

# Part II III IV - Maryland Poverty Level

Lebo Sango

2024-04-24

## Data Cleaning

```
maryland_raw = read.csv("/Users/lebsan/Documents/STAT 5084 - Time Series/County Level Project/Universe_1.csv")

maryland =maryland_raw %>%
  mutate( SAIPE = as.numeric(SAIPE),SNAP = as.numeric(SNAP),
          IRS_exempt_State = as.numeric(IRS_exempt_State),
          Poverty_Universe = as.numeric(Poverty_Universe)) %>%
  select(c(year,County, SAIPE, SNAP, IRS_exempt_State, Poverty_Universe)) %>%
  as_tsibble(index =year, key = County )

maryland %>% head(5)
```

```
## # A tsibble: 5 x 6 [1Y]
## # Key:      County [1]
##   year County      SAIPE  SNAP  IRS_exempt_State Poverty_Universe
##   <int> <chr>      <dbl> <dbl>      <dbl>      <dbl>
## 1  1998 Allegany County 10473  6650      472945      69532
## 2  1999 Allegany County  9270  6294      468976      69404
## 3  2000 Allegany County  9445  5922      465555      68408
## 4  2001 Allegany County  8954  6365      475208      68151
## 5  2002 Allegany County  9418  6864      487317      67632
```

## Linear models

```
## # A tibble: 7 x 4
##   Linear_models      AIC    CV    BIC
##   <chr>      <dbl> <dbl> <dbl>
## 1 SAIPE_c_IRS      -2547. 0.277 -2462.
## 2 SAIPE_c_PovUniverse -2559. 0.286 -2474.
## 3 SAIPE_c_IRS_PovUniverse -2665. 0.236 -2552.
## 4 SAIPE_c_SNAP_IRS    -2808. 0.184 -2695.
## 5 SAIPE_c_SNAP_IRS_PovUniverse -2861. 0.173 -2719.
## 6 SAIPE_c_SNAP      -2823. 0.174 -2738.
## 7 SAIPE_c_SNAP_PovUniverse -2878. 0.163 -2764.
```

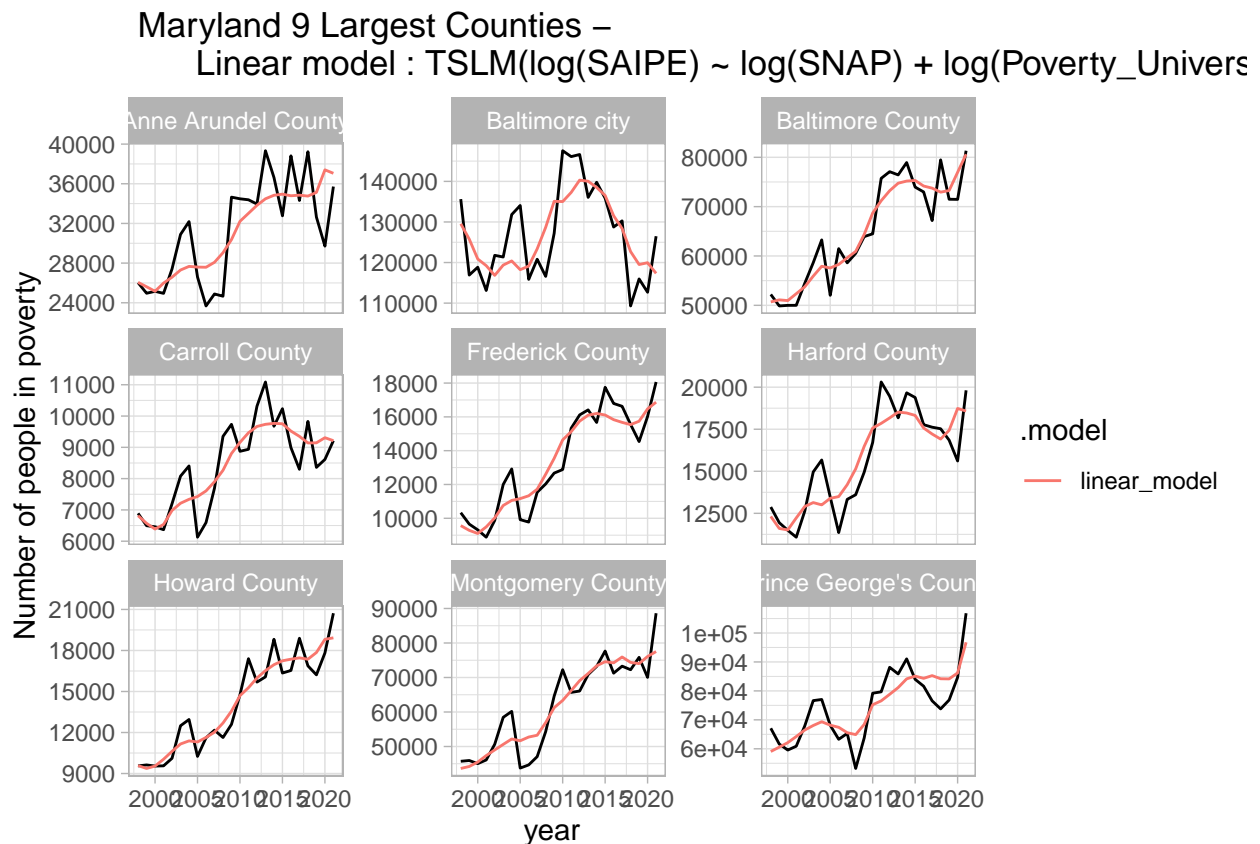
Lower cross validation and BIC model 6 (SAIPE = B0 + B1SNAP + B2PovUniverse) is the best model. SNAP and Poverty Universe displayed a strong correlation coefficient with SAIPE. In addition, model 1 including all the dependent variable had very close precision crieteria compared to model6. I decided to

