# baytrends—Long Term Water Quality Trend Analysis

**Rebecca Murphy[a], Elgin Perry[b], Jennifer Keisman[c], Jon Harcum[d], Erik W Leppo[d]**

**[a]University of Maryland Center for Environmental Science,**
**[b]Statistics Consultant, [c]U.S. Geological Survey, [d]Tetra Tech, Inc.**

## Contents

## Introduction

In general, we expect that users wanting to test baytrends will have a basic understanding of R and will use baytrends from within RStudio.

The package, baytrends, includes built in support files (stationMasterList, layerLukup, parameterList, and usgsGages) and a sample data set (dataCensored) that facilitate testing the functionality and flexibility of using the general additive model (GAM) approach for evaluating long-term trends.

If you are interested in digging into baytrends further, contact the package maintainer (see package documentation for email) for the following additional data files:

- **cbpDataAll_04_QW_baytrends.rda**. Similar to dataCensored, but more data.
- **CBPstatsLookupTables(v1.0.x).xlsm**. Template spreadsheets for creating user defined support files.
- **1984-2016 seasonally detrended salinity data.rda**. Used for modeling with salinity as an independent variable.
- **1983-2016 seasonally detrended flow data.rda**. Used for modeling with flow as an independent variable.
- **template.docx.** An MS Word file that includes pre-defined header styles for outputting R markdown scripts to MS Word.

Vignettes have been created for the following specialized topics and can be found by reviewing the package documentation:

- **Vignette_QW**. Instructions for creating a qw-formatted data set similar to the built-in data set dataCensored.
- **Vignette_Create_Seasonally_Detrended_Flow_and_Salinity_Data_Sets**. Instructions for creating seasonally-detrended flow and salinity data sets similar to **1983-2016 seasonally detrended flow data.rda** and **1984-2016 seasonally detrended salinity data.rda**.

# The Basics

The most basic analysis can be demonstrated with the below code chunk. The first three lines of the below code chunk load the baytrends package and specify the working directory. The baytrends' function analysisOrganizeData processes an input data frame (dataCensored) and sets up overall analysis parameters. The processed data frame (df) and overall analysis specifications (analySpec) can be extracted from the returned list (dfr). By default, analysisOrganizeData will output some basic record count statistics and tables listing information about parameters, layers, and stations. The function, gamTest, is the workhorse of baytrends and is called after selecting surface (layer='S') total nitrogen (dep='tn') at station CB5.4 (stat='CB5.4') for analysis.

*Code Chunk 1. Minimal baytrends example.*

```
library(baytrends)
ProjRoot <-  "E:/Dropbox/CBP/cbpTrends/work10"
setwd(ProjRoot)

dfr       <- analysisOrganizeData(dataCensored)
df        <- dfr[["df"]]
analySpec <- dfr[["analySpec"]]

stat = 'CB5.4'; dep = 'tn'; layer = 'S';
gamResult <- gamTest(df, dep, stat, layer, analySpec)
```

The values for station (stat), dependent variable (dep), and layer (layer) can be changed to other available values in the processed data frame (df) which can be determined by inspecting the data frame or reviewing the output tables from analysisOrganizeData.

In this basic mode, gamTest will output a series of tables and plots for the following three models:

- **Linear Trend with Seasonality (gam0).** The dependent variable is modeled with a <u>linear</u> term as a function of year and a stationary seasonal term.
- **Non-linear Trend with Seasonality (gam1).** The dependent variable is modeled with a <u>non-linear</u> term as a function of year and a stationary seasonal term.
- **Non-linear trend with Seas+Int (gam2).** The dependent variable is modeled with a <u>non-linear</u> term as a function of year. Seasonality includes an interaction term which allows seasonality to vary over the period of record.

In baytrends, the above models are referred to as gam0, gam1, and gam2. The following two additional models, gam3 and gam4, have been included with baytrends as well.

- **Non-linear trend with Seas+Int. & Intervention (gam3).** Same as gam2 but includes an intervention term which allows for a data shift (i.e., step change) that could be caused by changing analytical methods or implementing a source control that might be expected to cause a step change in the time series.
- **Non-linear trend with Seas+Int. & Hydro Adj (gam4).** Same as gam2 but includes a hydrologic term that allows for factoring wet/dry conditions in the model.

