

Math 320 Research Project: Forecasting Birth Rates by Race

Rachel Hong and Joshulyne Park

December 15, 2016

Contents

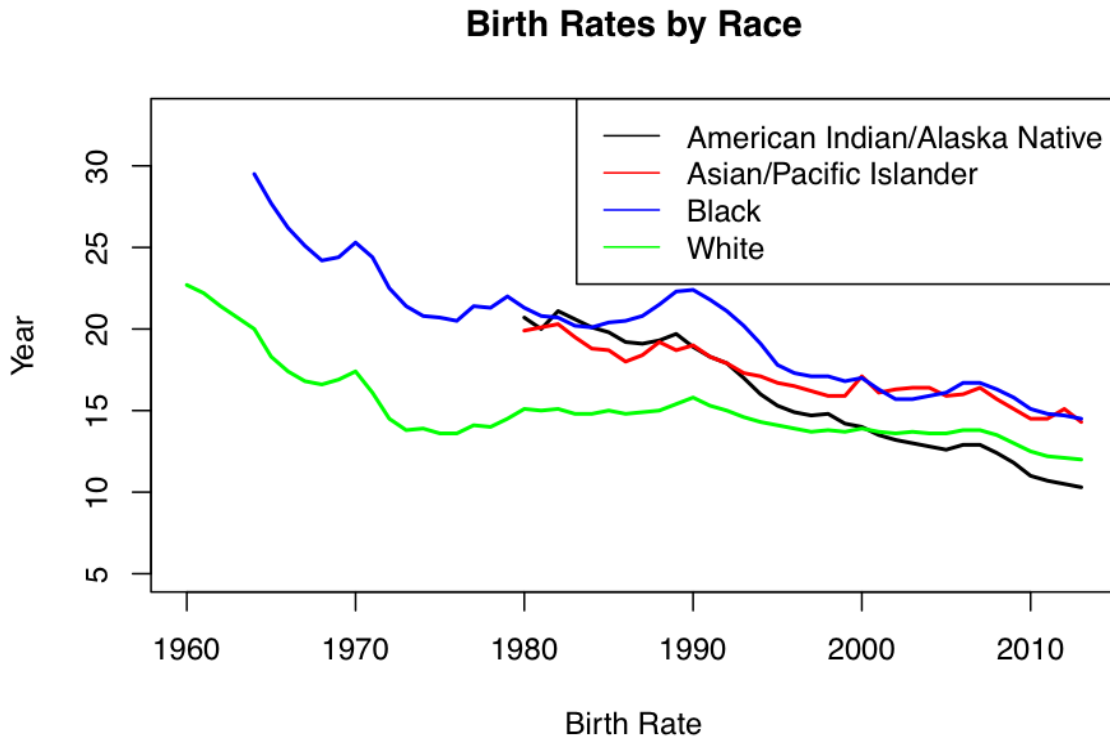
Summary	1
Data Prep	2
Time Series Models	2
Linear Model	2
Exponential Model	4
Autoregressive Integrated Moving Average Model	6
Forecast Comparison / Model Selection	8
Conclusion	11
Models: In-Depth Analysis	11
R Code	11

Summary

Our project focuses on a dataset from the National Center for Health Statistics, “NCHS - Births, Birth Rates, and Fertility Rates, by Race of Mother: United States, 1960-2013.” It contains data regarding birth rates, categorized by race of the mother from the time period 1960-2013. The general perceived notion that we are exploring is that birth rates have been in a decline post the “baby boom” years. To confirm this idea, we will perform experiments to develop a more detailed understanding of this trend. Using R, we will build time series models that forecast birth rates for each race. The three main models we will examine and implement are: linear, exponential, and autoregressive integrated moving average (ARIMA). All three are commonly used time series models, starting from the most simple to a more sophisticated model that takes into consideration important time series elements such as trend and seasonality. While some of these models are already functions within R, we will break down the mathematical formulas used to generate these models and the forecasts. For each dataset, we will find the best model by analyzing the model fit and the residuals and ultimately select the best model with the smallest error, or the root mean squared error. Once our model is chosen, we will forecast birth rates for each subset of data for the next 10 years.

Data Prep

While the data was generally clean, there were a few missing values that needed to be taken out. After the data consisted of only relevant values, we organized the data into four subsets, one for each race in order to run our time series models. Below is a plot of birth rates by race.



Time Series Models

Linear Model

10 Year Horizon Forecast Values for Race: American Indian/Alaska Native

##	Point Forecast	Lo 95	Hi 95
## 2014	9.672193	8.326344	11.018041
## 2015	9.328755	7.976330	10.681179
## 2016	8.985317	7.625986	10.344648
## 2017	8.641879	7.275316	10.008443
## 2018	8.298442	6.924325	9.672558
## 2019	7.955004	6.573018	9.336989
## 2020	7.611566	6.221402	9.001730
## 2021	7.268128	5.869481	8.666775
## 2022	6.924691	5.517261	8.332120
## 2023	6.581253	5.164748	7.997758

10 Year Horizon Forecast Values for Race: Asian/Pacific Islander

##	Point Forecast	Lo 95	Hi 95
----	----------------	-------	-------

## 2014	14.29893	13.12239	15.47547
## 2015	14.13752	12.95524	15.31981
## 2016	13.97612	12.78780	15.16444
## 2017	13.81471	12.62007	15.00936
## 2018	13.65331	12.45206	14.85456
## 2019	13.49190	12.28377	14.70003
## 2020	13.33050	12.11522	14.54577
## 2021	13.16909	11.94640	14.39178
## 2022	13.00769	11.77731	14.23806
## 2023	12.84628	11.60797	14.08458

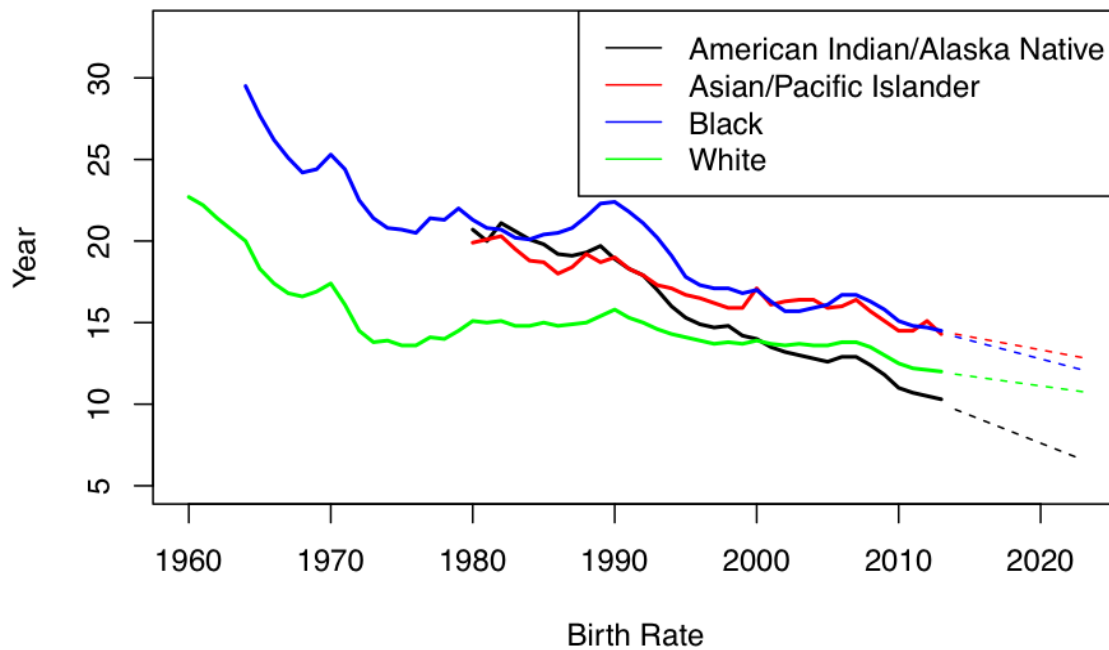
10 Year Horizon Forecast Values for Race: Black