proceed with model 6 because two of the precision criteria were the smallest. This distinction can be caused by the low correlation coefficient of 0.118 between SAIPE and IRS\_exempt\_State which hinder the model.

*# Plot of the fitted predictions of the nine biggest counties with the best linear model.*  
maryland %>%

```
filter(County %in% c("Montgomery County", "Prince George's County",
  "Baltimore County", "Anne Arundel County",
  "Baltimore city", "Howard County",
  "Frederick County", "Harford County",
  "Carroll County" ) ) %>%

model(linear_model = TSLM(log(SAIPE) ~ log(SNAP) + log(Poverty_Universe))) %>%
augment() %>% ggplot(aes(x=year))+
geom_line(aes(y=SAIPE)) +
geom_line(aes(y=.fitted, color=.model)) +
facet_wrap(~County, scales = "free_y")+
labs(title = " Maryland 9 Largest Counties -
  Linear model : TSLM(log(SAIPE) ~ log(SNAP) + log(Poverty_Universe))",
  y=" Number of people in poverty")+
theme_light()
```

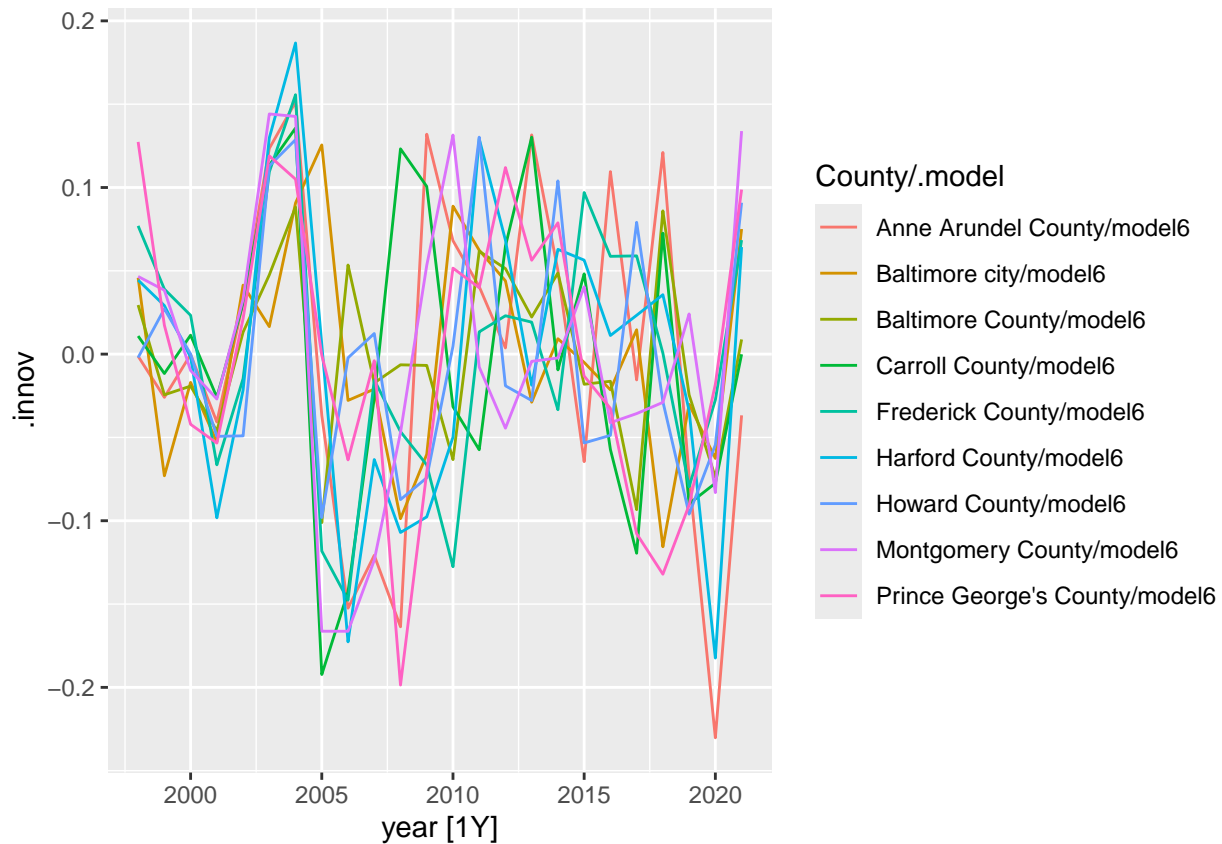


*# Residual plot of the nine largest counties*  
MD\_resid = maryland %>% filter(County %in% c("Montgomery County", "Prince George's County",  
"Baltimore County", "Anne Arundel County",  
"Baltimore city", "Howard County",  
"Frederick County", "Harford County",  
"Carroll County" ) ) %>%

```

model(model6 = TSLM(log(SAIBE) ~ log(SNAP) + log(Poverty_Universe))) %>%
augment()
MD_resid %>% autoplot(.innov)

```



*# LjungBox test on every county of Maryland state*

```

MD_resid2 = maryland %>%
  model(model6 = TSLM(log(SAIBE) ~ log(SNAP) + log(Poverty_Universe))) %>%
  augment()

```

```

MD_resid2 %>% select(County, .model, .innov) %>% group_by(County) %>%
  features(.innov, lbjung_box) %>% filter(lb_pvalue <= 0.05)

```

```

## # A tibble: 4 x 4
##   County      .model lb_stat lb_pvalue
##   <chr>      <chr>    <dbl>    <dbl>
## 1 Cecil County model6      8.87  0.00291
## 2 Dorchester County model6      5.53  0.0187
## 3 Prince George's County model6      5.32  0.0211
## 4 Talbot County model6      4.59  0.0322

```

The only counties that do have white noise are Prince George's county, Talbot County, Dorchester County and Cecil county while the rest does not exhibits autocorrelation. Overall the model does better at capturing the trend but fails to capture cyclicalities. Furthermore, I expect to employ more sophisticated models that can capture the cyclicalities and fluctuations of SAIBE.