To demonstrate gam3 and gam4, it is necessary to include some supplemental data by updating the previous code as shown below (Code Chunk 2). Two R data files (*.rda) are loaded and a variable methodsList, is artificially created to simulate a method change on May 1, 1998 for total suspended solids (dep='tss'). An additional gamTest call is also added and expanded to pass the detrended salinity and detrended flow data. Selected figures from running this code are presented next.

***Code Chunk 2. Example code chunk that includes loading supplemental salinity and flow data; and establishes an example intervention term. (Gray lines of code previously discussed.)***

```r
library(baytrends)
ProjRoot <-  "E:/Dropbox/CBP/cbpTrends/work10"
setwd(ProjRoot)

dfr       <- analysisOrganizeData(dataCensored)
df        <- dfr[["df"]]
analySpec <- dfr[["analySpec"]]

load('../CBP_data/1984-2016 seasonally detrended salinity data.rda')
load('../CBP_data/1983-2016 seasonally detrended flow data.rda')
methodsList <- data.frame(stationMethodGroup = 'MD-Potomac',
                          parameter          = 'tss',
                          beginDate          = as.POSIXct('1998-05-01'),
                          intervention       = TRUE,
                          labels             = 'demonstration',
                          stringsAsFactors = FALSE)

stat = 'CB5.4'; dep = 'tn'; layer = 'S';
gamResult <- gamTest(df, dep, stat, layer, analySpec)
                    , salinity.detrended = salinity.detrended
                    , flow.detrended = flow.detrended)
stat = 'TF2.2'; dep = 'tss'; layer = 'S';
gamResult <- gamTest(df, dep, stat, layer, analySpec
                    , salinity.detrended = salinity.detrended
                    , flow.detrended = flow.detrended)
```

## Figure Descriptions

Figures 1 and 2 display the gam2 and gam4 models for surface total nitrogen at CB5.4. Figures 3 and 4 display the gam2 and gam3 models for surface total suspended solids at TF2.2. Figures for gam0 and gam1 have less complexity and their content can be inferred from the discussion related to Figures 1-4.
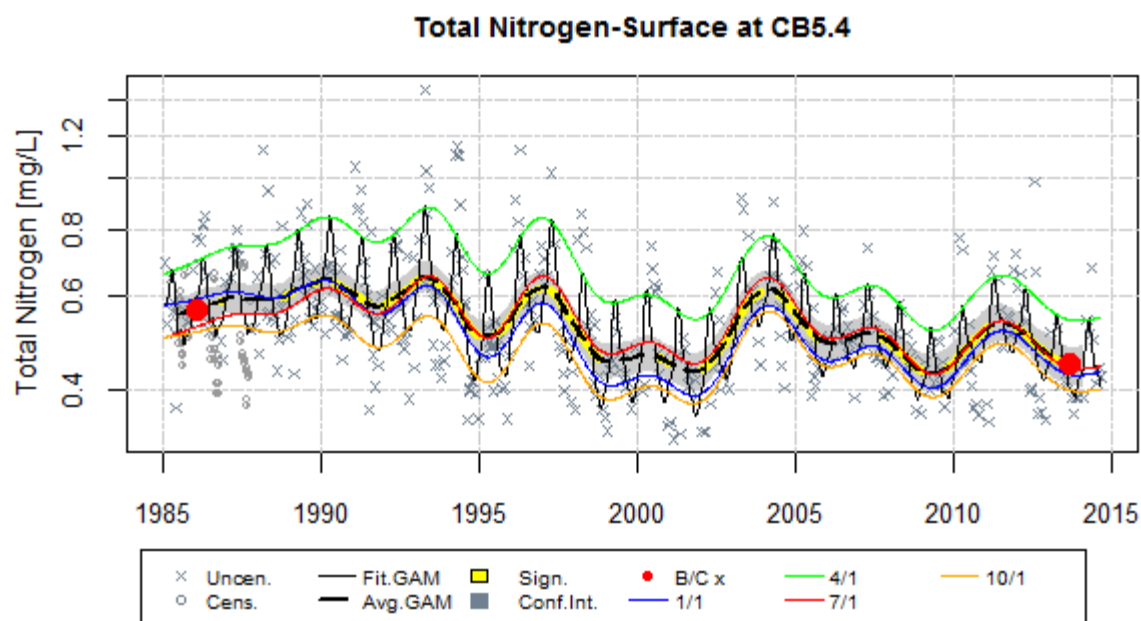
**Observations.** All figures present the observed data using two primary data symbols. Uncensored (detected) observations are presented as "x" while censored observations are presented as "o". Often censored values are thought of as non-detects and are represented as "less-than" the detection limit. The concept of "less-than" works well for direct measures; however, may lead to difficulties for computed variables such as species of nitrogen or phosphorus which are often sums or differences of other variables. Averaging multiple observations can also introduce complexities, i.e., what is the average of "<0.5" and "1"? In baytrends these values are represented as interval censored observations with a lower and upper bound. For example, the average of "<0.5" and "1" is represented as the range 0.5-0.75. Graphically, both the 0.5 and 0.75 are plotted as an "o" with a line connecting the points. This approach adds some additional 'upfront' burden in preparing data for baytrends, but is preferred over

