

# A tutorial on Generalised Additive Mixed Effects Models for bilingualism research

## Motivation

...

## Generalised additive (mixed) models

...

## Pre-requisites

This tutorial assumes readers are already familiar with R and have at least some experience fitting linear models, including models with random effects (variably known as mixed, hierarchical, nested, multi-level models. The following packages need to be installed:

- tidyverse
- mgcv
- tidygam

You can find the tutorial code and data here: [XXX](#)

## Case study 1: U-shaped learning

### U-shaped learning

[Lit review]

## The data

We have simulated data of learning scores from 200 subjects, taken at 10 time points. A proficiency score was also included for each subject at each time point. The data is intended to simulate a study in which participants perform a learning task and take a proficiency test at 10 consecutive time points. This is what the data looks like.

```
dat1
```

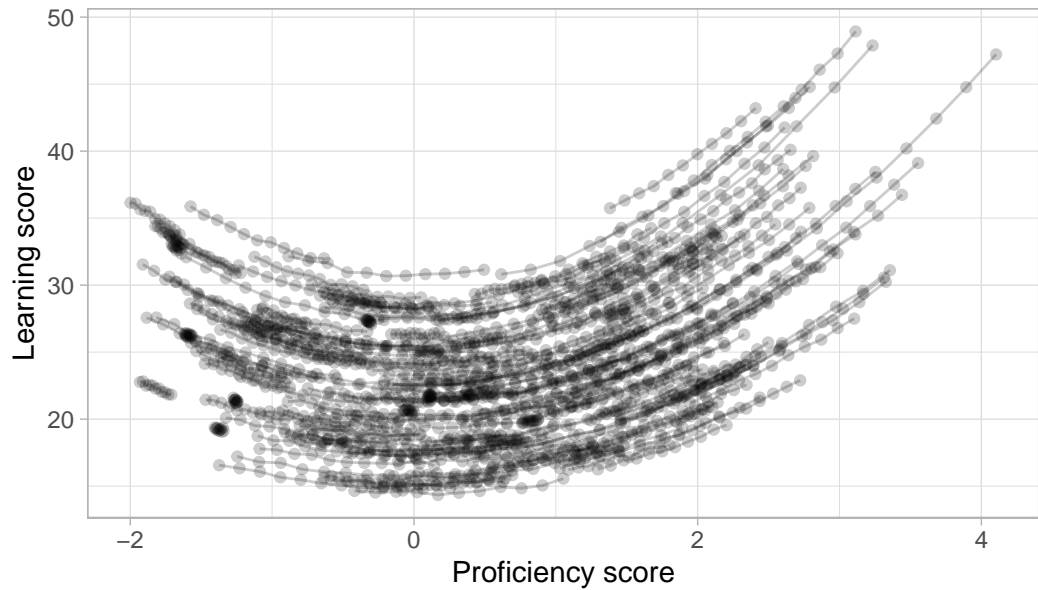
```
# A tibble: 2,200 x 4
  score proficiency subj   time
  <dbl>         <dbl> <fct> <int>
1  26.7         -1.37  s1     0
2  26.4         -1.29  s1     1
3  26.0         -1.21  s1     2
4  25.7         -1.14  s1     3
5  25.5         -1.06  s1     4
6  25.2         -0.983 s1     5
7  25.0         -0.906 s1     6
8  24.9         -0.830 s1     7
9  24.6         -0.753 s1     8
10 24.5         -0.677 s1     9
# i 2,190 more rows
```

The **score** column contains the learning scores, while the **proficiency** column the proficiency scores. The subject ID is given in **subj**. time point 0 to 9 is in **time**. Figure 1 shows the relationship between proficiency and learning scores for individual subjects. From the figure, it is clear that such relationship has a U-shape, by which learning scores initially decrease, plateau, then increase again with higher proficiency.

In the following sections we will analyse this data using GAMMs. For pedagogical purposes, we will first focus on the effect of proficiency scores on learning scores at a single time point, time point 5. Subsequently, we will analyse the entire data, to illustrate how to conduct a time-series analysis.