## Part 3 - Stochastic Models

### Single County Forecasts

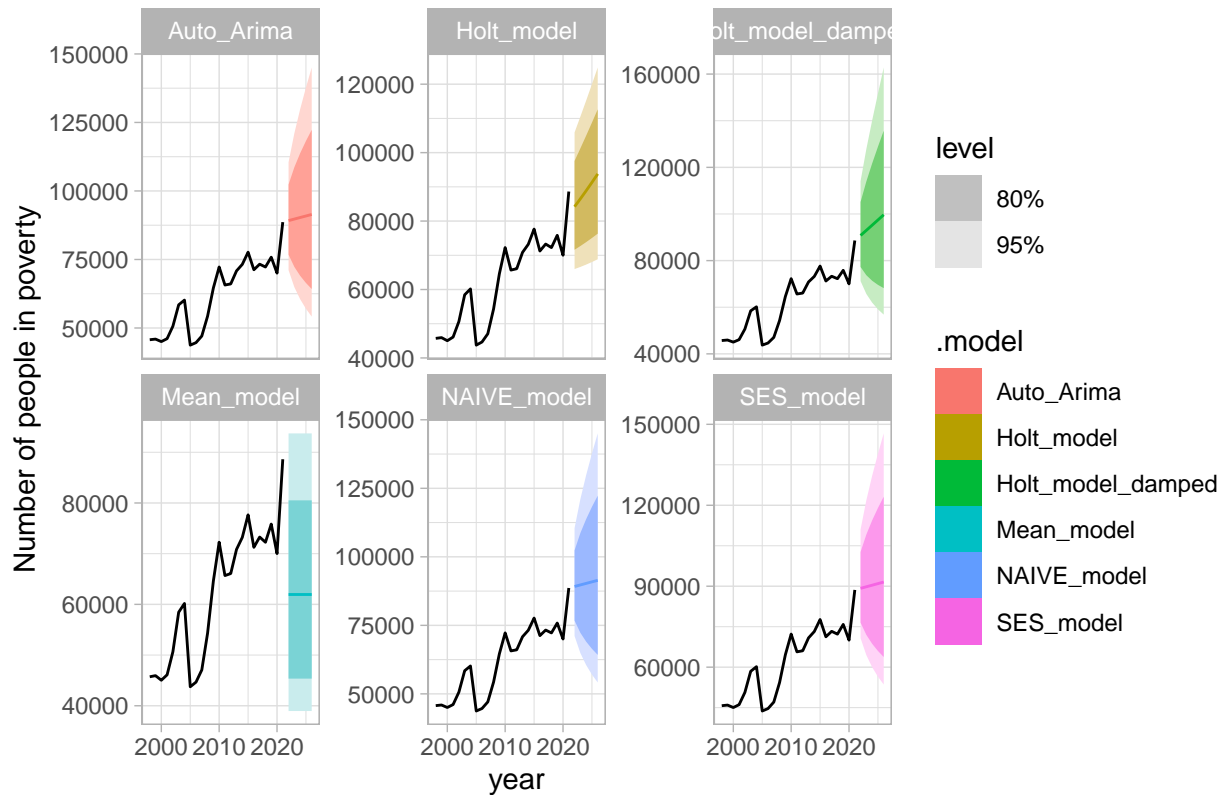
```
stochastic_model = maryland %>% filter(County %in% "Montgomery County") %>%  
  model(  
    NAIVE_model = NAIVE(log(SAIPe)),  
    Mean_model = MEAN(log(SAIPe)),  
    SES_model = ETS(log(SAIPe) ~ error("A")+trend("N")+  
      season("N")),  
    Holt_model = ETS(log(SAIPe) ~ error("A")+trend("A")+  
      season("N")),  
    Holt_model_damped = ETS(log(SAIPe) ~ error("A")+trend("Ad")+  
      season("N")),  
    Auto_Arima = ARIMA(log(SAIPe)))  
  
stochastic_model
```

```
## # A mable: 1 x 7  
## # Key:      County [1]  
##   County      NAIVE_model Mean_model    SES_model    Holt_model Holt_model_damped  
##   <chr>         <model>      <model>      <model>      <model>      <model>  
## 1 Montgomery~  <NAIVE>      <MEAN>    <ETS(A,N,N)> <ETS(A,A,N)> <ETS(A,Ad,N)>  
## # i 1 more variable: Auto_Arima <model>
```

### Plotting the number in poverty data along with a five-year forecast

```
stochastic_model %>% forecast(h="5 years") %>% autoplot(maryland)+  
  facet_wrap(~.model, scales = "free_y")+  
  theme_light()+  
  labs(title = " Montgomery County - Forecast of Number of inhabitants in poverty",  
    y = "Number of people in poverty")
```

## Montgomery County – Forecast of Number of inhabitants in poverty



The best model for this county is the Auto Arima( evaluated at difference) with a low root mean square error and mean average percentage error.

```
# The auto arima is the model that exhibits the smallest RMSE accross Maryland counties.
stochastic_model %>% accuracy() %>%
  group_by( stochastic_models = .model, Maryland_County = County) %>%
  summarise(RMSE = sum(RMSE), MAPE = sum(MAPE)) %>%
  arrange(min(RMSE))
```

```
## 'summarise()' has grouped output by 'stochastic_models'. You can override using
## the '.groups' argument.
```

```
## # A tibble: 6 x 4
## # Groups:   stochastic_models [6]
##   stochastic_models Maryland_County RMSE MAPE
##   <chr>           <chr>      <dbl> <dbl>
## 1 Auto_Arima      Montgomery County 6804.  7.88
## 2 Holt_model      Montgomery County 6561.  8.48
## 3 Holt_model_damped Montgomery County 6640.  7.42
## 4 Mean_model      Montgomery County 12979. 19.7
## 5 NAIVE_model     Montgomery County 6950.  8.18
## 6 SES_model       Montgomery County 6804.  7.84
```

