

시계열자료분석팀

5팀

김민주

박준영

곽동길

강서진

황호성

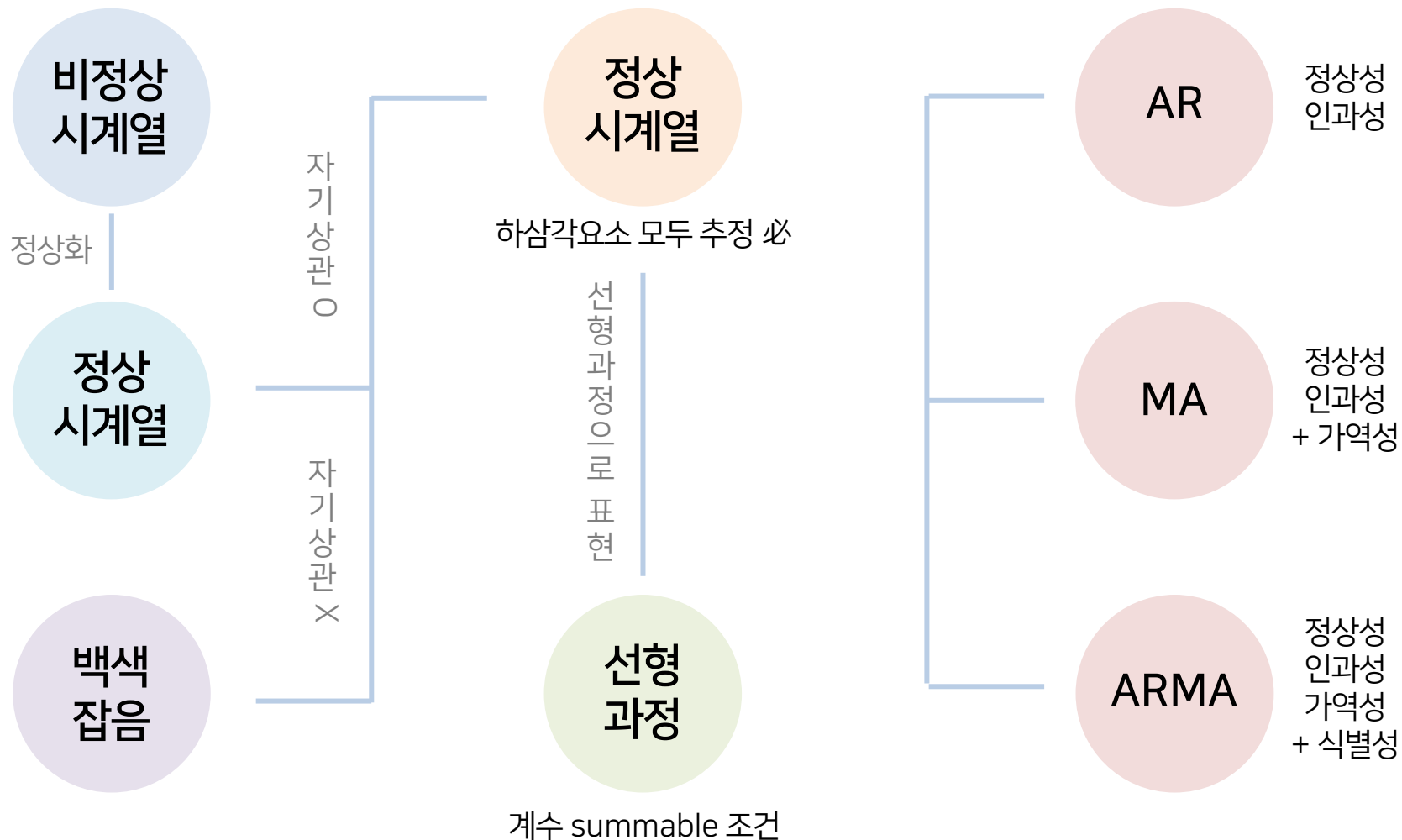
INDEX

- | | |
|-----------|--------------|
| 0. 2주차 복습 | 4. 이분산모형 |
| 1. ARIMA | 5. ARMAX |
| 2. SARIMA | 6. VAR |
| 3. ARFIMA | 7. 시계열과 머신러닝 |

0

2주차 복습

시계열 자료 분석 흐름 정리



정상 시계열 모형 적합 절차

Step 1) 모형 식별

AR, MA는 ACF 또는 PACF를 확인하여 차수 결정

ARMA는 IC(Information Criteria)를 최소로 하는 모형 선택

	AR(p)	MA(q)	ARMA(p,q)
ACF	지수적으로 감소	q 이후로 절단	지수적으로 감소
PACF	p 이후에 절단	0으로 수렴	지수적으로 감소

Step 2) 모수 추정

MLE, LSE, MME 등의 추정량으로 모수 추정

정상 시계열 모형 적합 절차

Step 3) 모형 진단

모형 진단은 모수에 대한 검정과 잔차에 대한 검정으로 나뉨

모수 검정

정상성, 가역성, 식별성 만족 여부 확인
모수 $\neq 0$ 인지 확인

잔차 검정

추세, 계절성, 이상치 확인
백색잡음 여부 확인
정규성 만족 여부 확인

Step 4) 예측

MSPE를 최소화하는 방향으로 예측 진행

1

ARIMA

비정상 시계열 모형



현실에서의 대부분의 시계열 데이터는 비정상 시계열 데이터

⋮

- ① 시계열이 약정상성을 만족하지 않음
- ② 평균 또는 분산이 일정하지 않아 추세 혹은 계절성 존재
- ③ 공분산이 시점에 의존

비정상 시계열 모형



현실에서의 대부분의 시계열 데이터는 비정상 시계열 데이터

정상화 과정을 사전에 매번 거쳐야 해서 번거로움



- ① 시계열이 약정상성을 만족하지 않음
- ② 평균 또는 분산이 일정하지 않아 추세 혹은 계절성 존재
- ③ 공분산이 시점에 의존

비정상 시계열 모형



현실에서의 대부분의 시계열 데이터는 비정상 시계열 데이터

정상화 과정을 사전에 매번 거쳐야 해서 번거로움

⋮

⋮

① 시계열이 약정상성을 만족하지 않음

사전에 정상화할 필요 없이

② 평균 또는 분산이 일정하지 않아 추세 혹은 계절성 존재
모델 자체에서 정상화를 이루는 **비정상 시계열 모형** 이용!

③ 공분산이 시점에 의존

ARIMA (자기회귀누적이동평균 모형)

ARIMA의 정의

d차 차분 후, 오차 y_t 가 ARMA(p,q)를 만족하는 시계열 데이터

→ ARIMA(p,d,q)를 따른다고 표현

⋮

- > ARMA 모형에 차분이 결합된 형태
- > 추세가 있어 정상성을 만족하지 않는 시계열 자료에 적용 가능

ARIMA (자기회귀누적이동평균 모형)

ARIMA 수식

$$\phi(B)(1 - B)^d X_t = \theta(B)Z_t$$

$$(1 - \phi_1 B - \dots - \phi_p B^p)(1 - B)^d X_t = (1 + \theta_1 B + \dots + \theta_q B^q)Z_t$$

$\phi(B)$ 는 AR의 특성방정식, $\theta(B)$ 는 MA의 특성방정식, $(1 - B)^d$ 는 d 차 차분

⋮

① $d = 0$: ARIMA(p,0,q) = ARMA(p,q)

② $d \geq 1$: ARIMA(p,d,q)

→ d 차 추세 존재, 오차가 ARMA(p,q)를 따를 때 적용 가능

ARIMA (자기회귀누적이동평균 모형)

ARIMA 수식

$$\phi(B)(1 - B)^d X_t = \theta(B)Z_t$$

$$(1 - \phi_1 B - \dots - \phi_p B^p)(1 - B)^d X_t = (1 + \theta_1 B + \dots + \theta_q B^q)Z_t$$

$\phi(B)$ 는 AR의 특성방정식, $\theta(B)$ 는 MA의 특성방정식, $(1 - B)^d$ 는 d 차 차분

⋮

① $d = 0$: $\text{ARIMA}(p, 0, q) = \text{ARMA}(p, q)$

② $d \geq 1$: $\text{ARIMA}(p, d, q)$

→ d 차 추세 존재, 오차가 $\text{ARMA}(p, q)$ 를 따를 때 적용 가능

ARIMA (자기회귀누적이동평균 모형)



ARIMA 수식

ARIMA 모형의 장점

$$\phi(B)(1-B)^d X_t = \theta(B)Z_t$$

$$(1 - \phi_1 B - \dots - \phi_p B^p)(1 - B)^d X_t = (1 + \theta_1 B + \dots + \theta_q B^q)Z_t$$

① 모형 자체가 차분을 포함하고 있어 간결하게 표현 가능

$\phi(B)$ 는 AR의 특성방정식, $\theta(B)$ 는 MA의 특성방정식, $(1-B)^d$ 는 d 차 차분

② 여전히 선형 과정을 따르기 때문에 예측이 용이

① $d = 0$: $\text{ARIMA}(p, 0, q) = \text{ARMA}(p, q)$

② $d \geq 1$: $\text{ARIMA}(p, d, q)$

→ d 차 추세 존재, 오차가 $\text{ARMA}(p, q)$ 를 따를 때 적용 가능



ARDMA가 아닌 ARIMA라고 하는 이유?



ARIMA에서 **I**는 **누적(integration)**을 의미

ARIMA 모형의 상식

⋮

$$\phi(B)(1 - B)X_t = \theta(B)Z_t$$

① 모형 자체가 차분을 포함하고 있어 간결하게 표현 가능

$$Y_t = (1 - B)X_t = X_t - X_{t-1}$$

② $X_t = X_{t-1} + Y_t = (X_{t-2} + Y_{t-1}) + Y_t = \dots = X_0 + \sum_{j=1}^t Y_j$

