

C3M4_autograded

April 1, 2022

1 C3M4 Autograded Assignment

1.0.1 Outline:

Here are the objectives of this assignment:

1. Differentiate between additive and multiplicative models.
2. Apply GAMs to different types of responses, including Binary and Poisson.
3. Explore how effective degrees of freedom relates to parameter linearity.
4. Evaluate the performance of GAM models.

Here are some general tips:

1. Read the questions carefully to understand what is being asked.
2. When you feel that your work is completed, feel free to hit the **Validate** button to see your results on the *visible* unit tests. If you have questions about unit testing, please refer to the “Module 0: Introduction” notebook provided as an optional resource for this course. In this assignment, there are hidden unit tests that check your code. You will not receive any feedback for failed hidden unit tests until the assignment is submitted. **Do not misinterpret the feedback from visible unit tests as all possible tests for a given question—write your code carefully!**
3. Before submitting, we recommend restarting the kernel and running all the cells in order that they appear to make sure that there are no additional bugs in your code.

```
[1]: # Load Required Libraries
library(testthat)
library(tidyverse)
library(mgcv) # For GAM models
library(ggplot2)
```

```
Attaching packages: tidyverse
1.3.0

ggplot2 3.3.0    purrr   0.3.4
tibble  3.0.1    dplyr   0.8.5
tidyr   1.0.2    stringr 1.4.0
readr   1.3.1    forcats 0.5.0
```

Conflicts

```
tidyverse_conflicts()
  dplyr::filter() masks stats::filter()
  purrr::is_null() masks
testthat::is_null()
  dplyr::lag() masks stats::lag()
  dplyr::matches() masks
tidyr::matches(), testthat::matches()
```

Loading required package: nlme

Attaching package: 'nlme'

The following object is masked from 'package:dplyr':

collapse

This is mgcv 1.8-31. For overview type 'help("mgcv-package")'.

2 Problem 1: Identifying Additive Models

GAMs are specified to be the sum of separate functions of predictors. In math terms, that means

$$Y_i = f(x_{i,1} + \dots + f_{i,p}) + \epsilon_i = f_1(x_{i,1}) + \dots + f_p(x_{i,p}) + \epsilon_i$$

. For the following models, answer TRUE if the model can be expressed as a GAM and FALSE if it can not.

1. $f(x_1, x_2, x_3) = x_1 + x_2 + (x_1 x_2)^2 + \sqrt{x_3}$
2. $f(x_1, x_2, x_3) = \frac{x_1 x_2 x_3 + x_2 x_3 + x_3}{x_2 x_3}$
3. $f(x_1, x_2) = \log(\sqrt{x_1 x_2})$
4. $f(x_1, x_2) = \beta_0 + \exp(x_1 x_2)$
5. $f(x_1) = 0$

[6]: *# Code your answers as TRUE or FALSE.*

```
prob.1.1 = FALSE
```

```
prob.1.2 = TRUE
```

```
prob.1.3 = TRUE
```

```
prob.1.4 = FALSE
```

```
prob.1.5 = TRUE

# your code here
```

```
[ ]: # Test Cell
# Make sure your answers are booleans!
# This cell has hidden test cases that will run after submission.
```

```
[ ]: # Test Cell
# This cell has hidden test cases that will run after submission.
```

```
[ ]: # Test Cell
# This cell has hidden test cases that will run after submission.
```

```
[ ]: # Test Cell
# This cell has hidden test cases that will run after submission.
```

```
[ ]: # Test Cell
# This cell has hidden test cases that will run after submission.
```

3 Problem 2: GAMs with Binary Response

In 1911, the Titanic sailed off on its maiden voyage from Southampton, on its way towards New York City. Unfortunately, the ship would eventually collide with an iceberg and sink to the bottom of the ocean. As the ship was sinking, it is said that lifeboats (and floating doors) were prioritized for women and children, and thus many of them were able to survive until rescue arrived. Although a bit grizzly, it does pose an interesting statistical question. If we have the list of passengers, can we predict who will survive the sinking of the Titanic?