## Exponential Smoothing Models

```
ES_maryland = maryland %>% model(

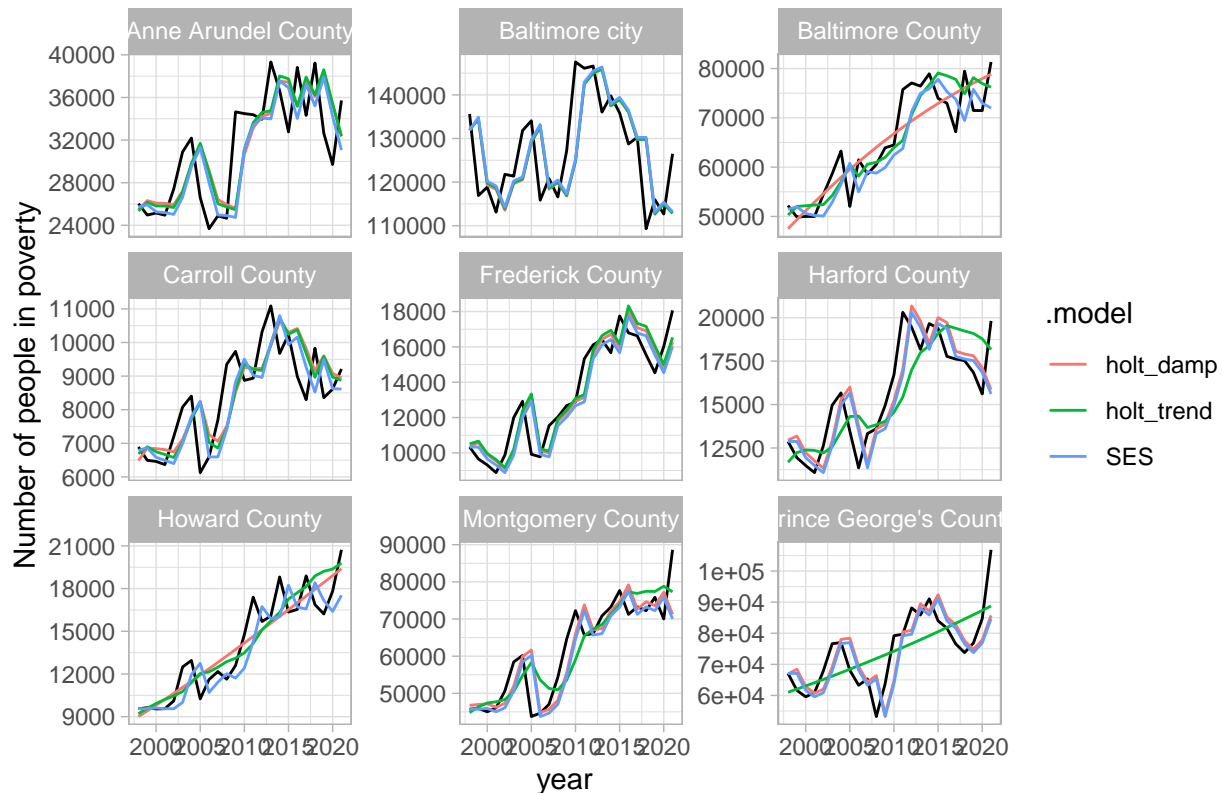
  SES = ETS(log(SAYPE)~error("A") + trend("N") +season("N")),

  holt_trend = ETS(log(SAYPE)~error("A")+trend("A")+season("N")),

  holt_damp = ETS(log(SAYPE)~error("A")+trend("Ad")+season("N")))

# Fitting the ETS models into every country of the state of Maryland
# I proceeded to graph the top 9 largest counties by population as graphing
# 24 counties would not properly fit.
ES_maryland %>%
  filter(County %in% c("Montgomery County", "Prince George's County"
    , "Baltimore County", "Anne Arundel County",
    "Baltimore city", "Howard County",
    "Frederick County", "Harford County",
    "Carroll County")) %>%
  augment() %>% ggplot(aes(x=year))+
  geom_line(aes(y=SAIPE)) +
  geom_line(aes(y=.fitted, color=.model)) +
  facet_wrap(~County, scales = "free_y") +
  labs(title = " ETS Models - All Maryland Counties",
    y = " Number of people in poverty")+
  theme_light()
```

## ETS Models – All Maryland Counties



### Best Performing ETS Model

I selected the Holt\_damped model because per the results below, it shows the smallest RSME and MAPE values against the other models. One particularity of the holt damped model is its additive trend feature which implies that long run forecast as  $h$  approaches infinity, the damping parameter will be constant, while in the short- forecast will be trended.

```
ES_maryland %>% accuracy() %>%
  group_by(Exponential_smoothing_model = .model) %>%
  summarise(RMSE = sum(RMSE),
            MAPE = sum(MAPE)) %>%
  arrange(min(RMSE))
```

```
## # A tibble: 3 x 3
##   Exponential_smoothing_model  RMSE  MAPE
##   <chr>                        <dbl> <dbl>
## 1 SES                        48970.  175.
## 2 holt_damp                 46967.  174.
## 3 holt_trend                 48702.  180.
```

