

**Course:** Stats 4999Z - Project in Statistical Sciences

**Project:** Final Report

**Title:** Utilizing Seasonal-Trend Decomposition using LOESS and Generalized Additive  
Models to Understand the Seasonality and Long-Term Trend of Area Burned by  
Wildfire in Alberta

**Name:** Matthew Jones

**Supervisor:** Douglas G. Woolford

**Date:** April 30, 2024

## **Abstract**

The main goal of this paper is to understand the seasonality and long-term trend of area burned by wildfire in the Canadian Province of Alberta. To understand the seasonality and long-term trend of area burned, two smoothing techniques were utilized: Generalized Additive Models (GAMs) and Seasonal-Trend Decomposition using LOESS (STL). This paper also set to determine the advantages and disadvantages of using each smoothing technique to understand patterns in area burned.

Understanding the seasonality and long-term trend of area burned is especially important from a fire management perspective, as these patterns can inform short-term and long-term decision-making. It was discovered that the area burned in Alberta typically peaks during May or June. Furthermore, it was found that in Alberta, between 1959 and 2020, months May through September had strictly increasing trends in area burned, while months October through April had neither strictly increasing nor strictly decreasing trends in area burned. This paper also identified GAMs as a more advantageous approach when a smoothed long-term trend of area burned is of interest. In contrast, it was determined the STL approach is more suitable for situations where shorter-term fluctuations in area burned are of interest. In addition, GAMs were identified as a more appropriate option when users lack a statistical background.

## **1 Introduction**

Wildfires in Canada are common. Between 1970 and 2017, 8000 wildfires occurred, and on average, these wildfires burnt 2.25 million hectares (ha) annually (Tymstra et al., 2020). The two main causes of these wildfires were human activity and lighting. Human activity accounts for 49% of wildfires in Canada, while lighting accounts for 47% (Tymstra et al., 2020). Canadian wildfires have a wide range of effects. One common effect of Canadian wildfires is economic burden. For example, the 2016

Horse River Wildfire in northeastern Alberta caused \$3.84 billion in insurance claims and was the costliest natural disaster in Canadian history (2017 \$) (Tymstra et al., 2020). A second common effect of Canadian wildfires is human health impacts. Canadian wildfires impact human health as their smoke is associated with worse birth outcomes and cardiovascular events (D'Evelyn et al., 2022). Current research suggests that fire activity in Canada, measured as the number of fires and total area burned, has increased, and will continue to increase because of climate change (Coops et al., 2018; Wotton et al., 2010). Understanding fire activity trends and seasonality is essential from a fire management perspective. This is because seasonal patterns can help inform decisions like when to hire contract firefighters or lease aerial firefighting equipment like helicopters, and the long-term trends can help inform decisions on whether to invest in more firefighting equipment, such as water bombers. This paper employs two smoothing techniques to understand the seasonality and long-term trend of area burned by wildfire in the Canadian province of Alberta. The first smoothing technique is Seasonal-Trend Decomposition using LOESS (STL), and the second is Generalized Additive Models (GAMs). While the primary goal of this paper is to understand the trends and seasonality of area burned in Alberta, this paper also intends to accomplish the following sub-objectives: (1) Provide an explanation of what STL and GAMs are and how they can be used to understand the trends and seasonality fire activity, and (2) Determine the advantages and disadvantages of using each smoothing technique to analyze area burned data.

## **2 Literature Review**

### **2.1 Fire Activity Trends in Canada**

Many studies have examined the trends of fire activity in Canada. Wotton et al. (2010) utilized general circulation models (GCMs) to predict that fire occurrence in Canada will increase 25% by 2030 and 75% by the end of the century. Similarly, Flannigan et al.

(2005) utilized GCMs in conjunction with historical relationships between area burned, fire danger, and weather to project a 74-118% increase in area burned in Canada by the end of the century. Coops et al. (2018) examined the trend of area burned by wildfire in Canada, and within Canada's ecozones. An ecozone is an area with a particular type of natural environment. They found that area burned by wildfire in Canada increased 11% per year between the years 2006 and 2015. They also found that increases in area burned were particularly evident in the Taiga Plains and Taiga Shield West ecozones. Area burned increased by 26.1% per year in the Taiga Plains ecozone, and area burned increased by 12.7% per year in the Taiga Shield West ecozone. The only ecoregion they found to have a negative trend in area burned was the Atlantic maritime region.

## **2.2 Seasonality of Fire Activity in Canada**

Other studies have investigated the seasonality of fire activity in Canada. Jain et al. (2017) analyzed the trend of fire season length in North America between 1979 and 2015. They found evidence that fire season length increased over large areas of North America, with increases being particularly evident in Eastern Canada. Ahmed and Hassan (2023) investigated the trend and seasonality of fire occurrences between 1959 and 2021 within the subregions of Alberta, Canada, using non-parametric statistical tests, specifically, the Mann-Kendall test and Sen's slope estimator. They found significant increases in fire occurrences across subregions. In addition, they found that the fire season in some subregions increased in length by 4 to 5 days. Coogan et al. (2020) examined the seasonality of human- and lighting-caused wildfires in Canada. They found that lightning-caused fires peaked between June and August, whereas human-caused fires peaked in May. In addition, they found that the seasonality of human- and lighting-caused fires varied across ecozones.

### 3 Data

The data used in this study was provided by the Canadian Wildland Fire Information System (CWFIS) (Natural Resources Canada, n.d.). The dataset contains records of forest fires occurring in Canada between 1946 and 2021. Various variables were recorded for each forest fire. These variables include fire date, location, and area burned in hectares (ha). Fire management agencies in Canada, such as provinces, territories, and Parks Canada, inputted these variables. The variables can be used to create a monthly area burned time series in Alberta from 1959 to 2020. This time series excludes areas burned within the national parks found in Alberta. The time series of monthly area burned in Alberta had outliers: months with exceptionally high amounts of area burned. A log base ten transformation was applied to the time series to remove the outliers and normalize the time series. The two smoothing techniques were applied to the log transformed area burned time series to understand the seasonality and long-term trend of area burned in Alberta. The differences between the transformed and non-transformed time series can be seen in Figure 1 of the Appendix.

## 4 Methodology

### 4.1 STL

#### 4.1.1 Ordinary Least Squares Regression

Suppose the relationship between an independent variable  $x_i$  and a dependent variable  $Y_i$  for  $i = 1, 2, \dots, n$  is of form,

$$Y_i = g(x_i) + \epsilon_i$$

Where  $g(x_i)$  is a function capturing the relationship between the variables and  $\epsilon_i$  is the random error term. The random error terms are assumed to be normally distributed, independent, have zero mean, and have constant variance.

When the relationship between the variables is linear, a regression analysis technique called ordinary least squares (OLS) is commonly used. OLS represents  $g$  through the formula:

$$g(x) = \beta_1 + \beta_2 x$$

Where  $\beta_1$  and  $\beta_2$  are the unknown coefficients of the model, which are to be estimated by minimizing the residual sum of squares (RSS),

$$RSS = \sum_{i=1}^n (Y_i - g(x_i))^2$$

The RSS is the squared difference between the observed and fitted values.

#### 4.1.2 Locally Estimated Scatterplot Smoothing

In many situations however,  $g$  is non-linear. In such situations, alternative techniques to OLS are needed. One commonly used non-parametric technique is locally estimated scatterplot smoothing (LOESS). LOESS estimates a smooth function  $g$  through local regression.