artificial substitution practices such as replacing censored values with zero, one-half detection limit, or detection limit.
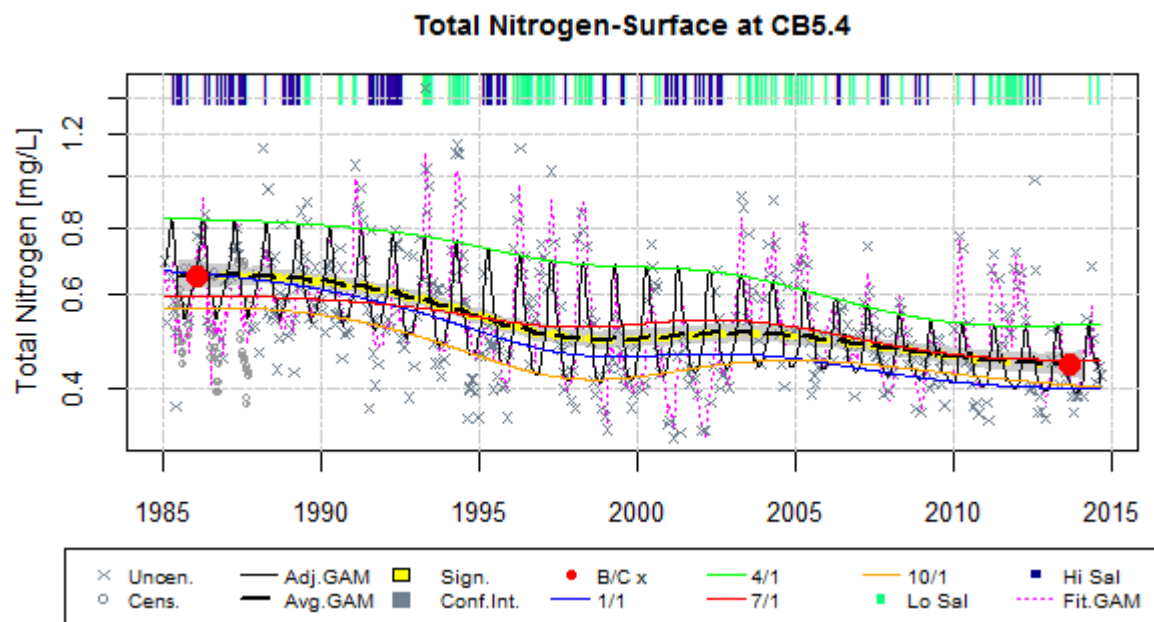
**Non-linear trend with Seas+Int (gam2).** The non-linear trend with seasonality and interaction term (gam2) model is displayed in Figures 1 and 3 for surface total nitrogen at CB5.4 and total suspended solids at TF2.2. Both figures include the fitted GAM (Fit. GAM) as a thin black line. The fitted GAM shows a clear seasonal pattern for both figures. The average GAM (Avg. GAM) is an annually-smoothed concentration based on a one-year sliding window and is shown as a thick black dashed line. The gray shaded region is the 95 percent confidence interval for the average GAM and periods of statistically significant increases or decreases are highlighted in yellow. (By default, the nominal significance level is set to 0.05 but can be set by the user.) Seasonal models representing January 1, April 1, July 1, and October 1 are shown as blue, green, red, and orange lines, respectively. There is a clear distinction between the seasonal models for total nitrogen that is only present for the January 1 model for total suspended solids. The solid red dot represents the average baseline and current conditions. As shown in these figures the average baseline and current conditions are based on the first two and last two years of the modeled record, respectively.

