . Therefore, the variance of a Random Walk will increase over time, and hence the prediction intervals will increase over time. This is a normal feature you see when you do ARIMA forecasts, and this is the basic reason why.

How to decide in practice on degree of differencing d :

- (1) Eye-ball Test.
- (2) Formal Stationary Tests [Dickey-Fuller, Kwiatkowski-Phillips-Schmidt-Shin (KPSS)].

(3) Cross-validation.

5.6 ARIMA Modelling Example

[\[R Code\] ARIMA Modelling Example](#)

Chapter 6

Week 6

6.1 SARIMA Models

Frequently, time series exhibit “seasonality.”

Rough Definition of Seasonality

A time series X_t is said to be “seasonal” if it exhibits regular variation so that for some lag s , X_t is “similar” to X_{t-s} . Some sources of seasonality are weather or scheduled events. These typically lead to yearly, weekly, monthly, or quarterly cycles.

REMARK 6.1.1

ARIMA models are not ideal for modelling seasonality.

ARIMA Models \implies Random Walk with Stationary Errors

Random walks do not seasonality.

DEFINITION 6.1.2: Seasonal ARIMA

X_t is said to follow a **Seasonal ARIMA** model (SARIMA) of orders p, d, q and P, D, Q and seasonal period s if

$$\Phi_P(B^s)\phi_p(B)(1 - B^s)^D(1 - B)^d = \Theta_Q(B^s)\theta_q(B)W_t$$

We abbreviate the SARIMA p, d, q, P, D, Q model with seasonal period s as SARIMA(p, d, q) \times (P, D, Q) $_s$.

$$\begin{aligned}\Phi_P(B) &= 1 - \Phi_1 B - \dots - \Phi_P B^P \\ \Phi_P(B^s) &= 1 - \Phi_1 B^s - \dots - \Phi_P B^{Ps} \\ \phi_p(B) &= 1 - \phi_1 B - \dots - \phi_p B^p \\ \Theta_Q(B) &= 1 + \Theta_1 B + \dots + \Theta_Q B^Q \\ \Theta_Q(B^s) &= 1 + \Theta_1 B^s + \dots + \Theta_Q B^{Qs} \\ \theta_q(B) &= 1 + \theta_1 B + \dots + \theta_q B^q\end{aligned}$$

DEFINITION 6.1.3

The **seasonal** autoregressive and moving average polynomials are defined by

$$\begin{aligned}\Phi(z) &= 1 - \Phi_1 z - \dots - \Phi_P z^P \\ \Theta(z) &= 1 + \Theta_1 z + \dots + \Theta_Q z^Q\end{aligned}$$

EXAMPLE 6.1.4

Let $X_t \sim \text{SARIMA}(1, 1, 1) \times (1, 1, 1)_{13}$.

$$\Phi(z) = 1 - \Phi_1 z$$

$$\phi(z) = 1 - \phi_1 z$$

$$\Theta(z) = 1 + \Theta_1 z$$

$$\theta(z) = 1 + \theta_1 z$$

Therefore,

$$(1 - \Phi_1 B^{13})(1 - \phi_1 B) \underbrace{(1 - B^{13})(1 - B)X_t}_{Y_t} = \Theta(B^{13})\theta(B)W_t$$

$$Y_t - \Phi_1 Y_{t-13} - \phi_1 Y_{t-1} - \phi_1 \Phi_1 Y_{t-14} = \text{MA term}$$

$$Y_t = f(Y_{t-13}, Y_{t-1}, \text{MA noise}, Y_{t-14})$$

where Y_{t-13} is the seasonal lag.

REMARK 6.1.5

- (1) $Y_t = (1 - B^s)^D (1 - B)^d X_t$, a SARIMA model is just one big ARMA model for Y_t .
- (2) Advantage over ARMA and ARIMA models is **parsimony**. Since seasonal series have the feature that X_t is similar to X_{t-s} , we introduce just a few additional terms to model X_t as a function of X_{t-s} .

Fitting SARIMA Models

- (1) Usually the seasonal lag s is known.
- (2) Differencing and seasonal differencing can be decided upon by:
 - (a) Eye-ball test and/or examining the ACF and PACF.
 - (b) Stationarity tests.
 - (c) Cross-validation.

We will discuss (b) and (c).
- (3) Choosing the order and estimating the components of $\Phi, \phi, \Theta, \theta$ can be done in the same way as with ARMA models.