⋮

① $d = 0$: $ARIMA(p, 0, q) = ARMA(p, q)$

X_t 는 Y_t 의 누적합으로 볼 수 있으며 이는 random walk의 확장으로 해석 가능

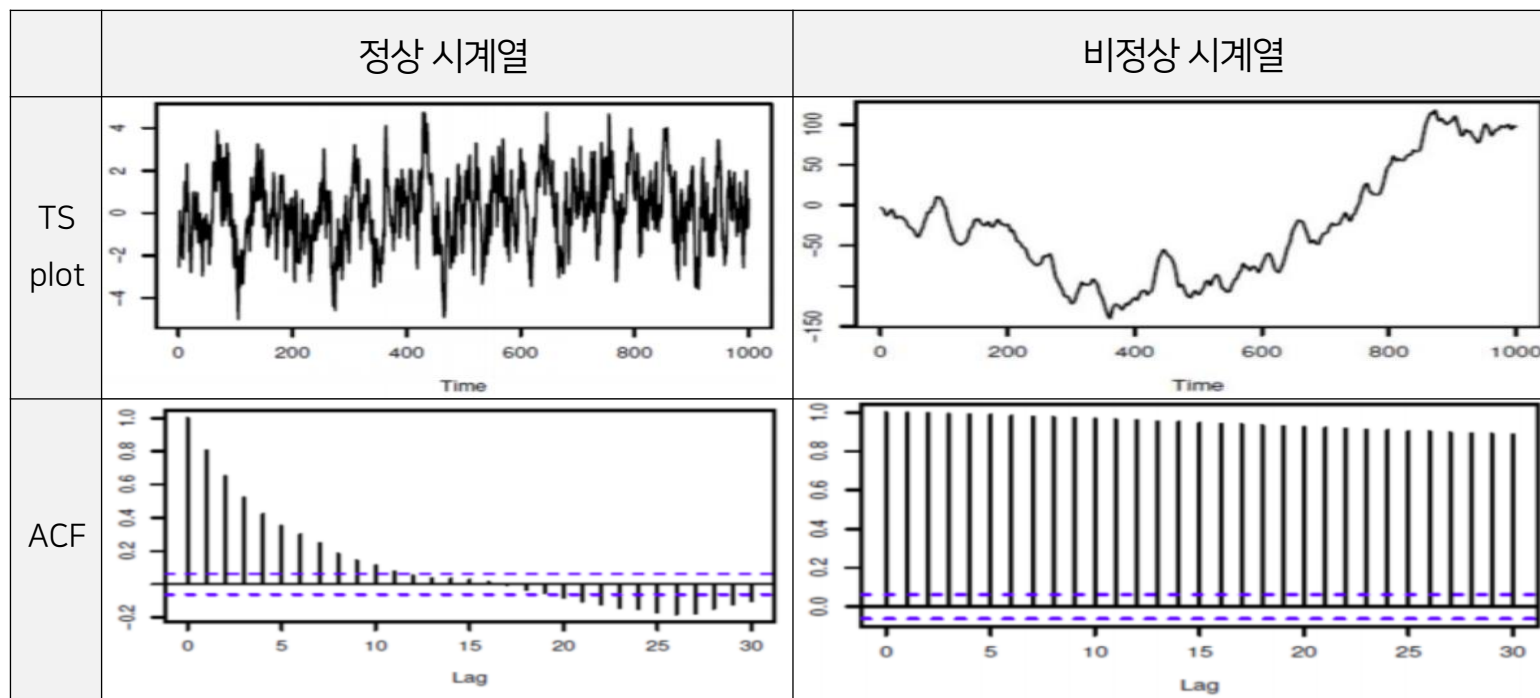
→ ARIMA 즉, 자기회귀**누적**이동평균모형이라고 불림

1

ARIMA

ARIMA | ARIMA 적합 절차

[1] TS plot과 ACF를 통해 정상/비정상 시계열 여부 판단

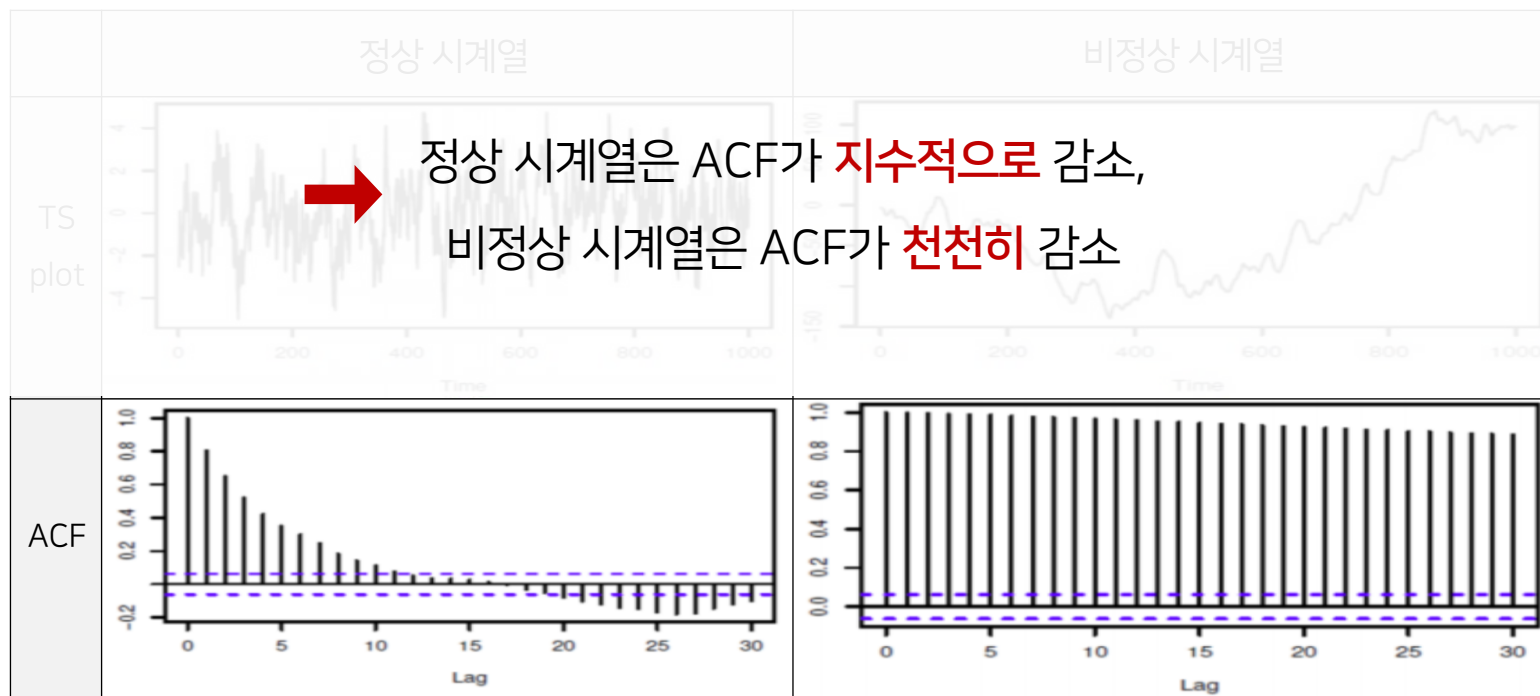


1

ARIMA

ARIMA | ARIMA 적합 절차

[1] TS plot과 ACF를 통해 정상/비정상 시계열 여부 판단



ARIMA | ARIMA 적합 절차

[2] 비정상 시계열에서 추세가 관측된다면, 차분을 통해 정상화 진행

⋮

비정상 시계열 확인 및 추세 관측



d차 차분 적용 (보통 1차 또는 2차 차분을 적용)



ACF와 PACF 그래프가 지수적으로 감소하면 적절한 정상화

ARIMA | ARIMA 적합 절차

[2] 비정상 시계열에서 추세가 관측된다면, 차분을 통해 정상화 진행

⋮

비정상 시계열 확인 및 추세 관측



d차 차분 적용 (보통 1차 또는 2차 차분을 적용)



ACF와 PACF 그래프가 지수적으로 감소하면 적절한 정상화

ARIMA | ARIMA 적합 절차

[2] 비정상 시계열에서 추세가 관측된다면, 차분을 통해 정상화 진행



d차 차분이 아닌 1, 2차 차분을 적용하는 이유?

d차 차분 시 필요 이상으로 차분 되는 **과대차분** 발생 가능성 존재

(ACF가 복잡해지거나 분산이 커지며, 불필요한 상관관계 생성)

d차 차분 적용 (보통 1차 또는 2차 차분을 적용)



과대차분을 방지하기 위해 1, 2차 차분 사용

ACF와 PACF 그래프가 지속적으로 감소하면 적절한 정상화

ARIMA | ARIMA 적합 절차

[3] p, q 차수 결정 및 모수 추정 후 진단을 마친 모형을 통해 예측 진행





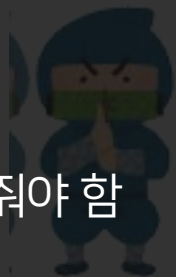
ARIMA | ARIMA 적합 **정상 vs 비정상 시계열 모형**

① 정상 시계열 적용

입력값(X)은 정상화를 거친 오차항
결과값(Y)은 오차항에 대한 예측값



최종 예측을 위해서는
추정한 추세와 계절성을 더해줘야 함

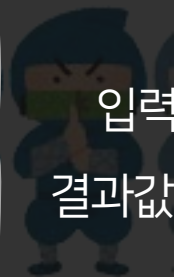


② 비정상 시계열 적용

연산 과정에 차분을 통한 정상화 포함!



입력값(X)은 원본 시계열 데이터
결과값(Y)은 원본 데이터 최종 예측값



2

SARIMA

SARIMA (Seasonal ARIMA)

정상화에서의 계절성은 **결정적**, 즉 모든 주기에서 계절성이 **동일함**을 가정

⋮

하지만 현실에서는 모든 계절성분이 **비결정적**, 즉 다른 성분들과의 **상관관계** 有

SARIMA (Seasonal ARIMA)

정상화에서의 계절성은 **결정적**, 즉 모든 주기에서 계절성이 **동일함**을 가정

⋮

하지만 현실에서는 모든 계절성분이 **비결정적**, 즉 다른 성분들과의 **상관관계** 有

➡ **주기가 변함에 따라 계절성도 변화**할 수 있음을 반영하는 SARIMA 모형

SARIMA (Seasonal ARIMA)

SARIMA의 정의

Seasonal ARIMA의 줄임말

추세와 계절성이 모두 존재하는 비정상 데이터에 적용할 수 있는 모형

⋮

> 가장 확장된 형태의 ARMA 모형

> 계절성 사이의 상관관계에 대해 모델링하는 확률적 접근법

SARIMA (Seasonal ARIMA)

SARIMA의 정의

Seasonal ARIMA의 줄임말



추세와 계절성이 모두 존재하는 비정상 데이터에 적용할 수 있는 모형
주기가 12인($s = 12$) 예시 데이터를 가지고

? SARIMA 모형이 계절성을 어떻게 다루는지 알아보자!



> 가장 확장된 형태의 ARMA 모형

> 계절성 사이의 상관관계에 대해 모델링하는 확률적 접근법

SARIMA (Seasonal ARIMA)

[Step 1] Between year model

⋮

	Jan	Feb	...	Dec
Year 1	Y_1	Y_2	...	Y_{12}
Year 2	Y_{13}	Y_{14}	...	Y_{24}
⋮	⋮	⋮	⋮	⋮
Year r	$Y_{1+12(r-1)}$	$Y_{2+12(r-1)}$...	$Y_{12+12(r-1)}$

SARIMA (Seasonal ARIMA)

	Jan	Feb	...	Dec
Year 1	Y_1	Y_2	...	Y_{12}
Year 2	Y_{13}	Y_{14}	...	Y_{24}
\vdots	\vdots	\vdots	\vdots	\vdots
Year r	$Y_{1+12(r-1)}$	$Y_{2+12(r-1)}$...	$Y_{12+12(r-1)}$
	\downarrow	\downarrow	\downarrow	\downarrow
	$ARMA(P, Q)$	$ARMA(P, Q)$	$ARMA(P, Q)$	$ARMA(P, Q)$



$Y_{j+12t} + \phi_1 Y_{j+12(t-1)} + \dots + \phi_p Y_{j+12(t-p)} = \theta_1 U_{j+12(t-1)} + \dots + \theta_p U_{j+12(t-p)} + U_{j+12t}$

> 각 열(month)을 계절성분으로 정의

(where $t = 0, 1, \dots, r$, $U_t \sim WN(0, \sigma_u^2)$)

> 각 열을 고정한 후 모든 열은 동일한 ARMA 모델을 따른다고 가정

> Y_j 끼리 상관관계 존재

$$\Phi(B^{12})Y_t = \Theta(B^{12})U_t$$

SARIMA (Seasonal ARIMA)

	Jan	Feb	...	Dec
Year 1	Y_1	Y_2	...	Y_{12}
Year 2	Y_{13}	Y_{14}	...	Y_{24}
\vdots	\vdots	\vdots	\vdots	\vdots
Year r	$Y_{1+12(r-1)}$	$Y_{2+12(r-1)}$...	$Y_{12+12(r-1)}$
	\downarrow	\downarrow	\downarrow	\downarrow
	$ARMA(P, Q)$	$ARMA(P, Q)$	$ARMA(P, Q)$	$ARMA(P, Q)$

