

Compass TB1 Group Project: Lab Notebook

Euan Enticott, Daniel Ward, Shannon Williams

25/01/2021

1 Introduction

In this lab notebook we present the analyses undertaken as part of the Compass TB1 group project. We focus on a large, open source dataset provided by the City of Chicago consisting of reported incidents of crime from 2001 to present day, as reported by the Chicago Police Department. The dataset can be obtained from the website <https://data.cityofchicago.org/Public-Safety/Crimes-2001-to-Present/ijzp-q8t2>.

The notebook will be broken up into three sections:

1. Workflow and Computational Techniques: A summary of the approach used in the development of the **chigcrim** package.
2. Classification: Predicting whether a reported crime led to the arrest of a suspect. We consider three different discriminative classifiers: logistic regression, support vector machines, and random forests.
3. Prediction: Predicting the number of reported crimes which occur in a given region over a specified time period (of a given type, if desired), through interpolation. We consider kernel ridge regression and generalised additive models (GAMs).

2 Workflow and Computational Techniques

In the development of the package **chigcrim**, we utilised a number of computational techniques which are documented below. The package can be accessed on GitHub: <https://github.com/shannon-wms/chicago-crime>.

2.1 Object-Oriented Programming

We decided to use an object oriented approach to constructing models. In particular, we made use of the **R6Class** from the **R6** package. **R6** objects are mutable and use an object oriented approach similar to most programming languages (unlike the more basic **S3** and **S4** classes). **R6** classes are very similar to reference classes, although generally **R6** classes are preferred. This is for a variety of reasons, including the fact that **R6** classes are much faster, and they handle fields in a neater way by putting them in a separate environment. More information on this can be found here. Notably, **R6** classes support both public (available to the user) and private (not visible to the user) members.

Each of the models considered can be initialised with a call to `model$new(...)`, and then fitted using `model$fit(X, y)` where **X** and **y** are our training data matrix and response vector, respectively. Predictions can then be made with `model$predict(X)`. The datasets, fitted model and predictions, along with other objects unique to the model, are then stored as fields accessed through `model$field`.

It is worth noting that this use of `fit` and `predict` methods in object-oriented machine learning packages is used elsewhere, particularly in the **Python** package **scikit-learn**. The key advantage of this approach is that it abstracts away the underlying complexities, and creates a consistent structure that makes it easier to compare different models.

2.2 Parallel Programming

Parallel programming was used to speed up slow code, utilising the package **doParallel** which calls the packages **parallel** and **foreach**. In particular, we implemented a parallel cross-validation function **kfold_cv**, which facilitated more accurate and faster assessment of model performance.

2.3 Documentation

In order to document functions in the package, we made use of **roxygen2**. This package automatically creates the **.Rd** documentation files for functions, classes and data, from comments added above the function definitions. In addition to being an efficient way to generate documentation, this allows the code and documentation to coexist, making it easier to remember to update the documentation.

2.4 API Querying

The Chicago Crime dataset is updated on a daily basis and, at the time of writing, contains over 7 million rows of data. For each reported incident, a total of 22 features are recorded (for a full list of features and their descriptions see the City of Chicago website). We elected not to download the dataset in full and include this in the package as it would consume a significant amount of storage, and quickly become out of date.

Instead, we utilised the **RSocrata** package which allows easy interaction with the online open data portal. The wrapper **load_data** allows the end user to directly download the dataset into the R environment, via the function **read.socrata**. **load_data** provides functionality for filtering by year, limiting the number of data points received, and omitting missing/NA values.

2.5 Testing and Continuous Integration

For tests, we used the **testthat** package. This simply creates a new directory **./tests/testthat/** in which the user can define tests for the package functions. Although the tests can be run alone, generally, we made use of the package **rcmdcheck**. This contains the function **rcmdcheck**, which not only runs the tests, but also carries out a more sophisticated check. Some examples of what is checked are listed below:

- Checks the package installs correctly.
- Checks for missing documentation.
- Automatically runs examples in documentation, to check they run successfully.
- Checks for undefined variables.
- Checks for missing dependencies.

In order to ensure that tests and checks were run on a regular basis, continuous integration was set up using Github Actions. This runs **rcmdcheck()** for each pull request and push to the main branch, thus limiting the possibility of unintentionally introducing any breaking changes.

3 Classification

In this section we consider the binary classification task of predicting whether a suspect is arrested, given data on the reported crime. Three models will be investigated: logistic regression, support vector machines and random forest classifiers.

3.1 Assessing Performance

To assess performance, three metrics will be used:

- Overall accuracy: The proportion of correct predictions.
- Sensitivity: The proportion of positive observations correctly predicted.
- Specificity: The proportion of negative observations correctly predicted.

These metrics can be computed using the functions `classification_accuracy`, `classification_sensitivity`, and `classification_specifity`, respectively, from the **chigcrim** package.

Note that logistic regression is a probabilistic classifier whereas support vector machines and random forests are non-probabilistic. As such, the output from logistic regression is a probability in $[0, 1]$ while the outputs from the other models are class predictions in $\{0, 1\}$. In our analysis, the probabilities given by logistic regression are rounded to yield predictions that can be assessed using the metrics above. This facilitates the comparison of these results with the other two classifiers.

```
# Load required packages
library(chigcrim)

# Use seed to ensure reproducibility
set.seed(1)

# List of metrics used to evaluate classification models
metrics <- list(accuracy = classification_accuracy,
                 sensitivity = classification_sensitivity,
                 specificity = classification_specifity)
```

For each classifier, we will use repeated k -fold cross validation to obtain mean values of the metrics. Specifically, 5-fold cross validation is repeated five times (run in parallel using the function `kfold_cv`), and the metrics are averaged across the folds and repeats.

3.2 Feature Selection

Due to the time and computational limitations involved in downloading, storing, and training models on the entire dataset, we shall train the models only on a subset of the data from the year 2019.

```
# Load 2019 data to build classification models
df <- load_data(year = 2019, strings_as_factors = FALSE)

print(as.data.frame(head(df)))
```

##	id	case_number	date	block	iucr
## 1	11552577	JC100040	2019-01-01 00:31:00	032XX W LAWRENCE AVE	1310
## 2	11552587	JC100034	2019-01-01 00:05:00	006XX E 83RD PL	1310
## 3	11552596	JC100045	2019-01-01 00:03:00	001XX W HURON ST	0430
## 4	11552605	JC100030	2019-01-01 00:01:00	004XX N MONTICELLO AVE	143A
## 5	11552609	JC100028	2019-01-01 00:05:00	013XX S CENTRAL PARK AVE	0486
## 6	11552610	JC100044	2019-01-01 00:02:00	031XX W 25TH ST	1310

##	primary_type	description	location	description	arrest
## 1	CRIMINAL DAMAGE	TO PROPERTY	RESTAURANT		FALSE
## 2	CRIMINAL DAMAGE	TO PROPERTY	RESIDENCE		FALSE
## 3	BATTERY AGGRAVATED: OTHER DANG WEAPON		HOTEL/MOTEL		FALSE
## 4	WEAPONS VIOLATION	UNLAWFUL POSS OF HANDGUN	ALLEY		TRUE
## 5	BATTERY	DOMESTIC BATTERY SIMPLE	APARTMENT		FALSE
## 6	CRIMINAL DAMAGE	TO PROPERTY	RESIDENCE		FALSE

##	domestic	beat	district	ward	community_area	fbi_code	x_coordinate	y_coordinate
## 1	FALSE	1713	17	33	14	14	1153943	1931709
## 2	FALSE	632	6	6	44	14	1182085	1849762
## 3	FALSE	1832	18	42	8	04B	1175159	1905043
## 4	FALSE	1122	11	27	23	15	1151958	1902815
## 5	TRUE	1011	10	24	29	08B	1152589	1893426
## 6	FALSE	1033	10	12	30	14	1155950	1887234

##	year	updated_on	latitude	longitude
## 1	2019	2019-01-10 15:16:50	41.96844	-87.70934
## 2	2019	2019-01-10 15:16:50	41.74297	-87.60841
## 3	2019	2019-01-10 15:16:50	41.89482	-87.63213
## 4	2019	2019-01-10 15:16:50	41.88920	-87.71740
## 5	2019	2019-01-10 15:16:50	41.86342	-87.71533
## 6	2019	2019-01-10 15:16:50	41.84636	-87.70316

Below we outline our feature selection and data pre-processing choices. These decisions involved incorporating information from conducting exploratory data analysis (EDA), as well as considering what is computationally feasible.

- Encode the date as the day of the year, `yday` (an integer between 1-365, or 1-366 for leap years).
- Encode time of day as a floating point in 0-24 (`time`), as this showed promising results in our EDA.
- Drop `year` as we are only using data from 2019.
- Drop `id` as this is the unique key for each data point.
- We drop `case_number`, as the values are almost all unique and are not informative.
- Drop `primary_type`, `description` and `iucr` code. Keep `fbi_code` as an indicator of crime type.
- Drop `updated_on` as this only refers to when the data was last updated and is uninformative.
- Drop `latitude` and `longitude`, and keep `x_coordinate` and `y_coordinate`. Note that these co-ordinates are likely not particularly useful for linear classifiers, but should be useful for non-linear methods.

