

시계열 분석

(출생아 수 예측 모형)



응용통계학과
201552001
유승우

목 차

I. 분석 개요	1
II. 시계열 분석	2
1. 데이터 설명	2
2. ETS 모델 분석	5
3. ARIMA 모델 분석	7
III. 결론	14
[부록] R 코드	16

I 분석 개요

한국의 출생율이 바닥을 치고 있다. 정부는 지난 19년간 점진적으로 예산을 늘려가며 약 230조를 투입하고, 정책의 방향성까지 바꿔가며 저출생 문제를 극복하고자 했지만, 여전히 2020년 출생율이 OECD 회원국 중 한국이 꼴찌이며, 평균 0.84명으로 유일하게 0명대를 기록하고 있다.

우리나라의 출생율 하락세는 계속해서 진행 중이므로, 이를 극복하기 위해 출생율 예측으로 앞으로의 정책의 방향성과 예산집행의 기초자료로 활용하기 위해 분석을 실시한다.



[그림 1.1] 최근 19년 출생율 변화¹⁾



[그림 1.2] 최근 저출생·고령화 예산 규모²⁾

1) <http://news.tf.co.kr/read/ptoday/1849187.htm>
2) <https://m.news.zum.com/articles/54972035>

II. 시계열 분석

1. 데이터 설명

○ 데이터 설명

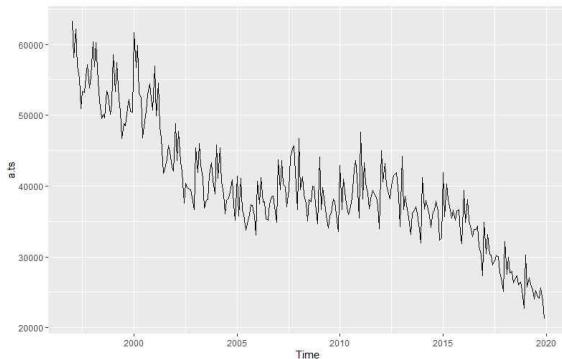
⑩ 시군구/성/월별 출생	
▣ 자료갱신일: 2020-08-26 / 수록기간: 월, 년 1997.01 ~ 2019.12 / 자료문의처: 02-2012-9114, 042-481~	
임금설정 *	
시점	전국 계(명)
2019. 12	21,228
2019. 11	23,727
2019. 10	25,613
2019. 09	24,090
2019. 08	24,371
2019. 07	25,222
2019. 06	23,992
2019. 05	25,299
2019. 04	26,104
2019. 03	27,049
2019. 02	25,710
2019. 01	30,271
2018. 12	22,767
2018. 11	25,301
2018. 10	26,474
2018. 09	26,066
2018. 08	27,381
2018. 07	27,033
2018. 06	26,357
2018. 05	27,949
2018. 04	27,734
2018. 03	29,987
2018. 02	27,576
2018. 01	27,576
2017. 12	27,576
2017. 11	27,576
2017. 10	27,576
2017. 09	27,576
2017. 08	27,576
2017. 07	27,576
2017. 06	27,576
2017. 05	27,576
2017. 04	27,576
2017. 03	27,576
2017. 02	27,576
2017. 01	27,576
2016. 12	27,576
2016. 11	27,576
2016. 10	27,576
2016. 09	27,576
2016. 08	27,576
2016. 07	27,576
2016. 06	27,576
2016. 05	27,576
2016. 04	27,576
2016. 03	27,576
2016. 02	27,576
2016. 01	27,576
2015. 12	27,576
2015. 11	27,576
2015. 10	27,576
2015. 09	27,576
2015. 08	27,576
2015. 07	27,576
2015. 06	27,576
2015. 05	27,576
2015. 04	27,576
2015. 03	27,576
2015. 02	27,576
2015. 01	27,576
2014. 12	27,576
2014. 11	27,576
2014. 10	27,576
2014. 09	27,576
2014. 08	27,576
2014. 07	27,576
2014. 06	27,576
2014. 05	27,576
2014. 04	27,576
2014. 03	27,576
2014. 02	27,576
2014. 01	27,576
2013. 12	27,576
2013. 11	27,576
2013. 10	27,576
2013. 09	27,576
2013. 08	27,576
2013. 07	27,576
2013. 06	27,576
2013. 05	27,576
2013. 04	27,576
2013. 03	27,576
2013. 02	27,576
2013. 01	27,576
2012. 12	27,576
2012. 11	27,576
2012. 10	27,576
2012. 09	27,576
2012. 08	27,576
2012. 07	27,576
2012. 06	27,576
2012. 05	27,576
2012. 04	27,576
2012. 03	27,576
2012. 02	27,576
2012. 01	27,576
2011. 12	27,576
2011. 11	27,576
2011. 10	27,576
2011. 09	27,576
2011. 08	27,576
2011. 07	27,576
2011. 06	27,576
2011. 05	27,576
2011. 04	27,576
2011. 03	27,576
2011. 02	27,576
2011. 01	27,576
2010. 12	27,576
2010. 11	27,576
2010. 10	27,576
2010. 09	27,576
2010. 08	27,576
2010. 07	27,576
2010. 06	27,576
2010. 05	27,576
2010. 04	27,576