ARMA(P,Q)

$$Y_{j+12t} - \Phi_1 Y_{j+12(t-1)} - \dots - \Phi_P Y_{j+12(t-P)} = U_{j+12t} + \Theta_1 U_{j+12(t-1)} + \dots + \Theta_Q U_{j+12(t-Q)}$$

$$(where\ t = 0, 1, \dots, r, \quad U_t \sim WN(0, \sigma_U^2))$$

$$\Updownarrow$$

$$\Phi(B^{12})Y_t = \Theta(B^{12})U_t$$

$j \in [1, 12]$ 는 month에 해당하며, j 가 달라져도 Φ 와 Θ 동일

SARIMA (Seasonal ARIMA)

	Jan	Feb	...	Dec
Year 1	Y_1	Y_2	...	Y_{12}
Year 2	Y_{13}		...	Y_{24}
⋮				⋮
Year r	$Y_{1+12(r-1)}$	$Y_{2+12(r-1)}$...	$Y_{12+12(r-1)}$
	ARMA(P,Q)	ARMA(P,Q)	ARMA(P,Q)	ARMA(P,Q)

SARIMA에서
같은 월에 해당하는 시계열끼리 상관관계를 가지고
ARMA(P,Q)를 따르는 계절성을 갖는다고 가정

ARMA(P,Q)

$$Y_{j+12t} - \Phi_1 Y_{j+12(t-1)} - \cdots - \Phi_P Y_{j+12(t-P)} = U_{j+12t} + \Theta_1 U_{j+12(t-1)} + \cdots + \Theta_Q U_{j+12(t-Q)}$$

$$(where\ t = 0, 1, \dots, r, \quad U_t \sim WN(0, \sigma_U^2))$$

$$\Updownarrow$$

$$\Phi(B^{12})Y_t = \Theta(B^{12})U_t$$

$j \in [1, 12]$ 는 month에 해당하며, j 가 달라져도 Φ 와 Θ 동일

SARIMA (Seasonal ARIMA)

SARIMA에서

같은 월에 해당하는 시계열끼리 상관관계를 가지고
ARMA(P,Q)를 따르는 계절성을 갖는다고 가정

ARMA(P,Q)

$Y_{j+12t} - \phi_1 Y_{j+12(t-1)} - \dots - \phi_p Y_{j+12(t-p)} = \theta_1 U_{j+12(t-1)} + \dots + \theta_q U_{j+12(t-q)}$ 즉, 주기마다 계절성이 동일하지 않지만



Year의 계절성분 간 상관관계를 모델에 의해 설명 가능



$$\Phi(B^{12})Y_t = \Theta(B^{12})U_t$$

$j \in [1,12]$ 는 month에 해당하며, j 가 달라져도 Φ 와 Θ 동일

SARIMA (Seasonal ARIMA)

[Step 2] Between month model

■
■
■

	Jan	Feb	...	Dec		
Year 1	Y_1	Y_2	...	Y_{12}	→	$ARMA(p, q)$
Year 2	Y_{13}	Y_{14}	...	Y_{24}	→	$ARMA(p, q)$
⋮	⋮	⋮	⋮	⋮	→	$ARMA(p, q)$
Year r	$Y_{1+12(r-1)}$	$Y_{2+12(r-1)}$...	$Y_{12+12(r-1)}$	→	$ARMA(p, q)$
	↓	↓	↓	↓		
	$ARMA(P, Q)$	$ARMA(P, Q)$	$ARMA(P, Q)$	$ARMA(P, Q)$		

SARIMA (Seasonal ARIMA)

	Jan	Feb	...	Dec		
Year 1	Y_1	Y_2	...	Y_{12}	→	$ARMA(p, q)$
Year 2	Y_{13}	Y_{14}	...	Y_{24}	→	$ARMA(p, q)$
⋮	⋮	⋮	⋮	⋮	→	$ARMA(p, q)$
Year r	$Y_{1+12(r-1)}$	$Y_{2+12(r-1)}$...	$Y_{12+12(r-1)}$	→	$ARMA(p, q)$
	↓	↓	↓	↓		
	$ARMA(P, Q)$	$ARMA(P, Q)$	$ARMA(P, Q)$	$ARMA(P, Q)$		

⋮

행을 고정하여 동일 연도의 연속된 시계열 데이터끼리의 상관관계를 파악

SARIMA (Seasonal ARIMA)

	Jan	Feb	...	Dec		
Year 1	Y_1	Y_2	...	Y_{12}	→	$ARMA(p, q)$
Year 2	Y_{13}	Y_{14}	...	Y_{24}	→	$ARMA(p, q)$
⋮	⋮	⋮	⋮	⋮	→	$ARMA(p, q)$
Year r	$Y_{1+12(r-1)}$	$Y_{2+12(r-1)}$...	$Y_{12+12(r-1)}$	→	$ARMA(p, q)$
	↓	↓	↓	↓		
	$ARMA(P, Q)$	$ARMA(P, Q)$	$ARMA(P, Q)$	$ARMA(P, Q)$		

동일한 월에 속하는 시계열끼리의 상관관계 : $ARMA(P, Q)$

한 주기 내에서 연속된 시계열끼리의 상관관계 : $ARMA(p, q)$

SARIMA (Seasonal ARIMA)



Between month model 과정

ARMA(P,Q) 와 ARMA(p,q) 구분					
Year	Jan	Feb	...	Dec	ARMA(p,q)
Year 1	Y_1	Y_2	\dots	Y_{12}	\rightarrow ARMA(p,q)
Year 2	Y_{13}	Y_{14}	\dots	Y_{24}	\rightarrow ARMA(p,q)
\vdots	\vdots	\vdots	\vdots	\vdots	\rightarrow ARMA(p,q)
Year r	$(\{Y_{13}, \dots, Y_{12}\})$	와 같이 행의 상관관계 존재 가능성 반영을 위해			
	\downarrow	\downarrow	\downarrow	\downarrow	
2) 위의 과정을 식으로 나타내면, $\phi(B)U_t = \theta(B)Z_t, \quad Z_t \sim WN(0, \sigma^2)$					

행을 고정하면 동일한 월에 속하는 시계열끼리의 상관관계 : ARMA(P,Q)

한 주기 내에서 연속된 시계열끼리의 상관관계 : ARMA(p,q)

SARIMA (Seasonal ARIMA)

[Step 3] Seasonal ARMA

⋮

Seasonal ARMA

$$\Phi(B^{12})Y_t = \Theta(B^{12}) \phi^{-1}(B)\theta(B)Z_t$$

$$\phi(B)\Phi(B^{12})Y_t = \theta(B)\Theta(B^{12})Z_t, \quad Z_t \sim WN(0, \sigma^2)$$

파란색 부분은 12시점 전과의 상관관계(주기마다의 계절성분)
 빨간색 부분은 바로 전 시점과의 상관관계(연속된 시계열 상관관계)

SARIMA (Seasonal ARIMA)

Step1과 Step2 과정을 합쳐서 하나의 식으로 나타냄
 Seasonal ARMA = SARMA(p,q)x(P,Q)

Seasonal ARMA

$$\Phi(B^{12})Y_t = \Theta(B^{12}) \phi^{-1}(B)\theta(B)Z_t$$

$$\phi(B)\Phi(B^{12})Y_t = \theta(B)\Theta(B^{12})Z_t, \quad Z_t \sim WN(0, \sigma^2)$$

파란색 부분은 12시점 전과의 상관관계(주기마다의 계절성분)
 빨간색 부분은 바로 전 시점과의 상관관계(연속된 시계열 상관관계)

SARIMA (Seasonal ARIMA)

[Step 4] SARIMA(p,d,q)×(P,D,Q)

SARIMA에 차분을 더해서 **SARIMA** 모형 완성

⋮

SARIMA(p,d,q)×(P,D,Q)

$$\phi(B)\Phi(B^{12})(1-B)^d(1-B^{12})^D X_t = \theta(B)\Theta(B^{12}) Z_t$$

$$(Z_t \sim WN(0, \sigma^2))$$



SARIMA(p,d,q)×(P,D,Q) 식의 의미

SARIMA (Seasonal ARIMA)

© SARIMA(p,d,q)×(P,D,Q)