**Non-linear trend with Seas+Int. & Hydro Adj (gam4).** Including a hydrologic term in the total nitrogen model is shown in Figure 2. This model uses salinity as the independent variable used as the hydrology term. Here the fitted gam (Fit.GAM) is shown as a dashed magenta line. The solid black line is referred to as the adjusted GAM (Adj.GAM) and represents the model after removing the hydrology effect. Periods of low salinity (wet conditions) and high salinity (dry conditions) are indicated along the top horizontal axis with salinity values less than the $20^{th}$ percentile and greater than the $80^{th}$ percentile shown as green and blue striping, respectively. (By default, total nitrogen at this station is modeled using salinity; however, the modeler can change the hydrology term to flow.)

**Non-linear trend with Seas+Int. & Intervention (gam3).** The gam3 model includes an intervention term which allows for a data shift that could be caused by changing analytical methods. The effect of including an intervention term is shown in Figure 4. The vertical blue line represents the date when there was a laboratory change and more consistent rinsing techniques were implemented. Visually, there is a decrease in concentration. This figure also includes an adjusted baseline and current condition estimate which is represented by the solid blue dot. As shown in this figure, the blue dot at the beginning of the record represents the modeled baseline condition that would have been observed if the current method could have been retroactively applied to the earlier samples.

**Figure 1. Surface Total Nitrogen at Station CB5.4 modeled using Non-linear trend with Seas+Int (gam2).**



**Figure 2. Surface Total Nitrogen at Station CB5.4 modeled using Non-linear trend with Seas+Int. & Hydro Adj (gam4).**

*Figure 3. Surface Total Suspended Solids at Station TF2.2 modeled using Non-linear trend with Seas+Int (gam2).*



*Figure 4. Surface Total Suspended Solids at Station TF2.2 modeled using Non-linear trend with Seas+Int. & Intervention (gam3).*

## Tabular Output

The gamTest function also outputs a series of tables summarizing the developed model.

- GAM Analysis of Variance
- GAM Parameter Coefficients
- GAM Diagnostics
- Estimates of Change

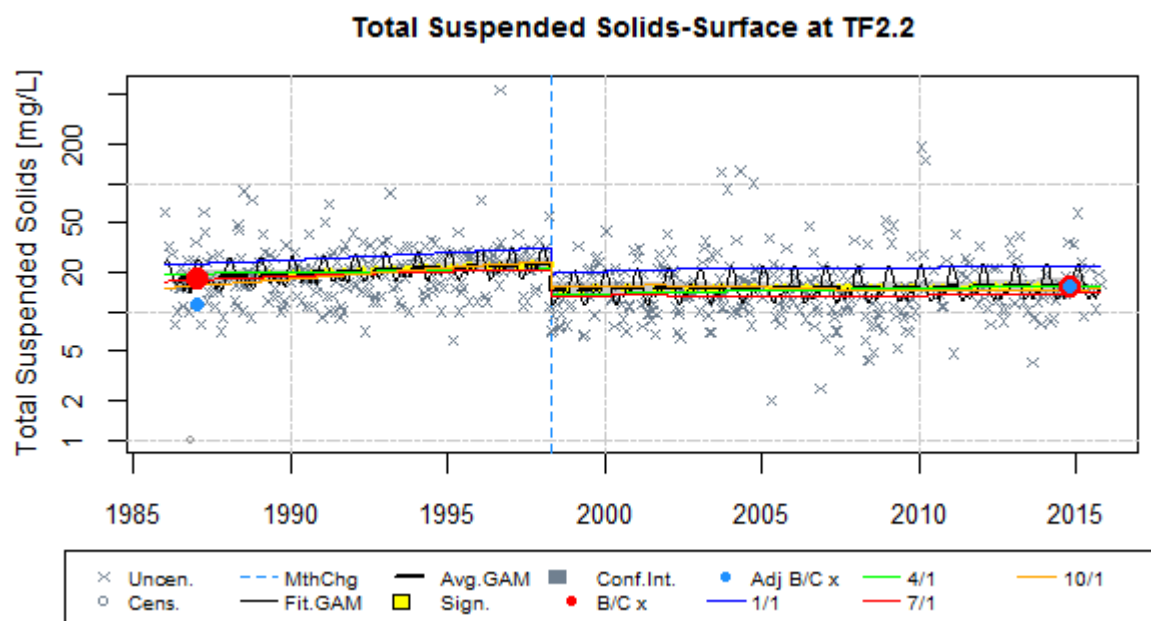Examples of these tables are provided in Table 1 for surface total nitrogen at station CB5.4 for gam2 and gam4. Table 2 displays the results for surface total suspended solids at station TF2.2 for gam2 and gam3.