The first step before estimating  $g$  using LOESS is to choose a window size. The window size determines the number of points around  $x_i$  to include in the local regression. For example, a window size of 5 would include  $x_i$  and its two nearest neighbors on either side. After selecting the window size, the following steps are followed for each  $x_i$ :

1. Centre the window around the current  $x_i$ .

2. Assign weight to the points within the window. Weights are assigned using the tricube weight function. The tricube weight functions gives more weight to points close to the center of the window and less weight to points further from the center. The tricube weight function is defined as:

$$w(x_k) = \left( 1 - \left( \frac{|x_k - x_i|}{|x_i - x_{max}|} \right)^3 \right)^3$$

Where  $x_i$  is the center of the window,  $x_k$  represents a point in the window,  $w(x_k)$  is the weight assigned to  $x_k$ , and  $x_{max}$  is the point furthest from  $x_i$  in the window.

3. Perform weighted OLS on the points within the window. Points with higher weight have greater influence on the fitted line. Points with lower weight have less influence on the fitted line.
4. Use the weighted OLS regression to obtain a fitted value at  $x_i$ . This fitted value becomes a point of the LOESS curve.

After completing the above steps for each  $x_i$ , an initial smooth curve  $g$  is obtained. However, this curve can be influenced by outliers. To combat this, LOESS refines the curve iteratively using robustness weights. These weights are designed to reduce the influence of outliers on the fitted curve. The iterative process is as follows:

1. Calculate the residuals of the initially fit curve. These are the differences between the observed ( $Y_i$ ) and predicted values ( $\hat{Y}_i$ ). The residual for the  $i^{th}$  data point is defined as:

$$e_i = Y_i - \hat{Y}_i$$

2. Calculate the robustness weight for each observation. The robustness weights are defined as,

$$\delta_i = \begin{cases} \left(1 - \left(\frac{e_i}{6s}\right)^2\right)^2, & \left|\frac{e_i}{6s}\right| < 1 \\ 0, & \left|\frac{e_i}{6s}\right| \geq 1 \end{cases}$$

Where  $\delta_i$  is the robustness weight for the  $i^{th}$  observation and  $s$  is the median of the absolute residuals ( $|e_i|$ ). Observations with smaller residuals are given larger robustness weights, and observations with larger residuals are given smaller robustness weights.

3. Fit new LOESS curve considering both the original weights and the robustness weights. Weights assigned to points within a window are the product of the original weight ( $w(x_k)$ ) and the robustness weight ( $\delta_k$ ). The resulting curve will be less influenced by outliers.
4. Repeat steps 1-3 iteratively until a stopping criterion. 2 robust iterations are generally considered an appropriate stopping criterion (Cleveland, 1979).

It should be mentioned that sometimes polynomial regression of degree 2 is used instead of OLS when obtaining a LOESS curve. An explanation of polynomial regression can be found in Section 4.2.2.

#### 4.1.3 Seasonal Trend Decomposition using LOESS

Let  $Y_v$  for  $v = 1, \dots, N$  be a time series.  $Y_v$  can be decomposed in the following way,

$$Y_v = T_v + S_v + R_v$$



Where  $T_v$  is the trend component,  $S_v$  is the seasonal component, and  $R_v$  is the remainder component. The trend component represents the long-term change in the data, the seasonal component represents the repeated pattern at regular intervals in the data, and the remainder component represents the remaining variation in the data after the trend and seasonal components are removed. STL is a decomposition method that can extract the trend, seasonal, and remainder components from a time series. It decomposes the time series in the manner detailed above. That is, it extracts the components in a way so that when they are added together, they equal the original value of the time series.

When decomposing a time series, STL uses two loops: an inner loop, and an outer loop. Each run of the outer loop is followed up with multiple passes of the inner loop. The outer loop calculates the robustness weight for the inner loop, and the inner loop calculates the trend, seasonal, and remainder components through iteration.

When performing an STL decomposition, the following steps can be followed:

**Step 1:** Input the time series  $Y_n$  to be decomposed and initialize the value of the trend component ( $T_v^{(0)}$ ) and the robustness weight ( $\rho_v$ ).  $T_v^{(0)}$  is initialized to 0 and  $\rho_v$  is initialized to 1.

**Step 2:** Enter the inner loop. Each pass of the inner loop consists of seasonal smoothing that updates the seasonal component, followed by trend smoothing that updates the trend component. Let  $S_v^{(k)}$  be the seasonal component at the end of the  $k^{th}$  pass and let  $T_v^{(k)}$  be the trend component at the end of the  $k^{th}$  pass. The inner loops consist of six sub-steps to obtain  $S_v^{(k+1)}$  and  $T_v^{(k+1)}$ .

(1) Sub-step 1: Detrend  $Y_n$ . A detrended series  $Y_v - T_v^{(k)}$  is computed.

- (2) Sub-step 2: Smooth the cycle-subseries of the detrended series. A cycle-subseries refers to the subseries of values at one position in a seasonal cycle. Suppose we have a monthly time series between the years 1959 and 2021. January 1959, January 1960,..., January 2021 would be a cycle-subseries, February 1959, February 1960,..., February 2021 would be a cycle-subseries, and so on. Each cycle-subseries is smoothed using LOESS with robustness weight  $\rho_v$ . The smoothed values for all the cycle-subseries are combined to form a temporary series denoted  $C_v^{(k+1)}$ .
- (3) Sub-step 3: Perform low-pass filtering on  $C_v^{(k+1)}$ . The filter consists of a moving average whose length is equal to the seasonal period, followed by a second moving average whose length is also equal to the seasonal period, followed by a moving average of length three, followed by a LOESS implementation with robustness weight  $\rho_v$ . The obtained series after performing the low-pass filtering is denoted  $L_v^{(k+1)}$ .
- (4) Sub-step 4: Detrend  $C_v^{(k+1)}$ . The seasonal component  $S_v^{(k+1)}$  is calculated as
- $$S_v^{(k+1)} = C_v^{(k+1)} - L_v^{(k+1)}.$$
- (5) Sub-step 5: Deseasonalize  $Y_n$ . A deseasonalized series  $Y_v - S_v^{(k+1)}$  is computed.
- (6) Sub-step 6: Perform LOESS on  $Y_v - S_v^{(k+1)}$  with robustness weight  $\rho_v$ . The result of the LOESS implementation is the trend component  $T_v^{(k+1)}$ .

**Step 3:** Check whether a predetermined number of passes through the inner loop have been reached. If the predetermined number of passes has not been reached, let  $T_v^{(k)} = T_v^{(k+1)}$ , and return to Step 2. If the predetermined number of passes has been reached, the inner loop is exited, and the outer loop is entered.

**Step 4:** Enter the outer loop. Let  $T_v = T_v^{(k+1)}$ ,  $S_v = S_v^{(k+1)}$ , and  $R_v = Y_v - S_v - T_v$ . If the predetermined number of iterations of the outer loop has been reached, the trend component  $T_v$ , the seasonal component  $S_v$ , and the remainder component  $R_v$  are plotted. Otherwise, the robustness weight  $\rho_v$  is updated, and the inner loop is reentered.  $\rho_v$  is updated using the following formula,

$$\rho_v = B\left(\frac{|R_v|}{6 \cdot \text{median}(|R_v|)}\right)$$

Where  $B$  is the bisquare weight function defined as,

$$B(u) = \begin{cases} (1 - u^2)^2, & 0 < u \leq 1 \\ 0, & u > 1 \end{cases}$$

STL has 6 parameters that must be chosen.  $n_{(p)}$  is the number of observations in the cycle of the seasonal component ( $n_{(p)} = 12$  for monthly data),  $n_{(i)}$  is the number of passes through the inner loop,  $n_{(o)}$  is the number of times the outer loop is run,  $n_{(l)}$  is the smoothing parameter of the low-pass filter,  $n_{(t)}$  is the smoothing parameter for the trend component, and  $n_{(s)}$  is the smoothing parameter for the cycle-subseries. A smoothing parameter sets the window size when performing LOESS. Guidelines exist for the selection of all parameters except  $n_{(s)}$  (Cleveland et al., 1990).  $n_{(s)}$  is to be chosen using diagnostic methods, such as a seasonal-diagnostic plot (Cleveland et al., 1990). The larger the value of  $n_{(s)}$  the smoother the seasonal component will be.