$$\phi(B)\Phi(B^{12})(1-B)^d(1-B^{12})^D X_t = \theta(B)\Theta(B^{12}) Z_t$$

$$(Z_t \sim WN(0, \sigma^2))$$

SARIMA에 차분을 더해서 SARIMA 모형 완성

소문자 (p,d,q)는 전체 시계열에 대한 차수 의미

대문자 (P,D,Q)는 주기 패턴에 대한 차수 의미

$$\phi(B)\Phi(B^{12})(1-B)^d(1-B^{12})^D X_t = \theta(B)\Theta(B^{12}) Z_t$$

> $(1-B)^d$ 는 전체 시계열의 추세에 대한 d 차 차분

> $(1-B^{12})^D$ 는 계절성분이 갖는 D 차 추세에 대해 lag 12 차분 D 번 적용

SARIMA | 적합절차

[1] 분산 안정화

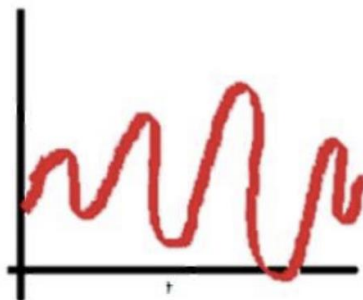
TS plot과 잔차를 확인하여 이분산성이 나타나는지 파악

⋮

이분산성 발견 시 분산 안정화 변환

→ Log 변환, 제곱근 변환, Box-Cox 변환 등

자세한 내용은 1주차 클린업 참고

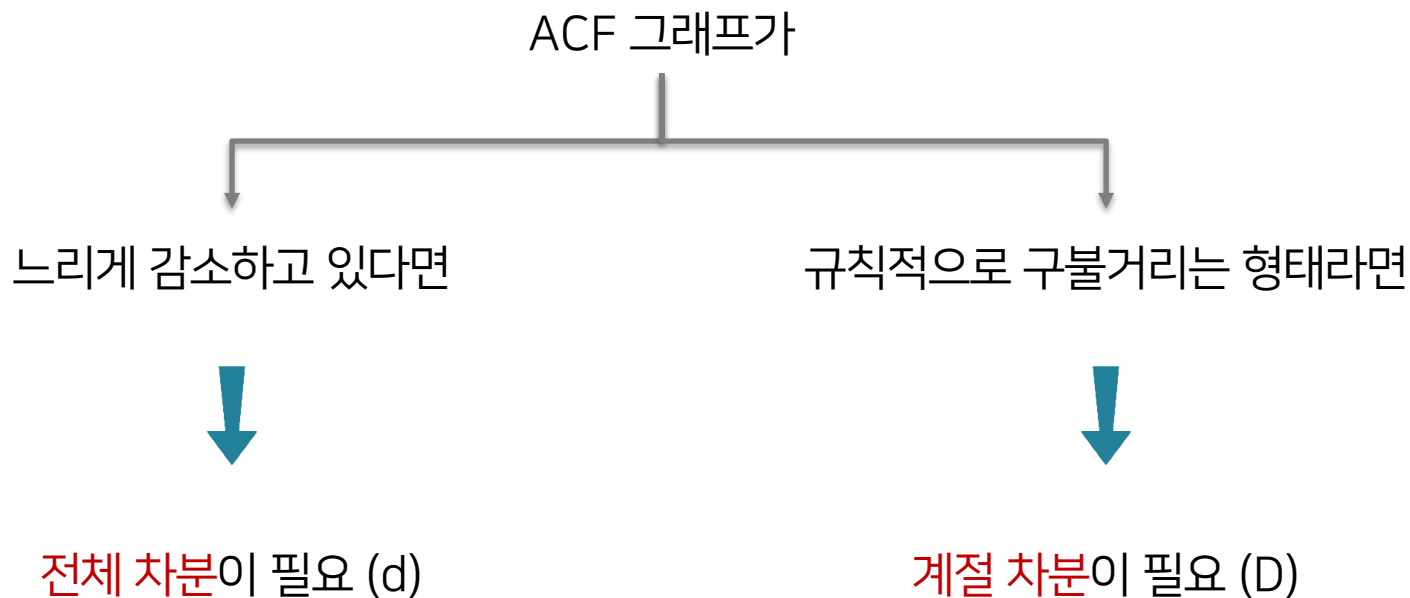


SARIMA | 적합절차

[2] d차 또는 lag-d 차분 진행 여부 결정

과대차분 되지 않도록 주의

ACF 그래프의 개형에 따라 d와 D의 차수를 결정



SARIMA | 적합절차

[2] d차 또는 lag-d 차분 진행 여부 결정

과대차분 되지 않도록 주의

ACF 그래프의 개형에 따라 d와 D의 차수를 결정

ACF 그래프가

느리게 감소하면서 동시에 규칙적으로 구불거리는 형태인 경우

느리게 감소하고 있다면



규칙적으로 구불거리는 형태라면

계절 차분 후 전체 차분 진행

전체 차분 시 (d-1)차 차분을 진행해야 함을 유념

1주차 클린업 참고

SARIMA | 적합절차

[3] p , q 와 P , Q 의 차수 결정

p 와 q 는 ACF, PACF 그래프 또는 모형 식별 방법으로 구할 수 있음

P 와 Q 도 마찬가지로, 계절성분의 모수이므로 그래프 해석 방식에 차이가 있음

→ **주기**마다의 그래프를 확인할 것!



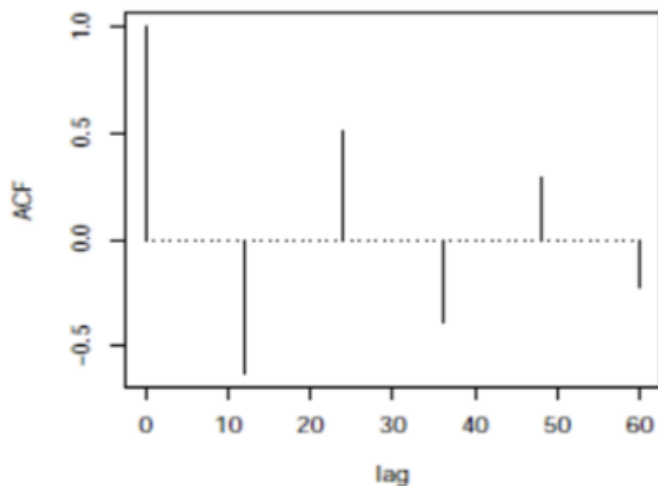
예시를 통해 확인해보자



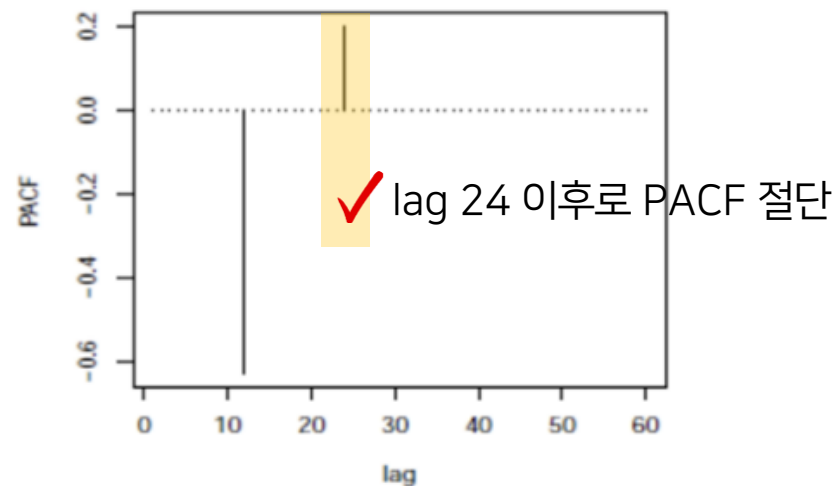
SARIMA | 적합절차

예시 1) 주기가 12인 시계열 데이터 A

SARIMA(0,0,0)(2,0,0): ACF



SARIMA(0,0,0)(2,0,0): PACF



SARIMA(0,0,0)(2,0,0)의 의미 : 계절성분이 ARMA(2,0), 즉 AR(2)를 따름

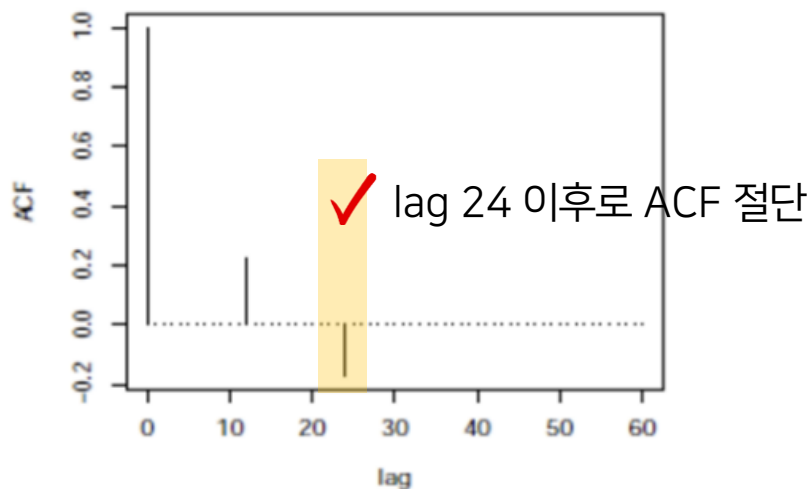
∴ ACF는 지수적으로 감소, PACF는 시차 이후 절단되는 양상

→ 여기서 시차는 주기의 배수 단위로 해석 !

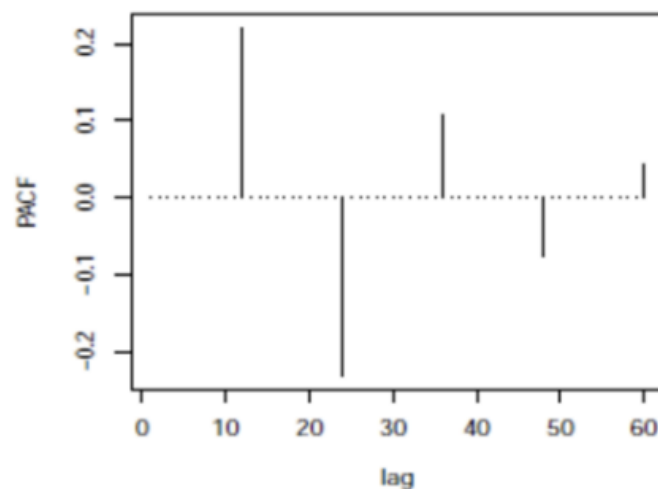
SARIMA | 적합절차

예시 2) 주기가 12인 시계열 데이터 B

SARIMA(0,0,0)(0,0,2): ACF



SARIMA(0,0,0)(0,0,2): PACF



SARIMA(0,0,0)(0,0,2)의 의미 : 계절성분이 ARMA(0,2), 즉 MA(2)를 따름

∴ ACF는 시차 이후로 절단, PACF는 지수적으로 감소하는 양상

→ 여기서 시차는 주기의 배수 단위로 해석 !

SARIMA | 적합절차

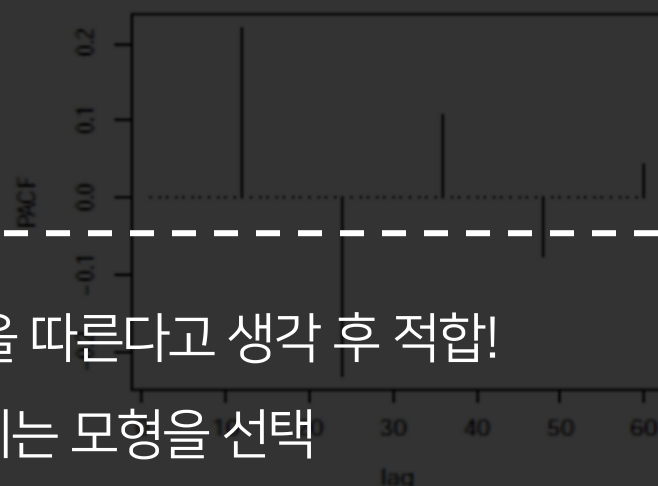
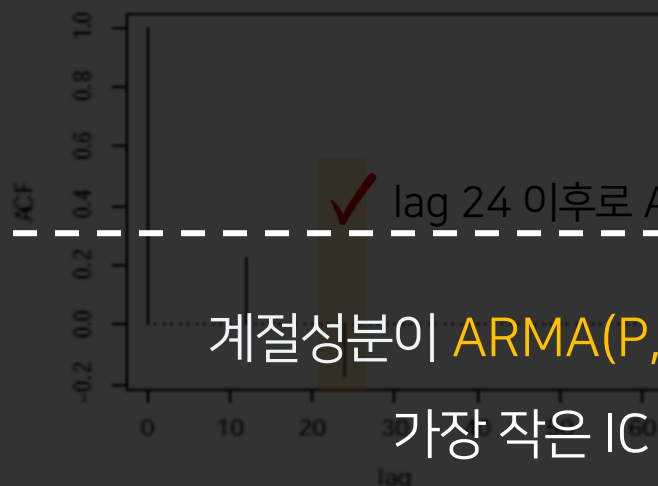


예시 2) 주기가 12인 시계열 데이터 B

ACF와 PACF 모두 **지수적으로 감소**하는 양상이라면?

SARIMA(0,0,0)(0,0,2): ACF

SARIMA(0,0,0)(0,0,2): PACF



계절성분이 **ARMA(P, Q)모형**을 따른다고 생각 후 적합!

가장 작은 IC 값을 가지는 모형을 선택

cf) 현실적으로 모든 모수의 IC 비교 불가 → 범위를 정해놓고 해당 범위 내에서 찾을

SARIMA(0,0,0)(0,0,2)의 의미 : 계절성분이 ARMA(0,2), 즉 MA(2)를 따름

∴ ACF는 시차 이후로 절단, PACF는 지수적으로 감소하는 양상

→ 여기서 시차는 주기의 배수 단위로 해석!

SARIMA | 적합절차

[4] 모수 추정 후 예측 진행

SARIMA 모형은 정상화를 따로 거치지 않는 모형
→ 적합한 모형의 결과값이 원본데이터에 대한 최종 예측값



별도의 후처리 과정이 필요 없음 !



3

ARFIMA

ARFIMA 모형

ARFIMA (AR Fractionally IMA)

ARIMA 모델에서 차분의 차수를 양의 정수가 아닌 실수까지 허용하여,
시계열이 가진 **장기 기억**을 보존하는 모형



정수 차원의 차분은 과거 관측치를 빼주는 의미이므로
장기 기억을 분석하는 것이 불가능하여 ARFIMA가 등장!