##	Point Forecast	Lo 95	Hi 95
## 2014	14.14220	11.285594	16.99881
## 2015	13.91406	11.050863	16.77725
## 2016	13.68591	10.815895	16.55592
## 2017	13.45776	10.580691	16.33482
## 2018	13.22961	10.345254	16.11396
## 2019	13.00146	10.109584	15.89334
## 2020	12.77331	9.873685	15.67294
## 2021	12.54516	9.637557	15.45277
## 2022	12.31701	9.401203	15.23282
## 2023	12.08886	9.164624	15.01311

10 Year Horizon Forecast Values for Race: White

##	Point Forecast	Lo 95	Hi 95
## 2014	11.83990	8.663654	15.01615
## 2015	11.71987	8.537330	14.90241
## 2016	11.59984	8.410794	14.78889
## 2017	11.47981	8.284048	14.67557
## 2018	11.35978	8.157092	14.56247
## 2019	11.23975	8.029929	14.44957
## 2020	11.11972	7.902560	14.33688
## 2021	10.99969	7.774986	14.22439
## 2022	10.87966	7.647208	14.11211
## 2023	10.75963	7.519228	14.00003

Linear Model: Birth Rates by Race



Exponential Model

10 Year Horizon Forecast Values for Race: American Indian/Alaska Native

##	Point Forecast	Lo.95	Hi.95
## 2014	10.349156	9.512041	11.259943
## 2015	10.120373	9.297930	11.015563
## 2016	9.896646	9.088451	10.776711
## 2017	9.677866	8.883509	10.543254
## 2018	9.463922	8.683014	10.315061
## 2019	9.254708	8.486877	10.092006
## 2020	9.050118	8.295009	9.873967
## 2021	8.850052	8.107324	9.660823
## 2022	8.654408	7.923737	9.452457
## 2023	8.463089	7.744164	9.248755

10 Year Horizon Forecast Values for Race: Asian/Pacific Islander

##	Point Forecast	Lo.95	Hi.95
## 2014	14.45608	13.51625	15.46126
## 2015	14.32073	13.38531	15.32153
## 2016	14.18666	13.25541	15.18332
## 2017	14.05384	13.12657	15.04660
## 2018	13.92226	12.99877	14.91136
## 2019	13.79191	12.87201	14.77756
## 2020	13.66279	12.74629	14.64519
## 2021	13.53487	12.62160	14.51422

```
## 2022      13.40815 12.49795 14.38464
## 2023      13.28262 12.37533 14.25643
```

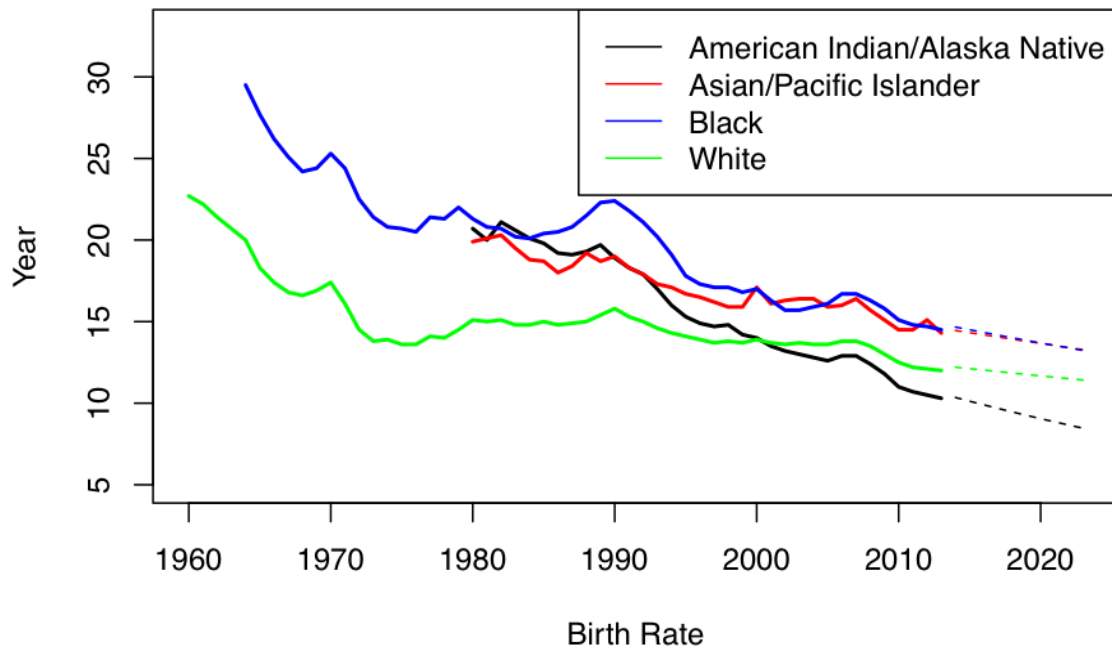
10 Year Horizon Forecast Values for Race: Black

```
##      Point Forecast    Lo.95    Hi.95
## 2014      14.66684 12.87004 16.71450
## 2015      14.49956 12.71942 16.52883
## 2016      14.33418 12.57042 16.34541
## 2017      14.17069 12.42304 16.16420
## 2018      14.00906 12.27725 15.98516
## 2019      13.84928 12.13304 15.80828
## 2020      13.69132 11.99040 15.63352
## 2021      13.53516 11.84932 15.46085
## 2022      13.38078 11.70977 15.29024
## 2023      13.22816 11.57175 15.12168
```

10 Year Horizon Forecast Values for Race: White

```
##      Point Forecast    Lo.95    Hi.95
## 2014      12.20612 10.191368 14.61917
## 2015      12.11570 10.112260 14.51606
## 2016      12.02595 10.033646 14.41386
## 2017      11.93687  9.955523 14.31255
## 2018      11.84845  9.877892 14.21211
## 2019      11.76068  9.800750 14.11255
## 2020      11.67356  9.724097 14.01386
## 2021      11.58709  9.647931 13.91601
## 2022      11.50126  9.572251 13.81900
## 2023      11.41606  9.497056 13.72283
```