## Modelling a non-linear effect

Let's first focus on how to model a non-linear effect: here, we can look at the effect of the proficiency score on the learning score. As we have seen in Figure 1, the effect is U-shaped, i.e. it is not linear. To simplify things, let's look only at data from time point. We will extend the analysis to also include time point as a predictor later. We can model a non-linear effect of proficiency on the learning score with the following code.



Connected points are measurements that belong to a single subject, taken at each of the 10 time points.

Figure 1: ?(caption)

```
dat15 <- filter(dat1, time == 5)

# attach the mgcv package
library(mgcv)

# fit the model
gam_1 <- gam(
  score ~ s(proficiency),
  data = dat15
)
```

The formula states that `score`, the outcome variable (also known as the dependent variable) should be modelled as a function of `proficiency`, but we use the `s()` to indicate that we want to estimate a (potentially) non-linear effect. The name of the function, `s`, stands for “smooth term”. Smooth terms (aka smoothers) are mathematical objects that allow GAMs to fit non-linear effects. A detailed treatment of smoothers is beyond the scope of this tutorial, so we refer the readers to XXX.

Before looking at the summary of the `gam_1` model, let’s plot the predicted effect of proficiency. We will use the `tidygam` package, which provides users with utility functions that make extracting and plotting predictions from GAMs easier.

```
# attach tidygam
library(tidygam)

# extract model predictions
gam_1_preds <- predict_gam(gam_1)

# plot predictions
plot(gam_1_preds, series = "proficiency")
```

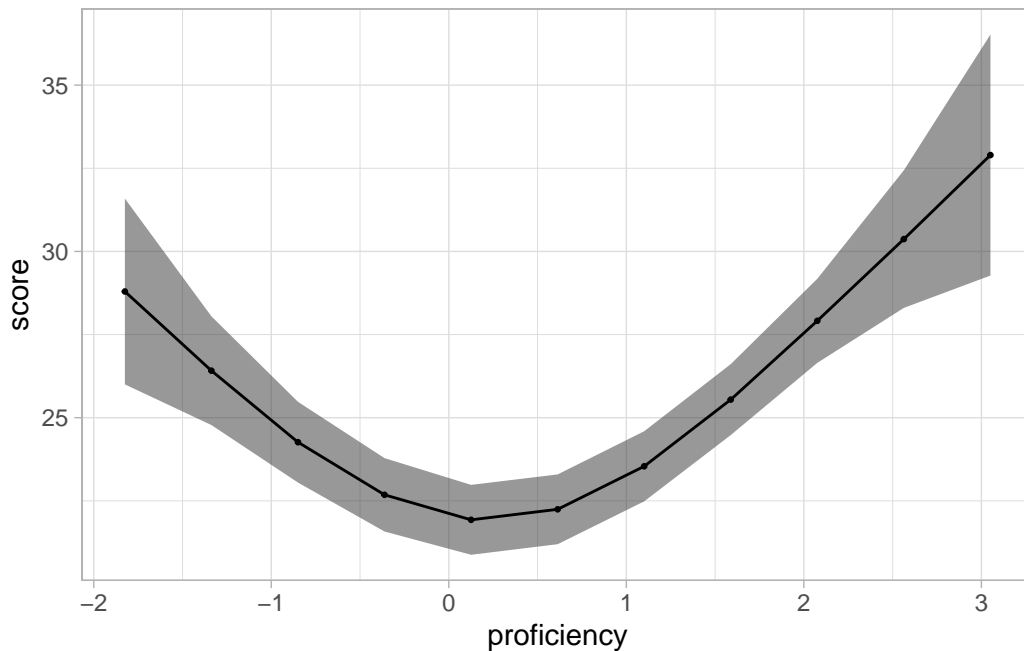


Figure 2: ?(caption)