6.2 SARIMA Cardiovascular Mortality Example

[R Code] SARIMA Cardiovascular Mortality Example

6.3 Time Series Cross-Validation

DEFINITION 6.3.1: Cross-validation

Cross-validation is a data driven model evaluation and selection tool for predictive models that entails the following.

- (1) Splitting the available data into training and testing sets.
- (2) Fitting models on the training sets.
- (3) Evaluating predictions of the model on the tests sets as an overall evaluation of model quality.

Standard Cross-Validation

Suppose (Y_i, X_i) for $1 \leq i \leq n$ satisfy $Y_i = f(X_i) + \varepsilon_i$. Let M be a model used to estimate f using \hat{f} , with the goal of minimizing $L(Y_i, \hat{f}(X_i))$.

K -fold Cross-Validation

- (1) Split (Y_i, X_i) for $1 \leq i \leq n$ randomly into K -groups G_1, \dots, G_k .
- (2) For each $1 \leq i \leq K$, use M to estimate $\hat{f}^{(-j)}$ when the data G_i is left out.
- (3) Evaluate error on G_i with

$$CV_j = \sum_{(Y_i, X_i) \in G_j} L(Y_i, \hat{f}^{(-j)}(X_i))$$

- (4) The total cross-validation error of the model is:

$$CV(M) = \sum_{j=1}^k CV_j$$

REMARK 6.3.2

- K is often called the number of **folds**.
- If $K = n$, the procedure is often called the “leave-one-out” cross-validation.
- $K = 10$ is called “10-fold cross validation.”

Problems with Time Series Cross-Validation

- (1) Randomly splitting the data scrambles up any serial dependence relationships.
- (2) In time series forecasting, it is often most natural to use the past (recent past) to predict future values.

Time Series Cross-Validation Algorithm

- (1) Split the data into training and testing ranges $1 \leq t_r \leq T$ where $t_r \approx 0.75T$ is 75% of the training sample. The test sample is X_{t_r+1}, \dots, X_T .
- (2) For each j in $t_r + 1, \dots, T$, use model to forecast $\hat{X}_{j+1|j}$ based on X_1, \dots, X_j . Calculate loss

$$L(\hat{X}_{j+1|j}; X_{j+1}) = L_j$$

- (3) Cross-validation score of model

$$CV(M) = \sum_{j=t_r+1}^T L_j$$

REMARK 6.3.3

- (1) If interested in longer horizon forecasting, you can compare

$$\hat{X}_{j+1|j}, \dots, \hat{X}_{j+h|j} \quad \text{to} \quad X_{j+1}, \dots, X_{j+h}$$

in the loss calculation step.

- (2) Stationarity is *crucial* in time series cross validation since the model errors in the present must be similar to errors in the future.
 (3) One normally cannot cross-validate everything as this is computationally infeasible.

6.4 Cross-Validation Example

[R Code] Cross-Validation Example

6.5 Simulated and Bootstrapped Prediction Intervals

Usually forecasts are of the form

$$\hat{X}_{T+1|T} = g(X_T, X_{T-1}, \dots, X_1, W_{T+1})$$

where W_{T+1} is a strong white noise innovation.

Often, even models are additive so that

$$\hat{X}_{T+1|T} = g(X_T, \dots, X_1) + W_{T+1}$$

Simple and powerful models to produce prediction intervals use simulation!

Simulated Prediction Intervals

- (1) Choose a distribution for $\{W_t\}$. A common choice is $W_t \sim \mathcal{N}(0, \hat{\sigma}_W^2)$.
- (2) For $b = 1, \dots, B$ where B is a large number, simulate $\{W_{T+1}^{(b)}\}$.
- (3) Compute $\hat{X}_{T+1|T}^{(b)} = g(X_T, \dots, X_1) + W_{T+1}^{(b)}$ for $b = 1, \dots, B$.
- (4) Denote the empirical q^{th} quantile of $\{\hat{X}_{T+1}^{(b)} : b = 1, \dots, B\}$ by $\hat{Q}_{T+1}(q)$. We set the $(1 - \alpha)$ prediction interval as

$$\left(\hat{Q}_{T+1}\left(\frac{\alpha}{2}\right), \hat{Q}_{T+1}\left(1 - \frac{\alpha}{2}\right) \right)$$

REMARK 6.5.1

For longer horizon forecasts, prediction intervals can be obtained by iteration:

$$\hat{X}_{T+h|T}^{(b)} = g(\hat{X}_{T+h-1|T}^{(b)}, \dots, \hat{X}_{T+1|T}^{(b)}, X_T, \dots, X_1) + W_{T+h}^{(b)}$$

The prediction interval is

$$\left(\hat{Q}_{T+h}\left(\frac{\alpha}{2}\right), \hat{Q}_{T+h}\left(1 - \frac{\alpha}{2}\right) \right)$$

where $\hat{Q}_{T+h}(q)$ the empirical q^{th} quantile of $\hat{X}_{T+h}^{(b)}$.