Exponential Model: Birth Rates by Race



Autoregressive Integrated Moving Average Model

10 Year Horizon Forecast Values for Race: American Indian/Alaska Native

##	Point Forecast	Lo 95	Hi 95
## 2014	9.984848	9.163265	10.806432
## 2015	9.669697	8.507802	10.831592
## 2016	9.354545	7.931520	10.777571
## 2017	9.039394	7.396226	10.682562
## 2018	8.724242	6.887125	10.561360
## 2019	8.409091	6.396629	10.421552
## 2020	8.093939	5.920232	10.267646
## 2021	7.778788	5.454997	10.102578
## 2022	7.463636	4.998884	9.928388
## 2023	7.148485	4.550408	9.746562

10 Year Horizon Forecast Values for Race: Asian/Pacific Islander

##	Point Forecast	Lo 95	Hi 95
## 2014	14.13030	13.123130	15.13748
## 2015	13.96061	12.536249	15.38496
## 2016	13.79091	12.046435	15.53538
## 2017	13.62121	11.606866	15.63556
## 2018	13.45152	11.199408	15.70362
## 2019	13.28182	10.814759	15.74888
## 2020	13.11212	10.447392	15.77685
## 2021	12.94242	10.093709	15.79114
## 2022	12.77273	9.751209	15.79425
## 2023	12.60303	9.418070	15.78799

10 Year Horizon Forecast Values for Race: Black

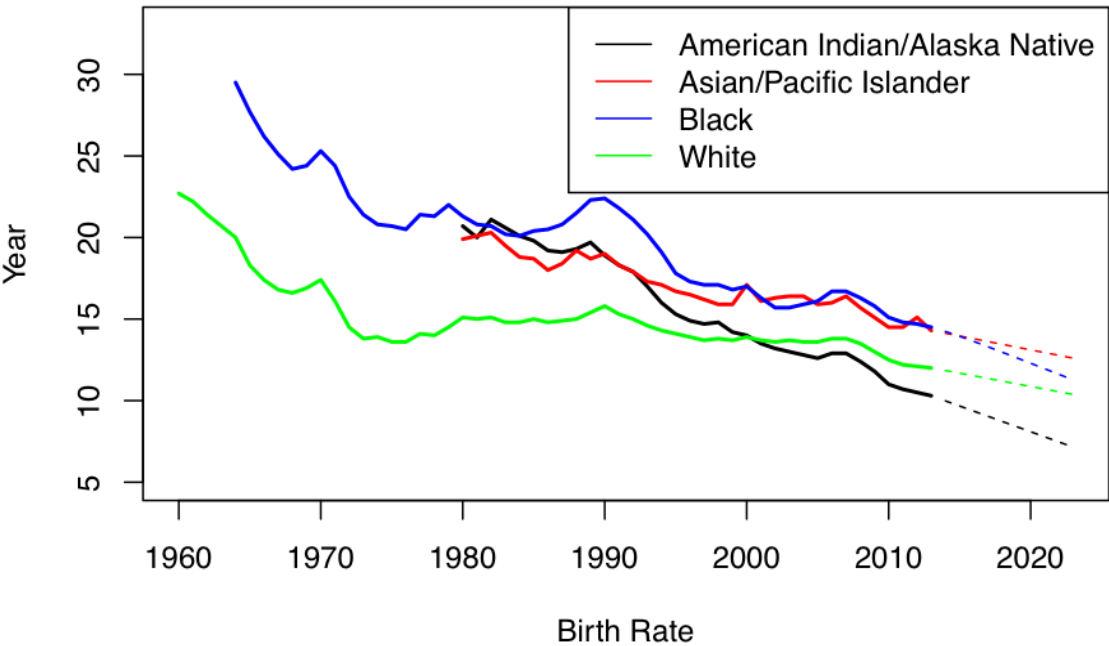
##	Point Forecast	Lo 95	Hi 95
## 2014	14.23952	13.156519	15.32251
## 2015	13.94421	11.923109	15.96532
## 2016	13.62886	10.739130	16.51860
## 2017	13.30198	9.626053	16.97790
## 2018	12.96845	8.583371	17.35353
## 2019	12.63109	7.603422	17.65877
## 2020	12.29154	6.677137	17.90594
## 2021	11.95072	5.796088	18.10535
## 2022	11.60916	4.953081	18.26525
## 2023	11.26719	4.142192	18.39219

10 Year Horizon Forecast Values for Race: White

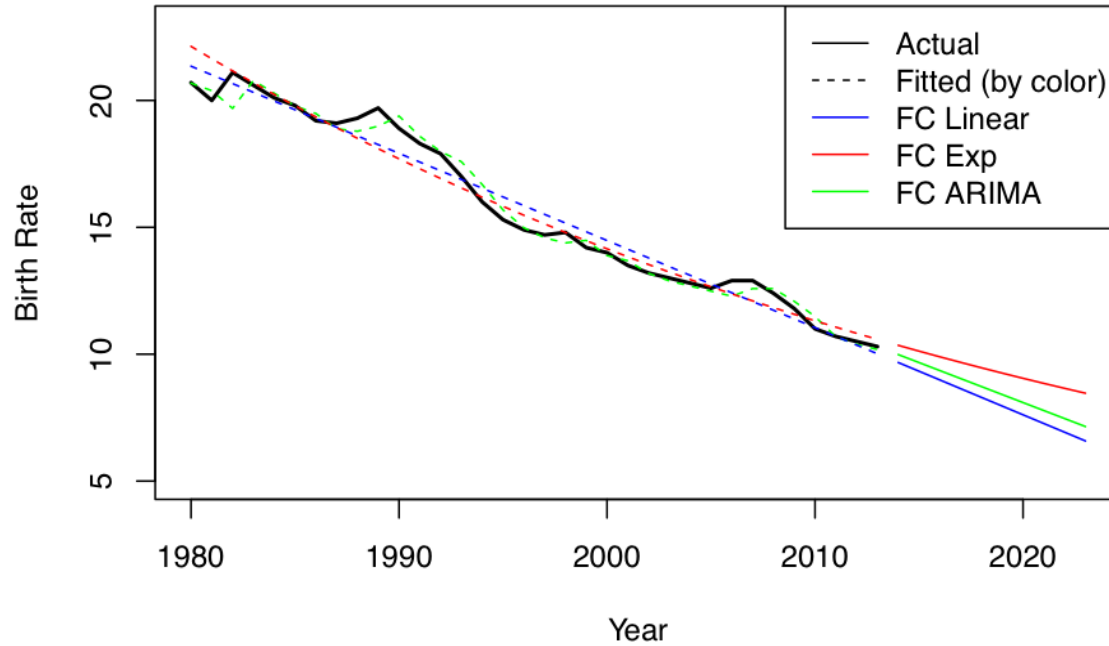
##	Point Forecast	Lo 95	Hi 95
## 2014	11.84784	11.043678	12.65201
## 2015	11.68438	10.126049	13.24271
## 2016	11.52092	9.357491	13.68434
## 2017	11.35745	8.626565	14.08834
## 2018	11.19399	7.905307	14.48267

## 2019	11.03053	7.182348	14.87871
## 2020	10.86706	6.452099	15.28203
## 2021	10.70360	5.711531	15.69567
## 2022	10.54014	4.958916	16.12136
## 2023	10.37667	4.193249	16.56010