- Keep community_area, but other areas (district, beat, ward and block) are dropped.
- Omit all NA values.
- Particularly rare factors will be grouped into a variable other (see ?otherise).

```
# Convert rarely used fbi categories to "OTHER"
```

```
df$fbi_code <- otherise(df$fbi_code, 500)
```

```
## 7 out of 26 categories were converted to OTHER corresponding to 0.51%
```

```
## of observations.
```

```
df$location_description <- otherise(df$location_description, 1000)
```

```
## 126 out of 155 categories were converted to OTHER corresponding to 10.82%
```

```
## of observations.
```

```
# List of features for removal
```

```
remove_features <- c("id", "year", "case_number", "primary_type", "description",  
                     "iucr", "updated_on", "latitude", "longitude",  
                     "date", "district", "beat", "ward", "block")
```

```
# Add more time features, remove missing values and convert columns to factor
```

```
df %<>% mutate(day = yday(df$date),  
               time = hour(df$date) + minute(df$date) / 60) %>%  
  select(-all_of(remove_features)) %>%  
  filter(complete.cases(df)) %>%  
  mutate(location_description = as.factor(location_description),  
         fbi_code = as.factor(fbi_code),  
         community_area = as.factor(community_area))
```

```
print(as.data.frame(head(df)))
```

```
##   location_description arrest domestic community_area fbi_code x_coordinate  
## 1      RESTAURANT    FALSE      FALSE           14      14      1153943  
## 2      RESIDENCE    FALSE      FALSE           44      14      1182085  
## 3    HOTEL/MOTEL    FALSE      FALSE            8     04B      1175159  
## 4         ALLEY     TRUE      FALSE           23      15      1151958  
## 5    APARTMENT    FALSE      TRUE           29     08B      1152589  
## 6      RESIDENCE    FALSE      FALSE           30      14      1155950  
##   y_coordinate day      time  
## 1      1931709   1 0.51666667  
## 2      1849762   1 0.08333333  
## 3      1905043   1 0.05000000  
## 4      1902815   1 0.01666667  
## 5      1893426   1 0.08333333  
## 6      1887234   1 0.03333333
```

3.3 Logistic Regression

We will first consider a logistic regression classifier, which is a linear classifier for binary data. We use our R6 class `LogisticRegression` which has methods `fit` and `predict`. `fit` calls the function `optim()` internally to minimize the cross-entropy/log-loss function using gradient-based optimisation methods (either BFGS, L-BFGS-B or CG); in our analysis we use BFGS as we don't run into memory issues.

```
# Split into data matrix and response vector
X <- df %>% select(-arrest)
y <- df$arrest

# Initialise new LogisticRegression object
lr <- LogisticRegression$new(solver = "BFGS",
                             control = list(maxit = 1000, reltol = 1e-4),
                             rounding = TRUE)

# Find repeated cross-validation error
lr_results <- kfold_cv(lr, X, y, metrics, k = 5, n_reps = 5,
                      parallel = TRUE, n_threads = 5)

# Export results for use in report
tab <- xtable::xtable(as_tibble(lr_results), digits = 4)
print(tab, type = "latex", comment = FALSE)
```

	accuracy	sensitivity	specificity
1	0.8614	0.4347	0.9788
2	0.8613	0.4351	0.9786
3	0.8613	0.4345	0.9787
4	0.8613	0.4341	0.9789
5	0.8612	0.4346	0.9786

3.4 Support Vector Machines

As the support vector machine (SVM) allows use of a kernel function, it should be able to capture non-linear relationships in the original feature space (given an appropriate kernel). As the relationship between the features and the class is unknown, a radial basis function kernel is used. The R6 class `SupportVectorMachine`, again with `fit` and `predict` methods, is a wrapper of the function `svm` from the dedicated SVM package `e1071`. Unfortunately, support vector machines do not scale well to large data sets, so here we will only use 30,000 rows.

```
# Subset the data randomly
index <- sample(1:nrow(df), 30000)
X_subset <- X[index, ]
y_subset <- y[index]
```

In order to obtain optimal values for `cost` and `gamma` to be provided to `svm` for training the SVM, we utilise the built-in generic function `tune` from `e1071` which performs a grid search over given parameter ranges.

```
svm_tuned <- tune(e1071::svm, y_subset ~ ., data = cbind(X_subset, y_subset),
  ranges = list(cost = 2^(-2:5), gamma = 2^(-15:-4)),
  tunecontrol = tune.control(nrepeat = 3,
    sampling = "cross",
    cross = 5),
  kernel = "radial", type = "C-classification")
```

```
svm_tuned
```

```
##
## Parameter tuning of 'e1071::svm':
##
## - sampling method: 5-fold cross validation
##
## - best parameters:
##   cost   gamma
##    16 0.03125
##
## - best performance: 0.128
```

We obtain optimal values of 2^4 for `cost` and 2^{-5} for `gamma`, which we use to train our SVM.

```
# Initialise new SupportVectorMachine object
svm <- SupportVectorMachine$new(kernel = "radial")

# Run repeated 5-fold cross-validation to evaluate performance
svm_results <- kfold_cv(svm, X_subset, y_subset, metrics, k = 5,
  n_reps = 5, parallel = FALSE,
  n_threads = 5, gamma = 2^-5, cost = 2^4)

# Export results for use in report
tab <- xtable::xtable(as_tibble(svm_results), digits = 4)
print(tab, type = "latex", comment = FALSE)
```

	accuracy	sensitivity	specificity
1	0.8624	0.4076	0.9870
2	0.8627	0.4091	0.9870
3	0.8626	0.4087	0.9870
4	0.8622	0.4067	0.9870
5	0.8619	0.4052	0.9870

3.5 Random Forests

The random forest trains multiple decision trees using a subset of features to create each tree. Decision trees are an ordered set of rules defined on the feature values that attempt to split the data into classes. As decision trees have a tendency to overfit, the random forest attempts to overcome this by growing multiple trees on different features within the dataset. In the classification case, a prediction is made by predicting from each decision tree separately and finding the most common predicted class. This is a black-box method but has the potential for good predictive performance.

Our R6 class `RandomForest` is a wrapper for the function `randomForest` from the `randomForest` package. Due to the high computational cost involved in fitting random forest models we again use a smaller subset of the data. Here we use the same dataset used to fit the SVM, with the factor `community_area` removed as the algorithm for fitting the model cannot handle such a large number of factors.

```
X_subset <- select(X_subset, - "community_area")

# Initialise new RandomForest object
rf <- RandomForest$new()

# Run repeated 5-fold cross-validation to evaluate performance
rf_results <- kfold_cv(rf, X_subset, y_subset, metrics, k = 5,
                      n_reps = 5, parallel = TRUE, n_threads = 5)

# Export results for use in report
tab <- xtable::xtable(as_tibble(rf_results), digits = 4)
print(tab, type = "latex", comment = FALSE)
```

	accuracy	sensitivity	specificity
1	0.8762	0.5258	0.9722
2	0.8780	0.5293	0.9736
3	0.8762	0.5299	0.9712
4	0.8775	0.5293	0.9729
5	0.8761	0.5259	0.9721

3.6 Model Comparisons

The models all performed quite similarly. Unsurprisingly, the support vector machine and the random forest classifier performed best, due to their ability to capture non-linear relationships. It is worth noting that this difference is surprisingly small, although this could be related to the fact that the logistic regression model was trained on a larger dataset. Alternatively, the SVM results could suggest that encoding the community areas as factors was sufficient to capture much of the spatial aspect of the crimes in logistic regression.

Logistic regression had higher sensitivity than the support vector machine, and thus could be viewed as a preferable option if false negatives are particularly detrimental. It is worth noting however that both support vector machines and logistic regression can be adjusted to reflect unbalanced costs associated with incorrect predictions.

Being an ensemble model, the random forest unsurprisingly performed the best based on all three metrics. It is likely that the random forest could be improved by training on a larger subset of the data and optimising the hyperparameters. The downside of using a random forest is that, as a black-box method, it does not provide any inference as to how the features are used within the model. We are at least able to look at how many times each feature was chosen as a splitting point in all of the trees, providing us with some information as to the importance of each of the features.

```
# Fit model with optimal hyper parameters
rf$fit(X_subset, y_subset)
```



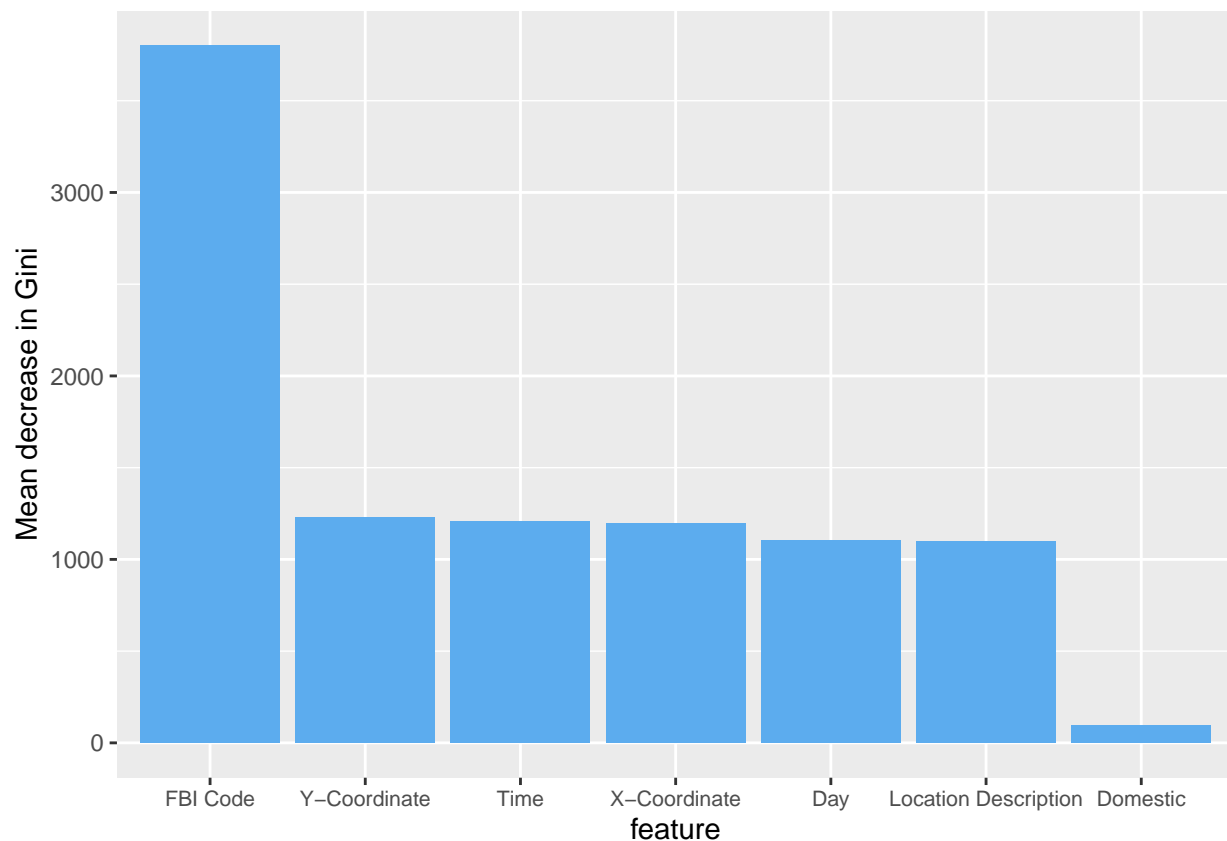
```

# Get importance of features from model
imp <- randomForest::importance(rf$fitted_model)
imp <- tibble(feature = row.names(imp), mean_decrease_gini = as.vector(imp))

imp <- imp %>%
  arrange(desc(mean_decrease_gini)) %>%
  mutate(feature = factor(feature, levels = .$feature))

# Plot importance
imp_plot <- ggplot(aes(x = feature, y = mean_decrease_gini), data = imp) +
  geom_col(fill = "steelblue2") + ylab("Mean decrease in Gini") +
  scale_x_discrete(labels = c("FBI Code", "Y-Coordinate", "Time", "X-Coordinate",
                             "Day", "Location Description", "Domestic")) +
  theme(axis.text.x = element_text(size = rel(0.9)))
imp_plot

```



The plot suggests that the most “important” feature is `fbi_code`, which is somewhat unsurprising as this tells us what type of crime was reported. This would intuitively have a significant impact on whether the incident report led to a suspect being arrested. Note that none of the models achieved an accuracy greater than 90%.