Distributions to Choose for W_t

- (1) $W_t \sim \mathcal{N}(0, \hat{\sigma}_W^2)$ where $\hat{\sigma}_W^2$ is estimated from residuals which leads to approximately the same “well known” prediction intervals.
- (2) A distribution fit to the estimated residuals \hat{W}_t ; e.g., a t -distribution, Pareto, etc.
- (3) The empirical distribution of the residuals \hat{W}_t ; that is, randomly drawing $\{\hat{W}_1, \dots, \hat{W}_T\}$ which is commonly known as **bootstrapping**.

Note: An important consideration of the bootstrap is that the residuals should be white! We can check the whiteness of the residuals using the ACF or a white noise test.

6.6 Bootstrap Prediction Intervals Example

[R Code] [Bootstrap Prediction Intervals Example](#)

Chapter 7

Week 7

7.1 Exponential Smoothing Models Introduction

- **ARIMA Models:** Model a time series, potentially after differencing towards stationarity, in terms of its autocorrelation (linear process).
- **Exponential Smoothing:** Flexibly model the trend and seasonality observed in a time series.

Simple Exponential Smoothing

Suppose we wish to forecast a time series X_1, \dots, X_T . Two extreme forecasts are

$$\hat{X}_{T+1|T} = X_T \quad [\text{Random Walk}]$$

$$\hat{X}_{T+1|T} = \bar{X} = \frac{1}{T} \sum_{t=1}^T X_t \quad [\text{IID Sequence}]$$

Compromise: *Exponential Smoothing*.

$$\hat{X}_{T+1|T} = \alpha X_T + \alpha(1 - \alpha)X_{T-1} + \dots + \alpha(1 - \alpha)^{T-1}X_1 \quad (\alpha \in [0, 1])$$

Weights applied to past observations decrease exponentially quickly.

Simple exponential smoothing can be stated as a recursive system of equations.

- Prediction Equation: $\hat{X}_{T+1} = \ell_T$.
- Smoothing/Level Equation: $\ell_T = \alpha X_T + (1 - \alpha)\ell_{T-1} = \ell_T(\alpha, \ell_0)$ which is a convex combination of last observed value and last prediction or “level.”
- Initial Condition: ℓ_0 .
- Parameters Defining Model are $\alpha \in [0, 1]$ and ℓ_0 .

Estimation may be conducted using MLE (later) or OLS. For OLS,

$$(\hat{\alpha}, \hat{\ell}_0) = \arg \min_{0 \leq \alpha \leq 1, \ell_0 \in \mathbf{R}} \sum_{i=2}^T [X_i - \ell_i(\alpha, \ell_0)]^2$$

$$\hat{X}_{T+1} = \hat{\alpha}X_T + (1 - \hat{\alpha})\ell_T(\hat{\alpha}, \hat{\ell}_0)$$

which can be calculated by iterating the level equation back to ℓ_0 .

Linear Trend Exponential Smoothing

- Prediction Equation: $\hat{X}_{T+h} = \ell_T + hb_T$ where ℓ_T is the **level** and b_T is the **slope**.
- Level Equation: $\ell_T = \alpha X_T + (1 - \alpha)(\ell_{T-1} + b_{T-1})$ which is the convex combination of last observation and last “level” or prediction.
- Trend/Slope Equation: $b_T = \beta(\ell_T - \ell_{T-1}) + (1 - \beta)b_{T-1}$ where $\ell_T - \ell_{T-1}$ is the last “observed” slope or change in level.
- Initial Conditions: ℓ_0 and b_0 .
- Parameters: $\alpha, \beta \in [0, 1]$, $\ell_0, \beta_0 \in \mathbf{R}$ which are estimated using MLE/OLS.

Trend + Seasonal Exponential Smoothing (Holt Winters ES, 1960s)

Suppose h is the forecast horizon of interest and time series has seasonal period p . Set $k = \lfloor (h-1)/p \rfloor$.