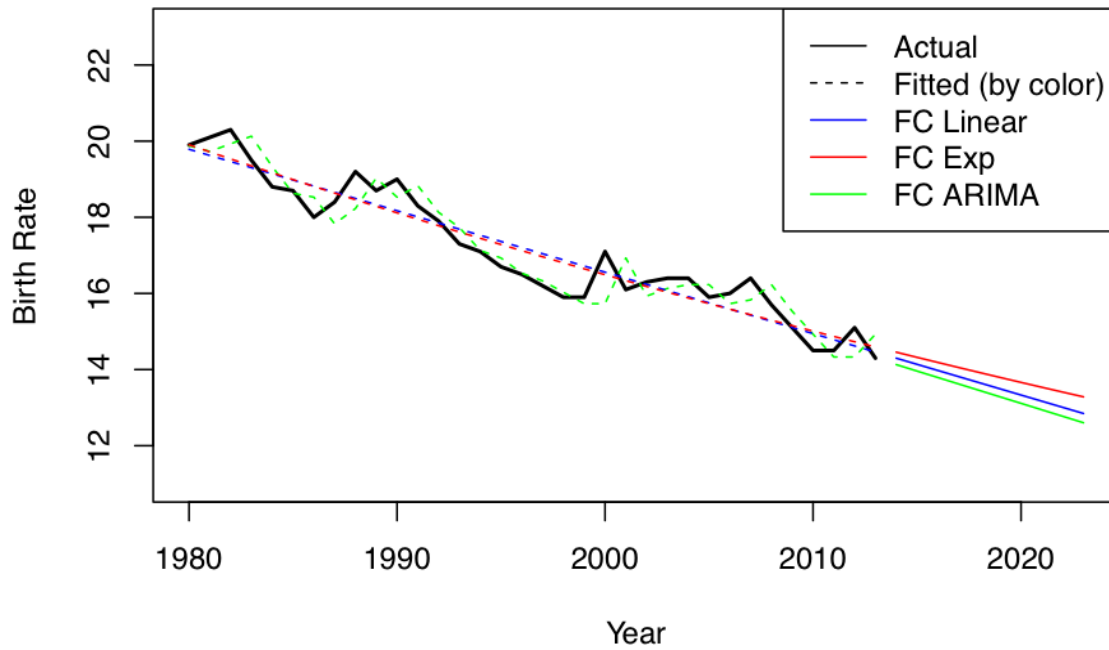
ARIMA Model: Birth Rates by Race



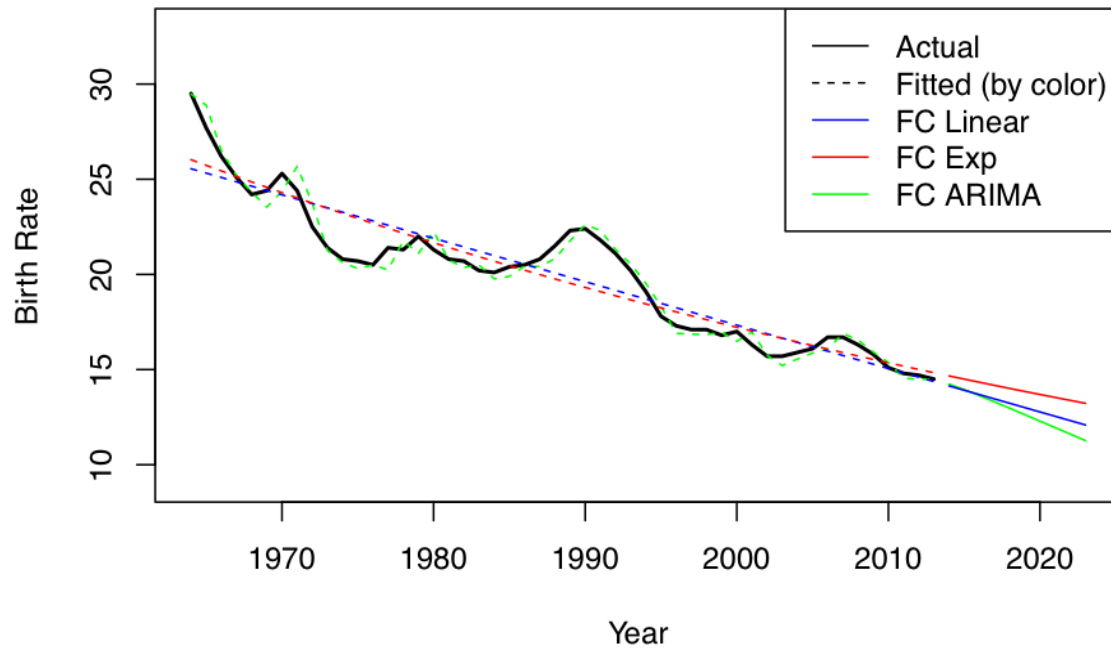
Forecast Comparison / Model Selection

Fitted and Forecast Values for Race: American Indian/Alaska Native

Model	RMSE Value
Linear	0.6048757
Exponential	0.6541382
ARIMA	0.4066674

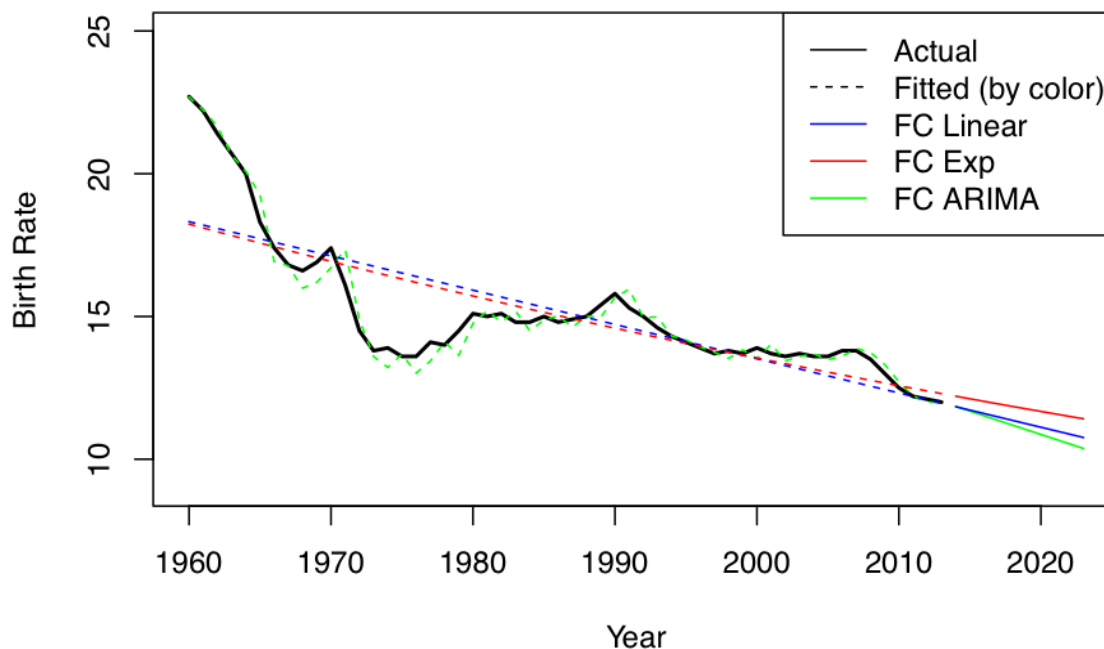
Fitted and Forecast Values for Race: Asian/Pacific Islander