Inspection of p-values in the analysis of variance and parameter coefficients table can provide insight on which dependent variables significantly contribute to the overall model. As shown in Table 1 (and as one might expect based on inspection of Figures 1 and 2), seasonal (doy) and a hydrologic (flw_sal) terms are statistically significant for modeling surface total nitrogen at station CB5.4. While the seasonal term (doy) is also significant for total suspended solids, it is perhaps more interesting to examine the statistical significance associated with the method change represented by the blue vertical line (intervention) (Table 2). The magnitude of interventionB (-0.46) represents the change in total suspended solids in log-space relative to the intercept term.

***Table 1. Surface Total Nitrogen at Station CB5.4 model diagnostics for gam2 and gam4.***

| Non-linear trend with Seas+Int (gam2) (see Figure 1) | Non-linear trend with Seas+Int. & Hydro Adj (gam4) (see Figure 2) |
|---|---|

**Table: GAM Analysis of Variance.**

| Type | Source | edf | F-stat | p-value |
|---|---|---|---|---|
| parametric terms | cyear | 1.00 | 0.5286 | 0.4676 |
| smoothed terms | s(cyear) | 17.31 | 6.3767 | <0.0001 |
| " " | s(doy) | 6.93 | 23.9563 | <0.0001 |
| " " | ti(cyear,doy) | 5.38 | 1.7958 | 0.0001 |

**Table: GAM Analysis of Variance.**

| Type | Source | edf | F-stat | p-value |
|---|---|---|---|---|
| parametric terms | cyear | 1.00 | 1.8602 | 0.1734 |
| smoothed terms | s(cyear) | 5.09 | 6.6042 | <0.0001 |
| " " | s(doy) | 7.10 | 35.3265 | <0.0001 |
| " " | ti(cyear,doy) | 5.44 | 3.1072 | <0.0001 |
| " " | s(flw_sal) | 2.35 | 72.9011 | <0.0001 |
| " " | ti(flw_sal,doy) | 4.72 | 4.8515 | <0.0001 |
| " " | ti(flw_sal,cyear) | 5.53 | 3.1203 | 0.0032 |
| " " | ti(flw_sal,doy,cyear) | 2.76 | 0.2089 | 0.0028 |

**Table: GAM Parameter Coefficients.**

| Parameter | Estimate | Std. Err. | t value | p-value |
|---|---|---|---|---|
| (Intercept) | -0.528750 | 0.126936 | -4.1655 | <0.0001 |
| cyear | 0.044284 | 0.060909 | 0.7270 | 0.4676 |

**Table: GAM Parameter Coefficients.**

| Parameter | Estimate | Std. Err. | t value | p-value |
|---|---|---|---|---|
| (Intercept) | -0.583830 | 0.022755 | -25.6570 | <0.0001 |
| cyear | 0.015011 | 0.011006 | 1.3639 | 0.1734 |

**Table: GAM Diagnostics.**

| AIC | RMSE | Adj. R-squared |
|---|---|---|
| -181.49 | 0.1892 | 0.4768 |

**Table: GAM Diagnostics.**

| AIC | RMSE | Adj. R-squared |
|---|---|---|
| -420.05 | 0.143 | 0.7004 |

**Table: Estimates of Change from 1985-2014.**

| Calculation | Estimate |
|---|---|
| Baseline log mean (geometric mean) | -0.5669 (0.5673) |
| Current log mean (geometric mean) | -0.8093 (0.4452) |
| Estimated log difference | -0.2424 |
| Std. Err. log difference | 0.0543 |
| 95% Confidence interval for log difference | (-0.3488 , -0.136) |
| Difference p-value | <0.0001 |
| Period of Record Percent Change Estimate (%) | -21.53% |

**Table: Estimates of Change from 1985-2014.**

| Calculation | Estimate |
|---|---|
| Baseline log mean (geometric mean) | -0.4301 (0.6504) |
| Current log mean (geometric mean) | -0.8169 (0.4418) |
| Estimated log difference | -0.3868 |
| Std. Err. log difference | 0.0443 |
| 95% Confidence interval for log difference | (-0.4737 , -0.2999) |
| Difference p-value | <0.0001 |
| Period of Record Percent Change Estimate (%) | -32.08% |