#### 4.1.4 Implementation in R

This report will use STL to extract the trend, seasonal, and remainder components from Alberta's log-transformed monthly area burned. There are two options when performing STL. The first option is to extract a constant seasonal component, and the

second is to extract a non-constant seasonal component. When a constant seasonal component is extracted, the monthly values on the curve are the same for each year. For example, January 1959, January 1960, and January 2021 all take the same value on the curve. When a non-constant seasonal component is extracted, monthly values on the curve can take on different values, so January 1959, January 1960, and January 2021 can take on different values. The advantage of extracting a constant seasonal component over time is that it gives insight into the general interannual pattern of area burned across years. The disadvantage of assuming constant seasonality is that it doesn't allow one to see changes in the seasonality of wildfires over time. It should be noted that extracting a constant seasonal component is a simplification, as the seasonality of area burned in Alberta isn't constant due to climatic patterns.

The 'stl()' function in R's 'stats' package allows one to perform STL. The 's.window' parameter corresponds to the smoothing parameter for the cycle-subseries ( $n_{(s)}$ ). As mentioned, it is to be chosen by the data analyst via diagnostic methods. This paper will use a seasonal-diagnostic plot to select the appropriate size of the smoothing parameter. The smoothing parameter only needs to be chosen when a non-constant seasonal component is to be extracted. When a constant seasonal component is to be extracted, the 's.window' parameter is set to 'periodic.'

## 4.2 GAMs

### 4.2.1 Polynomial regression

Polynomial Regression is an alternative parametric technique to LOESS that can be used to estimate the function  $g$  when it is non-linear. It expresses  $g$  as a  $n^{th}$  degree polynomial. Mathematically,

$$g(x) = \beta_1 + \beta_2 x + \beta_3 x^2 + \dots + \beta_{n+1} x^n$$

Where  $\beta_1, \beta_2, \beta_3, \dots, \beta_{n+1}$  are the coefficients of the model, which are estimated by minimizing the RSS.

Alternatively, a polynomial regression model can be represented using basis functions. Mathematically,

$$g(x) = \beta_1 b_1(x) + \beta_2 b_2(x) + \beta_3 b_3(x) + \cdots + \beta_{n+1} b_{n+1}(x)$$

Where  $b_1(x) = 1$ ,  $b_2(x) = x$ ,  $b_3(x) = x^2$ , ...,  $b_{n+1}(x) = x^n$ .

Often polynomial regression is a poor choice for modelling  $g$ . This is because for large  $n$ , polynomial basis functions tend to exhibit multicollinearity, which can cause issues in the estimation of the coefficients. In addition, high order polynomials tend to fluctuate wildly if there are wide gaps in the  $x_i$ 's.

#### 4.2.2 Regression Splines

Regression splines are a more flexible alternative to polynomial regression that can be used to estimate the function  $g$ .

To build a regression spline the following steps can be followed:

1. Define a set of knots  $\tau_1 < \tau_2 < \cdots < \tau_K$  over the domain of  $g$ . Knots are values that split the domain of  $g$  into different intervals.
2. Perform polynomial regression within each interval under the restriction that adjacently fit polynomials are smooth and continuous at their respective knots. Continuity is obtained by requiring adjacent polynomials to connect at knots, and smoothness is obtained by requiring adjacent polynomials first and second derivatives to equal at knots.
3. The resulting piecewise continuous polynomial function is called a regression spline.

If the polynomials fit within each interval are of degree- $n$ , then we say the piecewise continuous polynomial function is a degree- $n$  regression spline. A degree-3 regression spline is referred to as a cubic spline. Cubic splines can be constructed using basis functions. The basis representation of a cubic spline with  $K$  knots is given by the following equation,

$$g(x_i) = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \beta_3 x_i^3 + \sum_{j=1}^K \beta_{3+j} (x_i - \tau_j)_+^3$$

Where  $\beta_0, \beta_1, \dots, \beta_{3+K}$  are the coefficients of the model, which are estimated through the minimization of the RSS. The term  $(x_i - \tau_j)_+^3$  is called a truncated power basis function. Mathematically, it is defined as,

$$(x_i - \tau_j)_+^3 = \begin{cases} (x_i - \tau_j)^3, & \text{if } x_i > \tau_j \\ 0, & \text{otherwise} \end{cases}$$

Under the constraint that  $g(x_i)$  is linear outside the interval  $[\tau_1, \tau_K]$ , we say that  $g(x_i)$  is a natural cubic spline.

#### 4.2.3 Penalized Regression Splines

One issue with regression splines is the requirement to choose knot locations. Knot locations can dramatically influence modelling results, making their selection both in terms of the number of knots and their locations complicated. Penalized regression splines offer a solution to this problem. They do so by utilizing a relatively large number of knots, making the fitted model less sensitive to the knot locations. They then introduce a penalty, ensuring the fitted model is not excessively wiggly. The penalty term avoids the issue of overfitting, which is a common problem when many knots are used.

As mentioned, a cubic spline is fit by minimizing the RSS. Penalized cubic splines, however, are fit by minimizing the RSS plus a penalty term. That is, they are fit by minimizing the following equation,

$$\sum_{i=1}^n (y_i - g(x_i))^2 + \lambda \int g''(t)^2 dt$$

Where  $g$  is the basis representation of a cubic spline, the first term in the equation is the RSS, and the second term in the equation is called the roughness penalty, and is controlled by  $\lambda$ , which is a nonnegative tuning parameter.

The roughness penalty is defined as the second derivative of the function  $g$  squared, integrated over the domain of  $g$ . The second derivative is a measure of how wiggly  $g$  is. Functions that have large second derivatives are wiggly, and functions with small second derivatives are not wiggly. For example, constant and linear functions have second derivatives equal to 0. So, the roughness penalty can be thought of as the sum of the wiggleness of  $g$  over its domain. The square in the penalty term ensures that the sign of the second derivative is always positive, and the integration does the summing over the domain.

When  $\lambda$  is small, the roughness penalty has less weight, and the estimated  $g$  will be wigglier. When  $\lambda$  is large, the roughness penalty has more weight, and the estimated  $g$  will be less wiggly. The optimal  $\lambda$  is chosen through cross-validation.

#### 4.2.4 Multiple Linear Regression

Suppose the relationship between a dependent variable  $Y_i$  and  $p$  independent variables  $x_{i1}, x_{i2}, \dots, x_{ip}$   $i = 1, 2, \dots, n$  is of form,

$$Y_i = h(x_{i1}, x_{i2}, \dots, x_{ip}) + \epsilon_i$$

Where  $h$  is a function capturing the relationship between the multiple independent variables and the dependent variables, and  $\epsilon_i$  is the random error term. The random error terms are assumed to be normally distributed, independent, have zero mean, and have constant variance.

When the relationship between each independent variable and the dependent variable is linear, multiple linear regression (MLR) is commonly used to represent  $h$ . MLR represents  $h$  as,

$$h(x_{i1}, x_{i2}, \dots, x_{ip}) = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip}$$

The unknown coefficients of the model are estimated by minimizing the RSS, similarly to the estimation process in OLS.