2010. 03	27,576
2010. 02	27,576
2010. 01	27,576
2009. 12	27,576
2009. 11	27,576
2009. 10	27,576
2009. 09	27,576
2009. 08	27,576
2009. 07	27,576
2009. 06	27,576
2009. 05	27,576
2009. 04	27,576
2009. 03	27,576
2009. 02	27,576
2009. 01	27,576
2008. 12	27,576
2008. 11	27,576
2008. 10	27,576
2008. 09	27,576
2008. 08	27,576
2008. 07	27,576
2008. 06	27,576
2008. 05	27,576
2008. 04	27,576
2008. 03	27,576
2008. 02	27,576
2008. 01	27,576
2007. 12	27,576
2007. 11	27,576
2007. 10	27,576
2007. 09	27,576
2007. 08	27,576
2007. 07	27,576
2007. 06	27,576
2007. 05	27,576
2007. 04	27,576
2007. 03	27,576
2007. 02	27,576
2007. 01	27,576
2006. 12	27,576
2006. 11	27,576
2006. 10	27,576
2006. 09	27,576
2006. 08	27,576
2006. 07	27,576
2006. 06	27,576
2006. 05	27,576
2006. 04	27,576
2006. 03	27,576
2006. 02	27,576
2006. 01	27,576
2005. 12	27,576
2005. 11	27,576
2005. 10	27,576
2005. 09	27,576
2005. 08	27,576
2005. 07	27,576
2005. 06	27,576
2005. 05	27,576
2005. 04	27,576
2005. 03	27,576
2005. 02	27,576
2005. 01	27,576
2004. 12	27,576
2004. 11	27,576
2004. 10	27,576
2004. 09	27,576
2004. 08	27,576
2004. 07	27,576
2004. 06	27,576
2004. 05	27,576
2004. 04	27,576
2004. 03	27,576
2004. 02	27,576
2004. 01	27,576
2003. 12	27,576
2003. 11	27,576
2003. 10	27,576
2003. 09	27,576
2003. 08	27,576
2003. 07	27,576
2003. 06	27,576
2003. 05	27,576
2003. 04	27,576
2003. 03	27,576
2003. 02	27,576
2003. 01	27,576
2002. 12	27,576
2002. 11	27,576
2002. 10	27,576
2002. 09	27,576
2002. 08	27,576
2002. 07	27,576
2002. 06	27,576
2002. 05	27,576
2002. 04	27,576
2002. 03	27,576
2002. 02	27,576
2002. 01	27,576
2001. 12	27,576
2001. 11	27,576
2001. 10	27,576
2001. 09	27,576
2001. 08	27,576
2001. 07	27,576
2001. 06	27,576
2001. 05	27,576
2001. 04	27,576
2001. 03	27,576
2001. 02	27,576
2001. 01	27,576
2000. 12	27,576
2000. 11	27,576
2000. 10	27,576
2000. 09	27,576
2000. 08	27,576
2000. 07	27,576
2000. 06	27,576
2000. 05	27,576
2000. 04	27,576
2000. 03	27,576
2000. 02	27,576
2000. 01	27,576
1999. 12	27,576
1999. 11	27,576
1999. 10	27,576
1999. 09	27,576
1999. 08	27,576
1999. 07	27,576
1999. 06	27,576
1999. 05	27,576
1999. 04	27,576
1999. 03	27,576
1999. 02	27,576
1999. 01	27,576
1998. 12	27,576
1998. 11	27,576
1998. 10	27,576
1998. 09	27,576
1998. 08	27,576
1998. 07	27,576
1998. 06	27,576
1998. 05	27,576
1998. 04	27,576
1998. 03	27,576
1998. 02	27,576
1998. 01	27,576
1997. 12	27,576
1997. 11	27,576
1997. 10	27,576
1997. 09	27,576
1997. 08	27,576
1997. 07	27,576
1997. 06	27,576
1997. 05	27,576
1997. 04	27,576
1997. 03	27,576
1997. 02	27,576
1997. 01	27,576

[그림 2.1.1] 월별 출생아 수³⁾

KOSIS 국가 통계 포털에서 1997. 01 ~ 2019. 12 기간의 월별 출생아 수를 사용하였다.