*Table 2. Surface Total Suspended Solids at Station TF2.2 model diagnostics for gam3.*

## Total Suspended Solids - Non-linear trend with Seas+Int (gam2) (see Figure 3)

### Table: GAM Analysis of Variance.

| Type | Source | edf | F-stat | p-value |
|---|---|---|---|---|
| parametric terms | cyear | 1.00 | 0.5585 | 0.4552 |
| smoothed terms | s(cyear) | 11.94 | 2.5671 | 0.0014 |
| " " | s(doy) | 6.73 | 5.6356 | <0.0001 |
| " " | ti(cyear,doy) | 3.43 | 0.6734 | 0.0220 |

### Table: GAM Parameter Coefficients.

| Parameter | Estimate | Std. Err. | t value | p-value |
|---|---|---|---|---|
| (Intercept) | 2.898533 | 0.108590 | 26.6924 | <0.0001 |
| cyear | 0.085368 | 0.114228 | 0.7473 | 0.4552 |

### Table: GAM Diagnostics.

| AIC | RMSE | Adj. R-squared |
|---|---|---|
| 895.57 | 0.5498 | 0.1617 |

### Table: Estimates of Change from 1986-2015.

| Calculation | Estimate |
|---|---|
| Baseline log mean (geometric mean) | 2.881 (17.8326) |
| Current log mean (geometric mean) | 2.8255 (16.8693) |
| Estimated log difference | -0.0555 |
| Std. Err. log difference | 0.1445 |
| 95% Confidence interval for log difference | (-0.3387 , 0.2277) |
| Difference p-value | 0.7009 |
| Period of Record Percent Change Estimate (%) | -5.4% |

## Non-linear trend with Seas+Int. & Intervention (gam3) (see Figure 4)

### Table: GAM Analysis of Variance.

| Type | Source | edf | F-stat | p-value |
|---|---|---|---|---|
| parametric terms | intervention | 1.00 | 16.8562 | <0.0001 |
| " " | cyear | 1.00 | 3.7880 | 0.0522 |
| smoothed terms | s(cyear) | 1.48 | 4.9537 | 0.0171 |
| " " | s(doy) | 6.67 | 5.2260 | <0.0001 |
| " " | ti(cyear,doy) | 3.32 | 0.6419 | 0.0235 |

### Table: GAM Parameter Coefficients.

| Parameter | Estimate | Std. Err. | t value | p-value | Int Chg p-val | Int Chg est |
|---|---|---|---|---|---|---|
| (Intercept) | 3.050024 | 0.068317 | 44.6452 | <0.0001 | - | NA |
| interventionB | -0.460427 | 0.112145 | -4.1056 | <0.0001 | <0.0001 | -0.460427 |
| cyear | -0.025664 | 0.013186 | -1.9463 | 0.0522 | - | NA |

### Table: GAM Diagnostics.

| AIC | RMSE | Adj. R-squared |
|---|---|---|
| 896.47 | 0.5551 | 0.1454 |

### Table: Estimates of Change from 1986-2015.

| Calculation | Estimate | Adj. Estimate |
|---|---|---|
| Baseline log mean (geometric mean) | 2.9052 (18.2689) | 2.4448 (11.5279) |
| Current log mean (geometric mean) | 2.7611 (15.8179) | 2.7611 (15.8179) |
| Estimated log difference | -0.1441 | 0.3164 |
| Std. Err. log difference | 0.1051 | 0.1681 |
| 95% Confidence interval for log difference | (-0.3501 , 0.062) | (-0.0132 , 0.6459) |
| Difference p-value | 0.1712 | 0.0604 |
| Period of Record Percent Change Estimate (%) | -13.42% | 37.21% |

## Production Analyses

To maximize the efficiency in creating reports that involve multiple stations and parameters, baytrends was developed with the intention that users would incorporate baytrends in looping structures to direct baytrends on which combinations of stations and parameters to evaluate. In Code Chunk 3, the names of two comma delimited files are established for storage of tabular output (i.e., information from Table 1).

The list, analySpec, returned from analysisOrganizeData contains data frames of the dependent variables, stations, and layers. The variables, deps, layers, and stations, are extracted from analySpec and used as lists with for loops. The if structure inside the for loops is used to compile tabular output from each call of the gamTest function. After completing the for loops, the variables, stat.gam and chng.gam are written to the comma delimited files. The gamDiffModel argument specifies which gam models are used for evaluating changes that are stored in the variable chng.gam.