Model	RMSE Value
Linear	0.5287808
Exponential	0.5139744
ARIMA	0.4985302

Fitted and Forecast Values for Race: Black

Model	RMSE Value
Linear	1.33798
Exponential	1.309761
ARIMA	0.5357265

Fitted and Forecast Values for Race: Black



Model	RMSE Value
Linear	1.497296
Exponential	1.455345
ARIMA	0.3948074

Conclusion

For all of our models, the RMSE values revealed that the ARIMA model was the best fit for each of subset of data. This is unsurprising as the ARIMA model really takes into consideration the historical data in creating its forecasts unlike the linear and exponential model. In conclusion, our assumption that birth rates are in a decline were confirmed to be true. Although some datasets had periods of increase, overall the general declining trend overwhelmed such periods.

Models: In-Depth Analysis

Please see the attached document named “In-Depth Analysis on Linear, Exponential, and ARIMA models”.

R Code

Please see the attached appendix at the end.

In-Depth Analysis on Linear, Exponential, and ARIMA models

Linear Model

Linear least squared regression is a basic method for finding the relationship between two variables. It is used in time series modeling to find a trend line that represents the observed values (the dependent variable) as a linear function of time (the independent variable) so that we can see how it changes over the years.

The result is a linear equation in the form $y_t = \beta_t x_t + \epsilon$, where y_t is the dependent variable and x_t is the vector of independent variables, which are also called the regressors or the predictor variables. For our model, x_t represents time t . β_i is the vector of regression coefficients, which represents the change in y for a one-unit change in x_t – in other words, it is the slope; β_0 is the intercept ($x_0 = 1$); ϵ is the error term. The residual is the difference between the approximate value of y_t and the true value of y , and with the least-squared method, the linear equation is determined by finding the values for β_t that minimize the sum of the squared residuals.

For our data set, we conduct an analysis of birth rates over time for each race. Since our data is time series, our independent variable is time in units of years and our dependent variable is the birth rate by rate. β_1 represents the amount that we expect the dependent variable (birth rate) to change each year, and β_0 represents the birth rate during the first year in our analysis.

To find the coefficients that minimize the sum of the squared residuals, we can write the equation $y_t = \beta_t x_t + \epsilon$ in matrix form:

$$\{y\} = \{\beta\}[X] + \{\epsilon\}$$

For our model, we have:

$$[X] = \begin{bmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_2 \end{bmatrix}$$

$\{\beta\}$ contains the coefficients, $\{\epsilon\}$ contains the residuals, and $\{y\}$ contains the observed values of the dependent variables. They are all column variables. To solve for the coefficients, we solve the system of equations:

$$\{\beta\} = ([X]^T [X]) \setminus ([X]^T * \{y\})$$

To solve for the residuals between our fitted linear equation and the actual values of the data, we use the equation:

$$\{\epsilon\} = \{y\} - \{y_i\} = \{y\} - [X]\{\beta\}$$

For our analysis of this data set, we used the *tslm* function in R, which is a wrapper for *lm* that is specific to time series. We pass in a vector regressed on trend, where the vector is numeric responses, and the function uses QR decomposition to solve the matrix system of equations. The output includes coefficients, residuals, and fitted values.

Exponential Trend Model

Exponential models are useful and common in time series because it allows us to look at growth rates over time. Often, the observed data does not follow a linear trend, but the change in the observed data does. Exponential modeling is especially applicable when looking at the growth of a population, which is often non-linear when we look at it in intervals but linear once we take the logarithm.

The exponential equation can be written as $y_t = \beta_0 e^{\beta_1 x_t}$, where the birth rate has a constant rate of growth at rate β . If we take the logarithm of this, then our equation becomes:

$$\ln(y_t) = \ln(\beta_0) + \beta_1 x_t$$

Thus, $\ln(y_t)$ is a linear function of x_t . For our analysis, x_t is time, so the interpretation is that $\ln(y_t)$ is a linear function of time, and thus we can solve for it with linear least squares regression and then use the exponentiation of our result to find β .

To find the coefficients, we write the logarithm of our original exponential equation $y_t = \beta_0 e^{\beta_1 x_t}$ and put it in matrix form:

$$\{\ln(y_t)\} = \{\beta\}[X] + \{\epsilon\}$$

Thus we can solve the linearized equation of the exponential function in the same way as the linear function.

For our exponential model that was built for this data set, we again used *tslm*. However, we took the logarithm of the values before passing them into the function, and then exponentiated the output values in order to view the forecast and graph in our original units.

Autoregressive Integrated Moving Average

The autoregressive integrated moving average (ARIMA) model is a combination of three different time series components designed to create the best fit model for a time series data set. The autoregressive (AR) component can be simplified to a stochastic difference equation. It is a model in which the current value of the series is linearly related to its past values with an additive shock value. The moving average (MA) component takes another approach to time series modeling. Utilizing the fact that variation in time series data is driven by past shocks, the approach takes distributed lags of current and past shocks to model the current value of the series. Lastly, the integrated (I) aspect of the ARIMA model takes into consideration a very important concept in time series modeling, stationarity. A stationary time series is a data set whose properties do not depend on time. Therefore, a dataset that exhibits trend or seasonality is not stationary, however, a white noise series is stationary – random and independent of when it is observed.

To stabilize the variance dependent on time, the integrated component, I(d) calculates the differences between consecutive values in the time series, a process known as differencing.

$$y'_t = y_t - y_{t-1}$$

If the data set does not appear stationary after differencing, it may be necessary to difference the data a second time, or second-order differencing.

$$\begin{aligned} y''_t &= y'_t - y'_{t-1} \\ &= (y_t - y_{t-1}) - (y_{t-1} - y_{t-2}) \\ &= y_t - 2y_{t-1} + y_{t-2} \end{aligned}$$

In lag operator form or backshift notation, differencing of the the time series data set x_t is represented by:

$$y_t = (1 - L)^d x_t$$

The AR(p) model is represented by:

$$y_t = c + \phi_1 y_{t-1} + \phi_2 y_{t-2} + \dots + \phi_p y_{t-p} + e_t = c + \sum_{i=1}^p \phi_i y_{t-i} + e_t$$

where c is a constant, ϕ_1, \dots, ϕ_p are the parameters of the model, and e_t is white noise.

In lag operator form or backshift notation:

$$\Phi(L)y_t = (1 - \phi_1 L - \phi_2 L^2 - \dots - \phi_p L^p)y_t = e_t$$

The MA(q) model is represented by:

$$y_t = c + e_t + \theta_1 e_{t-1} + \theta_2 e_{t-2} + \dots + \theta_q e_{t-q} = c + \sum_{i=1}^q \theta_i e_{t-i} + e_t$$

where c is a constant, $\theta_1, \dots, \theta_q$ are the parameters of the model, and e_t is white noise.