The issue with MLR lies in its assumption that the relationship between each predictor and the response is linear. In many real-world scenarios, this is not the case. In such situations, alternative modelling techniques are needed.

#### 4.2.5 Generalized Additive Models

Generalized additive models (GAMs) address MLR's limitation by allowing for non-linear relationships between the dependent variable and each independent variable. GAMs allow for non-linear relationships by replacing the linear components  $\beta_j x_{ij}$  in MLR with smooth non-linear functions  $f_j(x_{ij})$ . A GAM for a response variable  $Y_i$  with  $p$  predictors is defined as,

$$Y_i = \beta_0 + f_1(x_{i1}) + f_2(x_{i2}) + \dots + f_p(x_{ip}) + \epsilon_i$$

Where  $\beta_0$  is the intercept,  $f_j(x_{ij})$  is a smooth function for the  $j^{th}$  predictor  $x_{ij}$ , and  $\epsilon_i$  is the random error term. The random error terms are assumed to be normally distributed, independent, have zero mean, and have constant variance. Often the smooth functions  $f$  in GAMs are represented using penalized cubic splines.

To understand the fitting procedure of a GAM, consider the case in which we are modelling data  $(Y_i, x_{1i}, x_{2i})$ . A GAM for such data would be,

$$Y_i = \beta_0 + f_1(x_{1i}) + f_2(x_{2i}) + \epsilon_i$$

As mentioned, the smooth functions  $f$  are often represented using penalized cubic splines. So,  $f_1$  and  $f_2$  would each be represented using the cubic spline basis outlined



in section 4.2.2. The coefficients of the GAM are then estimated through the minimization of the RSS plus a penalty term for each function's roughness. Mathematically, this would be the minimization problem,

$$\text{minimize } \sum_{i=1}^n (Y_i - \beta_0 - f(x_{1i}) - f(x_{2i}))^2 + \lambda_1 \int \frac{\partial^2 f_1}{\partial x_1^2} dx_1 + \lambda_2 \int \frac{\partial^2 f_2}{\partial x_2^2} dx_2$$

#### 4.2.6 Time Series Decomposition using GAMs

GAMs can be used to decompose a time series in the same way STL can. However, rather than decomposing a time series into a trend, seasonal, and remainder components, GAMs decompose a time series into a trend, seasonal, and intercept component. The GAM to perform this type of decomposition on log area burned in Alberta is as follows,

$$Y_i = \beta_0 + f(m_i) + f(t_i) + \epsilon_i$$

Where  $Y_i$  represents log transformed area burned in Alberta,  $m_i$  represents the month, and  $t_i$  represents year.

The term  $f(t_i)$  represents the trend component and  $f(m_i)$  represents the seasonal component. The extracted seasonal component is constant over time. The trend component can be interpreted as the pattern log area burned follows as years progress when month is held constant, and the seasonal component can be interpreted as the pattern log area burned follows within a year when year is held constant.

To perform this decomposition the 'mgcv' package in R programming language can be used (Wood, 2006). When a GAM is fit using the 'mgcv' package, metrics that evaluate its performance such as adjusted R-squared are calculated (Wood, 2006).

#### 4.2.7 GAM Surface

GAMs can be used to fit a surface to a time series. This surface can then be used to understand the trends and seasonality of the time series. The GAM to fit a surface to the log area burned time series in Alberta is as follows,

$$Y_i = \beta_0 + f(m_i, t_i) + \epsilon_i$$

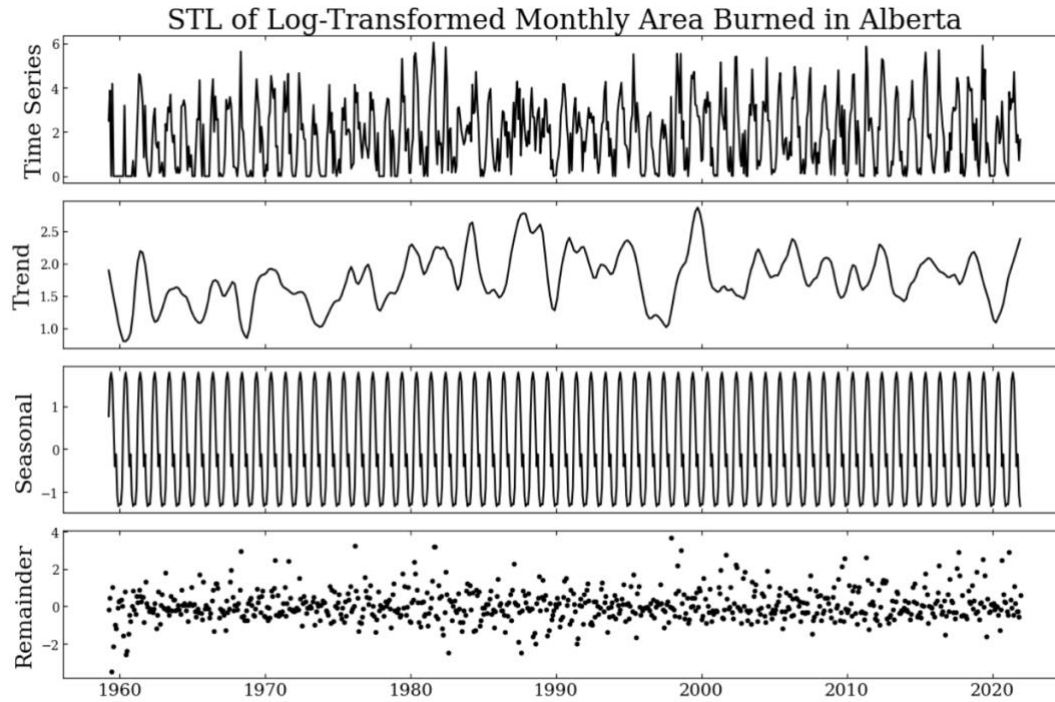
Where  $f(m_i, t_i)$  is a smooth bivariate function capturing the seasonality and trend of log area burned. The function is fit to the data plotted with month on the x-axis, year on the y-axis, and log area burned on the z-axis. Again, this GAM can be fit using R programming languages 'mgcv' package.