## ARIMA Models

The most commonly selected model is the ARIMA(0,1,0) evaluated at difference, the ARIMA(1,0,0) with a mean constant, and ARIMA(0,1,1) with a drift.

```
# The most selected ARIMA model is the model evaluated at difference
maryland %>% model(ARIMA(log(SAYPE))) %>% print()
```

```
## # A mable: 24 x 2
## # Key:      County [24]
##   County      'ARIMA(log(SAYPE))'
##   <chr>        <model>
## 1 Allegany County    <ARIMA(1,0,0) w/ mean>
## 2 Anne Arundel County <ARIMA(0,1,0)>
## 3 Baltimore County  <ARIMA(0,1,1) w/ drift>
## 4 Baltimore city    <ARIMA(1,0,0) w/ mean>
## 5 Calvert County    <ARIMA(0,1,0)>
## 6 Caroline County    <ARIMA(0,1,0)>
## 7 Carroll County    <ARIMA(0,1,0)>
## 8 Cecil County      <ARIMA(1,1,0)>
## 9 Charles County    <ARIMA(0,1,0)>
## 10 Dorchester County <ARIMA(0,1,0)>
## # i 14 more rows
```

## Best Performing Arima model

The best ARIMA model is the model with the drift with the smallest RMSE and MAPE. Moreover, the model forecast follows a straight line. That is, the forecast indicates the number of people in poverty is increasing as the trend is sloping upward, so the constant is non-zero and d is 1. As suggested by the RMSE, the best model is the ARIMA with drift

```
maryland %>% model(Difference = ARIMA(log(SAYPE) ~ pdq(0,1,0)),
                  ARIMA_Drift = ARIMA(log(SAYPE) ~ 1 + pdq(0,1,1)),