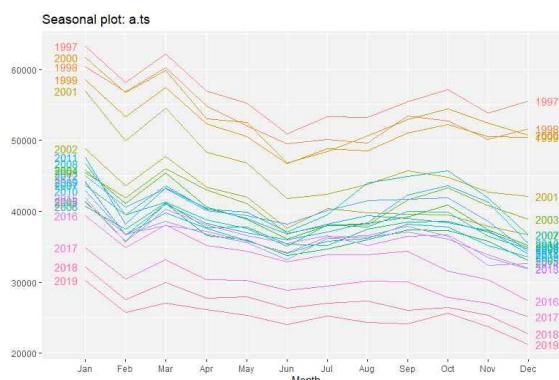
3) https://kosis.kr/statisticsList/statisticsListIndex.do?vwcd=MT_ZTITLE&menuId=M_01_01#content-group

○ 데이터 분포

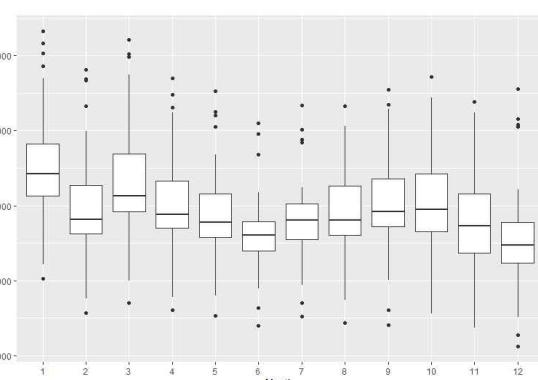


[그림 2.1.2] 출생아 수 분포

출생아 수 분포에서 출생아 수가 감소하는 추세와 계절성이 보이고, 변동 폭의 변화도 크지 않게 보인다.



[그림 2.1.3] 연도별 월별 출생아 수



[그림 2.1.4] 월별 출생아 수

월별 출생아 수는 주로 1~3월의 출생아 수가 많았고, 6~7월, 12월이 적게 나타났다

○ Train & Test 자료 분리

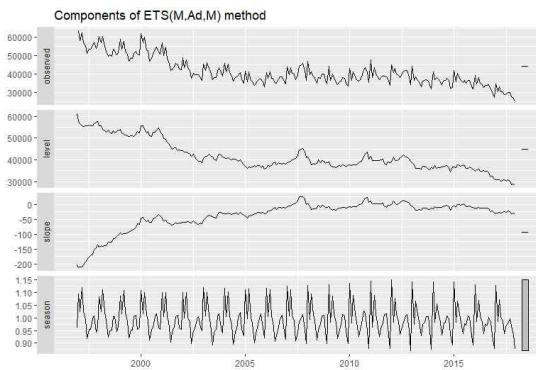
Train	1997. 01 ~ 2017. 12 기간의 출생아 수
Test	2018. 01 ~ 2019. 12 기간의 출생아 수

[표 2.1.1] Train & Test 자료 분리

시계열 예측 모형의 예측 오차를 비교하기 위해서 2년치 데이터를 Test data로 분리시켜 분석을 진행한다.

2. ETS 모델 분석

○ ETS 모형 적합



[그림 2.2.1] 분해 요소별 분포

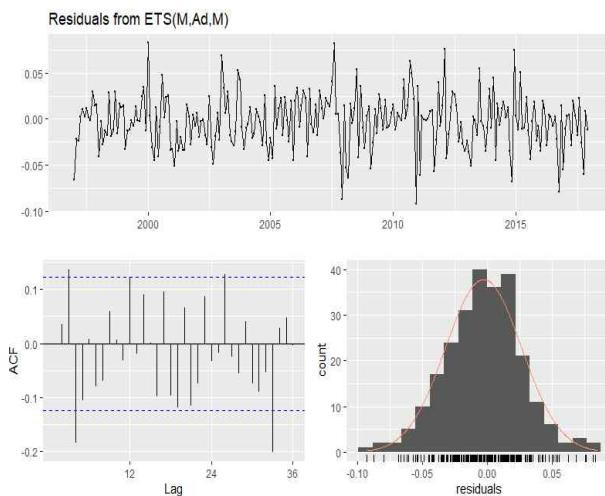
평활화 지수	
alpha	0.7617
beta	0.0042
gamma	0.2383
phi	0.9676

[표 2.2.1] 평활화 지수

분석을 시작하기에 앞서 Train data의 변동폭의 변화가 크지 않으므로 변환은 하지 않고 분석을 진행한다.