4 Prediction

In this section we consider the regression task of predicting the number of reported incidents of crime in a given region over a specified time period. We will look at a kernel ridge regression model before moving on to generalised additive models.

4.1 Assessing Performance

To assess performance, two metrics will be used: the root mean squared error (RMSE) and R-squared (R^2). These metrics can be computed using the functions `rmse_loss` and `r_squared`, respectively, from the `chigcrim` package.

```
# List of metrics used to evaluate regression models  
metrics <- list(rmse = rmse_loss, r2 = r_squared)
```

Repeated k -fold cross validation is used, and hyperparameters are chosen based on the RMSE metric. In particular, we will use 5 repetitions of 5-fold cross validation at each parameter value using the `kfold_cv` function.

4.2 Kernel Ridge Regression

The R6 class `KernelRegression` fits a kernel ridge regression model with a specified kernel: one of "linear" (equivalent to ridge regression), "polynomial", or "rbf" (radial basis kernel).

As the relationship between the data and response is complex and non-linear, only the radial basis kernel will be considered here. This implicitly induces an extremely flexible infinite-dimensional feature transform. Note that we fix the regularisation parameter `lambda` as 0.001 in the ensuing analysis.

```
# Load required packages
library(rgdal)
library(sf)
library(ggmap)
# Use seed to ensure reproducibility
set.seed(1)
```

4.2.1 Weekly Predictions

We first consider a simple, single-dimensional model, using only the number of crimes per week as the predictor.

```
df <- load_data(select = "date")

df %<>%
  mutate(date = date(date), week = week(date), year = year(date)) %>%
  filter(week < 53) %>%
  group_by(week, year) %>%
  mutate(week_start = min(date)) %>%
  ungroup() %>%
  count(week, week_start, year) %>%
  arrange(year, week) %>%
  mutate(cumulative_week = 1:nrow())
print(as.data.frame(head(df)))
```

```
##   week week_start year      n cumulative_week
## 1     1 2001-01-01 2001 8942              1
## 2     2 2001-01-08 2001 8589              2
## 3     3 2001-01-15 2001 8704              3
## 4     4 2001-01-22 2001 8300              4
## 5     5 2001-01-29 2001 8473              5
## 6     6 2001-02-05 2001 8418              6
```

We shall use a grid search with 5-fold cross validation, repeated 5 times to obtain a robust average for RMSE, in order to choose the bandwidth hyperparameter for the radial basis kernel.

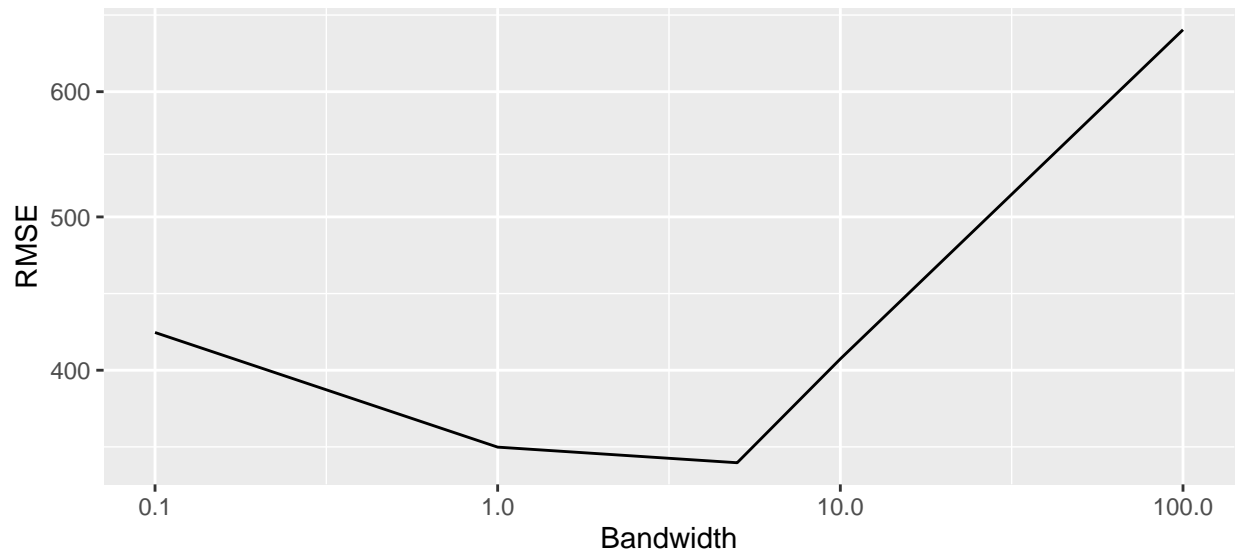
```
X <- as.matrix(df$cumulative_week)
X <- scale(X)
y <- df$n

# Candidate values for the bandwidth of the radial basis function
bandwidth <- c(0.1, 1, 5, 10, 100)
# Initialise vector of RMSE
rmse_vec <- c()
# Conduct a grid search with 5-fold CV
for (i in 1:length(bandwidth)){
  kr = KernelRidge$new(kernel = "rbf", lambda = 0.001, bandwidth[i])
  cv_results <- kfold_cv(kr, X, y, metrics, k = 5, n_reps = 5,
                        parallel = TRUE, n_threads = 5)
```

```

mean_rmse <- mean(cv_results$rmse)
rmse_vec <- c(rmse_vec, mean_rmse)
# If mean RMSE is better than for the last best model clone and treat as best
if (mean_rmse == min(rmse_vec)){
  # Clone the model and save
  kr_best <- kr$clone(deep = TRUE)
  # Update the best results
  cv_results_best <- cv_results
}
}
# Save as data frame and plot results
rbf_results <- data.frame(bandwidth, rmse = rmse_vec)
bandwidth_plot <- ggplot(rbf_results) +
  geom_line(aes(x = bandwidth, y = rmse)) +
  scale_x_log10() +
  scale_y_log10() +
  xlab("Bandwidth") + ylab("RMSE")
bandwidth_plot

```



We find that a bandwidth parameter value of 5 provided the lowest RMSE value. We can look at the results for the 5 repeats of 5-fold cross-validation for this model:

```

tab <- xtable::xtable(as_tibble(cv_results_best), digits = 4)
print(tab, type = "latex", comment = FALSE)

```

	rmse	r2
1	353.3775	0.9627
2	347.7463	0.9640
3	351.1616	0.9635
4	346.9327	0.9646
5	348.7730	0.9642

The model performs very well, explaining 96.38 % of the variance on average.

As this is single-dimensional, we can plot the prediction function for the model.

```

# Fit the final kernel ridge regression model
kr_best$fit(X, y)
# Obtain predictions on the dataset
y_hat <- kr_best$predict(X)
# Plot predictions
df$y_hat <- y_hat
week_plot <- ggplot(df) +
  geom_point(aes(x = week_start, y = n)) +
  geom_line(aes(x = week_start, y = y_hat), colour = "steelblue2", size = 1) +
  labs(x = "Year", y = "Number of reported incidents")
week_plot