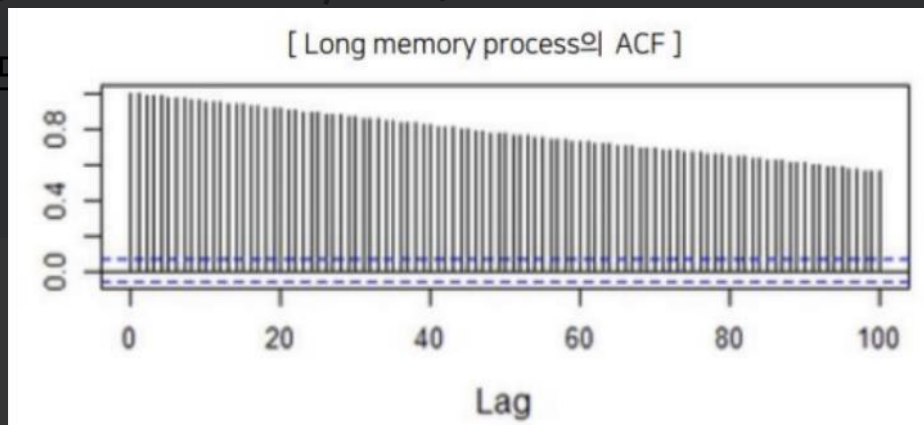


ARFIMA 모형

장기기억 확률과정이란?

ARFIMA (AR Fractionally IMA)

ARIMA 모형



허용하여,

 $\rho(k)$ 가 $0 < d < 0.5$ 인 어떤 실수 d 에 대해

 $\rho(k) \sim Ck^{2d-1}, k \rightarrow \infty (C > 0)$ 를 만족하는 확률과정 $\{Z_t\}$

정수 차원의 차분은 과거 관측치를 빼주는 의미이므로

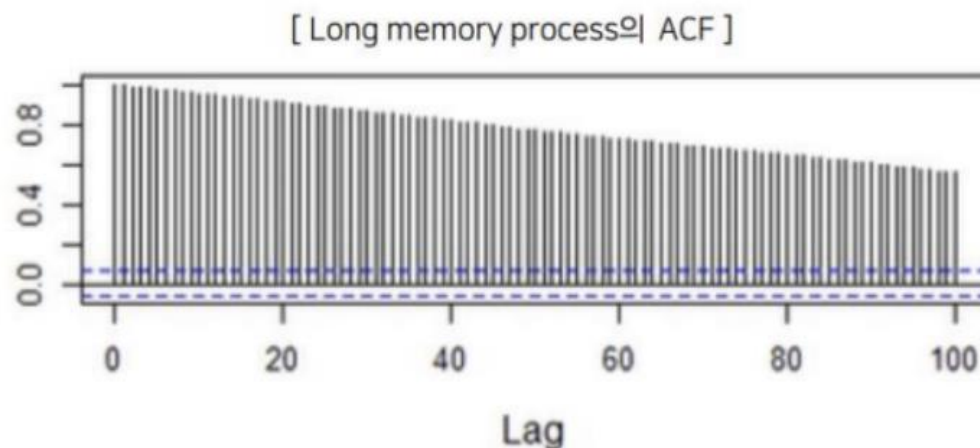
장기 기억을 분석하는 것이 불가능하여 ARFIMA가 등장

 \Rightarrow ACF의 합이 무한대로 발산 & ACF가 빠르게 0으로 수렴하지 X

ARFIMA 모형

ARFIMA

ARIMA



공하여,

[ACF가 천천히 감소하는 경우]

- ① 차분을 여러 번 진행해도 추세의 온전한 제거 불가
- ② 차분 횟수가 늘어나 추정해야 하는 모수 증가

→ ARFIMA 모델 사용이 적절 !

ARFIMA 모형

ARFIMA 수식

ARIMA와 동일

$$\phi(B)(1 - B)^d X_t = \theta(B)Z_t$$

d 의 값과 반영되는 과거 데이터의 양은 서로 반비례

⋮



정상성 조건을 만족하기 위해 $0 < d < 0.5$ 로 d 의 범위 제한
(0.5보다 커지면 분산이 발산 !)



이때 차분의 차수 d 는 LSE나 MLE로 추정

4

이분산모형

이분산모형

지금까지의 시계열 모형들

→ 분산에 변화가 없음을 가정 후, 평균의 움직임에만 관심을 가짐

⋮



수익률, 주가, 환율 등 금융 관련 시계열 자료는
분산이 과거 자료에 의존하는 특성 가정

이분산모형

지금까지의 시계열 모형들

→ 분산에 변화가 없음을 가정 후, 평균의 움직임에만 관심을 가짐

⋮



수익률, 주가, 환율 등 금융 관련 시계열 자료는
분산이 과거 자료에 의존하는 특성 가정



이분산 시계열 모형

관측치들 간 시간에 따라 변하는 분산에 관심을 갖는 모형

이분산모형

지금까지의 시계열 모형들

→ 분산에 변화가 없음을 가정, 평균의 움직임에만 관심을 가짐



다시 말해, 경제학 분야의 “변동성”을
통계적 용어인 “조건부 분산”을 통해 시간의 함수로 표현하는 시계열 모형

많은 금융 시계열 자료는 조건부 이분산성을 가짐

분산이 과거 자료에 의존하는 특성 가정

이분산 시계열 모형

관측치들 간 시간에 따라 변하는 분산에 관심을 갖는 모형

이분산모형 | 수익률

수익률 (Return)

t 시점의 가격 = P_t

- Simple return : $R_t = \frac{P_t - P_{t-1}}{P_{t-1}}$
- Log return : $r_t = \log P_t - \log P_{t-1} = \log(1 + R_t) \approx R_t$

⋮

Log return의 장점

1. 덧셈으로 표현 가능 2. log로 분산 안정화 효과

앞으로 다룰 수익률은 모두 **Log return !**

이분산모형 | 조건부 이분산성

조건부 등분산성

구조적인 변동성이 이전 기간 변동성의 영향을 받지 않는 것

$$Var(r_t|F_{t-1}) = constant$$

VS



조건부 이분산성

변동성이 시점에 의존하는 것
미래의 변동이 현재까지의 상황에 의존

$$Var(r_t|F_{t-1}) \neq constant$$

F_t = 현재까지의 모든 정보 집합 σ -field

ARCH 모형

ARCH: 자기회귀이분산모형

t 시점의 오차 변동성인 σ_t^2 가 AR(m) 모형을 따른다고 가정하여
과거 시점의 오차항으로 설명하는 모형

σ_t^2 은 조건부 분산

⋮

ARCH(1)

$$r_t = \sigma_t \varepsilon_t, \sigma_t \sim N(0, \sigma^2), \varepsilon_t \sim N(0, 1)$$

$$\sigma_t^2 = \text{Var}(r_t | F_{t-1}) = \alpha_0 + \alpha_1 r_{t-1}^2$$

$$r_t^2 = (\alpha_0 + \alpha_1 r_{t-1}^2) \varepsilon_t^2 = \sigma_t^2 \varepsilon_t^2$$

σ_t^2 이 r_{t-1}^2 에 의존함을 볼 수 있음 \Rightarrow 조건부 이분산

$r_t | F_{t-1}$: t-1까지의 모든 정보를 다 알고 있을 경우의 t 시점 수익률

ARCH 모형

ARCH: 자기회귀이분산모형

t 시점의 오차 변동성인 σ_t^2 가 AR(m) 모형을 따른다고 가정하여
과거 시점의 오차항으로 설명하는 모형

 σ_t^2 은 조건부 분산

⋮

ARCH(m)

$$r_t = \sigma_t \varepsilon_t$$

$$\sigma_t^2 = \alpha_0 + \alpha_1 r_{t-1}^2 + \cdots + \alpha_m r_{t-m}^2 \approx r_t^2$$

→ σ_t^2 에 AR(m) 가정한 것

ARCH 모형 | 비선형성

ARCH(1)

$$\begin{aligned}
 r_t^2 &= \sigma_t^2 \varepsilon_t^2 = (\alpha_0 + \alpha_1 r_{t-1}^2) \varepsilon_t^2 \\
 &= \alpha_0 \varepsilon_t^2 + \alpha_1 \varepsilon_t^2 \{(\alpha_0 + \alpha_1 r_{t-2}^2) \varepsilon_{t-1}^2\} \\
 &= \alpha_0 \varepsilon_t^2 + \alpha_0 \alpha_1 \varepsilon_t^2 \varepsilon_{t-1}^2 + \alpha_1^2 r_{t-2}^2 \varepsilon_t^2 \varepsilon_{t-1}^2 \\
 &\quad \vdots \\
 &= \alpha_0 \sum_{j=0}^n (\alpha_1^j \varepsilon_t^2 \varepsilon_{t-1}^2 \cdots \varepsilon_{t-j}^2) + \alpha_1^{n+1} r_{t-n-1}^2 \varepsilon_t^2 \varepsilon_{t-1}^2 \cdots \varepsilon_{t-j}^2
 \end{aligned}$$



곱셈식으로 표현되는 **비선형적 모델**

$0 < \alpha_1 < 1$ 일 때, 정상성 만족

ARCH 모형 | 비선형성



ARCH(1)

ARCH(m) 모형의 단점

$$r_t^2 = \sigma_t^2 \varepsilon_t^2 = (\alpha_0 + \alpha_1 r_{t-1}^2) \varepsilon_t^2$$

$$= \alpha_0 \varepsilon_t^2 + \alpha_1 \varepsilon_t^2 \{(\alpha_0 + \alpha_1 r_{t-2}^2) \varepsilon_{t-1}^2\}$$

$$= \alpha_0 \varepsilon_t^2 + \alpha_0 \alpha_1 \varepsilon_t^2 \varepsilon_{t-1}^2 + \alpha_1^2 r_{t-2}^2 \varepsilon_t^2 \varepsilon_{t-1}^2$$

ARCH(m)에서 m이 커지면

⋮

추정해야 할 **모수가 많아지고**, 추정량의 **정확도가 떨어짐**

$$= \alpha_0 \sum_{j=0}^n (\alpha_1^j \varepsilon_t^2 \varepsilon_{t-1}^2 \cdots \varepsilon_{t-j}^2) + \alpha_1^{n+1} r_{t-n-1}^2 \varepsilon_t^2 \varepsilon_{t-1}^2 \cdots \varepsilon_{t-j}^2$$

이를 해결하기 위해 **일반화된 모델인 GARCH 모형** 사용 $0 < \alpha_1 < 1$ 일 때, 정상성 만족

GARCH 모형

GARCH: 일반화자기회귀이분산모형

t 시점의 오차 변동성인 σ_t^2 가 ARMA(m, n) 모형을 따른다고 가정하여
과거 시점의 오차항으로 설명하는 모형

⋮

ARCH 모형에서 σ_t^2 에 AR 모형을 가정한 것처럼
GARCH 모형에서는 σ_t^2 에 ARMA 모형을 가정한 **확장 모델**

↓

두 모형 모두 **t 시점 수익률의 변동성**을 표현하는 것은 동일

GARCH 모형

r_t^2 는 σ_t^2 에 근사하며,
 둘 사이의 오차를 $\eta_t = r_t^2 - \sigma_t^2$ 라고 할 때,
 GARCH(1,1)을 식으로 나타내 보자