ETS모델 적합 결과 ETS(M, Ad, M) 모형이 적합되었다. alpha가 0.7617로 level에서는 큰 변동이 있고, gamma가 0.2383으로 계절에서는 약간의 변동이 있으며, beta가 0.0042로 기울기는 일정하게 나타났다.

○ ETS 모형 진단



[그림 2.2.2] 잔차 독립성 검정 결과

잔차의 독립성검정 (Ljung-Box test)	
통계량(Q)	47.462
df	7
P-value	4.535e-08

[표 2.2.2] 잔차 독립성 검정 결과

H_0 : 잔자는 독립이다.

H_1 : 잔자는 독립이 아니다.

ETS모델을 사용하기 전 오차의 가정을 만족 여부를 확인하기 위해 잔차의 독립성 검정을 실시한다. 검정 결과 P-value가 4.535e-08으로 귀무가설을 기각하여 잔차는 독립이 아니라는 충분한 근거가 있다.

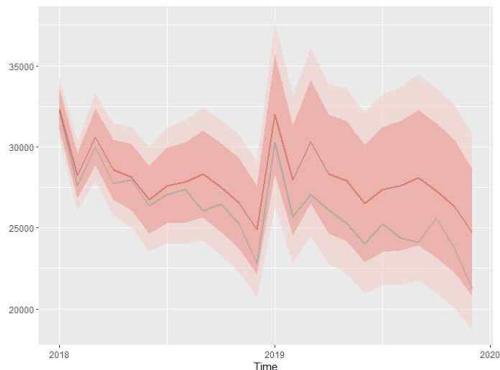
따라서, 분포가 정규성을 띠는 것으로 보이고, 독립성 가정이 위반되었다. 독립이 아니어도 예측값에는 큰 문제가 없지만, 예측 구간이 좁아질 수 있어 신뢰성이 떨어질 수 있다는 것을 인지하고 예측을 진행한다.

○ 예측 결과

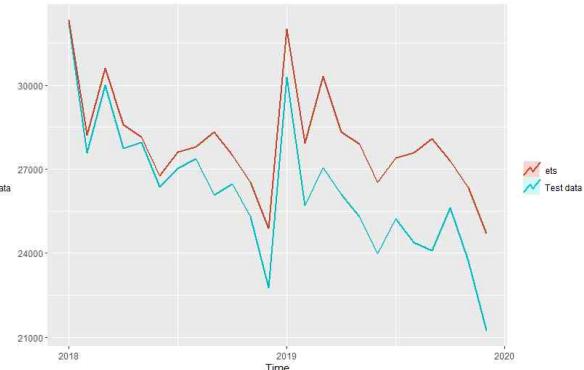
RMSE	MASE	MAPE
2077.466	0.654	7.009

[표 2.2.3] 예측 결과 주요 지표

ETS(M, Ad, M) 모형을 통한 예측과 Test data의 차이는 RMSE가 2077.466으로 나타났고, MASE가 0.653으로 1보다 작게 나타났다.



[그림 2.2.3] ETS모형 예측 1



[그림 2.2.4] ETS모형 예측 2

예측이 Test data와 약간의 차이가 있지만, 형태가 비슷하게 나타났고, 신뢰구간 안에 포함되어 있는 것을 확인할 수 있다.

3. ARIMA 모델 분석

○ 정상성 여부 확인

KPSS Unit Root Test					
critical values	10 pct	5 pct	2.5 pct	1 pct	test-statistic
	0.347	0.463	0.574	0.739	3.1342

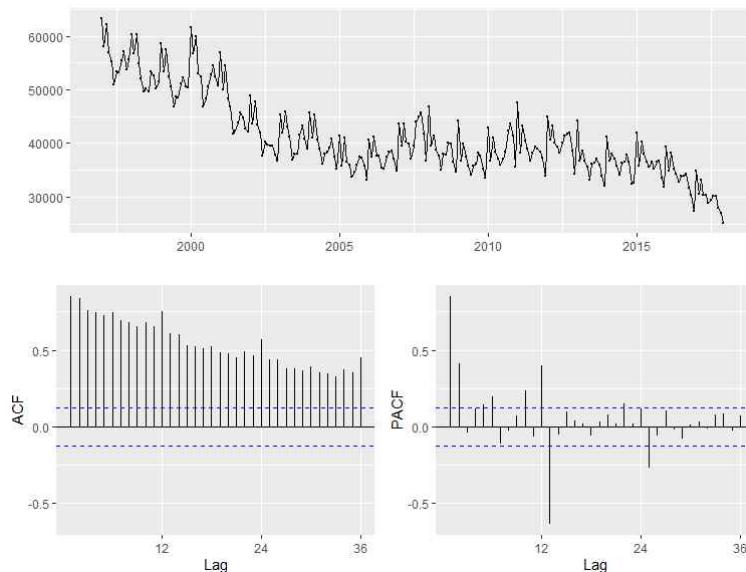
[표 2.3.1] KPSS 검정

H_0 : 시계열 데이터가 정상성을 가진다.