```



The plot suggests a cyclic pattern over years, so naturally we proceed to consider using a trigonometric transformation of the data (taking year as the period), in the hope that we can capture this periodicity with a simpler model. Note that, whilst kernel ridge regression is useful for interpolation, this model with a radial basis kernel actually has limited utility for forecasting reported crime counts into the future. We illustrate this point by attempting to extrapolate from this model.

```

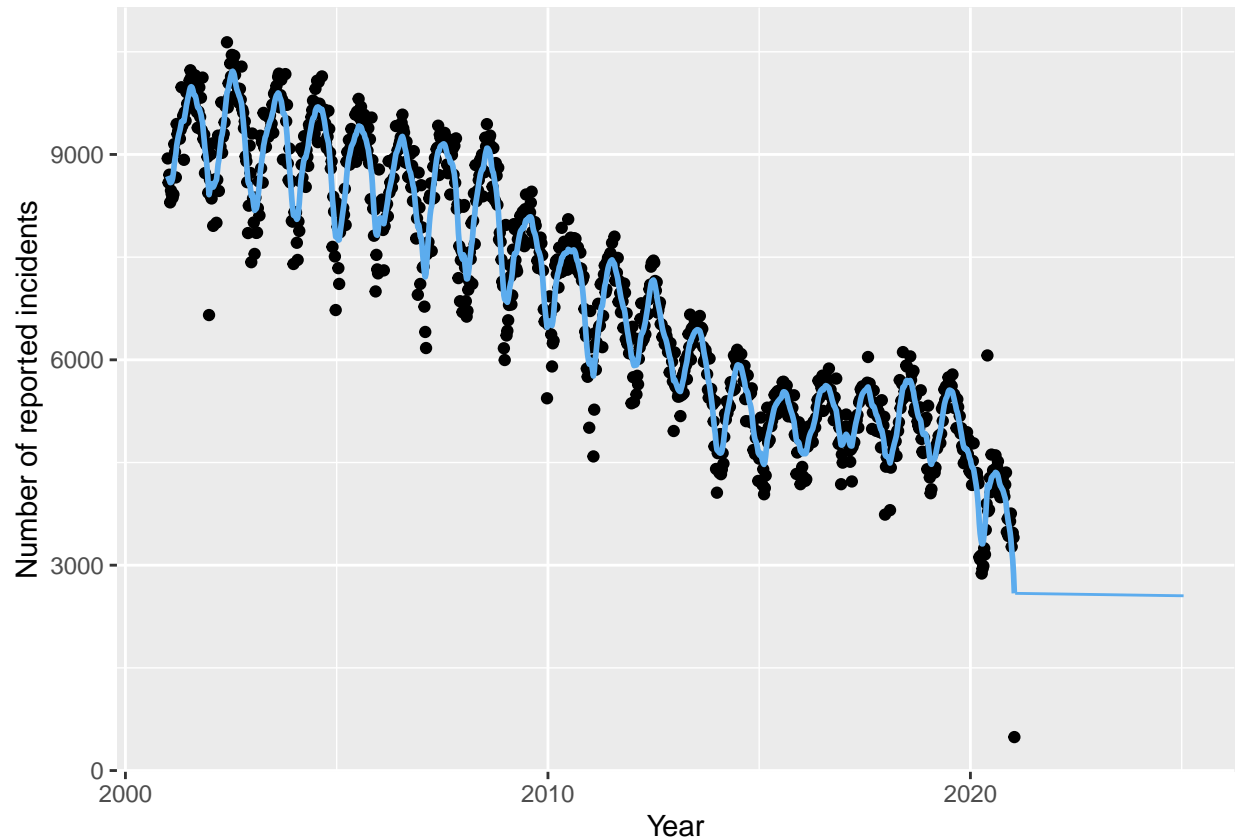
# Get the start of the next week not included in the data
start_week_date <- max(df$week_start) + 7
# Get the corresponding week number
start_week_int <- max(df$cumulative_week) + 1
# Add the dates
extra_dates <- seq(ymd(start_week_date), ymd(start_week_date + 365*4),
  by = '1 week')
extra_X <- start_week_int:(start_week_int + length(extra_dates) - 1)
extra_X <- (extra_X - attributes(X)$`scaled:center`) /

```

```

  attributes(X)$`scaled:scale`
# Predict using our fitted model
extra_y_hat <- kr_best$predict(as.matrix(extra_X))
# Attach these values to the data frame
extra_df <- tibble(extra_dates, extra_y_hat)
week_plot + geom_line(data = extra_df, aes(extra_dates, extra_y_hat), colour = "steelblue2")

```



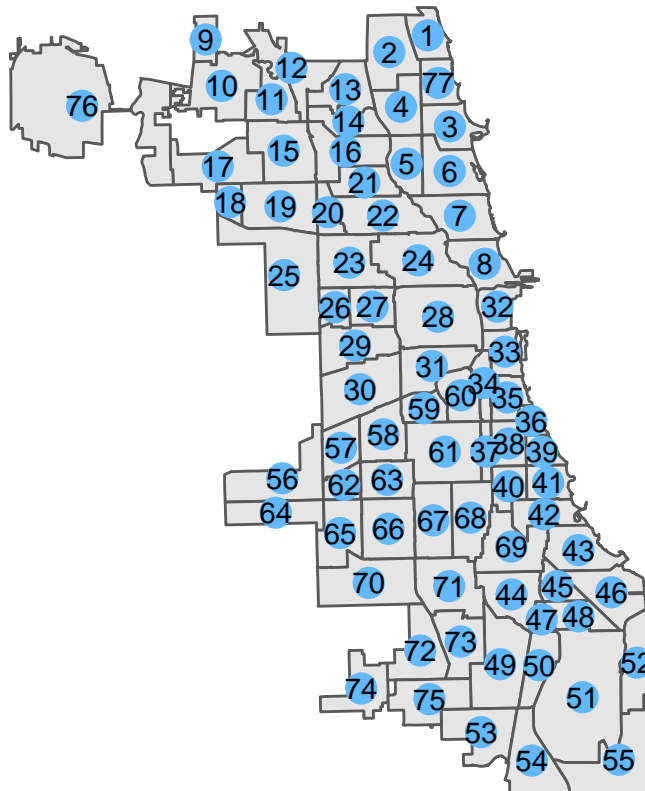
The cycles are captured solely due to being able to interpolate locally between points. When extrapolating with a radial basis function kernel, no cycles are predicted, so the out-of-sample accuracy drops substantially. The cross-validation error found is only applicable if the new data followed the same distribution, and this is simply not the case when considering future data.

4.2.2 Monthly Predictions with a Spatial Component

The previous example aggregated counts of reported crime over weeks, and as such used only a single-dimensional feature vector. We look at extending this analysis by aggregating over an additional spatial component. In the ensuing analysis, we consider predicting the reported crime counts over a month, aggregated over the community areas.

The positions of the community areas will be represented by the coordinates of their centroids, which are extracted from their shapefiles. These shapefiles are included in the package for ease.

```
data("community_bounds")
centroids <- community_bounds %>%
  arrange(as.integer(area_numbe)) %>%
  st_centroid() %>%
  st_coordinates() %>%
  as_tibble() %>%
  mutate(community_area = 1:nrow(.))
# Plot the locations of centroids on map
centroids_plot <- ggplot() +
  geom_sf(data = community_bounds) +
  geom_point(data = centroids, aes(X, Y), size = 5, colour = "steelblue1") +
  geom_text(data = centroids, aes(X, Y, label = community_area)) +
  theme_void()
centroids_plot
```



We shall consider data from the years 2016 to 2020.

```
df <- load_data(c(2016, 2020), strings_as_factors = FALSE, na_omit = TRUE)

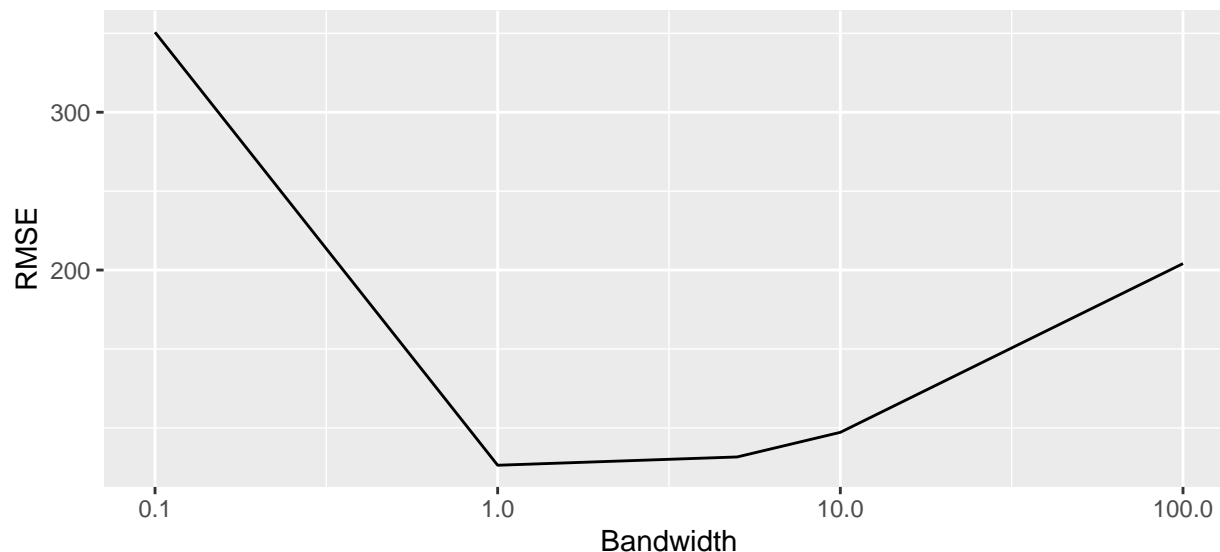
df %<>% select(date, community_area) %>%
  mutate(date = date(date), year = year(date), month = month(date)) %>%
  group_by(year, month) %>%
  mutate(date = min(date)) %>%
  group_by_all() %>%
  summarise(n = n()) %>%
  left_join(centroids, by = "community_area")
```

In order to choose a suitable bandwidth parameter for the radial basis function, we again perform repeated 5-fold cross-validation.

```
# Obtain data matrix for fitting the model
X <- df %>% ungroup() %>%
  select(-community_area, -date, -n) %>%
  as.matrix()
# Scale the parameters
X <- scale(X)
y <- df$n

# Bandwidth candidates
bandwidth <- c(0.1, 1, 5, 10, 100)
# Initialise vector for storing RMSE values
rmse_vec <- c()
# Set new seed for reproducibility
set.seed(3)
for (i in 1:length(bandwidth)){
  kr <- KernelRidge$new("rbf", lambda = 0.001, bandwidth[i])
  cv_results <- kfold_cv(kr, X, y, metrics, k = 5, n_reps = 5,
    parallel = TRUE, n_threads = 5)
  mean_rmse <- mean(cv_results$rmse)
  rmse_vec <- c(rmse_vec, mean_rmse)
  # Save best model
  if (mean_rmse == min(rmse_vec)){
    kr_best <- kr$clone(deep=TRUE)
    cv_results_best <- cv_results
  }
}

# Store as data frame and plot
rbf_results <- data.frame(bandwidth, rmse = rmse_vec)
band_plot <- ggplot(rbf_results) +
  geom_line(aes(x = bandwidth, y = rmse)) +
  scale_x_log10() +
  scale_y_log10()
band_plot
```

We obtain an optimal value of 1 which minimises the mean square error. Again, we can look at the results for the three repeats of 5-fold cross-validation for this model:

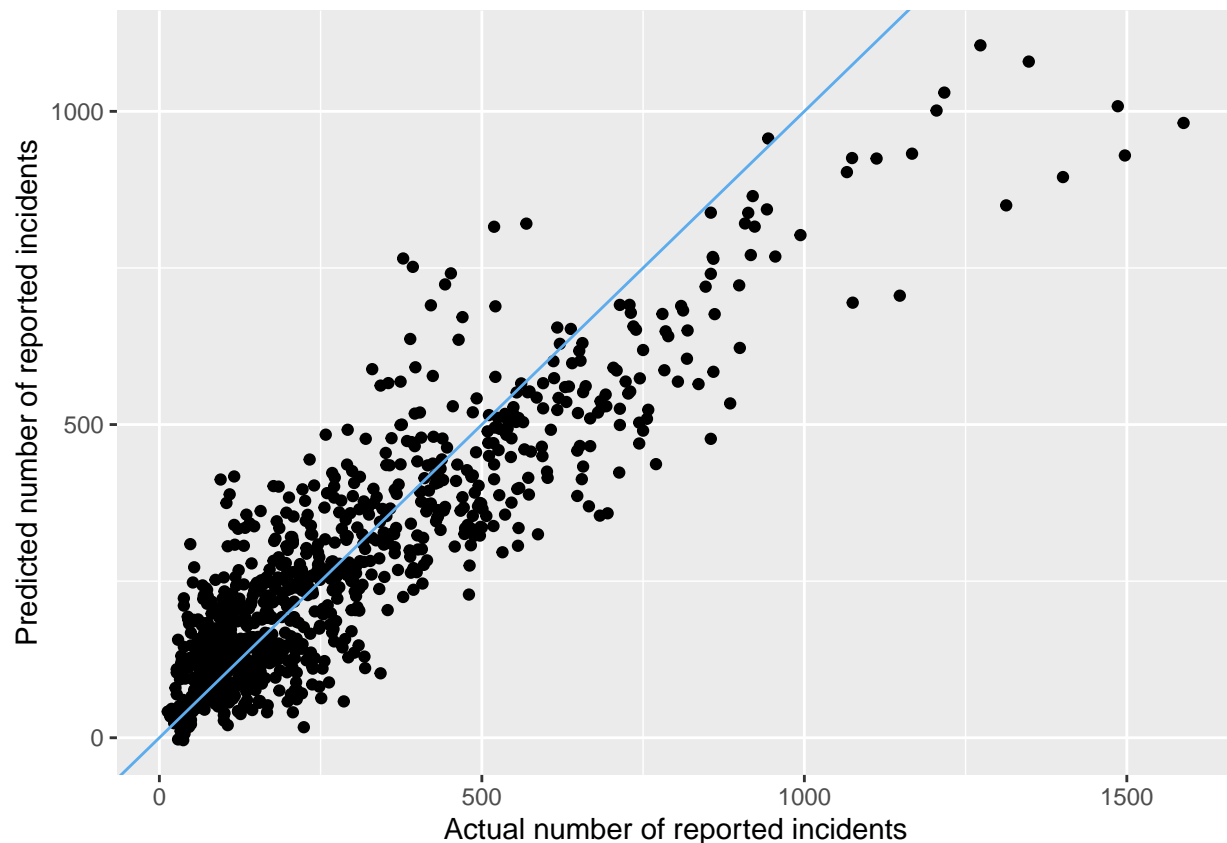
```
tab <- xtable::xtable(as_tibble(cv_results_best))
print(tab, type = "latex", comment = FALSE)
```

	rmse	r2
1	122.86	0.75
2	120.33	0.76
3	119.14	0.77
4	122.10	0.75
5	121.17	0.76

On average, this model explained 75.73% of the variance. Note the lower R^2 value: in contrast to the weekly model, observations were not averaged over the whole of Chicago.

We can plot the predicted values against the observed values on a held out data set to observe the predictive power of the model:

```
# Randomly generate training and test data
index <- sample(1:nrow(X), round(0.2 * nrow(X)))
# Split the data
X_train <- X[-index, , drop = FALSE]
y_train <- y[-index]
X_test <- X[index, ]
y_test <- y[index]
# Refit best model on full training dataset
kr_best$fit(X_train, y_train)
y_hat <- kr_best$predict(X_test)
# Plot predicted values against actual number of reports
pred_plot <- qplot(y_test, y_hat, geom = "point") +
  geom_abline(slope = 1, intercept = 0, colour = "steelblue2") +
  labs(x = "Actual number of reported incidents",
       y = "Predicted number of reported incidents") +
  scale_y_continuous(breaks = seq(0, 2000, by = 500))
pred_plot
```



The plot illustrates strong predictive power, particularly when the number of reported incidents is smaller. We do note the drop-off at the higher end of actual reported incidents, suggesting the model tends to under-predict in these cases. This is to be expected somewhat, as extremely high counts are rare and the model has not been trained on a dataset containing such high counts.

We conclude this section by noting that kernel ridge regression can provide prediction functions that explain a high proportion of the variance. However, a major limitation of this method in our case is that it would be primarily be useful for interpolation, rather than extrapolation, for the reasons outlined above.

In the subsequent section, we explore generalised additive models, which form a more principled approach for forecasting the number of crimes.

4.3 Generalised Additive Models

Our R6 class `PoissonGAM` is a wrapper for the Poisson family of GAMs, fitted via `gam` from the package `mgcv`. For full details on its use, see the package documentation.

For training the GAMs we will be using data from between 2015 and 2018. For testing we shall use data from 2019.

```
# Download training data
df_train <- load_data(c(2015, 2018), strings_as_factors = TRUE)
# Download test data
df_test <- load_data(2019, strings_as_factors = TRUE)
```

4.3.1 Daily Predictions

We first consider a GAM applied to daily counts of reported crime incidents, aggregated over the entire City of Chicago.

```
# Extract useful information from the date
df_train %<>% convert_dates(as_factors = FALSE, exclude = "hour")
df_test %<>% convert_dates(as_factors = FALSE, exclude = "hour")

# List of features to use in GAM
keep_features <- c("month", "week", "year", "day", "date", "community_area",
                  "beat", "fbi_code", "yday")

# Prep data for use in GAM
gam_train <- df_train %>% select(all_of(keep_features)) %>% na.omit()
gam_test <- df_test %>% select(all_of(keep_features)) %>% na.omit()
```

As each row in the datasets represent a single reported crime incident, we need to sum the number of rows for each day in order to get the number of reported crimes for each day. We use `count` from `dplyr` to achieve this.

```
count_train <- gam_train %>% count(yday, year, date) %>% arrange(year, yday)
count_test <- gam_test %>% count(yday, year, date) %>% arrange(year, yday)

print(as.data.frame(head(count_train)))
```

```
##   yday year      date    n
## 1     1 2015 2015-01-01 1221
## 2     2 2015 2015-01-02  673
## 3     3 2015 2015-01-03  648
## 4     4 2015 2015-01-04  514
## 5     5 2015 2015-01-05  523
## 6     6 2015 2015-01-06  503
```

We separate this dataset into a training set and a test set, and proceed to fit our baseline GAM which has the year and day of the year as predictors. We use a cyclic cubic smoother for the day of the year, and proceed to predict on both the training and test data.

```
# Convert data to matrix and vector format for model training
X_train <- count_train %>% select(-n)
y_train <- count_train$n
X_test <- count_test %>% select(-n)
y_test <- count_test$n

# Initialise new PoissonGAM object
gam1 <- PoissonGAM$new(time_period = "yday", region = NULL, crime_type = NULL,
```

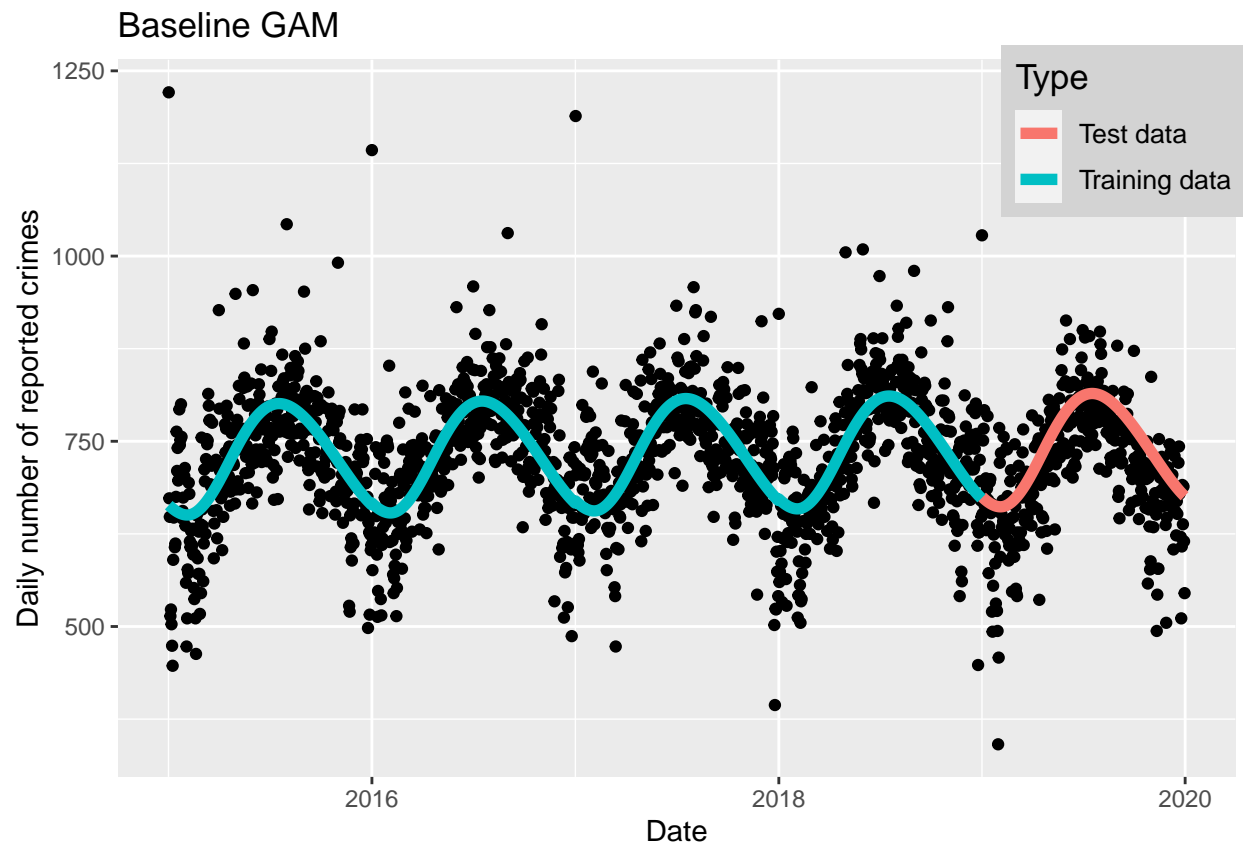
```

include_nb = FALSE, include_year = TRUE)
# Fit the baseline daily GAM
gam1$fit(X = X_train, y = y_train, n_threads = 3)
# Fit summary
gam1$fit_summary

##
## Family: poisson
## Link function: log
##
## Formula:
## n ~ s(as.numeric(yday), bs = "cc") + as.numeric(year)
##
## Parametric coefficients:
##               Estimate Std. Error z value Pr(>|z|)
## (Intercept)   -1.9118027  1.7432595  -1.097    0.273
## as.numeric(year)  0.0042184  0.0008645   4.880 1.06e-06 ***
## ---
## 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(as.numeric(yday)) 5.74      8   5605 <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## R-sq.(adj) =  0.382   Deviance explained = 38.2%
## UBRE = 5.3153   Scale est. = 1           n = 1461

# Predict on entire dataset
pred <- gam1$predict(X = bind_rows(X_train, X_test))
# Store type of each data
count_train$Type <- "Training data"
count_test$Type <- "Test data"
all_dat <- bind_rows(count_train, count_test)
all_dat$Type <- as.factor(all_dat$Type)
# Create plot for baseline model
base_gam_plot <- ggplot(all_dat) +
  geom_point(aes(date, n)) +
  xlab("Date") + ylab("Daily number of reported crimes") +
  geom_line(aes(date, pred, colour = Type), size = 2) +
  ggtitle("Baseline GAM") +
  theme(legend.position = c(0.9, 0.9),
        legend.background = element_rect(fill = "lightgrey", size = 0.5,
                                           linetype = "solid"),
        legend.title = element_text(size = 13),
        legend.text = element_text(size = 10))
base_gam_plot