- Prediction Equation: $\hat{X}_{T+1} = \ell_T + hb_T + s_{T+h-p(k+1)}$ where ℓ_T is the level, b_T is the slope, and $s_{T+1-p(k+1)}$ is the seasonal effect.
- Level Equation: $\ell_T = \alpha(X_t - s_{T-p}) + (1 - \alpha)(\ell_{T-1} + b_{T-1})$.
- Slope Equation: $b_T = \beta(\ell_T - \ell_{T-1}) + (1 - \beta)b_{T-1}$.
- Seasonal Equation: $s_T = \gamma(X_T - \ell_{T-1} + b_{T-1}) + (1 - \gamma)s_{T-p}$.
- Initial Conditions: $\ell_0, \beta_0, s_0, \dots, s_{-p+1}$.
- Parameters: $\alpha, \beta, \gamma \in [0, 1]$, $\ell_0, \beta_0, s_0, \dots, s_{-p+1} \in \mathbf{R}$.

7.2 Exponential Smoothing as a State Space Model

Consider Simple Exponential Smoothing:

- Prediction Equation: $X_{t|t-1} = \ell_{t-1}$.
- Level Equation: $\ell_t = \alpha X_t + (1 - \alpha)\ell_{t-1}$.

Re-arranging the level equation gives

$$\ell_t = \ell_{t-1} + \alpha \underbrace{(X_t - \ell_{t-1})}_{\text{residual } \varepsilon_t} = \ell_{t-1} + \alpha \varepsilon_t$$

Also, $X_t = \ell_{t-1} + \varepsilon_t$. Therefore, these equations can be reformulated as:

- Prediction Equation: $X_t = \ell_{t-1} + \varepsilon_t$.
- Level Equation: $\ell_t = \ell_{t-1} + \alpha \varepsilon_t$.

Why is this useful? If we make a parametric assumption on ε_t (e.g., $\varepsilon_t \sim \mathcal{N}(0, \sigma_\varepsilon^2)$), then we can use Likelihood techniques (MLE, AIC, simulation based Prediction Intervals).

Such equations are examples of “State Space” Models:

DEFINITION 7.2.1: State space model

We say X_T follows a general **state space model** if:

- Observation Equation: $X_t = A_t Y_t + \varepsilon_t$ where A_t is the **measurement matrix**, Y_t is the **state vector** (unobserved), and ε_t is an **observation error**.
- State Equation: $Y_t = \phi Y_{t-1} + W_t$.



ε_t and W_t are white noise terms that may depend on each other.

EXAMPLE 7.2.2: State Space Models

- AR(1): $X_t = Y_t$ where $Y_t = \phi Y_{t-1} + W_t$ where $W_t \sim$ strong white noise.
- Simple Exponential Smoothing:

$$X_t = Y_{t-1} + \varepsilon_t$$

$$Y_t = Y_{t-1} + \alpha \varepsilon_t$$

where $\varepsilon_t \sim$ strong white noise.

All ARMA and Exponential Smoothing models can be written in state-space form.

Parameter Estimation and Model Selection using State-Space Formulation

- $X_t = \ell_{t-1} + \varepsilon_t$.
- $\ell_t = \ell_{t-1} + \alpha \varepsilon_t$.
- $\varepsilon_t \sim \mathcal{N}(0, \sigma_\varepsilon^2)$.
- Initial Condition: ℓ_0 .

$$\mathcal{L}(X_1, \dots, X_T; \alpha, \ell_0, \sigma_\varepsilon^2) = \prod_{i=1}^T \frac{\mathcal{L}(X_i | X_{i-1}, \dots, X_1; \alpha, \ell_0, \sigma_\varepsilon^2)}{\mathcal{N}(\ell_{i-1}(\alpha, \ell_0), \sigma_\varepsilon^2)}$$

Likelihood can be maximized numerically, and we use this to calculate AIC/BIC.

7.3 Multiplicative Exponential Smoothing Models

Standard Exponential Smoothing has “additive” errors, in the sense that

$$X_t = \ell_{t-1} + \varepsilon_t$$

$$\ell_t = \alpha X_t + (1 - \alpha) \ell_{t-1}$$

Therefore, $\varepsilon_t = X_t - \ell_{t-1}$.