H_1 : 시계열 데이터가 정상성을 가지지 않는다.

ARIMA 모델을 분석에 사용하기 위해서는 시계열 데이터가 정상성을 만족해야 하는 가정이 필요하다. 따라서, 정상성 여부를 먼저 확인하며, ETS 분석과 마찬가지로 데이터의 변환은 하지 않고 진행한다.

Train data로 KPSS검정을 실시한 결과 통계량 값이 3.1342로 유의수준 5%에서의 0.463보다 크므로 귀무가설을 기각하여 출생아 수 데이터는 정상성을 가지지 않는다는 충분한 근거가 있다.



[그림 2.3.1] Train data의 ACF, PACF

또한, 그림[2.3.1]을 보면 Train data의 ACF가 감소하는 추세를 보이고 있으며, 12, 24, 36 시점에서 조금 증가한 것을 보아 계절 변동 요인도 강하게 있는 것으로 보인다.

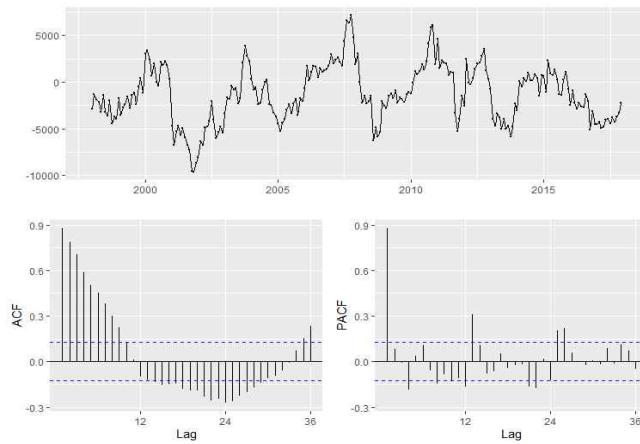
결과를 종합하면 출생아 수 데이터는 정상성을 나타내지 못하는 데이터로 보인다.

○ 일반 & 계절 차분 차수 결정

단위근 검정	
ndiffs	1
nsdiffs	1

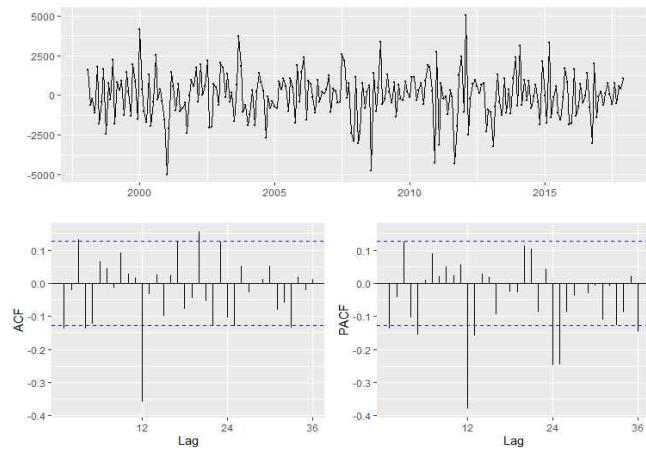
[표 2.3.2] 단위근 검정

Train data가 정상성을 나타내지 못하므로 차분 차수를 결정해야 한다. 차분 차수에 대한 검정으로 단위근 검정을 실시한 결과 ndiffs, nsdiffs가 각각 1씩 나타났고, 위의 그림[2.3.1]의 ACF를 보면 일반차분과 계절차분이 모두 필요한 것으로 보인다.



[그림 2.3.2] 계절 차분 실시 후 ACF, PACF

먼저, 계절 차분을 실시한 결과 여전히 ACF가 천천히 감소하는 추세이지만, 계절 변동 요인은 많이 감소한 것으로 보인다. 여전히 정상성을 만족하지 않는 것으로 보고, 1차 차분을 실시한다.



[그림 2.3.3] 계절 & 일반 차분 실시 후 ACF, PACF

1차 차분까지 실시한 결과 정상성을 만족하는 것으로 보이고, 일반 & 계절 차분의 차수 d와 D는 1이 적합하다고 판단된다.

○ ARIMA 모형 적합

AR과 MA의 차수를 결정하기 위한 첫 번째 방법으로 ACF, PACF를 근거로 선택하는 것과 두 번째 방법인 auto.arima를 사용하여 최소의 AICc를 가지는 모형을 선택하는 방법이 있다.