```



```
# Obtain metrics on test data  
pred_test <- gam1$predict(X = X_test)  
# Initialise vectors for RMSE and R squared values to compare subsequent models  
rmse <- c()  
rsq <- c()  
rmse[1] <- rmse_loss(pred_test, y_test)  
rsq[1] <- r_squared(pred_test, y_test)
```

4.3.2 Feature creation

There are several features that may be of interest. For example: whether a day falls on a weekend; the number of crime incidents reported on the previous day; or the number of crimes on the same day of the previous year.

We first summarise the dates on which the most and fewest crimes are reported, to gain some idea as to what factors might be influential.

```
most_crimes <- count_train %>% arrange(by = desc(n))
most_crimes$date <- as.character(most_crimes$date)
```

```
tab <- xtable::xtable(head(most_crimes[, -5]), )
print(tab, type = "latex", comment = FALSE)
```

	yday	year	date	n
1	1	2015	2015-01-01	1221
2	1	2017	2017-01-01	1189
3	1	2016	2016-01-01	1143
4	213	2015	2015-08-01	1043
5	245	2016	2016-09-01	1031
6	152	2018	2018-06-01	1009

```
least_crimes <- count_train %>% arrange(by = n)
least_crimes$date <- as.character(least_crimes$date)
```

```
tab <- xtable::xtable(head(least_crimes[, -5]), )
print(tab, type = "latex", comment = FALSE)
```

	yday	year	date	n
1	359	2017	2017-12-25	394
2	8	2015	2015-01-08	447
3	359	2018	2018-12-25	448
4	50	2015	2015-02-19	463
5	33	2015	2015-02-02	473
6	73	2017	2017-03-14	473

At first glance, we notice that the highest number of reported crimes often falls on the first day of the year. We also find that the lowest number appears to fall on Christmas day. We constructed the functions `add_prev_day`, `add_dow`, `add_is_fom`, `add_is_christmas` and `add_is_nyd` to attempt to capture some of this information. We proceed to fit a few different GAMs on the resulting data and observe their explanatory and predictive power.

```
# Add features to training data
count_train_new <- count_train %>%
  mutate(week = week(date), month = month(date)) %>%
  add_prev_day() %>% add_dow() %>% add_is_fom() %>%
  add_is_christmas() %>% add_is_nyd()
# Add features to test data
count_test_new <- count_test %>%
  mutate(week = week(date), month = month(date)) %>%
  add_prev_day() %>% add_dow() %>% add_is_fom() %>%
  add_is_christmas() %>% add_is_nyd()

count_full <- bind_rows(count_train_new, count_test_new)
```

```
data_2014 <- load_data(2014, strings_as_factors = TRUE)

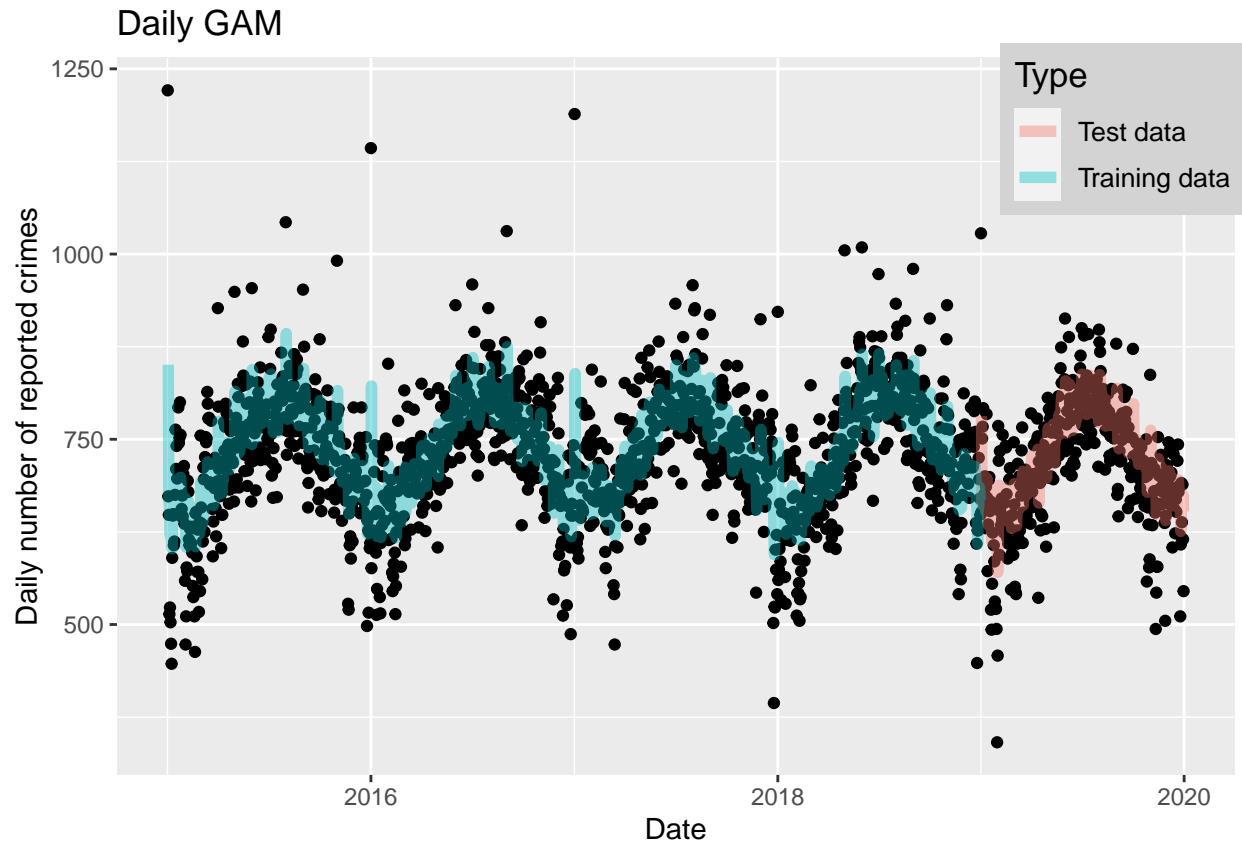
data_2014 <- data_2014 %>%
  convert_dates(exclude = "hour") %>%
  count(year, date) %>%
  arrange(year)
# Fill in previous year manually
count_train_new[1, "n_pre"] <- data_2014[data_2014$date == "2014-12-31", "n"][1,]
count_test_new[1, "n_pre"] <- count_train_new[count_train_new$date == "2018-12-31", "n"]
print(as.data.frame(head(count_train_new)))
```

```
##   yday year      date      n      Type week month n_pre dow is_fom
## 1     1 2015 2015-01-01 1221 Training data     1     1   585    5   TRUE
## 2     2 2015 2015-01-02   673 Training data     1     1  1221    6  FALSE
## 3     3 2015 2015-01-03   648 Training data     1     1   673    7  FALSE
## 4     4 2015 2015-01-04   514 Training data     1     1   648    1  FALSE
## 5     5 2015 2015-01-05   523 Training data     1     1   514    2  FALSE
## 6     6 2015 2015-01-06   503 Training data     1     1   523    3  FALSE
##   is_christmas is_nyd
## 1          FALSE  TRUE
## 2          FALSE FALSE
## 3          FALSE FALSE
## 4          FALSE FALSE
## 5          FALSE FALSE
## 6          FALSE FALSE
```