## 5 Results

### 5.1 Constant Seasonality

#### 5.1.1 STL

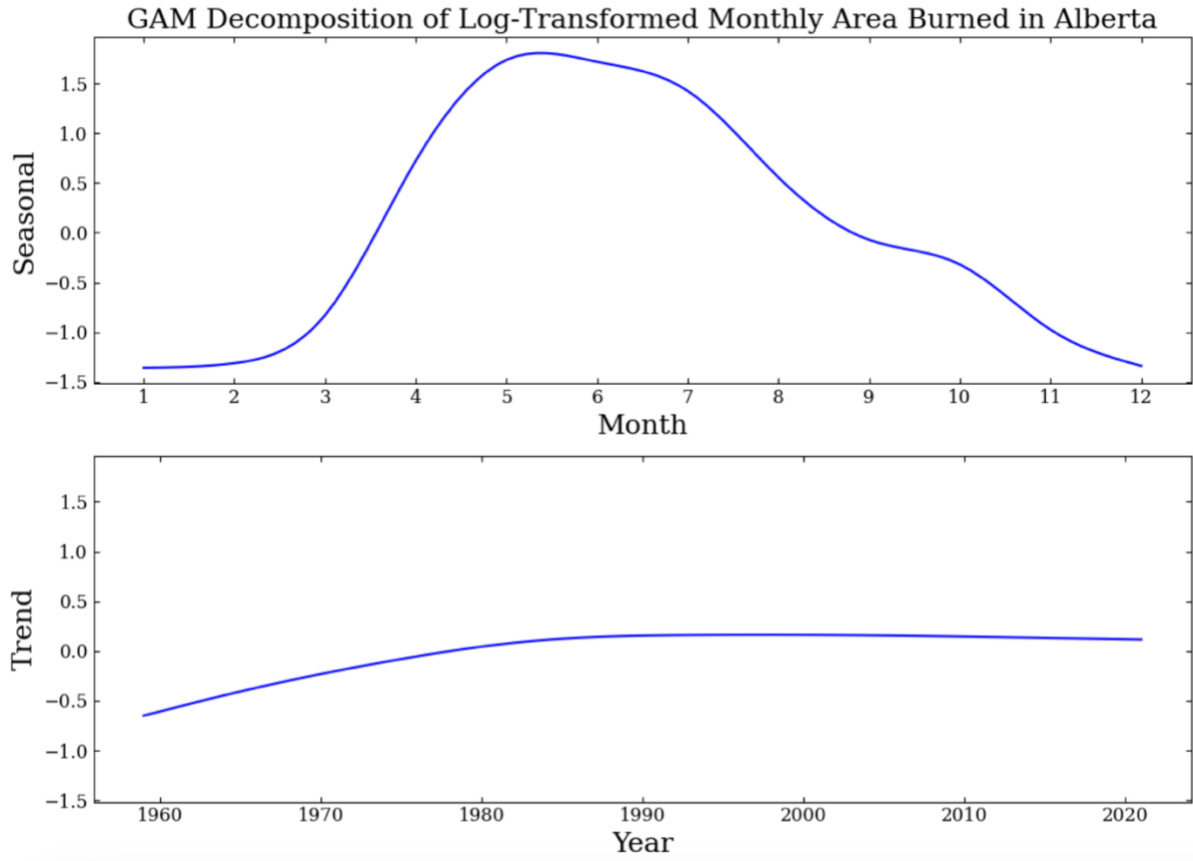
Figure 1 visualizes the STL when a constant seasonal component is extracted. The first row represents the Alberta monthly area burned time series log-transformed, the second row represents the trend component, the third represents the seasonal component, and the fourth represents the remainder component. From the figure, we can see that the trend component is not smooth. Instead, it fluctuates with time, with periods of increase and decrease.



**Figure 1:** STL of log-transformed monthly area burned in Alberta when a constant seasonal component is extracted.

### 5.1.2 GAM

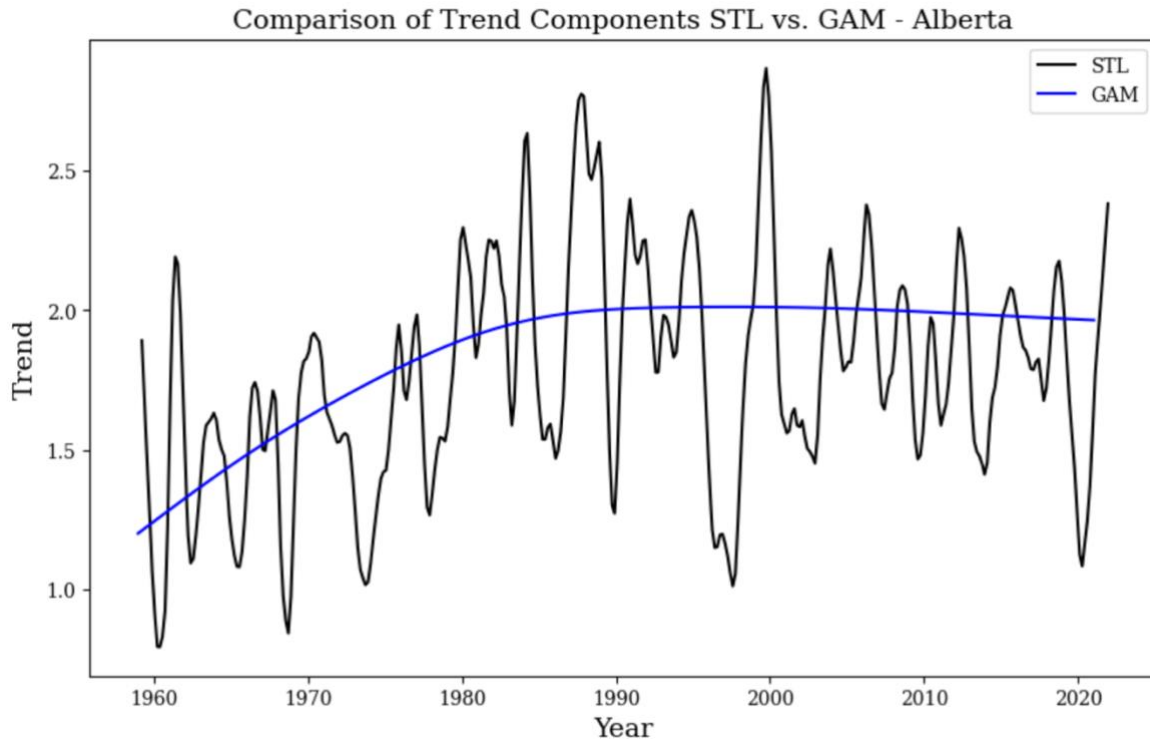
Figure 2 visualizes the GAM decomposition. The first row visualizes the seasonal component, and the second row visualizes the trend component. The seasonal component suggests that the area burned in Alberta peaks in May. The trend component suggests that the area burned in Alberta trended upwards between 1959 and 1990 but has since plateaued. The adjusted R-squared of the GAM was 0.61.