⋮

GARCH(1,1)

$$r_t = \sigma_t \varepsilon_t$$

$$\sigma_t^2 = \alpha_0 + \alpha_1 r_{t-1}^2 + \beta_1 \sigma_{t-1}^2$$

$$r_t^2 = \alpha_0 + (\alpha_1 + \beta_1) r_{t-1}^2 + \eta_t - \beta_1 \eta_{t-1}$$

AR(1)

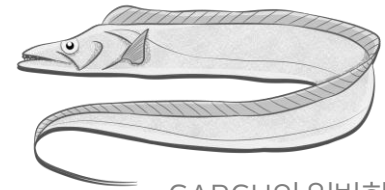
MA(1)

$$r_t^2 = \alpha_0 + \alpha_1 r_{t-1}^2 + \beta_1 r_{t-1}^2 + r_t^2 - \sigma_t^2 - \beta_1 (r_{t-1}^2 - \sigma_{t-1}^2)$$

GARCH 모형

이를 일반화하면 다음과 같이 나타낼 수 있음

⋮



GARCH의 일반화

GARCH(m,n)

$$r_t = \sigma_t \varepsilon_t$$

$$\sigma_t^2 = \alpha_0 + \sum_{i=1}^m \alpha_i r_{t-i}^2 + \sum_{j=1}^n \beta_j \sigma_{t-j}^2$$

5

ARMAX

정의

ARMAX 모형

ARMA 모형에 독립변수로 **외부요인**(eXogenous)을 추가한 모형

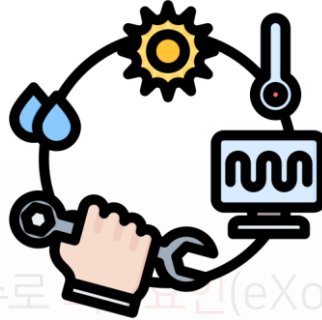


Y_t 와 X_t 의 관측값 수는 동일해야 함

$$\phi(B)Y_t = \theta(B)Z_t + \beta_0 + \beta_1 X_t$$

이 때, 추가된 독립변수는 연속형일 수도, 범주형일 수도 있음

정의



ARMAX 모형

ARMA 모형에 독립변수로 외생변수(exogenous)를 추가한 모형

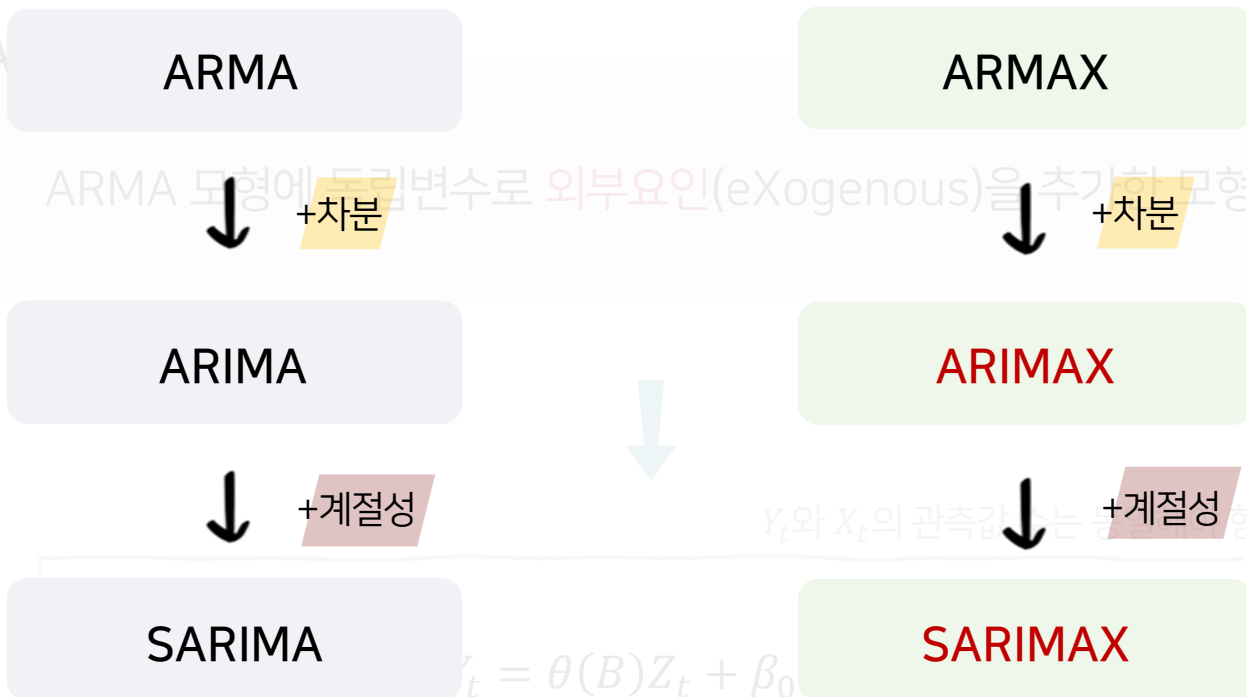
예를 들어 '날씨를 반영한 주가 예측' 모형식을 만들고자 한다면,
 위 식에서 Y_t 를 예측하고자 하는 주가, X_t 를 날씨 데이터라고 할 수 있음!

 Y_t 와 X_t 의 관측값 수는 동일해야 함

$$\phi(B)Y_t = \theta(B)Z_t + \beta_0 + \beta_1 X_t$$

이 때, 추가된 독립변수는 연속형일 수도, 범주형일 수도 있음

정의



$$\text{ARIMAX} : \phi(B)(1-B)^d Y_t = \theta(B)Z_t + \beta^T \underline{X}$$

$$\text{SARIMAX} : \phi(B)\Phi(B^s)(1-B)^d(1-B^s)^D Y_t = \theta(B)\Theta(B^s)Z_t + \beta^T \underline{X}$$

6

VAR

정의

VAR(Vector Auto Regressive)

AR 모형에 **Vector 구조를 결합**한 모형

현재 관측값을 자신과 다른 변수 각각의 과거 관측값으로 설명하는 모형

AR은 일변량 자기회귀모형, VAR은 이를 확장한 다변량 자기회귀모형



VAR(1) 모형

$$\begin{pmatrix} X_t \\ Y_t \end{pmatrix} = \begin{pmatrix} c_1 \\ c_2 \end{pmatrix} + \begin{pmatrix} \phi_{11} & \phi_{12} \\ \phi_{21} & \phi_{22} \end{pmatrix} \begin{pmatrix} X_{t-1} \\ Y_{t-1} \end{pmatrix} + \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \end{pmatrix}$$

정의

VAR(Vector Auto Regressive)

AR 모형에 **Vector 구조를 결합**한 모형

현재 관측값을 자신과 다른 변수 각각의 과거 관측값으로 설명하는 모형

AR은 일변량 자기회귀모형, VAR은 이를 확장한 다변량 자기회귀모형



VAR(1) 모형

아래와 같이 표현도 가능!

$$X_t = c_1 + \phi_{11}X_{t-1} + \phi_{12}Y_{t-1} + \varepsilon_1$$

$$Y_t = c_2 + \phi_{21}X_{t-1} + \phi_{22}Y_{t-1} + \varepsilon_2$$

정의

VAR(Vector Auto Regressive)



VAR모형은 단순 과거 데이터만을 고려하는 것이 아니라
여러 변수들의 의존성과 상호작용을 고려하는 모형!

회귀모형



VAR(1) 모형

아래와 같이 표현도 가능!

VAR은 상호연계성이 높은 경제구조를 분석할 때 용이하게 사용

$$X_t = c_1 + \phi_{11}X_{t-1} + \phi_{12}Y_{t-1} + \varepsilon_1$$
$$Y_t = c_2 + \phi_{21}X_{t-1} + \phi_{22}Y_{t-1} + \varepsilon_2$$

7

시계열과 머신러닝

시계열 자료에 사용되는 머신러닝



지금까지 알아본 전통적인 시계열 모형들이 아닌
머신러닝을 통해 시계열 자료를 다루는 방법을 알아보자!

(1) 결측치 보간

데이터 분석 과정에서 결측치 발견
→ 안정적인 분석을 위해 결측치 보간 必

평균, 최빈값 등의 일반적인 대표값으로 보간 시
시계열 데이터의 특성을 반영하지 못해 평균, 분산에 왜곡이 생길 수 있음

(1) 결측치 보간

데이터 분석 과정에서 결측치 발견
→ 안정적인 분석을 위해 결측치 보간 必



평균, 최빈값 등의 일반적인 대표값으로 보간 시,
시계열 데이터의 특성을 반영하지 못해 평균, 분산에 왜곡이 생길 수 있음

(1) 결측치 보간



시계열 데이터는 다른 방법을 통해 결측치 보간!

→ 안정적인 분석을 위해 결측치 보간 必



평균, 최빈값 등의 일반적인 대표값으로 보간 시,
시계열 데이터의 특성을 반영하지 못해 평균, 분산에 왜곡이 생길 수 있음

(1) 결측치 보간 | LOCF, NOCB, MA/MM

LOCF(Last Observation Carried Forward)

직전 관측치 값으로 결측치 대체

NOCB(Next Observation Carried Backward)

직후 관측치 값으로 결측치 대체

Moving Average / Moving Median

직전 time window 내 n개의 평균치 or 중앙값으로 대체

time window는 결측치 보간에 고려되는 데이터 포인트 범위

(1) 결측치 보간 | LOCF, NOCB, MA/MM

LOCF (Last Observation Carried Forward)

직전 관측치 값으로 결측치 대체

이 세 가지 방법을 일반적으로 사용하지만,

결측치를 기준으로 패턴이 급격하게 변화하는 경우에는

NOCB (Next Observation Carried Backward)

조금 더 복잡한 방법을 사용해야 함!

직후 관측치 값으로 결측치 대체

Moving Average / Moving Median

직전 N의 time window의 평균치 / 중앙값

Time window는 결측치 보간



(1) 결측치 보간 | 선형, 비선형, 스플라인 보간법

선형 보간법

근사 함수가 **선형(linear)** 함수임을 가정하여 보간

비선형 보간법

근사 함수가 **비선형(non-linear)** 함수임을 가정하여 보간

스플라인(Spline) 보간법

전체 구간을 **소구간으로 분할한 후**

저차수의 다항식으로 **매끄러운 함수를** 구하는 방법

(1) 결측치 보간 | 선형, 비선형, 스플라인 보간법

선형 보간법



근사 함수가 선형(linear) 함수임을 가정하여 보간

이 방법들로 보간하기 어려울 정도로 결측치가 많은 경우에는
결측치 보간 모델링을 통해 결측치를 예측하여 대체하는 방법도 고려 가능

근사 함수가 비선형(non-linear) 함수임을 가정하여 보간
But 모델링을 진행하기 위해서는 충분한 데이터가 있어야 함

스플라인(Spline) 보간법

전체 구간을 소구간으로 분할한 후

저차수의 다항식으로 매끄러운 함수를 구하는 방법

자세한 내용은 데마팀 2주차 클린업 참고

(2) 노이즈 처리 (Denoising)

노이즈(Noise)

다른 외부 요인의 간섭과 같이, **의도하지 않은 데이터의 왜곡**을 불러오는 모든 것



평균 대중교통 이용량



유명 가수 콘서트



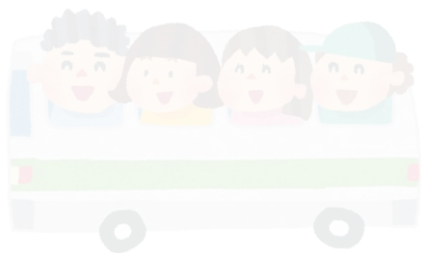
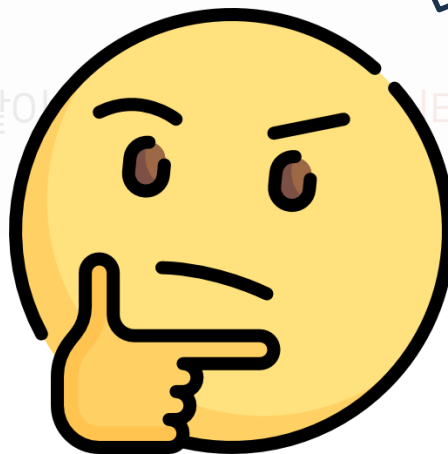
이용량 급격하게 증가

(2) 노이즈 처리 (Denoising)

노이즈(Noise)

다른 외부 요인의 간섭과 같이 데이터의 왜곡을 불러오는 모든 것

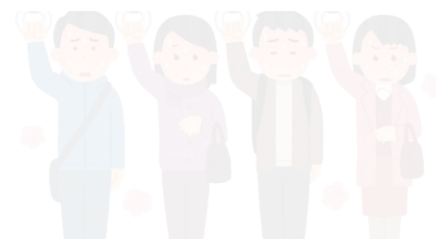
Noise!