Figure 2 is a plot of the predicted effect of proficiency on learning score, based on the `gam_1` model from above. We will look into the details of `predict_gam()`. For now, just notice that the learning score initially decreases with increasing proficiency until about 0.5 and then starts increasing, thus producing the typical U-shaped curve. Let's inspect the model summary now.

```
summary(gam_1)
```

```
Family: gaussian
Link function: identity
```

```

Formula:
score ~ s(proficiency)

Parametric coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)   24.559      0.331   74.19  <2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Approximate significance of smooth terms:
              edf Ref.df      F p-value
s(proficiency) 2.922  3.675 16.46  <2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

R-sq.(adj) =  0.242   Deviance explained = 25.3%
GCV = 22.352   Scale est. = 21.914      n = 200

```

The relevant parts of the summary for the time being are the **Parametric coefficients:** table and the **Approximate significance of smooth terms:** table. The first table contains the estimate of the intercept. This intercept is the same intercept you would get in a linear model: here, the estimate of the intercept is the predicted learning score when proficiency is 0. According to the summary, when proficiency is 0, the learning score is about 25. But what we are really interested in is the effect of proficiency on learning score, rather than the intercept per se.

Information on the effect of proficiency is found in the second table, which contains estimates of the smooth terms. Alas, these estimates don't say much about the effect per se. Rather, they just indicate if the effect is linear or not. More specifically, the **edf** estimate, which stands for Estimated Degrees of Freedom, would be 1 for perfectly linear effects and greater than 1 for non-linear effects. This number tells us nothing about the shape of the effect. The only way to assess this is to plot the model predictions, as we have done above.

In our **gam\_1** model, the **edf** for the smooth term over proficiency is 2.9. This is significantly different from 1 ( $p < 2e-16$ ), thus suggesting that the effect of proficiency is not linear (the **Ref.df**, reference degrees of freedom, and **F** values have the only function of being needed to get a  $p$ -value).

If we were to report this model and results, we could write something along the following lines: We fitted a Generalised Additive Model to learning score, with a smooth term over proficiency (to model non-linear effects). According to the model, the effect of proficiency is significantly non-linear ( $F = 16.46$ ,  $p < 0.0001$ ). Based on the prediction plot, we observe that

at lower proficiency, learning scores initially decrease and then start increasing again from about proficiency 0.5.

In the following section we will refit the data now including time point (note that this would be the model to fit in the first place and we have fitted a model only with proficiency for pedagogical purposes).

### Multiple smooth terms

The data `dat1` contains learning and proficiency scores from 200 subjects, taken at 10 time points. The data was simulated so that proficiency increased with time (at different degrees for different speakers). An interesting question is whether learning scores improve with time independently of proficiency, or if proficiency alone is causing learning scores to improve. We can approach with question by applying a statistical method called causal inference, based on directed acyclic graphs (DAGs) theory (a full treatment of this is beyond the scope of the paper. Readers are referred to XXX).

The DAGs in `?@fig-d1-dag` represent the causal relationships between learning scores, proficiency scores and time point in the two scenarios. In (a), time point affects proficiency and proficiency affects learning scores. In other words, time point has a direct effect on proficiency but not on learning scores; learning scores can be predicted from proficiency alone. In (b), on the other hand, time point affects proficiency as in (a) but also learning scores (and proficiency affects learning scores as in (a)). This means that time point affects learning scores in two ways: through its effect on proficiency and directly through its effect on learning scores.