**Figure 2:** GAM decomposition of log transformed monthly area burned in Alberta.

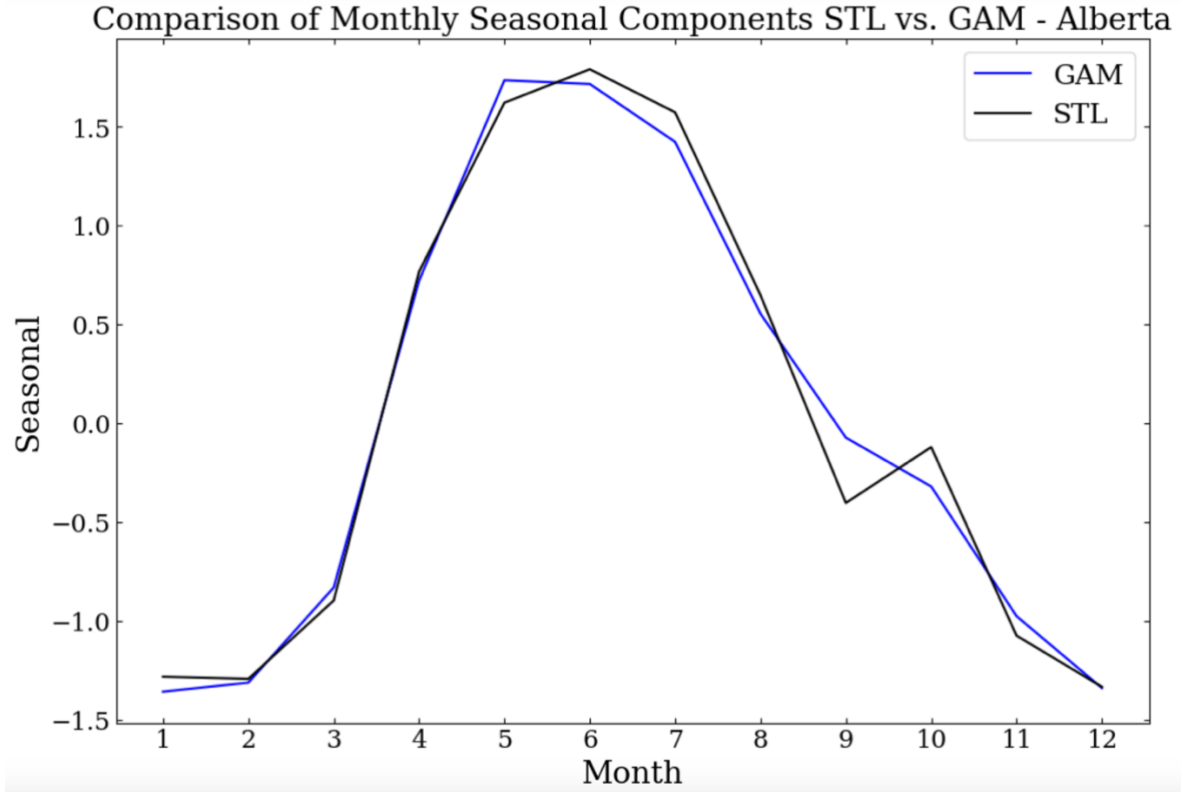
### 5.1.3 Comparison

The trend and constant seasonal component extracted using STL can be directly compared with the trend and constant seasonal component extracted using GAMs. Figure 3 compares the STL and GAM trend components. The GAM trend component is much smoother than the STL trend component. However, both components appear to follow a similar trajectory.



**Figure 3:** Comparison of GAM and STL trend components.

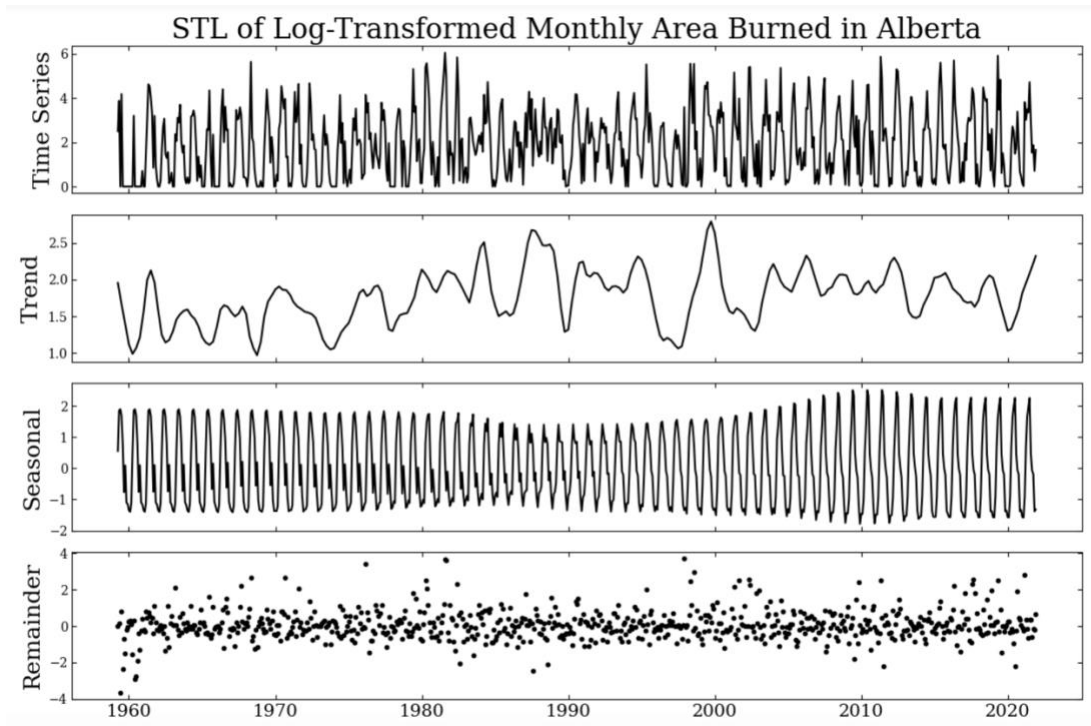
Figure 4 compares the STL and GAM seasonal components. In general, both components follow a similar pattern. However, the GAM seasonal component suggests that the area burned in Alberta peaks in May, whereas the STL seasonal component suggests that the area burned in Alberta peaks in June.



**Figure 4:** Comparison of STL and GAM seasonal components.

## 5.2 STL Non-Constant Seasonal Component

Figure 5 visualizes the STL when a non-constant seasonal component is extracted. Like Figure 1, the first row represents the monthly area burned time series in Alberta log-transformed, the second row represents the trend component, the third row represents the seasonal component, and the fourth row represents the remainder component. The size of the smoothing parameter for the cycle-subseries ( $n_s$ ) was chosen to be 17. This decision was facilitated using a seasonal-diagnostic plot, which can be found in Figure 2 of the appendix. From Figure 5 it can be seen that the seasonal component varies over time, suggesting that the seasonality of area burned in Alberta has evolved with time.