In lag operator form or backshift notation:

$$y_t = c + (1 + \theta_1 L + \theta_2 L^2 + \dots + \theta_q L^q) e_t = \Theta(L) e_t$$

Combining all three models results in the following autoregressive integrated moving average ARIMA(p,d,q) model and is represented by:

$$y_t = (1 - L)^d x_t \quad (I)$$

$$y_t = \phi_1 y_{t-1} + \phi_2 y_{t-2} + \dots + \phi_p y_{t-p} + e_t + \theta_1 e_{t-1} + \theta_2 e_{t-2} + \dots + \theta_q e_{t-q} \quad (AR \& MA)$$

where p is the order of the autoregressive model and q is the order of the moving average model.

In lag operator form or backshift notation:

$$y_t = (1 - L)^d x_t \quad (I)$$

$$\left(1 - \sum_{i=1}^p \phi_i L^i\right) y_t = \left(1 + \sum_{i=1}^q \theta_i L^i\right) e_t \quad (AR \& MA)$$

To determine the p , d , and q parameters of the ARIMA model, the Box-Jenkins approach utilizes the shape of the autocorrelation and the partial autocorrelation functions for model identification. The autocorrelation function measures autocorrelations within the dataset, or the relationship between y_t and y_{t-k} for various values of k . Autocorrelation between times t and $t + k$ can be defined as:

$$R(k) = \frac{E[(y_t - \mu_t)(y_{t+k} - \mu_{t+k})]}{\sigma_t \sigma_{t+k}}$$

where E is the expected value operator, μ is the mean and σ^2 is the variance. If the function $R(k)$ is well defined (no zero variance, existence of the mean, no infinite processes), then the value of $R(k)$ lies in the range between -1 and 1. The partial autocorrelation function also measures autocorrelations within the dataset, but takes into consideration the effects of other time lags and consequently removes them. The partial autocorrelation measures the direct correlation between times t and $t + k$ and removes any linear dependencies from $t + 1$ to $t + k - 1$ in which the autocorrelation function does take into effect.

For the AR(p) process, the autocorrelation function exhibits either an exponential decay to zero or an alternating positive and negative decay to zero shape. The partial autocorrelation function of the AR(p) process becomes zero at lag $p + 1$ and greater, allowing one to identify the p parameter. For the MA(q) process, the autocorrelation function exhibits one or more spikes and becomes zero at the lag $q + 1$ and greater, allowing one to identify the q parameter without the need to observe the partial autocorrelation function. To select the optimal combination of these p , d , and q parameters for the ARIMA model, the Akaike information criterion (AIC) is utilized. The AIC is calculated by the following:

$$AIC = -2 \log(L) + 2(p + q + k + 1)$$

where L is the likelihood of the data, p is the order of the autoregressive process, q is the order of the moving average process, and k is number of parameters in the model. The objective is to minimize the AIC values in order to obtain a good model fit. The AIC value, however, can only be used to compare ARIMA models with the same order of differencing.

The *auto.arima* function in R chooses the optimal p , d , and q parameters by running through the various combination of parameters and selecting the numbers that result in the lowest AIC value. This *auto.arima* function takes a time series dataset as an input and outputs a fit model of the dataset. This

function can also manually take on various constraints, such as maximum p and q values, or starting p and q values as well as binary arguments such as stationary, seasonal, or stepwise. Automatically, these values are set to be *False*, *True*, and *True* respectively. If stationary is set to be *True*, then the function restricts its search to stationary models. If seasonal is set to be *False*, then the function restricts its search once again to just non-seasonal models. If stepwise selection is set to be *False*, then the function runs a longer, exhaustive selection search that is generally slower, especially for seasonal models. All these parameters of the *auto.arima* function allow for greater control of creating a fit model according to what one is looking for.

All of the fit models generated from the *tslm* and *auto.arima* functions are used as inputs in the *forecast* function in R. This *forecast* function also takes in inputs such as forecast horizon, confidence interval levels and a robust parameter, which if set *True*, the function ignores missing values and outliers. The function has various outputs such as the forecasting method, point forecasts, high and low prediction intervals, residuals, and fitted values.

```

#=====

# MATH 320
#
# RESEARCH PROJECT - FORECASTING BIRTH RATES BY ETHNICITY
#
# FORECAST SCRIPT
#
#=====

# LOAD PACKAGES
library(sqldf)
library(forecast)
# Load necessary packages needed to run the script

# DATA PREP
#=====

# While the data was generally clean, there were a few missing values that needed to be taken out. The
#
# numeric values we needed were factor values and needed to be converted to numerical values. After the
#
# the data consisted of only relevant values, we organized the data into four subsets, one for each race
#
# in order to run our time series models. Lastly, the data needed to be placed in time series format in
#
# order to be able to run the time series models
#
#=====

data <- read.csv("~/Desktop/R Stuff/Math 320/NCHS_
_Births__Birth_Rates__and_Fertility_Rates__by_Race_of_Mother__United_States__1960-2013.csv",
  header = T)

# Read the data
data <- subset(data, Birth.Rate != "*" & Fertility.Rate != "*")
# Remove missing values
data_ai <- sqldf("select * from data where Race = 'American Indian/Alaska Native' order by Year")
data_ai$Birth.Rate <- as.numeric(levels(data_ai$Birth.Rate))[data_ai$Birth.Rate]
data_as <- sqldf("select * from data where Race = 'Asian/Pacific Islander' order by Year")
data_as$Birth.Rate <- as.numeric(levels(data_as$Birth.Rate))[data_as$Birth.Rate]
data_bl <- sqldf("select * from data where Race = 'Black' order by Year")
data_bl$Birth.Rate <- as.numeric(levels(data_bl$Birth.Rate))[data_bl$Birth.Rate]
data_wh <- sqldf("select * from data where Race = 'White' order by Year")
data_wh$Birth.Rate <- as.numeric(levels(data_wh$Birth.Rate))[data_wh$Birth.Rate]
# Subset the data while converting factor values to numeric
data_ai_ts <- ts(data_ai$Birth.Rate, start = 1980, frequency = 1)
data_as_ts <- ts(data_as$Birth.Rate, start = 1980, frequency = 1)
data_bl_ts <- ts(data_bl$Birth.Rate, start = 1964, frequency = 1)
data_wh_ts <- ts(data_wh$Birth.Rate, start = 1960, frequency = 1)
# Transform data into time series format
plot(data_ai$Year, data_ai$Birth.Rate, type = "l", xlim = c(1960, 2013), ylim = c(5, 33), lwd = 2, main
= "Birth Rates by Race",
  xlab = "Birth Rate", ylab = "Year")
lines(data_as$Year, data_as$Birth.Rate, col = "red", lwd = 2)
lines(data_bl$Year, data_bl$Birth.Rate, col = "blue", lwd = 2)
lines(data_wh$Year, data_wh$Birth.Rate, col = "green", lwd = 2)
legend(x="topright", legend=c("American Indian/Alaska Native", "Asian/Pacific Islander", "Black",
"White"),
  col=c("black", "red", "blue", "green"), lty = c(1,1,1,1), cex = 0.7)
# Plot the dataset

# MODEL
#=====

# For each of the subsets of data, we run a linear, exponential, and arima model and calculates forecast
#
# values. Then, we plot the data with its respective forecast values. We use the function tslm for
#
# linear and exponential models, the auto.arima function for the arima model, and the forecast function
#