As PoissonGAM does not currently have functionality to support additional custom features, we fit the following GAMs directly using `gam` from `mgcv`. The same analysis is performed on the following models.

```
library(mgcv)
gam2 <- gam(n ~ s(as.numeric(yday), bs = "cc") + n_pre,
            data = count_train_new, family = "poisson")
gam2_pred <- predict(gam2, type = "response", newdata = count_full)

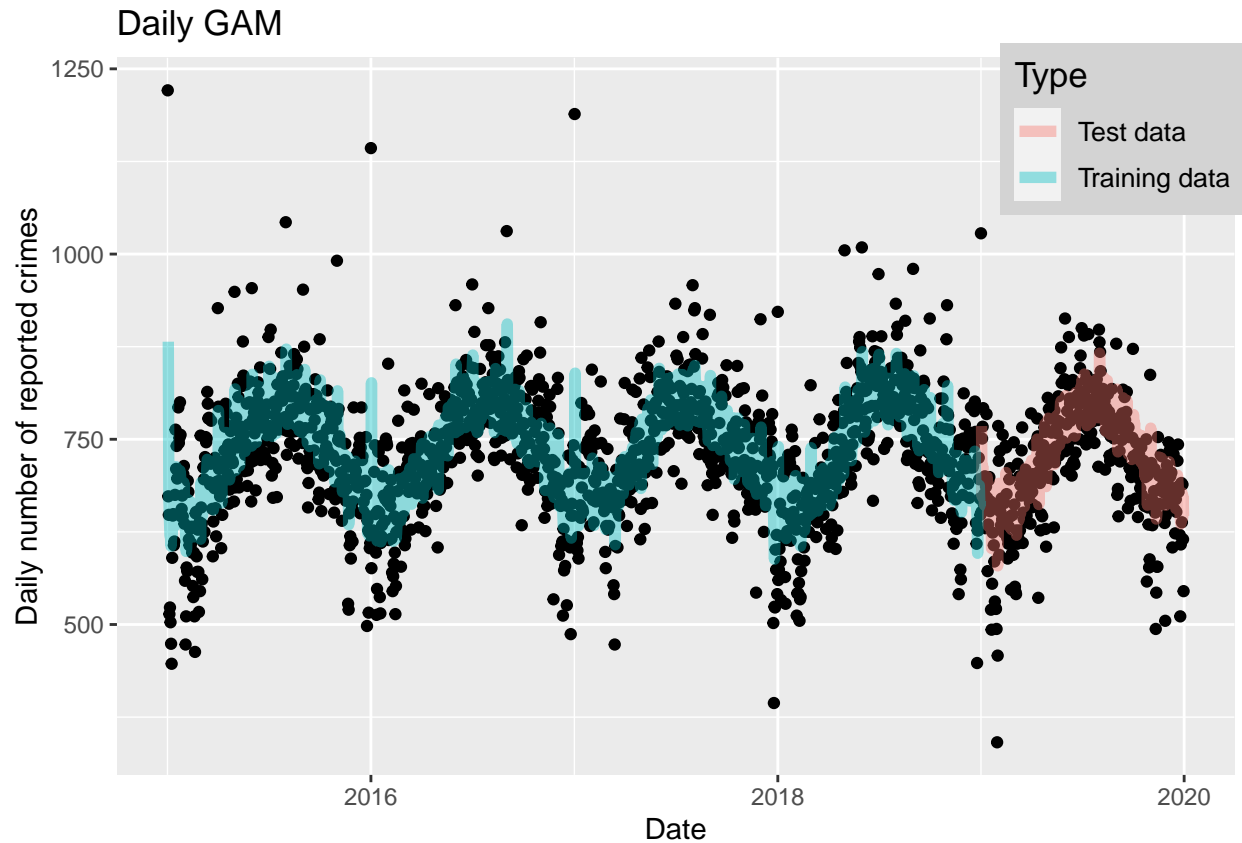
gam2_plot <- ggplot(all_dat) +
  geom_point(aes(date, n)) +
  xlab("Date") + ylab("Daily number of reported crimes") +
  geom_line(aes(date, gam2_pred, colour = Type), size = 2, alpha = 0.4) +
  ggtitle("Daily GAM") +
  theme(legend.position = c(0.9, 0.9),
        legend.background = element_rect(fill = "lightgrey", size = 0.5,
                                           linetype = "solid"),
        legend.title = element_text(size = 13),
        legend.text = element_text(size = 10))
gam2_plot
```



```
# Obtain metrics on test data
gam2_test <- predict(gam2, type = "response", newdata = count_test_new)
rmse[2] <- rmse_loss(gam2_test, count_test_new$n)
rsq[2] <- r_squared(gam2_test, count_test_new$n)

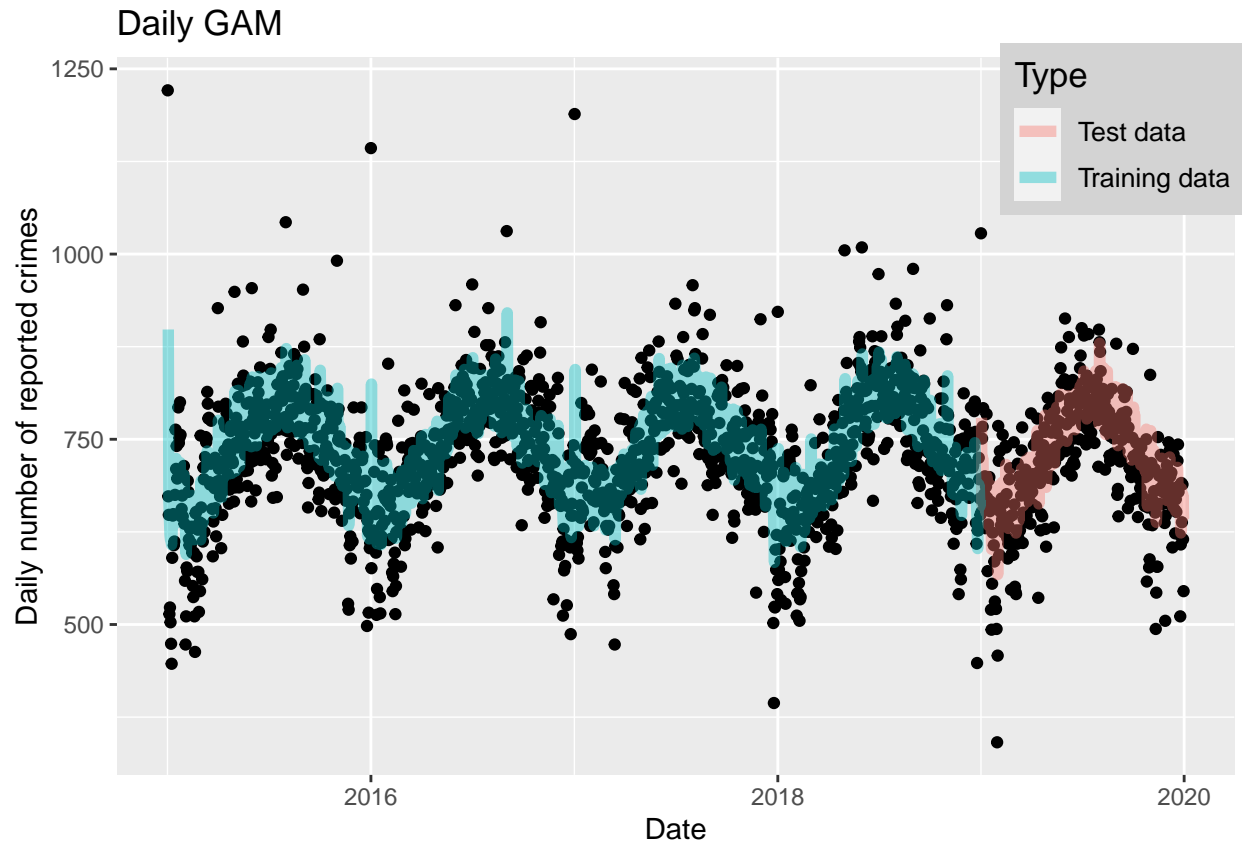
gam3 <- gam(n ~ s(as.numeric(yday), bs = "cc") +
             s(as.numeric(dow), bs = "cr", k = 5) + n_pre,
             data = count_train_new, family = "poisson")
gam3_pred <- predict(gam3, type = "response", newdata = count_full)

gam3_plot <- ggplot(all_dat) +
  geom_point(aes(date, n)) +
  xlab("Date") + ylab("Daily number of reported crimes") +
  geom_line(aes(date, gam3_pred, colour = Type), size = 2, alpha = 0.4) +
  ggtitle("Daily GAM") +
  theme(legend.position = c(0.9, 0.9),
        legend.background = element_rect(fill = "lightgrey", size = 0.5,
                                           linetype = "solid"),
        legend.title = element_text(size = 13),
        legend.text = element_text(size = 10))
gam3_plot
```