평균 대중교통 이용량



유명 가수 콘서트



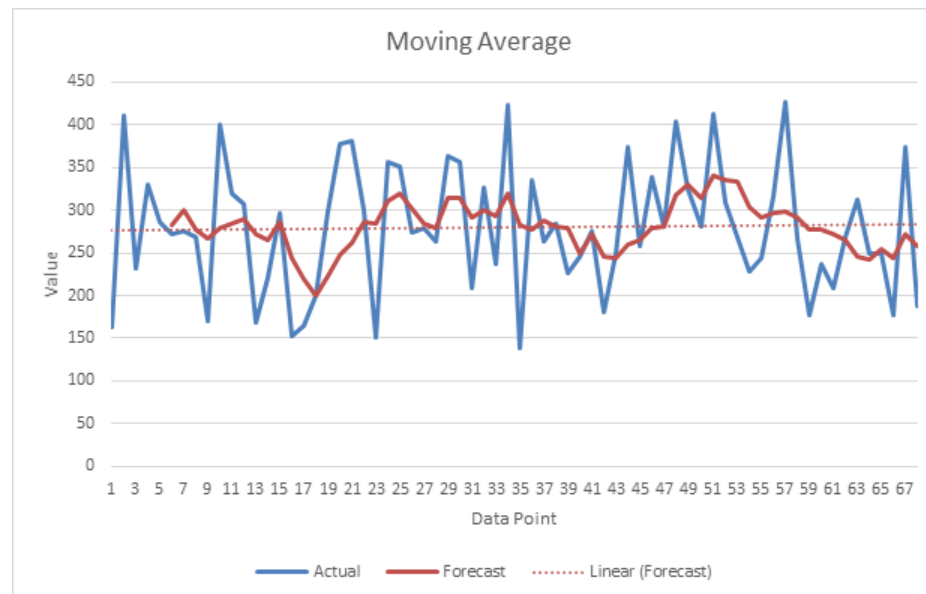
이용량 급격하게 증가

(2) 노이즈 처리 (Denoising)

Moving Average(MA)

평균값으로 관측치를 대체하여 평활하는 방법
주로 노이즈가 적은 데이터에 적용!

노이즈가 많은 경우 평균도 노이즈의 성격을 띠 수 있기 때문에



(2) 노이즈 처리 (Denoising)

Filtering

노이즈가 특정 분포를 따른다고 가정하고
해당 분포 값을 이용해 노이즈를 제거하는 과정

가우시안 필터링

(Gaussian Filtering)

노이즈가 정규분포를 따른다고 가정,
중심에 가까울 수록 큰 가중치를 부여

칼만 필터링

(Kalman Filtering)

잡음이 포함된 과거 측정값에서
현재 상태의 결합분포를 추정

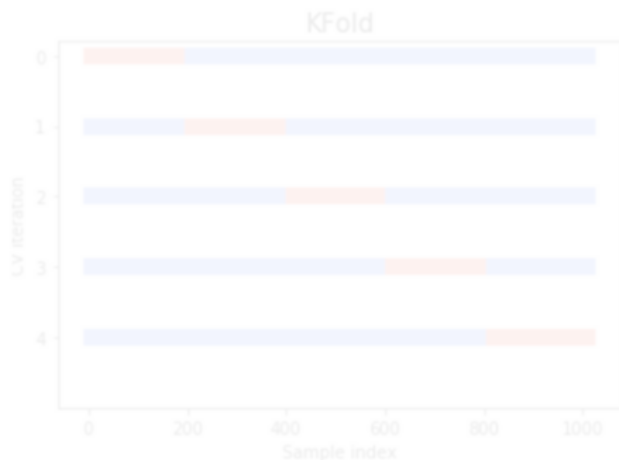


데이터가 다른 분포들의 결합이라 가정하고
데이터의 특성에 맞는 분포 모델링

(3) 시계열 데이터의 CV

교차검증(Cross Validation)

학습 데이터와 검증 데이터를 통해 모델의 성능을 측정하는 방법
모델의 과적합을 방지하기 위해 반드시 수행해야 함



일반적으로 사용하는 K-fold CV는 시간 순서를 고려 X
→ 시계열 자료의 특성을 반영하기 힘들

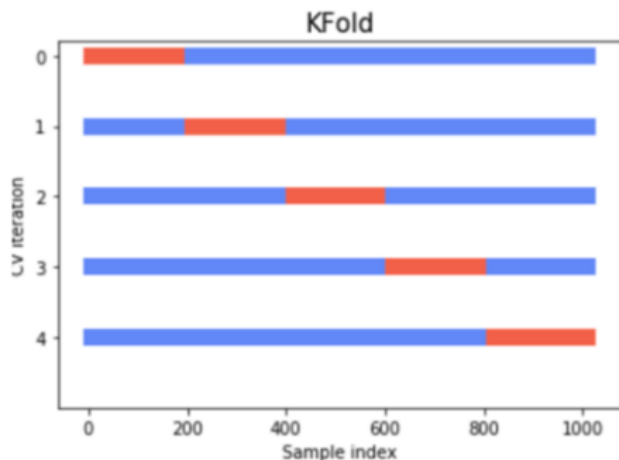
Time Series CV

Blocked Time Series CV

(3) 시계열 데이터의 CV

교차검증(Cross Validation)

학습 데이터와 검증 데이터를 통해 모델의 성능을 측정하는 방법
모델의 과적합을 방지하기 위해 반드시 수행해야 함



일반적으로 사용하는 K-fold CV는 시간 순서를 고려 X
→ 시계열 자료의 특성을 반영하기 힘들

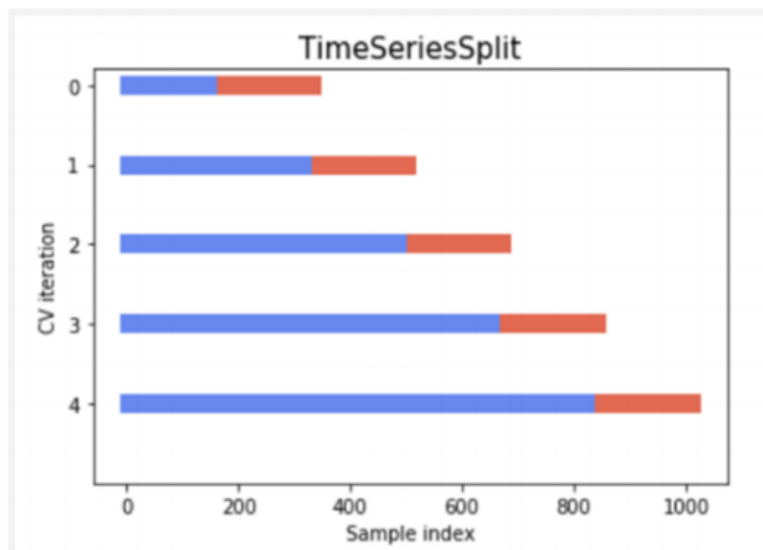
Time Series CV

Blocked Time Series CV

(3) 시계열 데이터의 CV

Time Series CV

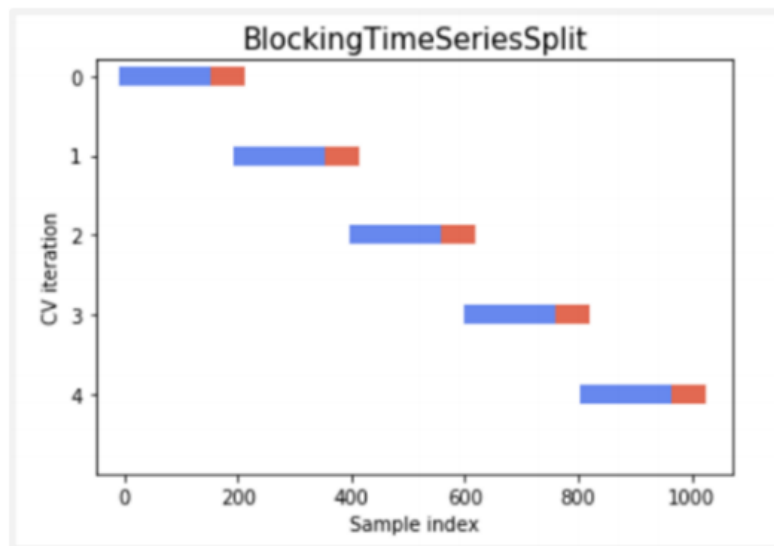
window를 누적하여 이전단계의 train set과 validation set을
다음 단계의 train set으로 사용해 교차검증 진행



(3) 시계열 데이터의 CV

Blocked Time Series CV

동일한 사이즈의 window를 옆으로 이동시키며
일정 비율로 train과 validation set으로 분할해 교차검증 진행



(4) 클래스 불균형

일반적인 데이터에서는 샘플링을 통해 클래스 불균형을 해결하지만,
시계열 데이터에서는 시간의 흐름을 반영해야 하므로 샘플링 사용이 제한됨



(4) 클래스 불균형

일반적인 데이터에서는 샘플링을 통해 클래스 불균형을 해결하지만,
시계열 데이터에서는 시간의 흐름을 반영해야 하므로 샘플링 사용이 제한됨



더 적은 수를 가지는 클래스에 **가중치를 부여**하는 방법
즉, 비용민감학습을 통해 클래스 불균형을 해결!