*Code Chunk 3. Example code chunk that uses for loops to evaluate multiple stations and parameters in specified order. Selected data are saved to comma delimited files. (Gray lines of code previously discussed.)*

```
library(baytrends)
ProjRoot <-  "E:/Dropbox/CBP/cbpTrends/work10"
setwd(ProjRoot)

statFile    <- file.path(ProjRoot,'statGAM_R.csv')
chngFile    <- file.path(ProjRoot,'chngGAM_R.csv')

dfr      <- analysisOrganizeData(dataCensored)
df       <- dfr[["df"]]
analySpec <- dfr[["analySpec"]]

deps     <- analySpec$depVarList$deps
stations <- analySpec$stationList$stations
layers   <- analySpec$layerList$layers

for (dep in deps)    {
  for (layer in layers) {
    for (stat in stations) {
      gamResult <- gamTest(df, dep, stat, layer, analySpec
                        , gamDiffModel = c(0,1,2))
      if (!is.na(gamResult[1])) {
        if(!exists("stat.gam")) stat.gam <- gamResult[["stat.gam.result"]] else
          stat.gam <- rbind(stat.gam,gamResult[["stat.gam.result"]])
        if(!exists("chng.gam")) chng.gam <- gamResult[["chng.gam.result"]] else
          chng.gam <- rbind(chng.gam,gamResult[["chng.gam.result"]])
      }
    }
  }
}

write.csv(stat.gam,file=statFile, row.names = FALSE)
write.csv(chng.gam,file=chngFile, row.names = FALSE)
```

## User defined support tables

Thus far, we have relied on built-in support and data files. Code Chunk 4 demonstrates one approach for loading user-defined support and data files. The baytrends package includes a loadExcel function that will read data from a spreadsheet. Users can contact the package maintainer (see package documentation for email information) for the example spreadsheet. The user can also study the structure of the built-in files and create their own files using an alternative approach.

***Code Chunk 4. Example code chunk that loads user-defined support files and data file. (Gray lines of code previously discussed.)***

```
library(baytrends)
ProjRoot <-  "E:/Dropbox/CBP/cbpTrends/work10"
setwd(ProjRoot)

statFile     <- file.path(ProjRoot,'statGAM_R.csv')
chngFile     <- file.path(ProjRoot,'chngGAM_R.csv')


fname <- "CBPstatsLookupTables(v1.0.x).xlsm"
stationMasterList <- loadExcel(folder='../CBP_data', file=fname,
                                sheet='stationMasterList', pk='station')
parameterList     <- loadExcel(folder='../CBP_data', file=fname,
                                sheet='parameterList',      pk='parm')
layerLukup        <- loadExcel(folder='../CBP_data', file=fname,
                                sheet='layerLukup',         pk='layers')


load('../CBP_data/cbpDataAll_04_QW_baytrends.rda')
dfr        <- analysisOrganizeData(cbpDataAll_04_QW
                                   , stationMasterList = stationMasterList
                                   , parameterList =  parameterList
                                   , layerLukup = layerLukup)
df         <- dfr[["df"]]
analySpec <- dfr[["analySpec"]]

deps       <- analySpec$depVarList$deps
stations  <- analySpec$stationList$stations
layers     <- analySpec$layerList$layers

for (dep in deps)    {
  for (layer in layers) {
    for (stat in stations) {
      gamResult <- gamTest(df, dep, stat, layer, analySpec, gamDiffModel=c(0,1,2))
      if (!is.na(gamResult[1])) {
        if(!exists("stat.gam")) stat.gam <- gamResult[["stat.gam.result"]] else
          stat.gam <- rbind(stat.gam,gamResult[["stat.gam.result"]])
        if(!exists("chng.gam")) chng.gam <- gamResult[["chng.gam.result"]] else
          chng.gam <- rbind(chng.gam,gamResult[["chng.gam.result"]])
      }
    }
  }
}

write.csv(stat.gam,file=statFile, row.names = FALSE)
write.csv(chng.gam,file=chngFile, row.names = FALSE)
```

The sample data set (dataCensored) included with baytrends includes selected 1985-2015 data for eight (8) stations from the Chesapeake Bay Monitoring Program. In this example, additional data are accessed by loading **cbpDataAll_04_QW_baytrends.rda** which is available from the package maintainer. Water quality variables are stored as class qw that allows for left- and interval-censored data. The reader is referred to the QW vignette for information on creating and manipulating variables stored as class qw. (Also, note that the argument for analysisOrganizeData needs to reflect the name of the data frame resulting from the load statement and the user-supplied stationMasterList, parameterList, and layerLukup are passed as arguments in the analysisOrganizeData function.)