첫 번째 방법으로 그림[2.3.3]을 비계절형 관점에서 보면 ACF, PACF 모두 1시점이 같은 상황에 2시점이 들어가있는 형태를 보인다. 이것을 절단으로 본다면 2시점에서 더 많이 들어가 있는 ACF를 절단으로 보고 MA1 혹은 MA2 모형으로 볼 수 있고, 둘 다 감소로 본다면 ARMA모형으로 볼 수도 있을 것이다.

계절형 관점에서 보면 ACF가 절단, PACF가 감소하는 형태를 띠고 있으므로, 계절형의 차수는 SMA1이 적합하다고 판단된다.

두 번째 방법인 auto.arima를 통하여 나타난 최소의 AICc 모형은 ARIMA(0,1,4)(0,1,1)[12]이고, 두 번째는 ARIMA(2,1,2)(0,1,1)[12] 이다.

ARIMA모형	AICc
ARIMA(0,1,4)(0,1,1)[12]	4096.654
ARIMA(2,1,2)(0,1,1)[12]	4098.021

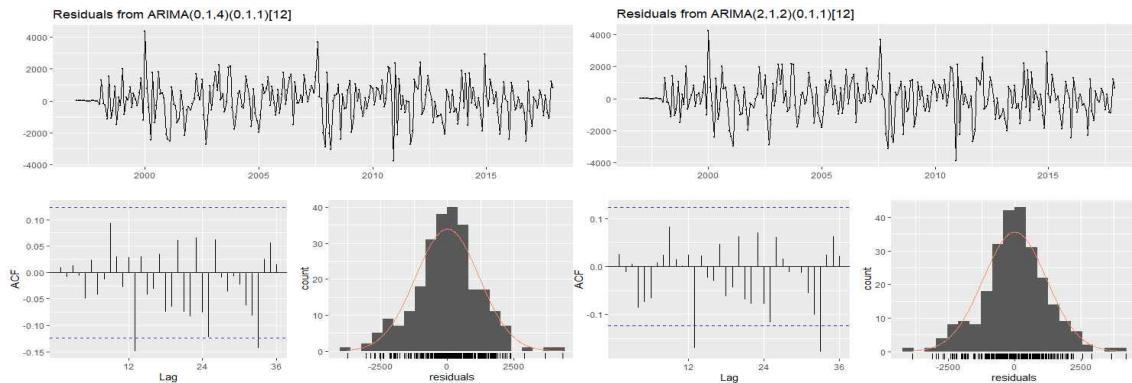
[표 2.3.3] 최소의 AICc 모델

선택된 모형들의 비계절형 차수가 ACF,PACF를 근거로 선택하려던 모형과 비슷한 형태이므로 더 나은 모형을 선택하기 위해 두 모형 모두 예측한 후 비교하여 모델을 선택해 보겠다.

○ ARIMA 모형 진단

잔차의 독립성검정 (Ljung-Box test)			
모형	통계량	df	P-value
ARIMA(0,1,4)(0,1,1)[12]	21.664	19	0.3013
ARIMA(2,1,2)(0,1,1)[12]	24.588	19	0.1745

[표 2.3.4] 모형별 잔차의 독립성검정



[그림 2.3.4] ARIMA(0,1,4)(0,1,1)[12]의 독립성 검정

[그림 2.3.5] ARIMA(2,1,2)(0,1,1)[12]의 독립성검정

H_0 : 잔차는 독립이다.

H_1 : 잔차는 독립이 아니다.

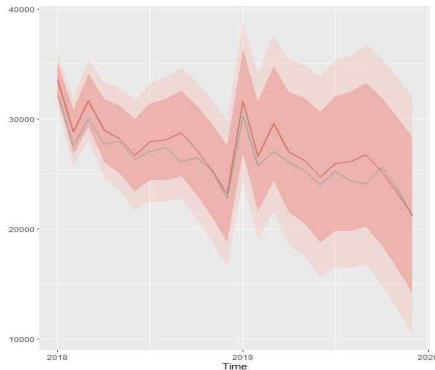
두 모형의 P-value가 각각 0.3013, 0.1745이고, 유의수준 5%에서 귀무가설을 기각하지 못하므로, 잔차는 독립이라는 충분한 근거가 있다. 각각의 ACF에서도 2시점에서 튀어나와 있지 만, 그렇게 큰 차이가 나지는 않다고 판단되므로 두 모형의 잔차는 백색잡음이라고 할 수 있다.