(4) 클래스 불균형

Scale_pos_weight

이진분류 문제에서 사용하는 파라미터로, 기본 설정은 0과 1로 라벨링 되었음을 가정
수가 더 적은 쪽을 1로 설정하며, 해당 클래스에 가중치를 부여하는 방식

Ex) 0, 1 예측 문제에서 1이 전체의 6% - `scale_pos_weight = 16`

Class_weight

샘플 수가 더 적은 쪽에 가중치를 부여하는 방법으로 클래스 별 가중치 제시

Ex) `model.fit(class_weight = {0:1, 1:2})`

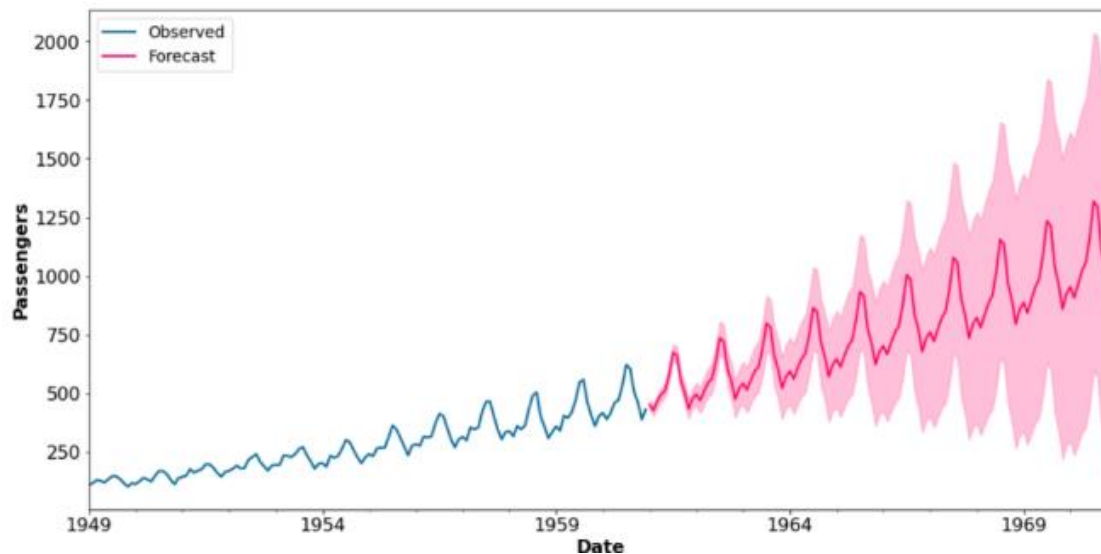
Sample weight

다중 분류에서 사용되며, 각 클래스 비율의 역수를 가중치로 계산하는 함수

Ex) `class_weight.compute_sample_weight(class_weight = "balanced")`

(5) 예측 및 평가지표

학습된 모델을 사용하여 미래 값을 예측할 때,
예측 구간과 신뢰구간을 고려하여 결과를 해석하는 것이 중요!



(5) 예측 및 평가지표 | MAPE

MAPE(Mean Absolute Percentage Error)

$$MAPE = \frac{100}{n} \sum_{t=1}^n \left| \frac{y_t - \hat{y}_t}{y_t} \right|$$

⋮

- MAPE는 $MAE(\frac{1}{n} \sum_{t=1}^n |y_t - \hat{y}_t|)$ 를 비율로 나타낸 것으로 실제값에 상대적으로 고려됨.
- Scale에 의존하지 않아 모델 간 성능을 비교할 때 용이하며, 퍼센트 값으로 나타낼 수 있기 때문에 직관적으로 이해할 수 있다는 장점

(5) 예측 및 평가지표 | MAPE

MAPE(Mean Absolute Percentage Error)

$$MAPE = \frac{100}{n} \sum_{t=1}^n \left| \frac{y_t - \hat{y}_t}{y_t} \right|$$

⋮

이를 수식으로 나타내면,

- MAPE는 $MAE(\frac{1}{n} \sum_{t=1}^n |y_t - \hat{y}_t|)$ 를 비율로 나타낸 것으로 실제값에 상대적으로 고려됨.

- Scale에 의존하지 않아 모델 간 성능 비교를 대응이하며, 퍼센트 값으로 나타낼 수 있기 때문에 직관적으로 이해할 수 있다는 장점

$$\frac{3 - 4}{3} \neq \frac{1012 - 1013}{1012}$$



(5) 예측 및 평가지표 | MAPE **MAPE의 한계**

MAPE(Mean Absolute Percentage Error)

1) 절대 오차를 실제 데이터로 나누기 때문에 0에 가까운 값이 있으면
MAPE가 상당히 크게 나타날 수 있으며, 0인 경우에는 계산이 불가능

2) 절대적인 값의 오차가 같더라도 실제값과 예측값과의 대소 관계에 따라
 과대 추정하는 **예측값에 페널티를 더 부여**하는 문제가 있음
 이를 수식으로 나타내면,

평균 오차가 20일 때 실제값과 예측값의 대소 관계에 따른 MAE, MAPE 비교

실제값이 예측값 보다 작을 때 (예측값이 과대추정)

MAE: 20.0

MAPE: 0.9133333333333333

실제값이 예측값 보다 클 때 (예측값이 과소추정)

MAE: 20.0

MAPE: 0.4371428571428571

있다는 장점
 (실제 코드 예시)

(5) 예측 및 평가지표 | SMAPE

SMAPE(Systematic Mean Absolute Percentage Error)

$$SMAPE = \frac{100}{n} \sum_{t=1}^n \frac{|y_t - \hat{y}_t|}{|y_t| + |\hat{y}_t|}$$

⋮

SMAPE는 **MAPE의 한계점을 보완**할 수 있음

평균 오차가 20일 때 실제값과 예측값의 대소 관계에 따른 MAE, SMAPE 비교

실제값이 예측값 보다 작을 때 (예측값이 과대추정)

MAE: 20.0

SMAPE: 0.29

실제값이 예측값 보다 클 때 (예측값이 과소추정)

MAE: 20.0

SMAPE: 0.29

(5) 예측 및 평가지표 | SMAPE

SMAPE(Systematic Mean Absolute Percentage Error)

$$SMAPE = \frac{100}{n} \sum_{t=1}^n \frac{|y_t - \hat{y}_t|}{|y_t| + |\hat{y}_t|}$$

그러나 SMAPE는 \hat{y}_t 에 의존하기 때문에,
예측값이 과소추정될 때 분모가 더 작아져 오차가 커짐

평균 오차가 20일 때 과소추정, 과대추정에 따른 MAE, SMAPE 비교

과대추정 시

MAE: 20.0

SMAPE: 0.14912698412698414

과소추정 시

MAE: 20.0

SMAPE: 0.21857142857142856

(5) 예측 및 평가지표 | SMAPE

SMAPE(Systematic Mean Absolute Percentage Error)



각각의 평가지표마다 분명한 특성이 존재하므로,
여러 평가지표들을 비교해보고 **분석과 예측의 목적에 따라**
적절한 평가지표를 선택하는 것이 중요!

예측값이 과소추정될 때 분모가 더 작아져 오차가 커짐

평균 오차가 20일 때 과소추정, 과대추정에 따른 MAE, SMAPE 비교

과대추정 시

MAE: 20.0

SMAPE: 0.14912698412698414

과소추정 시

MAE: 20.0

SMAPE: 0.21857142857142856



8

3주차 정리

ARIMA, SARIMA, ARFIMA

ARIMA(ARMA+차분)

$$\phi(B)(1 - B)^d X_t = \theta(B)Z_t$$

SARIMA(ARIMA+계절성)

$$\phi(B)\Phi(B^{12})(1 - B)^d(1 - B^{12})^D X_t = \theta(B)\Theta(B^{12})Z_t$$

ARFIMA(ARIMA+장기기억보존)

$$\phi(B)(1 - B)^d X_t = \theta(B)Z_t, \quad 0 < d < 0.5$$

조건부 이분산성, ARCH, GARCH

조건부 이분산성

미래의 변동이 현재에 의존하는 것

ARCH(자기회귀이분산모형)

조건부 분산 σ_t^2 를 과거시점의
오차항으로 설명

 σ_t^2 에 AR(M) 구조를 적용한 것

$$\sigma_t^2 = \alpha_0 + \alpha_1 r_{t-1}^2 + \cdots + \alpha_m r_{t-m}^2 \approx r_t^2$$

GARCH(일반화자기이분산모형)

 σ_t^2 에 ARMA모형을 가정하여 확장!

$$r_t^2 = \sigma_t^2 \varepsilon_t$$

$$\sigma_t^2 = \alpha_0 + \sum_{i=1}^p \alpha_i r_{t-i}^2 + \sum_{j=1}^q \beta_j \sigma_{t-j}^2$$

ARMAX

ARMAX(ARMA+외부요인)

$$\phi(B)Y_t = \theta(B)Z_t + \beta_0 + \beta_1 X_t$$

+ 차분

ARIMAX(ARMAX+차분)

$$\phi(B)(1 - B)^d Y_t = \theta(B)Z_t + B^t \underline{X}$$

+ 계절성

SARIMAX(ARMA+차분+계절성)

$$\phi(B)\Phi(B^s)(1 - B)^d Y_t = \theta(B)\Theta(B^s)Z_t + B^t \underline{X}$$

VAR

VAR(Vector Auto Regressive)

AR 모형에 **Vector** 구조를 결합한 모형

VAR(1) 모형

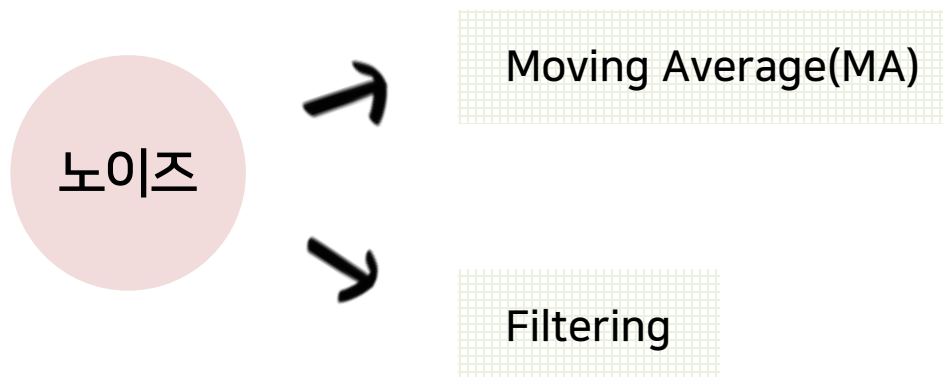
$$\begin{pmatrix} X_t \\ Y_t \end{pmatrix} = \begin{pmatrix} c_1 \\ c_2 \end{pmatrix} + \begin{pmatrix} \phi_{11} & \phi_{12} \\ \phi_{21} & \phi_{22} \end{pmatrix} \begin{pmatrix} X_{t-1} \\ Y_{t-1} \end{pmatrix} + \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \end{pmatrix}$$

시계열과 머신러닝 | 결측치 보간, 노이즈 처리

◎ 결측치 보간

시계열 데이터의 결측치 보간법	LOCF / NOCB / MA / MM
패턴이 급격하게 변하는 경우	선형, 비선형 보간 / Spline 보간

◎ 노이즈 처리



시계열과 머신러닝 | CV, 클래스 불균형

◎ 교차 검증

Time Series CV	이전단계의 train set과 validation set을 다음 단계의 train set으로 사용
Blocked Time Series CV	같은 사이즈의 윈도우 내에서 train, validation set으로 분할

◎ 클래스 불균형

Scale_force_weight	수가 더 적은 쪽을 1로 설정하고 가중치 부여
Class_weight	수가 더 적은 쪽에 가중치 부여
Sample_weight	각 클래스 비율의 역수를 가중치로 계산

시계열과 머신러닝 | 평가지표

MAPE(Mean Absolute Percentage Error)

$$MAPE = \frac{100}{n} \sum_{t=1}^n \left| \frac{y_t - \hat{y}_t}{y_t} \right|$$

SMAPE(Systematic Mean Absolute Percentage Error)

$$SMAPE = \frac{100}{n} \sum_{t=1}^n \frac{|y_t - \hat{y}_t|}{|y_t| + |\hat{y}_t|}$$

감사합니다