Turns out GAMs can be used for different kinds of response as well, including Binary responses. That means we can use GAMs to try to answer our question. We load in the data below. It may help to do some basic data analysis before doing too much modeling.

```
[7]: # Load in the Data
titanic = read.csv("titanic.csv", sep=",")
head(titanic)
```

A data.frame: 6 × 12

	PassengerId	Survived	Pclass	Name
	<int>	<int>	<int>	<fct>
1	1	0	3	Braund, Mr. Owen Harris
2	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Thayer)
3	3	1	3	Heikkinen, Miss. Laina
4	4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)
5	5	0	3	Allen, Mr. William Henry
6	6	0	3	Moran, Mr. James

3.0.1 2. (a) Data Cleaning

Before we get to work, we need to clean up this data a bit. If you look over it, you will notice a fair number of columns, many missing values and some NA values. We're going to need to clean our data before we're able to do any modelling. When you're done, you will have a cleaned dataset `titanic`, your training set `titanic.train` and your test set `titanic.test`. Here's what needs to be done:

- We have a lot of predictors, but we don't need them all. Restrict the data to on the `Survived`, `Pclass`, `Sex`, `Age` and `Fare` columns.
- There's still missing data in our dataframe. That won't do. Remove any rows that have at least one missing value in any column.
- If you look at the types of each column, you'll notice that some factors have been loaded as numeric. We should change that. Set `Survived` and `Pclass` to categorical.
- We will eventually want to analyze how well our model performs. Split the data into training and test sets. Do this by putting every fifth row into the test set, and use the rest for training. For example, the first 4 rows will be in the training set, and the 5th row will be in the test set. Repeat that pattern for the rest of the data.

```
[17]: titanic.train = NA
titanic.test = NA

# your code here
titanic = titanic %>%
  select(c('Survived', 'Pclass', 'Sex', 'Age', 'Fare')) %>%
  drop_na() %>%
  mutate(Survived = as.factor(Survived), Pclass = as.factor(Pclass))
```

```
[26]: index = seq(0, nrow(titanic), 5)
titanic.train = titanic[-index, ]
titanic.test = titanic[index, ]
```

```
[27]: nrow(titanic.train)
```

572

```
[28]: # Test Cell
# This cell has hidden test cases that will run after submission.
if(!test_that("Checking DataFrame Size", {expect_equal(nrow(titanic), 714)
  expect_equal(nrow(titanic.train), 572)})){
  print("Incorrect Dataset sizes. Make sure these are correct, or else your
  modelling could be incorrect.")
}
```

3.0.2 2. (b) Fit Your GAM

Now that our dataset is clean, we can fit our model. Fit your GAM as `titanic.gam` with `Survived` as your response and all other values as predictors. Make sure to smooth the necessary predictors!

Hint: The response is binary, so you will need to tell your model to expect that. How did we do that when we were fitting GLMs?

Look at the summary for your GAM model. Where any predictors insignificant? Save the string name of any/all insignificant predictors into the list `insig.predictors`.

```
[42]: titanic.gam = NA
      insig.predictors = c()

      # your code here
      titanic.gam = gam(Survived ~ Pclass + Sex + s(Age) + s(Fare), data=titanic,
        ↪train, family = binomial)
      summary(titanic.gam)
      insig.predictors = 'Fare'
```

Family: binomial

Link function: logit

Formula:

Survived ~ Pclass + Sex + s(Age) + s(Fare)

Parametric coefficients:

	Estimate	Std. Error	z value	Pr(> z)	
(Intercept)	2.6562	0.3394	7.826	5.03e-15	***
Pclass2	-1.4661	0.3756	-3.904	9.47e-05	***
Pclass3	-2.4575	0.4024	-6.107	1.01e-09	***
Sexmale	-2.4886	0.2315	-10.752	< 2e-16	***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Approximate significance of smooth terms:

	edf	Ref.df	Chi.sq	p-value	
s(Age)	3.995	4.959	19.28	0.00154	**
s(Fare)	2.024	2.524	1.26	0.52078	

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

R-sq.(adj) = 0.387 Deviance explained = 31.8%

UBRE = -0.043402 Scale est. = 1 n = 572

```
[ ]: # Test Cell
      # This cell has hidden test cases that will run after submission.
```

```
[ ]: # Test Cell
# This cell has hidden test cases that will run after submission.
```

3.0.3 2. (c) Effective Degrees of Freedom

Let's take a look at our continuous predictors and see if they appear linearly in our GAM model. There are two ways of doing this:

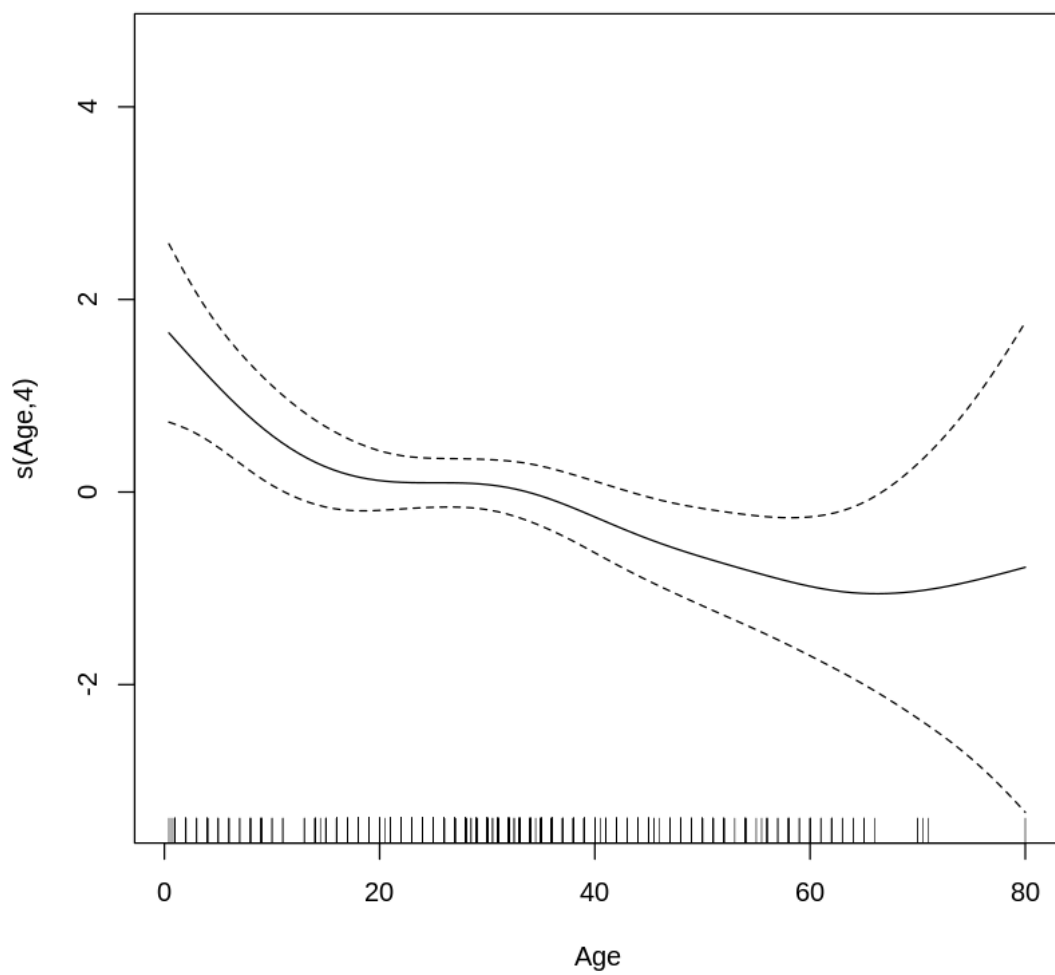
- Use the `plot.gam()` function to plot the curves of your continuous predictors.
- Look at the Effective Degrees of Freedom for the continuous variables.