○ 예측 결과

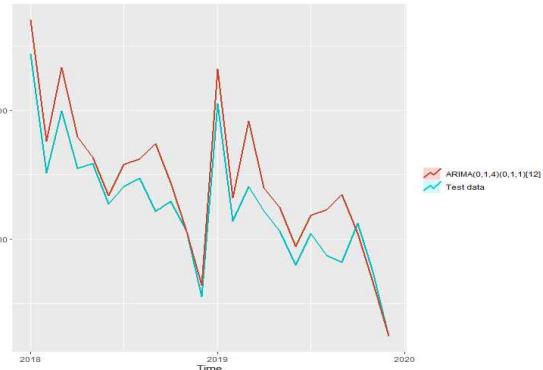
모형	RMSE	MASE	MAPE
ARIMA(0,1,4)(0,1,1)[12]	1274.239	0.383	3.876
ARIMA(2,1,2)(0,1,1)[12]	1173.846	0.352	3.583

[표 2.3.5] 예측 결과 주요 지표

ARIMA(2,1,2)(0,1,1)[12] 모형의 예측 오차가 ARIMA(0,1,4)(0,1,1)[12] 모형의 예측 오차보다 더 작다.



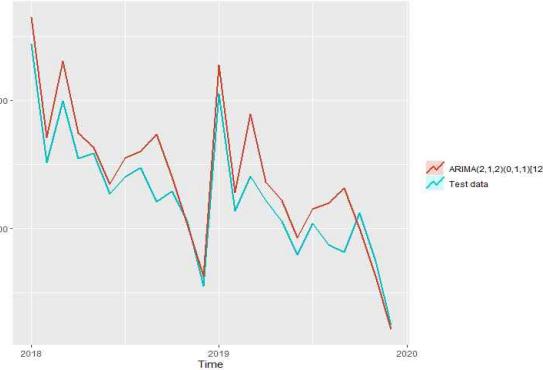
[그림 2.3.6] ARIMA(0,1,4)(0,1,1)[12] 예측 1



[그림 2.3.7] ARIMA(0,1,4)(0,1,1)[12] 예측 2



[그림 2.3.8] ARIMA(2,1,2)(0,1,1)[12] 예측 1



[그림 2.3.9] ARIMA(2,1,2)(0,1,1)[12] 예측 2

두 모형의 예측이 Test data와 약간의 차이가 있지만, 신뢰구간 안에 포함되어 있는 것을 확인할 수 있다. 모형의 예측 그래프들을 눈으로 봤을 때 차이점을 찾기 힘들고, 매우 비슷하게 나타난 것으로 보인다.

따라서, 최종 ARIMA 모형은 오차가 더 작은 ARIMA(2,1,2)(0,1,1)[12] 모형이고, 모형식은 $(1+0.614B+0.707B^2)(1-B^{12})(1-B)Y_t = (1+0.378B+0.642B^2)(1-0.65B^{12})\epsilon^t$ 이다.

ARIMA(2,1,2)(0,1,1)[12]	ar1	ar2	ma1	ma2	sma1
계 수	-0.614	-0.707	0.378	0.642	-0.650

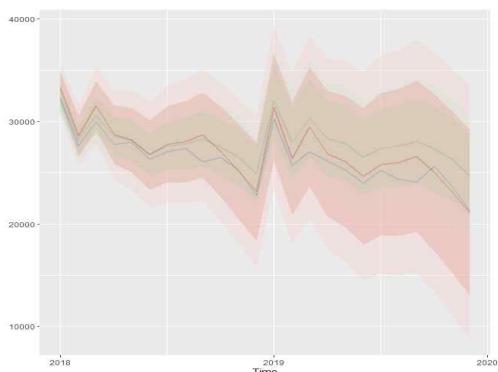
[표 2.3.6] 최종 ARIMA 모형의 모수

III. 결론

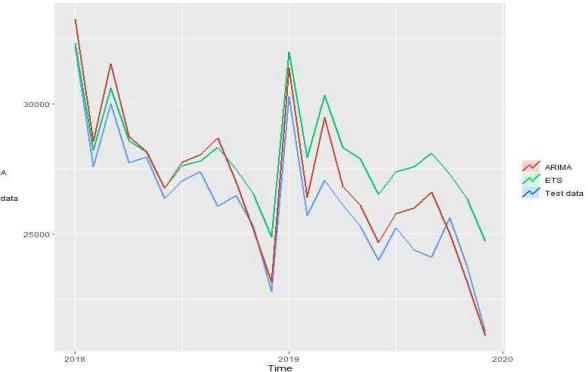
○ 최종 시계열 모형 선택

	ETS(M, Ad, M)	ARIMA(2,1,2)(0,1,1)[12]
RMSE	2077.466	1173.846
MASE	0.654	0.352
MAPE	7.009	3.583
잔차의 독립성검정	백색잡음 O	백색잡음 X

[표 3.1] 모형별 예측 비교



[그림 3.1] 모형별 예측 비교 1



[그림 3.2] 모형별 예측 비교 2

모형별 분석 결과 예측 오차가 더 작고 그래프상에서 Test data와 조금 더 가까운 ARIMA 모형을 분석에 사용하는 것이 적절하다.