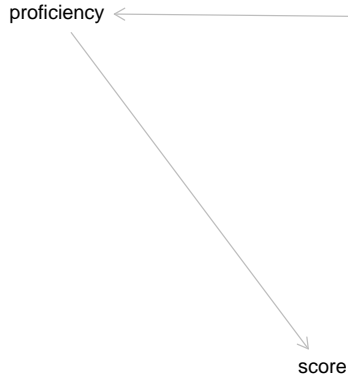


Figure 3: `?(caption)`

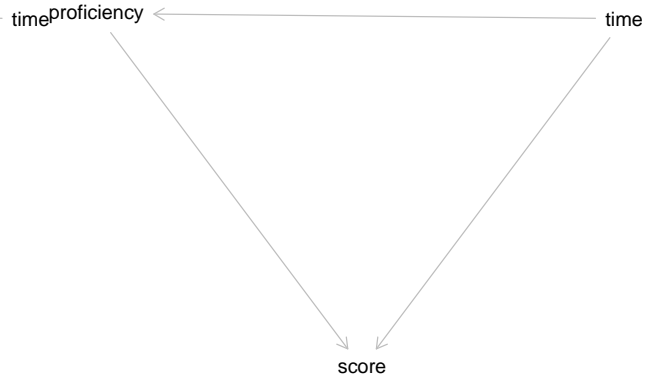


Figure 4: `?(caption)`

DAGs allow us to make causal statements based on statistical results. In this case, when including both time point and proficiency as predictors in a GAM, time should not have an effect on learning score if scenario (a) is correct (while it should have an effect if scenario (b) is correct). Let's model the data.

```
gam_2 <- gam(
  score ~ s(proficiency) + s(time) +
    s(subj, bs = "re"),
  data = dat1
)
```

In `gam_2` we fit a GAM to learning scores `score` with two predictors: a smooth term over proficiency and a smooth term over time point. We also include random effects to account for data from multiple participants. The syntax for random effects in GAMs is different from the syntax in `lme4`. With `gam()`, we can specify random effects using smooth terms and the `re` (for Random Effects) basis function. Random intercept are added with the syntax `s(ranint, bs = "re")` and random slopes with the syntax `s(ranint, ranslope, bs = "re")`. Here we just add a random intercept by subject for illustration. Let's inspect the model summary.

```
summary(gam_2)
```

Family: gaussian

Link function: identity

Formula:

`score ~ s(proficiency) + s(time) + s(subj, bs = "re")`

Parametric coefficients:

|             | Estimate | Std. Error | t value | Pr(> t )   |
|-------------|----------|------------|---------|------------|
| (Intercept) | 24.738   | 0.441      | 56.09   | <2e-16 *** |

---

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

Approximate significance of smooth terms:

|                | edf     | Ref.df  | F       | p-value      |
|----------------|---------|---------|---------|--------------|
| s(proficiency) | 8.699   | 8.976   | 4638.38 | < 2e-16 ***  |
| s(time)        | 1.000   | 1.000   | 20.27   | 7.36e-06 *** |
| s(subj)        | 198.943 | 200.000 | 2637.63 | < 2e-16 ***  |

---

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

R-sq.(adj) = 0.997 Deviance explained = 99.7%

GCV = 0.10023 Scale est. = 0.090681 n = 2200

...

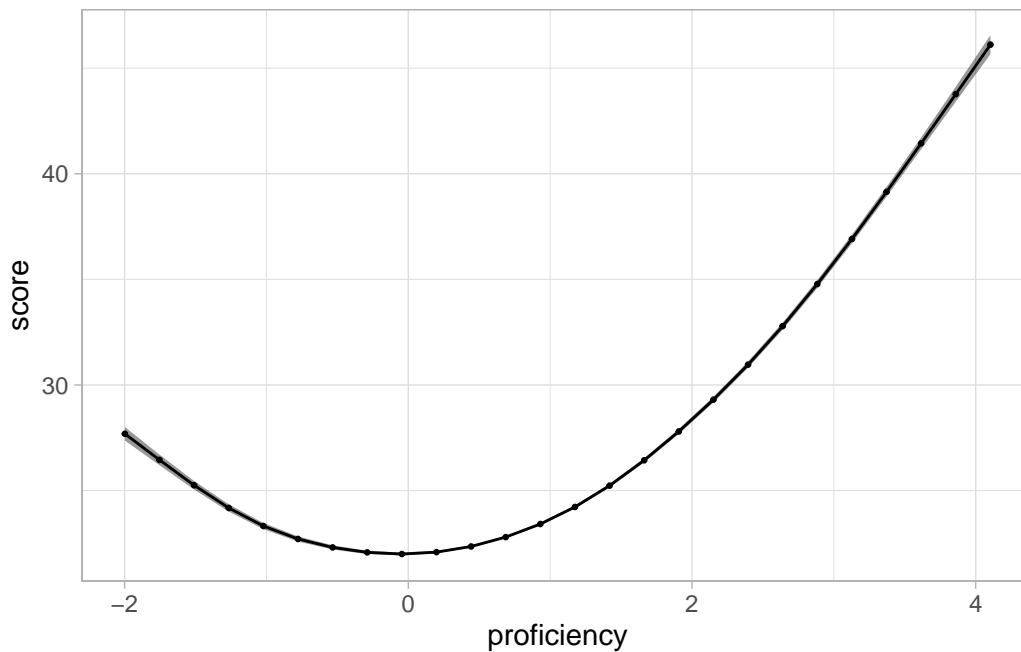
Let's plot the model predictions.

```
gam_2_preds <- predict_gam(  
  gam_2, length_out = 25,  
  series = "proficiency",  
  exclude_terms = c("s(subj)", "s(subj,proficiency)")  
)
```