                  ARIMA_mean = ARIMA(log(SAYPE) ~ 1 + pdq(1,0,0))) %>%
  accuracy() %>%
  group_by(Arima_models = .model) %>%
  summarise(RMSE = sum(RMSE),
            MAPE = sum(MAPE)) %>%
  arrange(min(RMSE))
```

```
## # A tibble: 3 x 3
##   Arima_models  RMSE  MAPE
##   <chr>        <dbl> <dbl>
## 1 ARIMA_Drift  47456.  174.
## 2 ARIMA_mean   48611.  178.
## 3 Difference   50380.  182.
```

```
# Fitting the ARIMA models to every county
```

```
maryland %>% filter(County %in% c("Montgomery County", "Prince George's County",
                                "Baltimore County", "Anne Arundel County",
                                "Baltimore city", "Howard County",
                                "Frederick County", "Harford County",
                                "Carroll County")) %>%
```



```

model(Difference_0_1_0 = ARIMA(log(SAIPe) ~ pdq(0,1,0)),

      ARIMA_Drift_0_1_1 = ARIMA(log(SAIPe) ~ 1 + pdq(0,1,1)),

      ARIMA_mean_1_0_0 = ARIMA(log(SAIPe) ~ 1 + pdq(1,0,0))) %>%

augment() %>% ggplot(aes(x=year))+

geom_line(aes(y=SAIPe)) +

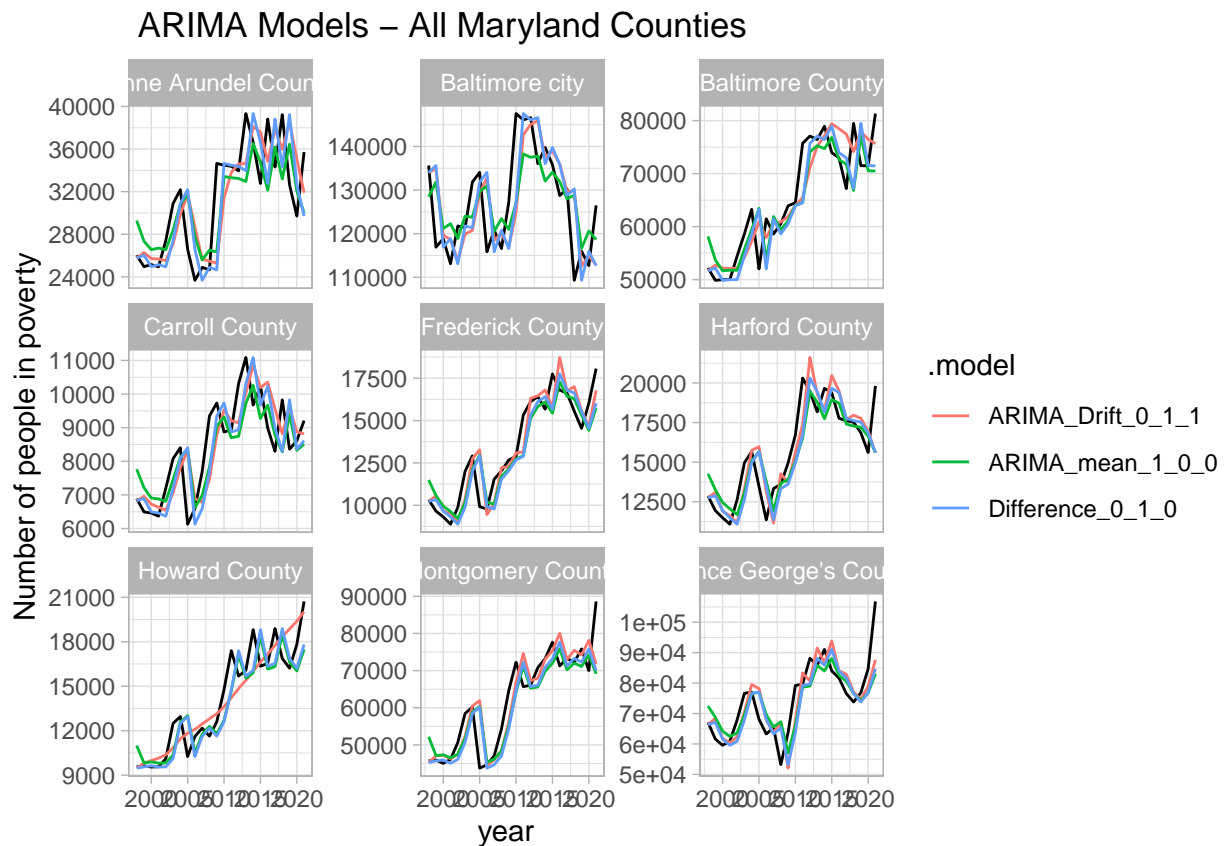
geom_line(aes(y=.fitted, color=.model)) +

facet_wrap(~County, scales = "free_y") +

labs(title = " ARIMA Models - All Maryland Counties",
      y = " Number of people in poverty")+

theme_light()