○ 결론

ARIMA모형으로 분석한 결과 출생아 수가 앞으로도 계속 줄어들 것으로 예측되며, 예측된 출생아 수를 이용하여 예산을 측정해 지원을 늘리는 등 저출생 문제에 대한 대책이 강구되어야 할 필요가 있다고 보인다.

[부록] R 코드

```
library(tidyverse)
library(forecast)
library(urca)

# 데이터 불러오기
birth <- read.csv("C:/Data/출생아 수.csv")
birth <- birth[-1,]
birth$전국 <- as.integer(birth$전국)
str(birth)
summary(birth)

# 시계열데이터 변환 & 분포확인
birth.ts <- ts(birth$전국,start=c(1997,1), freq=12)
birth.ts

autoplot(birth.ts)

ggseasonplot(birth.ts, year.labels=TRUE,year.labels.left = TRUE)

ggsubseriesplot(birth.ts)
data.frame(birth.ts=as.numeric(birth.ts), mon=as.factor(cycle(birth.ts))) %>%
  ggplot() +
  geom_boxplot(aes(x=mon,y=birth.ts)) +
  labs(x="Month")

# train & test 분할
train <- window(birth.ts , end=c(2017,12))
test <- window(birth.ts, start=c(2018,1))

## ETS 모델
# ETS 모델 적합 & 가정 검정
fit1 <- ets(train)
summary(fit1)

autoplot(fit1)
checkresiduals(fit1)
```

```

# 예측
fc1 <- forecast(fit1, h = length(test))
accuracy(fc1, test)

autoplot(train) +
  autolayer(test, series="Test data", size=1) +
  autolayer(fc1, series="ets", size=1, PI=FALSE) +
  labs(y= NULL, color=NULL)

autoplot(test, series="Test data", size=1) +
  autolayer(fc1, series="ets", size=1, alpha=0.5) +
  labs(y= NULL, color=NULL)

autoplot(test, series="Test data", size=1) +
  autolayer(fc1, series="ets", size=1, PI = FALSE) +
  labs(y= NULL, color=NULL)

## ARIMA모델
# 정상성만족 확인

train %>% ur.kpss() %>% summary()
ggrtsdisplay(train)

# 비계절,계절 차분 확인
ndiffs(train)
nsdiffs(train)

train_d <- diff(train,lag=12)
ggrtsdisplay(train_d)
ggrtsdisplay(diff(train_d))

```

```

# arima모형 적합
fit2 <- auto.arima(train,d=1,stepwise=FALSE,
                     approximation = FALSE,trace=TRUE)
fit3 <- auto.arima(train,d=1)

# 모델 가정 가설 검정
checkresiduals(fit2)
checkresiduals(fit3)

# 예측
fc2 <- forecast(fit2)
fc3 <- forecast(fit3)

# 정확도 비교
accuracy(fc2, test)
accuracy(fc3, test)

# 예측 그래프

autoplot(test, series="Test data", size=1) +
  autolayer(fc2, series="ARIMA(0,1,4)(0,1,1)[12]", size=1, alpha=0.5) +
  labs(y= NULL, color=NULL)

autoplot(test, series="Test data", size=1) +
  autolayer(fc2, series="ARIMA(0,1,4)(0,1,1)[12]", size=1, PI = FALSE) +
  labs(y= NULL, color=NULL)

autoplot(test, series="Test data", size=1) +
  autolayer(fc3, series="ARIMA(2,1,2)(0,1,1)[12]", size=1, alpha=0.5) +
  labs(y= NULL, color=NULL)

autoplot(test, series="Test data", size=1) +
  autolayer(fc3, series="ARIMA(2,1,2)(0,1,1)[12]", size=1, PI = FALSE) +
  labs(y= NULL, color=NULL)

```

```
# 최종 모형
coef(fit3)

confint(fit3)

# ETS & ARIMA 비교
autoplot(test, series="Test data", size=1) +
  autolayer(fc1, series="ETS", size=1, alpha=0.3) +
  autolayer(fc3, series="ARIMA", size=1, alpha=0.3) +
  labs(y= NULL, color=NULL)

autoplot(test, series="Test data", size=1) +
  autolayer(fc1, series="ETS", size=1, PI = FALSE) +
  autolayer(fc3, series="ARIMA", size=1, PI = FALSE) +
  labs(y= NULL, color=NULL)
```