We can also formulate exponential smoothing in terms of “multiplicative” errors, in the sense that

$$\varepsilon_t = \frac{X_t - \ell_{t-1}}{\ell_{t-1}}$$

where we note that the error is relative to the previous level. Therefore,

$$X_t = \ell_{t-1}(1 + \varepsilon_t)$$

$$\ell_t = \alpha X_t + (1 - \alpha)\ell_{t-1} = \alpha \varepsilon_t \ell_{t-1} + \alpha \ell_{t-1} + (1 - \alpha)\ell_{t-1} = \ell_{t-1}(1 + \alpha \varepsilon_t)$$

Why consider multiplicative errors? It is important to note that since the level follows the same exponential smoothing equation, the forecasts from multiplicative and additive error models will be the same. The difference arises from how uncertainty/error propagates in the model.

- Additive: $\hat{X}_{T+h} = \ell_T + \sum_{j=T+1}^{T+h} \varepsilon_j$ where we note that the MSE scales like h .
- Multiplicative: $\hat{X}_{T+h} = \ell_T \prod_{j=T+1}^{T+h} (1 + \varepsilon_j)$ where we note that the MSE (variance) is scaling like

$$\left(\mathbb{E}[(1 + \varepsilon_0)^2] \right)^h$$

which could grow very quickly as $h \rightarrow \infty$.

Multiplicative Linear + Trend and Holt Winters

Linear + Trend State Space Formulation:

$$\varepsilon_t = \frac{X_t - (\ell_{t-1} + b_{t-1})}{\ell_{t-1} + b_{t-1}}$$

$$X_t = (\ell_{t-1} + b_{t-1})(1 + \varepsilon_t)$$

$$\ell_t = (\ell_{t-1} + b_{t-1})(1 + \alpha \varepsilon_t)$$

$$b_t = b_{t-1} + \beta(\ell_{t-1} + b_{t-1})\varepsilon_t$$

where $\varepsilon_t \sim \mathcal{N}(0, \sigma_\varepsilon^2)$.

Multiplicative Seasonal Exponential Smoothing

Let p be the seasonal period.

$$X_t = (\ell_{t-1} + b_{t-1})s_{t-p}(1 + \varepsilon_t)$$

$$\ell_t = (\ell_{t-1} + b_{t-1})(1 + \alpha \varepsilon_t)$$

$$b_t = b_{t-1} + \beta(\ell_{t-1} + b_{t-1})\varepsilon_t$$

$$s_t = s_{t-p}(1 + \gamma \varepsilon_t)$$

When to use Additive versus Multiplicative

Seasonal Exponential Smoothing Models:

- (1) Multiplicative models imply that as the level increases (decreases) the seasonal fluctuations increase (decrease). Additive models suggest seasonal fluctuations remain constant as trend fluctuations.

Seasonal Fluctuations \uparrow as Level $\uparrow \implies$ Multiplicative.

- (2) Use AIC/BIC: The AIC can be evaluated for each state-space model and compared.

7.4 Exponential Smoothing Model Selection

Given the state-space formulation of exponential smoothing and the use of MLE to estimate the parameters, it is common to use AIC to choose among competing Exponential Smoothing (including additive versus multiplicative) models. Other options include:

- Cross-validation.
- Residual Analysis (white noise testing).

Prediction Intervals

Using the state-space formulation, valid prediction intervals may be computed using simulation.

EXAMPLE 7.4.1: Simple Exponential Smoothing

$$\hat{X}_{T+1|T} = \hat{\ell}_T$$

State-space formula:

$$\hat{X}_{T+1} \cong \hat{\ell}_T + \underbrace{\varepsilon_{T+1}}_{\mathcal{N}(0, \sigma_\varepsilon^2)}$$

(1) Estimate

$$\hat{\sigma}_\varepsilon^2 = \frac{1}{T-1} \sum_{j=2}^T (X_j - \hat{\ell}_{T-1})^2$$

(2) Simulate

$$\hat{X}_{T+1|T}^{(b)} = \hat{\ell}_T + \underbrace{\varepsilon_{T+1}^{(b)}}_{\mathcal{N}(0, \hat{\sigma}_\varepsilon^2)}$$

(3) Use 5% and 95% sample quantiles of $X_{T+1|T}^{(b)}$, $b = 1, \dots, B$ as prediction intervals.

REMARK 7.4.2

In many cases, the prediction MSE assuming $\varepsilon_t \sim \mathcal{N}(0, \sigma_\varepsilon^2)$ can be computed explicitly. See § 7.7 of HA.

An important consideration in applying this approach is that ε_t should behave like Gaussian white noise. We can check this using a residual analysis.

- White noise tests, ACF plots.
- Quantile-Quantile plot for Normality.

7.5 J and J Exponential Smoothing Forecast

[R Code] J and J Exponential Smoothing Forecast