```



## Cross validation

The best model is the ARIMA model evaluated at difference or ARIMA(0,1,0) depicting a random walk in which the changes in the level of poverty that oscillate up and down with an unpredictable patterns. The

results also demonstrate that SAIPE is non-stationary and that differencing is the suitable approach for SAIPE to be stationary.

## ETS

```
# Building Training Sets

ES_maryland_stretch = maryland %>% stretch_tsibble(.init = 20)

ES_maryland_training = ES_maryland_stretch %>% model(

  SES = ETS(log(SAIPE)~error("A") + trend("N") +season("N")),

  holt_trend = ETS(log(SAIPE)~error("A")+trend("A")+season("N")),

  holt_damp = ETS(log(SAIPE)~error("A")+trend("Ad")+season("N")))
```

```
# Accuracy Check
ES_maryland_training %>% accuracy() %>%
  group_by(ETS_Models = .model) %>%
  summarise( RMSE = sum(RMSE)) %>%
  arrange(RMSE)
```

```
## # A tibble: 3 x 2
##   ETS_Models    RMSE
##   <chr>        <dbl>
## 1 holt_damp    224481.
## 2 SES         232151.
## 3 holt_trend  233105.
```

```
# Cross validation between training and test data
ES_maryland_training %>%
forecast(h="5 years") %>%
  accuracy(maryland) %>%
  group_by(ETS_Models = .model) %>%
  summarise(RMSE = sum(RMSE)) %>%
  arrange(RMSE)
```

```
## # A tibble: 3 x 2
##   ETS_Models    RMSE
##   <chr>        <dbl>
## 1 holt_damp    70688.
## 2 SES         70717.
## 3 holt_trend  82792.
```

## ARIMA

```
ARIMA_maryland_strech = maryland %>% stretch_tsibble(.init = 20)

ARIMA_maryland_training = ARIMA_maryland_strech %>% model(Difference = ARIMA(log(SAIPe) ~ pdq(0,1,0)),
  ARIMA_Drift = ARIMA(log(SAIPe) ~ 1 + pdq(0,1,1)),
  ARIMA_mean = ARIMA(log(SAIPe) ~ 1 + pdq(1,0,0)))
```

*# Accuracy Check*

```
ARIMA_maryland_training %>% accuracy() %>%
  group_by(ARIMA_Models = .model) %>%
  summarise(RMSE = sum(RMSE)) %>%
  arrange(RMSE)
```

```
## # A tibble: 3 x 2
##   ARIMA_Models    RMSE
##   <chr>          <dbl>
## 1 ARIMA_Drift    225398.
## 2 ARIMA_mean     228739.
## 3 Difference     238381.
```

*# Cross validation between training and test data*

```
ARIMA_maryland_training %>%
  forecast(h="5 years") %>%
  accuracy(maryland) %>%
  group_by(ARIMA_Models = .model) %>%
  summarise(RMSE = sum(RMSE)) %>%
  arrange(RMSE)
```

```
## # A tibble: 3 x 2
##   ARIMA_Models    RMSE
##   <chr>          <dbl>
## 1 Difference      73027.
## 2 ARIMA_Drift     74115.
## 3 ARIMA_mean      77124.
```