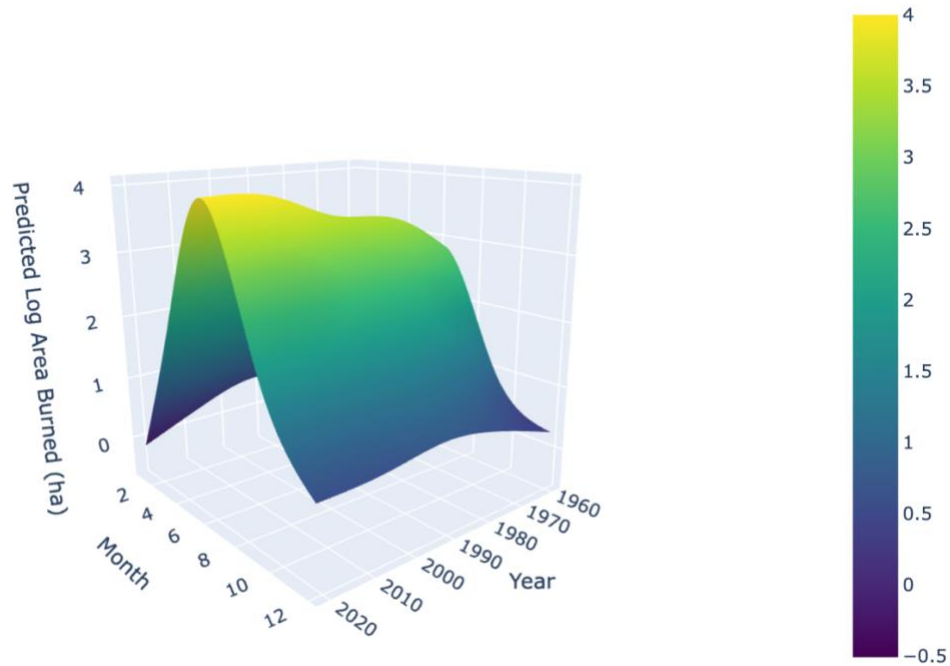


**Figure 5:** STL of log-transformed monthly area burned in Alberta when a non-constant seasonal component is extracted.

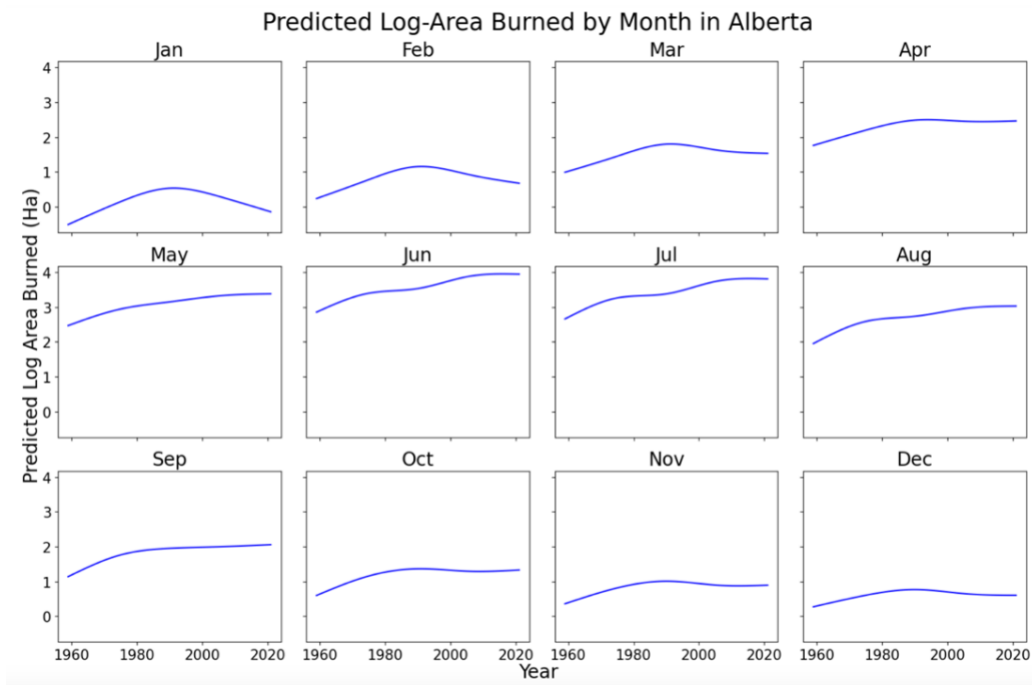
### 5.3 GAM Surface

Figure 6 shows the surface fit to the log area burned time series in Alberta. The surface represents the predicted log area burned. The x-axis is month, and the y-axis is year. Understanding the seasonality and trends by looking at the surface can be challenging. So, one can plot slices of the surface to make these patterns clearer. Figure 7 plots the surface sliced at each month. The approach shows how the predicted log area burned has varied across years within each month. From Figure 7, we can see that predicted area burned isn't strictly increasing or decreasing in months January, February, March, April, October, November, and December. However, in months May through September, predicted area burned is strictly increasing. The adjusted R-squared of the GAM was 0.574.

Predicted Log Area Burned in Alberta



**Figure 6:** Predicted log area burned in Alberta using GAM bivariate smoother.



**Figure 7:** Predicted log-area burned within each month across years using a GAM bivariate smoother.



## 6 Discussion

The GAM decomposition and the STL extracting a constant seasonal component suggest that area burned in Alberta peaks in May or June. The GAM decomposition suggested that the area burned in Alberta increased between 1959 and 1990; however, after the 1990s, it plateaued. Although the GAM decomposition suggested that the area burned in Alberta plateaued post-1990s, the GAM surface revealed a different story. The GAM surface indicated that the area burned in months May through September was strictly increasing from 1959 to 2020, while months October through April showed no clear trends in the area burned.

The advantage of using the GAM surface compared to the GAM decomposition is that it allows one to see the trends in area burned within each month. As discussed above, the GAM decomposition suggested that the area burned in Alberta was not strictly increasing due to a plateau post-1990s. However, the GAM surface revealed that the area burned in Alberta was strictly increasing, however, only in months May through June, rather than all months. The limitation of the GAM surface is when outliers are present. For example, if the area burned in May of one year is extreme, the predicted area burned in May of neighbouring years may increase. That is, the surface would be dragged up in neighbouring years even if those years experienced regular fire activity.

When both STL and GAMs are used, a smoothing parameter must be chosen. When STL is used, the smoothing parameter for the cycle-subseries must be chosen, and when GAMs are used, the non-negative tuning parameters of the penalties must be selected. An advantage of the GAM approach is that their smoothing parameters are algorithmically determined using R's 'mgcv' package using cross-validation. Meaning the user does not need to select the smoothing parameter manually. On the other hand, the STL approach requires the analyst to choose its smoothing parameter via

diagnostic methods or domain knowledge, which can be a challenging task without an understanding of STL's theory. Thus, it is suggested that if a user lacks a statistical background, GAMs should be used over STL to understand the trends and seasonality of area burned, as GAMs avoid the manual selection of smoothing parameters.

The GAM trend component is smoother than the STL trend component. A smoother trend component is advantageous when the overarching long-term trend of area burned is of interest. On the other hand, having a less smooth seasonal component is beneficial when shorter-term fluctuations in area burned are of interest. Therefore, it is suggested that the GAM approach be used if the goal of the analysis is obtaining a smoothed long-term trend, whereas the STL approach should be used in situations when shorter-term fluctuations are of interest.

## **7 Conclusion and Future Work**

The main purpose of this paper was to understand the long-term trend and seasonality of the area burned by wildfire in Alberta. It was found that the area burned in Alberta typically peaks during May or June. Furthermore, it was found that in Alberta, between 1959 and 2020, May through September had strictly increasing trends in area burned, while October through April had neither strictly increasing nor strictly decreasing trends in area burned.

This paper also sought to determine the advantages and disadvantages of using STL and GAMs to understand the seasonality and long-term trend of area burned by wildfire. It was found that GAMs should be used if a smoothed long-term trend of area burned is of interest, whereas STL should be used when shorter-term fluctuations in area burned are of interest. Furthermore, it is recommended that GAMs are used over STL to understand the trends and seasonality of area burned when the user doesn't have a strong statistical background.

Future work should revolve around determining factors that influence the trends and seasonality of area burned in Alberta. One thing that could be explored is the relationship between Pacific Ocean Sea surface temperature and area burned in Alberta. This relationship would be interesting to investigate as research has found a relationship between British Columbia lightning-caused fires and Pacific Ocean Sea surface temperatures (Wang et al., 2010). Future work should also investigate the trends and seasonality of area burned in provinces other than Alberta. Understanding each province's area burned trends is helpful as they can be used to develop fire management strategies specific to each province. Finally, future work should investigate the trends and seasonality of weekly or bi-weekly area burned in Alberta, as this paper only examined the trends and seasonality of monthly area burned.