```



```

# to calculate forecast values
#
#=====##

## Linear Model
lm_fit_ai <- tslm(data_ai_ts ~ trend)
lm_fc_ai <- forecast(lm_fit_ai, h = 10, level = 95)
lm_fit_as <- tslm(data_as_ts ~ trend)
lm_fc_as <- forecast(lm_fit_as, h = 10, level = 95)
lm_fit_bl <- tslm(data_bl_ts ~ trend)
lm_fc_bl <- forecast(lm_fit_bl, h = 10, level = 95)
lm_fit_wh <- tslm(data_wh_ts ~ trend)
lm_fc_wh <- forecast(lm_fit_wh, h = 10, level = 95)
# Create linear models for each subset of data and calculate forecast values
lm_fc_ai
# 10 Year Horizon Forecast Values for Race: American Indian/Alaska Native
lm_fc_as
# 10 Year Horizon Forecast Values for Race: Asian/Pacific Islander
lm_fc_bl
# 10 Year Horizon Forecast Values for Race: Black
lm_fc_wh
# 10 Year Horizon Forecast Values for Race: White
plot(data_ai$Year, data_ai$Birth.Rate, type = "l", xlim = c(1960, 2023), ylim = c(5, 33), lwd = 2,
     main = "Linear Model: Birth Rates by Race", xlab = "Birth Rate", ylab = "Year")
lines(data_ai$Year, data_ai$Birth.Rate, col = "red", lwd = 2)
lines(data_bl$Year, data_bl$Birth.Rate, col = "blue", lwd = 2)
lines(data_wh$Year, data_wh$Birth.Rate, col = "green", lwd = 2)
lines(2014:2023, lm_fc_ai$mean, col = "black", lwd = 1, lty = 2)
lines(2014:2023, lm_fc_as$mean, col = "red", lwd = 1, lty = 2)
lines(2014:2023, lm_fc_bl$mean, col = "blue", lwd = 1, lty = 2)
lines(2014:2023, lm_fc_wh$mean, col = "green", lwd = 1, lty = 2)
legend(x="topright", legend=c("American Indian/Alaska Native", "Asian/Pacific Islander", "Black",
"White"),
     col=c("black", "red", "blue", "green"), lty = c(1,1,1,1), cex = 0.7)
# Plot the dataset with its linear forecast values

## Exponential Model
ln_data_ai_ts <- ts(log(data_ai$Birth.Rate), start = 1980, frequency = 1)
ln_data_as_ts <- ts(log(data_as$Birth.Rate), start = 1980, frequency = 1)
ln_data_bl_ts <- ts(log(data_bl$Birth.Rate), start = 1964, frequency = 1)
ln_data_wh_ts <- ts(log(data_wh$Birth.Rate), start = 1960, frequency = 1)
# Log the time series dataset
ln_fit_ai <- tslm(ln_data_ai_ts ~ trend)
ln_fc_ai <- forecast(ln_fit_ai, h = 10, level = 95)
exp_ai <- exp(as.vector(ln_fc_ai$mean))
exp_ai1 <- as.data.frame(exp_ai)
colnames(exp_ai1) <- "Point Forecast"
exp_ai1$Lo.95 <- exp(as.vector(ln_fc_ai$lower))
exp_ai1$Hi.95 <- exp(as.vector(ln_fc_ai$upper))
rownames(exp_ai1) <- c("2014", "2015", "2016", "2017", "2018", "2019", "2020", "2021", "2022", "2023")
exp_ai_fit <- exp(as.vector(ln_fc_ai$fitted))
ln_fit_as <- tslm(ln_data_as_ts ~ trend)
ln_fc_as <- forecast(ln_fit_as, h = 10, level = 95)
exp_as <- exp(as.vector(ln_fc_as$mean))
exp_as1 <- as.data.frame(exp_as)
colnames(exp_as1) <- "Point Forecast"
exp_as1$Lo.95 <- exp(as.vector(ln_fc_as$lower))
exp_as1$Hi.95 <- exp(as.vector(ln_fc_as$upper))
rownames(exp_as1) <- c("2014", "2015", "2016", "2017", "2018", "2019", "2020", "2021", "2022", "2023")
exp_as_fit <- exp(as.vector(ln_fc_as$fitted))
ln_fit_bl <- tslm(ln_data_bl_ts ~ trend)
ln_fc_bl <- forecast(ln_fit_bl, h = 10, level = 95)
exp_bl <- exp(as.vector(ln_fc_bl$mean))
exp_bl1 <- as.data.frame(exp_bl)
colnames(exp_bl1) <- "Point Forecast"
exp_bl1$Lo.95 <- exp(as.vector(ln_fc_bl$lower))
exp_bl1$Hi.95 <- exp(as.vector(ln_fc_bl$upper))
rownames(exp_bl1) <- c("2014", "2015", "2016", "2017", "2018", "2019", "2020", "2021", "2022", "2023")
exp_bl_fit <- exp(as.vector(ln_fc_bl$fitted))
ln_fit_wh <- tslm(ln_data_wh_ts ~ trend)
ln_fc_wh <- forecast(ln_fit_wh, h = 10, level = 95)
exp_wh <- exp(as.vector(ln_fc_wh$mean))