## Forecasts

The 5 counties with the largest increase in poverty level in the next 5 years are Somerset, Baltimore city, Allegany, Dorechester, and Washington counties.

```
# Forecasting poverty
poverty_forecast = maryland %>% rename(county = County) %>%
  model(Arima_Diff = ARIMA(log(SAIPe) ~ pdq(0,1,0))) %>%
  forecast(h="5 years")

# Extracting current population which is in 2021 for every county
current_population = maryland %>% filter(year == 2021) %>%
  select(county =County, Poverty_Universe)
```

```
# Join current populatin in 2021 with forecast data
merge(current_population,poverty_forecast, by ="county" ) %>%
  mutate(Poverty_Percent_change = .mean/Poverty_Universe * 100) %>%
  group_by(county) %>% summarise(Poverty_level_Prct = max(Poverty_Percent_change)) %>%
  arrange(desc(Poverty_level_Prct))
```

```
## # A tibble: 24 x 2
##   county          Poverty_level_Prct
##   <chr>              <dbl>
## 1 Somerset County      24.2
## 2 Baltimore city       23.2
## 3 Allegany County      16.6
## 4 Dorchester County    15.1
## 5 Washington County    15.1
## 6 Wicomico County      14.5
## 7 Caroline County      13.7
## 8 Kent County          12.2
## 9 Prince George's County 11.8
## 10 Cecil County        11.5
## # i 14 more rows
```

```
# Mapping the forecast of poverty increases for the next 5 years
```

```
# Merging countypop with forecasted poverty level. Then, I substited county population with mean foreca
#-----
```

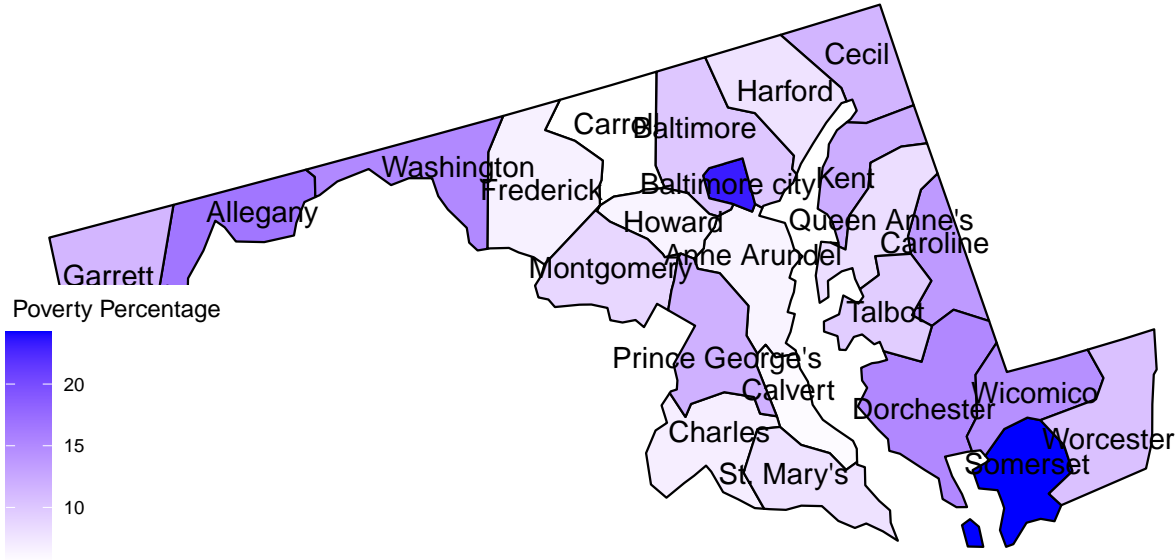
```
poverty_Prctchnge_forecast = merge(current_population,poverty_forecast, by ="county" ) %>%
  mutate(Poverty_Percent_change = .mean/Poverty_Universe * 100) %>%
  select(county,year = year.y, Poverty_Percent_change)
countypop = countypop
```

```
MD_countypop = merge(countypop,poverty_Prctchnge_forecast, by="county") %>%
  filter(abbr == "MD", year ==2026) %>%
  select(c(fips,county,abbr,Poverty_Percent_change))
```

```
#----- Map
```

```
plot_usmap(data = MD_countypop, values = "Poverty_Percent_change", include = "MD",
  labels = TRUE )+
  scale_fill_continuous(
    low = "white", high = "blue", name = " Poverty Percentage",
    label = scales::comma
  ) +
  labs(title = " Maryland Counties Poverty Level - 2026 Forecast ")
```

Maryland Counties Poverty Level – 2026 Forecast



The map reinforces the findings above and also provides a clear view of the projected poverty increase in 2026 in Somerset, Baltimore City, Allegany, Dorechester, and Washington counties.