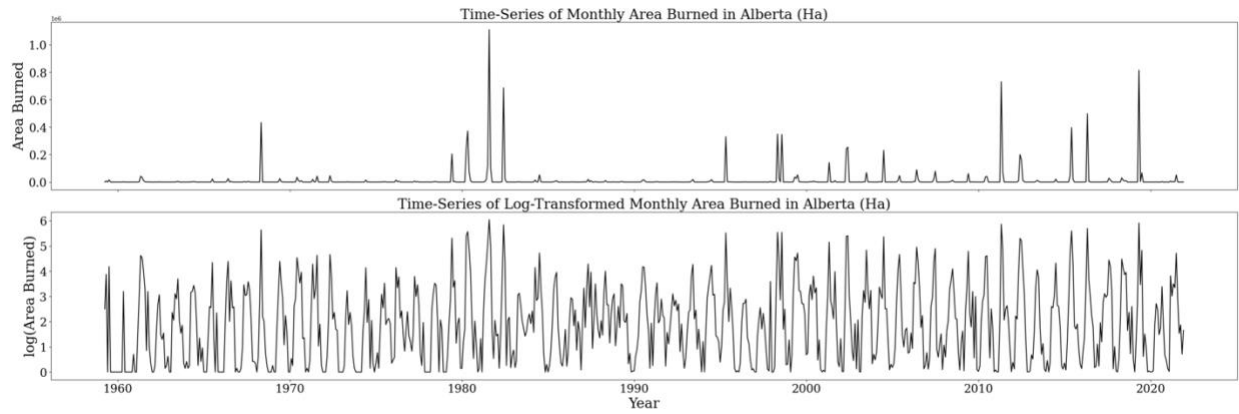
## References

- Ahmed, M. R., & Hassan, Q. K. (2023). Occurrence, Area Burned, and Seasonality Trends of Forest Fires in the Natural Subregions of Alberta over 1959–2021. *Fire*, 6(3), 96. <https://doi.org/10.3390/fire6030096>
- Cleveland, W. S. (1979). Robust Locally Weighted Regression and Smoothing Scatterplots. *Journal of the American Statistical Association*, 74(368), 829–836. <https://doi.org/10.1080/01621459.1979.10481038>
- Cleveland, R. B., Cleveland, W. S., McRae, J. E., & Terpenning, I. (1990). STL: A seasonal-trend decomposition procedure based on Loess. *Journal of Official Statistics*, 6(1), 3-73.
- Coogan, S. C. P., Cai, X., Jain, P., & Flannigan, M. D. (2020). Seasonality and trends in human- and lightning-caused wildfires  $\geq 2$  ha in Canada, 1959–2018. *International Journal of Wildland Fire*, 29(6), 473. <https://doi.org/10.1071/WF19129>
- Coops, N. C., Hermosilla, T., Wulder, M. A., White, J. C., & Bolton, D. K. (2018). A thirty year, fine-scale, characterization of area burned in Canadian forests shows evidence of regionally increasing trends in the last decade. *PLOS ONE*, 13(5), e0197218. <https://doi.org/10.1371/journal.pone.0197218>
- D'Evelyn, S. M., Jung, J., Alvarado, E., Baumgartner, J., Caligiuri, P., Hagmann, R. K., Henderson, S. B., Hessburg, P. F., Hopkins, S., Kasner, E. J., Krawchuk, M. A., Krenz, J. E., Lydersen, J. M., Marlier, M. E., Masuda, Y. J., Metlen, K., Mittelstaedt, G., Prichard, S. J., Schollaert, C. L., ... Spector, J. T. (2022). Wildfire, Smoke Exposure, Human Health, and Environmental Justice Need to

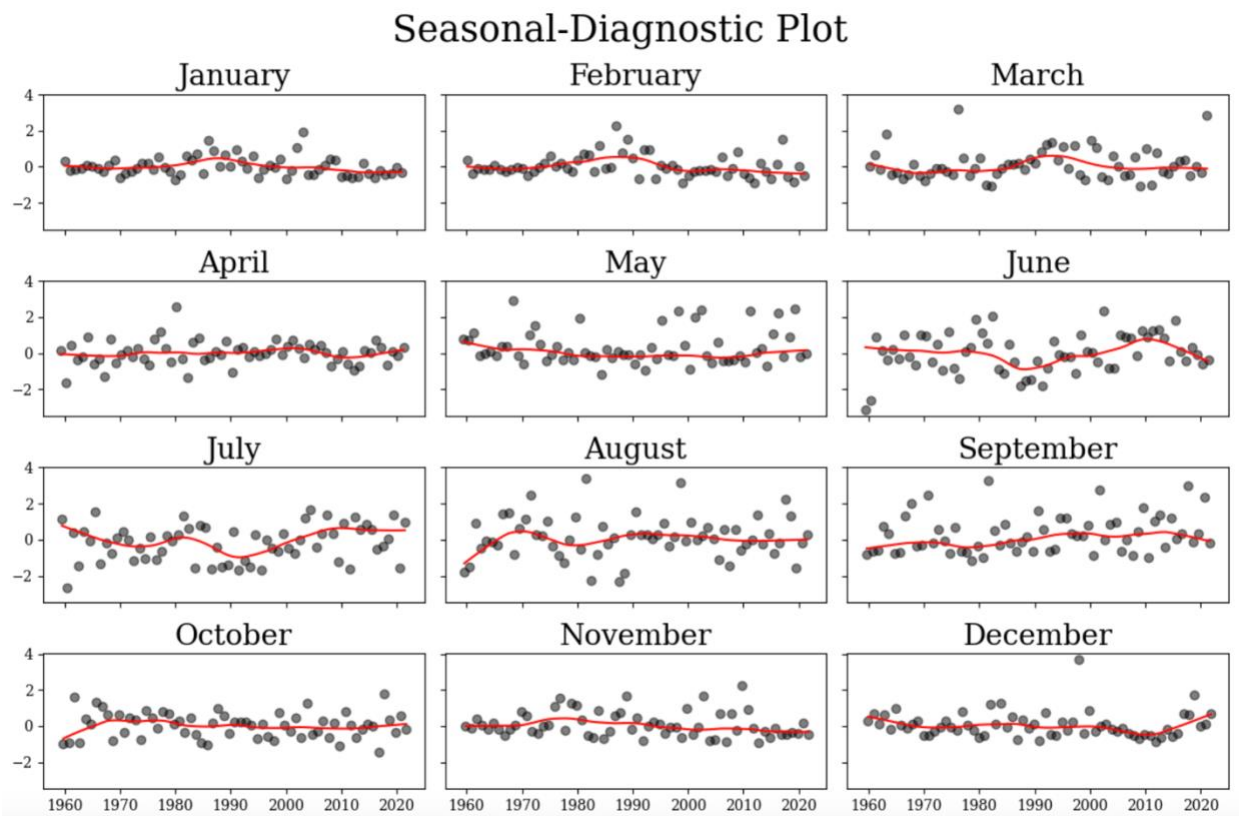
- be Integrated into Forest Restoration and Management. *Current Environmental Health Reports*, 9(3), 366–385. <https://doi.org/10.1007/s40572-022-00355-7>
- Flannigan, M. D., Logan, K. A., Amiro, B. D., Skinner, W. R., & Stocks, B. J. (2005). Future Area Burned in Canada. *Climatic Change*, 72(1–2), 1–16. <https://doi.org/10.1007/s10584-005-5935-y>
- Jain, P., Wang, X., & Flannigan, M. D. (2017). Trend analysis of fire season length and extreme fire weather in North America between 1979 and 2015. *International Journal of Wildland Fire*, 26, 1009–1020. <https://doi.org/10.1071/WF17008>
- Natural Resources Canada. (n.d.). CWFIS (Canadian Wildland Fire Information System) Datamart. Retrieved September 18, 2024, from <https://cwfis.cfs.nrcan.gc.ca/datamart>
- Tymstra, C., Stocks, B. J., Cai, X., & Flannigan, M. D. (2020). Wildfire management in Canada: Review, challenges and opportunities. *Progress in Disaster Science*, 5, 100045. <https://doi.org/10.1016/j.pdisas.2019.100045>
- Wang, Y., Flannigan, M., & Anderson, K. (2010). Correlations between forest fires in British Columbia, Canada, and sea surface temperature of the Pacific Ocean. *Ecological Modelling*, 221(1), 122–129. <https://doi.org/10.1016/j.ecolmodel.2008.12.007>
- Wood, S. N. (2006). *Generalized additive models: An introduction with R*. Chapman and Hall/CRC Press.

Wotton, B. M., Nock, C. A., & Flannigan, M. D. (2010). Forest fire occurrence and climate change in Canada. *International Journal of Wildland Fire*, 19(3), 253–271. <https://doi.org/10.1071/WF09002>

## Appendix



**Figure 1:** Difference between non-transformed and transformed monthly area burned time series in Alberta.



**Figure 2:** Seasonal diagnostic plot used to select the smoothing parameter for the cycle-subseries.