```

```

exp_wh1 <- as.data.frame(exp_wh)
colnames(exp_wh1) <- "Point Forecast"
exp_wh1$Lo.95 <- exp(as.vector(ln_fc_wh$lower))
exp_wh1$Hi.95 <- exp(as.vector(ln_fc_wh$upper))
rownames(exp_wh1) <- c("2014", "2015", "2016", "2017", "2018", "2019", "2020", "2021", "2022", "2023")
exp_wh_fit <- exp(as.vector(ln_fc_wh$fitted))
# Create exponential models for each subset of data and calculate forecast values
exp_ai1
# 10 Year Horizon Forecast Values for Race: American Indian/Alaska Native
exp_as1
# 10 Year Horizon Forecast Values for Race: Asian/Pacific Islander
exp_bl1
# 10 Year Horizon Forecast Values for Race: Black
exp_wh1
# 10 Year Horizon Forecast Values for Race: White
plot(data_ai$Year, data_ai$Birth.Rate, type = "l", xlim = c(1960, 2023), ylim = c(5, 33), lwd = 2,
     main = "Exponential Model: Birth Rates by Race", xlab = "Birth Rate", ylab = "Year")
lines(data_as$Year, data_as$Birth.Rate, col = "red", lwd = 2)
lines(data_bl$Year, data_bl$Birth.Rate, col = "blue", lwd = 2)
lines(data_wh$Year, data_wh$Birth.Rate, col = "green", lwd = 2)
lines(2014:2023, exp_ai, col = "black", lwd = 1, lty = 2)
lines(2014:2023, exp_as, col = "red", lwd = 1, lty = 2)
lines(2014:2023, exp_bl, col = "blue", lwd = 1, lty = 2)
lines(2014:2023, exp_wh, col = "green", lwd = 1, lty = 2)
legend(x="topright", legend=c("American Indian/Alaska Native", "Asian/Pacific Islander", "Black",
"White"),
     col=c("black", "red", "blue", "green"), lty = c(1,1,1,1), cex = 0.7)
# Plot the dataset with its exponential forecast values

## Autoregressive Integrated Moving Average Model
fit_ai <- auto.arima(data_ai_ts)
fc_ai <- forecast(fit_ai, h = 10, level = 95)
fit_as <- auto.arima(data_as_ts)
fc_as <- forecast(fit_as, h = 10, level = 95)
fit_bl <- auto.arima(data_bl_ts)
fc_bl <- forecast(fit_bl, h = 10, level = 95)
fit_wh <- auto.arima(data_wh_ts)
fc_wh <- forecast(fit_wh, h = 10, level = 95)
# Create arima models for each subset of data and calculate forecast values
fc_ai
# 10 Year Horizon Forecast Values for Race: American Indian/Alaska Native
fc_as
# 10 Year Horizon Forecast Values for Race: Asian/Pacific Islander
fc_bl
# 10 Year Horizon Forecast Values for Race: Black
fc_wh
# 10 Year Horizon Forecast Values for Race: White
plot(data_ai$Year, data_ai$Birth.Rate, type = "l", xlim = c(1960, 2023), ylim = c(5, 33),
     lwd = 2, main = "ARIMA Model: Birth Rates by Race", xlab = "Birth Rate", ylab = "Year")
lines(data_as$Year, data_as$Birth.Rate, col = "red", lwd = 2)
lines(data_bl$Year, data_bl$Birth.Rate, col = "blue", lwd = 2)
lines(data_wh$Year, data_wh$Birth.Rate, col = "green", lwd = 2)
lines(2014:2023, fc_ai$mean, col = "black", lwd = 1, lty = 2)
lines(2014:2023, fc_as$mean, col = "red", lwd = 1, lty = 2)
lines(2014:2023, fc_bl$mean, col = "blue", lwd = 1, lty = 2)
lines(2014:2023, fc_wh$mean, col = "green", lwd = 1, lty = 2)
legend(x="topright", legend=c("American Indian/Alaska Native", "Asian/Pacific Islander", "Black",
"White"),
     col=c("black", "red", "blue", "green"), lty = c(1,1,1,1), cex = 0.7)
# Plot the dataset with its arima forecast values

# FORECAST COMPARISON/MODEL SELECTION
#=====#

# After we run our different models, we also grab the fitted values and plot them to see how well they
#
# matched the actual dataset. We calculate the Root Mean Square Error (RMSE) values and select the model
#
# with the lowest RMSE value.
#
#=====#

```

```

plot(data_ai$Year, data_ai$Birth.Rate, type = "l", xlim = c(1980, 2023), ylim = c(5, 23), lwd = 2,
     main = "Fitted and Forecast Values for Race: American Indian/Alaska Native", xlab = "Year", ylab =
"Birth Rate")
lines(1980:2013, fc_ai$fitted, col = "green", lty = 2)
lines(2014:2023, fc_ai$mean, col = "green", lwd = 1, lty = 1)
lines(1980:2013, lm_fc_ai$fitted, col = "blue", lwd = 1, lty = 2)
lines(2014:2023, lm_fc_ai$mean, col = "blue", lwd = 1, lty = 1)
lines(1980:2013, exp_ai_fit, col = "red", lwd = 1, lty = 2)
lines(2014:2023, exp_ai, col = "red", lwd = 1, lty = 1)
legend(x="topright", legend=c("Actual", "Fitted (by color)", "FC Linear", "FC Exp", "FC ARIMA"),
     col=c("black", "black", "blue", "red", "green"), lty = c(1,2,1,1,1), cex = 0.7)
# Plot actual, fitted, and forecast values from the three different models for American Indian/Alaska
Natives race
plot(data_as$Year, data_as$Birth.Rate, type = "l", xlim = c(1980, 2023), ylim = c(11, 23), lwd = 2,
     main = "Fitted and Forecast Values for Race: Asian/Pacific Islander", xlab = "Year", ylab = "Birth
Rate")
lines(1980:2013, fc_as$fitted, col = "green", lty = 2)
lines(2014:2023, fc_as$mean, col = "green", lwd = 1, lty = 1)
lines(1980:2013, lm_fc_as$fitted, col = "blue", lwd = 1, lty = 2)
lines(2014:2023, lm_fc_as$mean, col = "blue", lwd = 1, lty = 1)
lines(1980:2013, exp_as_fit, col = "red", lwd = 1, lty = 2)
lines(2014:2023, exp_as, col = "red", lwd = 1, lty = 1)
legend(x="topright", legend=c("Actual", "Fitted (by color)", "FC Linear", "FC Exp", "FC ARIMA"),
     col=c("black", "black", "blue", "red", "green"), lty = c(1,2,1,1,1), cex = 0.7)
# Plot actual, fitted, and forecast values from the three different models for Asian/Pacific Islanders
race
plot(data_bl$Year, data_bl$Birth.Rate, type = "l", xlim = c(1964, 2023), ylim = c(9, 33), lwd = 2,
     main = "Fitted and Forecast Values for Race: Black", xlab = "Year", ylab = "Birth Rate")
lines(1964:2013, fc_bl$fitted, col = "green", lty = 2)
lines(2014:2023, fc_bl$mean, col = "green", lwd = 1, lty = 1)
lines(1964:2013, lm_fc_bl$fitted, col = "blue", lwd = 1, lty = 2)
lines(2014:2023, lm_fc_bl$mean, col = "blue", lwd = 1, lty = 1)
lines(1964:2013, exp_bl_fit, col = "red", lwd = 1, lty = 2)
lines(2014:2023, exp_bl, col = "red", lwd = 1, lty = 1)
legend(x="topright", legend=c("Actual", "Fitted (by color)", "FC Linear", "FC Exp", "FC ARIMA"),
     col=c("black", "black", "blue", "red", "green"), lty = c(1,2,1,1,1), cex = 0.7)
# Plot actual, fitted, and forecast values from the three different models for Black Americans race
plot(data_wh$Year, data_wh$Birth.Rate, type = "l", xlim = c(1960, 2023), ylim = c(9, 25), lwd = 2,
     main = "Fitted and Forecast Values for Race: Black", xlab = "Year", ylab = "Birth Rate")
lines(1960:2013, fc_wh$fitted, col = "green", lty = 2)
lines(2014:2023, fc_wh$mean, col = "green", lwd = 1, lty = 1)
lines(1960:2013, lm_fc_wh$fitted, col = "blue", lwd = 1, lty = 2)
lines(2014:2023, lm_fc_wh$mean, col = "blue", lwd = 1, lty = 1)
lines(1960:2013, exp_wh_fit, col = "red", lwd = 1, lty = 2)
lines(2014:2023, exp_wh, col = "red", lwd = 1, lty = 1)
legend(x="topright", legend=c("Actual", "Fitted (by color)", "FC Linear", "FC Exp", "FC ARIMA"),
     col=c("black", "black", "blue", "red", "green"), lty = c(1,2,1,1,1), cex = 0.7)
# Plot actual, fitted, and forecast values from the three different models for White Americans race

```