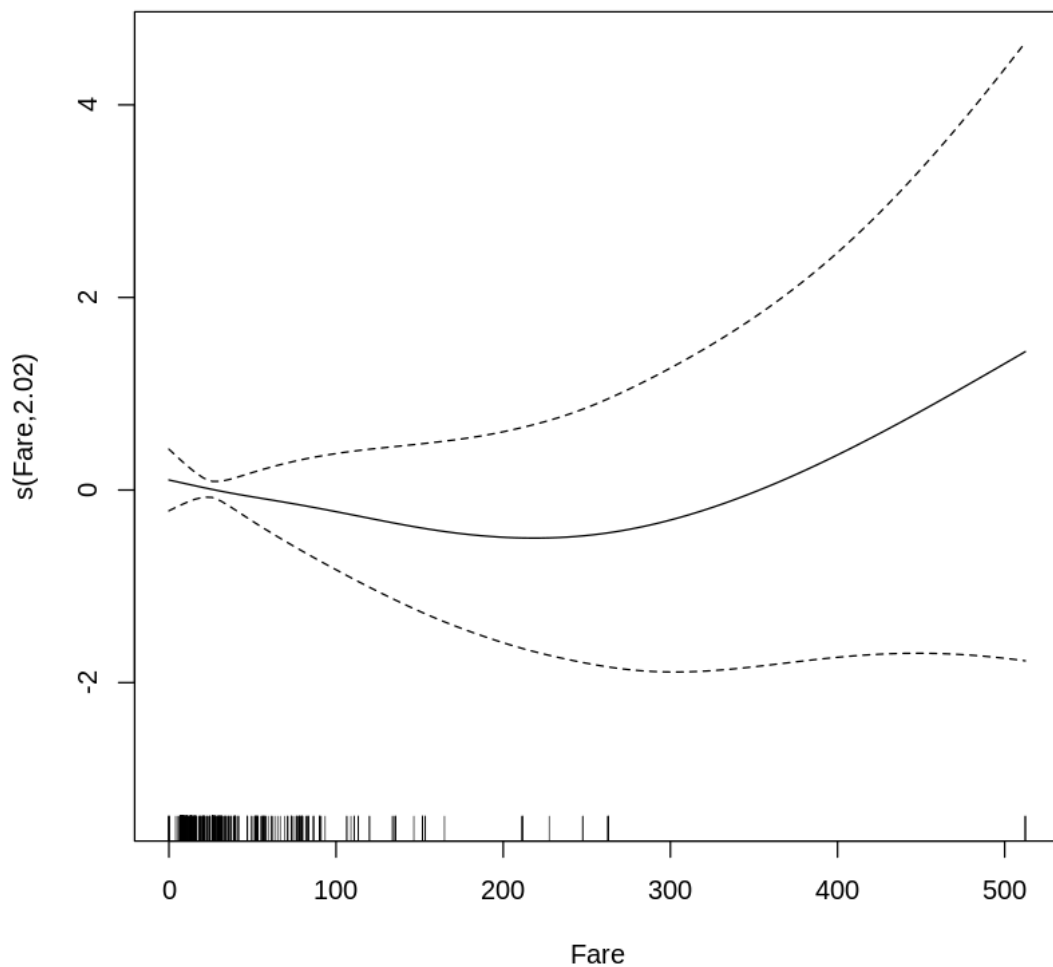
After conducting these analysis, determine whether each continuous predictor is linear or not. Remember, in statistical terms, a “smooth term” is linear if you can draw a line through the 95% confidence band.

Save your answer as `TRUE` if it is linear and `FALSE` if it is not. Use `age.is.linear` for `Age` and `Fare.is.linear` for `Fare`.

```
[46]: age.is.linear = NA
fare.is.linear = NA

# your code here
plot.gam(titanic.gam)
age.is.linear = TRUE
fare.is.linear = TRUE
```





```
[ ]: # Test Cell
      # This cell has hidden test cases that will run after submission.
```

```
[ ]: # Test Cell
      # This cell has hidden test cases that will run after submission.
```

3.0.4 2. (d) Predicting with GAMs

Let's use our Test set to determine how well our model performs on new data. Predict the **Survived** values for the data in your test set and compute the accuracy, precision, recall and F1 score. Save these values as `gam.acc`, `gam.prec`, `gam.rec` and `gam.f1`.

How well did the model do?


```
[66]: gam.acc = NA
gam.prec = NA
gam.rec = NA
gam.f1 = NA

# your code here
predictions = predict(titanic.gam, newdata=titanic.test, type="response")
predictions = factor(ifelse(predictions > 0.5, 1, 0))
table(predictions, titanic.test$Survived)
```

```
predictions  0  1
             0 69  5
             1 15 53
```

```
[71]: gam.acc = (53 + 69) / length(predictions)
gam.prec = 53 / (53 + 15)
gam.rec = 53 / (53 + 5)
gam.f1 = 2 * ((gam.prec * gam.rec) / (gam.prec + gam.rec))
gam.acc
gam.prec
gam.rec
gam.f1
```

```
0.859154929577465
```

```
0.779411764705882
```

```
0.913793103448276
```

```
0.841269841269841
```

```
[ ]: # Test Cell
# This cell has hidden test cases that will run after submission.
```