Warning in mgcv::predict.gam(model, newdata = pred\_grid, se.fit = TRUE, :  
non-existent exclude terms requested - ignoring

Warning: There was 1 warning in `dplyr::mutate()`.  
i In argument: `fit = rowSums(dplyr::across())`.  
Caused by warning:  
! Using `across()` without supplying `.cols` was deprecated in dplyr 1.1.0.  
i Please supply `.cols` instead.

```
plot(gam_2_preds)
```

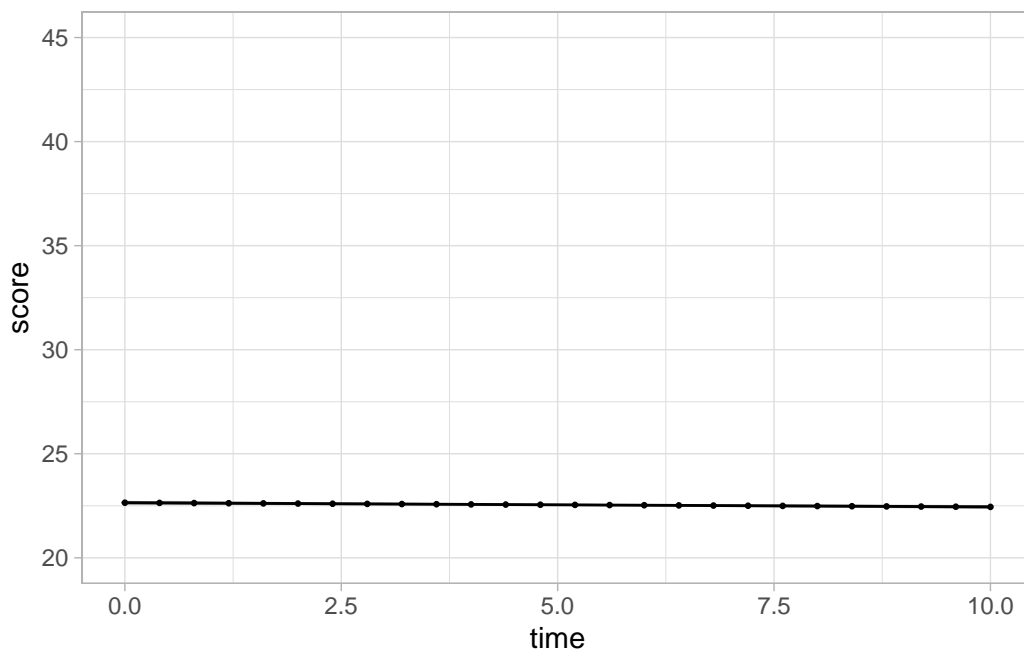




```
gam_2_preds_t <- predict_gam(
  gam_2, length_out = 25,
  series = "time",
  exclude_terms = c("s(subj)", "s(subj,proficiency)")
)
```

Warning in mgcv::predict.gam(model, newdata = pred\_grid, se.fit = TRUE, :  
non-existent exclude terms requested - ignoring

```
plot(gam_2_preds_t) +
  ylim(20, 45)
```



## Case study 2: simultaneous vs late bilinguals