## R Markdown

R Markdown enables directing tabular and graphical output to Microsoft Word files for further editing. Code Chunk 5 presents an R markdown (Rmd) script that would generate the same output as Code Chunk 3, except directed to an MS Word file. Unlike previous code chunks that would be stored as R scripts (.R files), Rmd scripts are stored as .Rmd files. The first portion of Code Chunk 5 (between the "---") is referred to as YAML (YAML Ain't Markup Language). The remainder of the code chunk is R markdown which is a format for writing reproducible, dynamic reports. In general, the code includes written content embedded with R code that can be rendered to a document. Syntax such as "#" and "###" will result in 1st and 3rd level headers in the MS Word document. Sets of three back ticks (```) signifies locations of R code.

We use a template MS word file, template.docx, that includes predefined styles for headers and font sizes. The design of the template.docx was made to best organize the output into a readable document. One trick associated with the template.docx is that the 2nd level header was defined to include a page break before the header. This feature helps with keeping the pages in a readable order when flipping through the pages. It is also important that template.docx be placed in the working directory. A second trick is the inclusion of two functions, .H2 and .H3, that can be integrated with the embedded code to force headers into the outputted MS word document. More information about YAML and Rmd files can be found here: https://www.rstudio.com/wp-content/uploads/2015/02/rmarkdown-cheatsheet.pdf

***Code Chunk 5. Example Rmd script for outputting to MS Word. (Gray lines of code previously discussed.)***

```
---
title: "General Additive Model Analysis"
output:
  word_document:
    fig_caption: yes
    fig_height: 3.5
    fig_width: 6.5
    reference_docx: template.docx
  html_document:
    fig_caption: yes
    fig_height: 3.5
    fig_width: 6.5
    keep_md: yes
    toc: yes
---

**Date: `r format(Sys.time(), "%B %d, %Y %X") # Month DD, YYYY H:MM:SS PM`
`r begin <- Sys.time()`**

# Initialization

### Load Library and Project Folder

```{r loadLibrary, results='hide'}
rm(list=ls())   # clear the global environment
cat("\014")     # clear the console
dev.off()       # clear any plots
```

```
library(baytrends)
ProjRoot <-   "E:/Dropbox/CBP/cbpTrends/work10"
setwd(ProjRoot)

statFile      <- file.path(ProjRoot,'statGAM_R.csv')
chngFile      <- file.path(ProjRoot,'chngGAM_R.csv')
```


### Prepare Data and Analysis Specifications

```{r dataPrep, results='asis', echo=FALSE}
dfr        <- analysisOrganizeData(dataCensored, reports = NA)
df         <- dfr[["df"]]
analySpec <- dfr[["analySpec"]]
depVarList<- analySpec$depVarList
layerList <- analySpec$layerList
deps       <- analySpec$depVarList$deps
layers     <- analySpec$layerList$layers
stations   <- analySpec$stationList$stations
```


## GAM Analysis

```{r stations_parameters, results='asis', echo=FALSE}

for (dep in deps)    {
  for (layer in layers) {
    .H2(paste(depVarList[deps==dep,"parmName"] , "--", layerList[layers==layer,"name"]))
    for (stat in stations) {
      .H3(stat)
      gamResult <- gamTest(df, dep, stat, layer, analySpec, gamDiffModel=c(0,1,2))
      if (!is.na(gamResult[1])) {
        if(!exists("stat.gam")) stat.gam <- gamResult[["stat.gam.result"]] else
          stat.gam <- rbind(stat.gam,gamResult[["stat.gam.result"]])
        if(!exists("chng.gam")) chng.gam <- gamResult[["chng.gam.result"]] else
          chng.gam <- rbind(chng.gam,gamResult[["chng.gam.result"]])
      }
    }
  }
}

write.csv(stat.gam,file=statFile, row.names = FALSE)
write.csv(chng.gam,file=chngFile, row.names = FALSE)
```
```