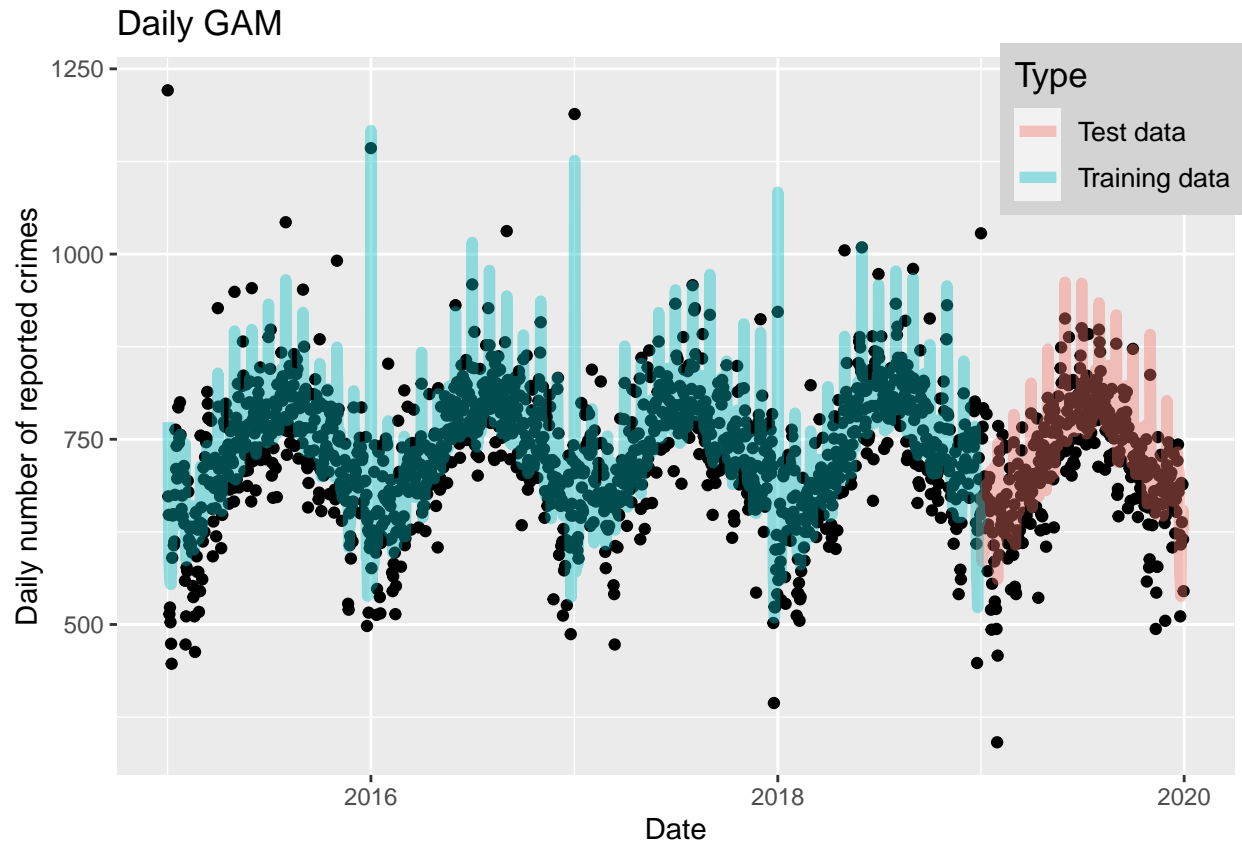
```
# Obtain metrics on test data
gam3_test <- predict(gam3, type = "response", newdata = count_test_new)
rmse[3] <- rmse_loss(gam3_test, count_test_new$n)
rsq[3] <- r_squared(gam3_test, count_test_new$n)

# Fit GAM
gam4 <- gam(n ~ s(as.numeric(yday), bs = "cc") +
            as.factor(dow) + n_pre,
            data = count_train_new, family = "poisson")
# Predict on entire dataset
gam4_pred <- predict(gam4, type = "response", newdata = count_full)
# Plot predictions
gam4_plot <- ggplot(all_dat) +
  geom_point(aes(date, n)) +
  xlab("Date") + ylab("Daily number of reported crimes") +
  geom_line(aes(date, gam4_pred, colour = Type), size = 2, alpha = 0.4) +
  ggtitle("Daily GAM") +
  theme(legend.position = c(0.9, 0.9),
        legend.background = element_rect(fill = "lightgrey", size = 0.5,
                                           linetype = "solid"),
        legend.title = element_text(size = 13),
        legend.text = element_text(size = 10))
gam4_plot
```



```
# Obtain metrics on test data
gam4_test <- predict(gam4, type = "response", newdata = count_test_new)
rmse[4] <- rmse_loss(gam4_test, count_test_new$n)
rsq[4] <- r_squared(gam4_test, count_test_new$n)

gam5 <- gam(n ~ s(as.numeric(yday), bs = "cc") + as.factor(dow) + n_pre +
            is_fom + is_christmas + is_nyd + as.factor(week),
            data = count_train_new, family = "poisson")
# Predict on entire dataset
gam5_pred <- predict(gam5, type = "response", newdata = count_full)
# Plot predictions
gam5_plot <- ggplot(all_dat) +
  geom_point(aes(date, n)) +
  xlab("Date") + ylab("Daily number of reported crimes") +
  geom_line(aes(date, gam5_pred, colour = Type), size = 2, alpha = 0.4) +
  ggtitle("Daily GAM") +
  theme(legend.position = c(0.9, 0.9),
        legend.background = element_rect(fill = "lightgrey", size = 0.5,
                                           linetype = "solid"),
        legend.title = element_text(size = 13),
        legend.text = element_text(size = 10))
gam5_plot
```



```
# Obtain metrics on test data
gam5_test <- predict(gam5, type = "response", newdata = count_test_new)
rmse[5] <- rmse_loss(gam5_test, count_test_new$n)
rsq[5] <- r_squared(gam5_test, count_test_new$n)
gam_metrics <- data.frame(RMSE = rmse, R2 = rsq)
```

```
tab <- xtable::xtable(gam_metrics)
print(tab, type = "latex", comment = FALSE)
```

	RMSE	R2
1	68.72	0.45
2	61.65	0.50
3	59.30	0.54
4	58.49	0.55
5	55.04	0.61

The final model appears to have the best performance, with the lowest RMSE and highest R^2 value. This model takes as input the most features, so is likely able to capture the variance to a greater extent.

4.3.3 Adding Spatial Data

In order to compare our model with the kernel ridge regression model, we aggregate over months and add a spatial component: the community area. `PoissonGAM` has functionality for taking a neighbourhood list along with the community area data, and fits a Markov random field smoother to this factor via `gam`.

```
# Construct training data set
spatial_train <- gam_train %>%
  get_count_data(time_period = "month", region = "community_area")
# Construct test data set
spatial_test <- gam_test %>%
  get_count_data(time_period = "month", region = "community_area")
# Split into data matrix and response vector
X_train <- spatial_train %>% select(-n)
y_train <- spatial_train$n
X_test <- spatial_test %>% select(-n)
metrics <- list(rmse = rmse_loss, r2 = r_squared)

# Initialise PoissonGAM object
mgam <- PoissonGAM$new(time_period = "month", region = "community_area",
  crime_type = NULL, include_nb = TRUE)

# Find CV error for given metrics
mgam_results <- kfold_cv(mgam, X = X_train, y = y_train,
  error_funcs = metrics, k = 5, n_reps = 5,
  parallel = TRUE, n_threads = 5)

names(mgam_results) <- c("RMSE", "R2")
tab <- xtable::xtable(as_tibble(mgam_results), digits = 3)
print(tab, type = "latex", comment = FALSE)
```

	RMSE	R2
1	40.991	0.974
2	40.795	0.974
3	40.599	0.974
4	41.220	0.974
5	40.934	0.974

```
# Re-fit model
mgam$fit(X_train, y_train, n_threads = 7)

# Summary of the fit
mgam$fit_summary

##
## Family: poisson
## Link function: log
##
## Formula:
## n ~ s(as.numeric(month), bs = "cc") + s(community_area, bs = "mrf",
##     xt = list(nb = self$nb_list))
##
## Parametric coefficients:
##             Estimate Std. Error z value Pr(>|z|)
## (Intercept)  5.313327   0.001396   3806   <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(as.numeric(month))  7.974      8   7145  <2e-16 ***
## s(community_area)    75.984     76 585373  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## R-sq.(adj) =  0.975   Deviance explained = 98.1%
## UBRE = 2.5478   Scale est. = 1           n = 3696
```

Notice that this model explains 98.1% of the variance in the training data.

```
# Predict on test data
spatial_test$predicted <- mgam$predict(X = X_test)

# Compute evaluation metrics
rmse <- rmse_loss(spatial_test$n, spatial_test$predicted)
rsq <- r_squared(spatial_test$n, spatial_test$predicted)
tab <- xtable::xtable(data.frame("RMSE" = rmse, "R2" = rsq))
print(tab, type = "latex", comment = FALSE)
```

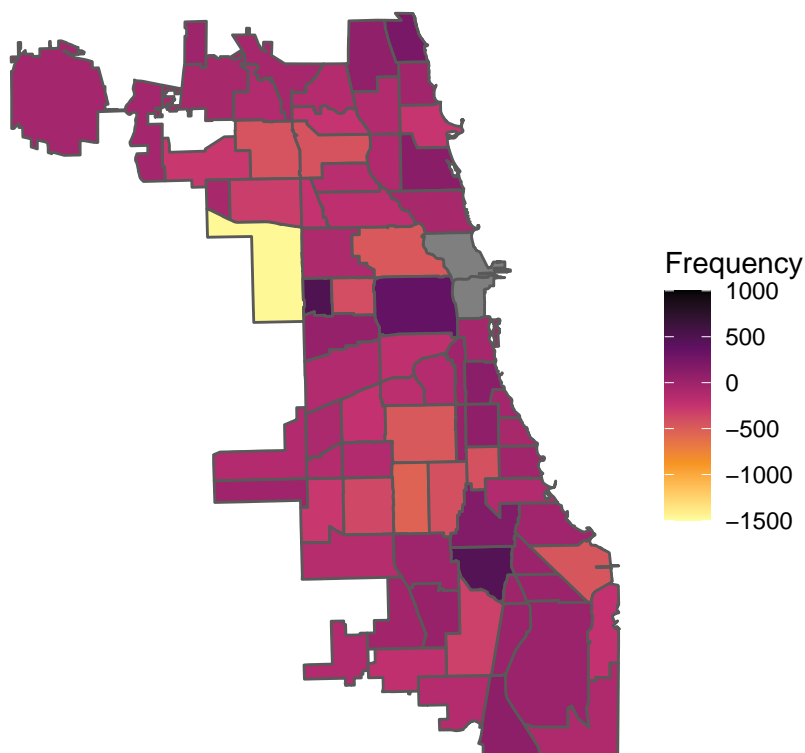
	RMSE	R2
1	38.83	0.98

To illustrate the predictive power of the fitted model over the community areas, we sum the true number of crimes reported in the test dataset and the predicted number generated by the model. Below we produce a map of the residuals.

```
data("community_bounds")
com_area <- spatial_test %>%
  group_by(community_area) %>%
  summarise(n = sum(n), n_pred = sum(predicted))
com_area$residual <- (com_area$n - com_area$n_pred)
community_bounds$residual <- com_area$residual

by_area <- ggplot(community_bounds) + geom_sf(aes(fill = residual)) +
  ggtitle("Residuals Produced by Community Area") +
  colorspace::scale_fill_continuous_sequential(5, palette = "Inferno") +
  theme_void() +
  guides(fill = guide_colorbar(title = "Frequency")) +
  scale_fill_gradientn(colors = colorspace::sequential_hcl(5, palette = "Inferno")[5:1],
    breaks = seq(-1500, 1000, 500), limits = c(-1500, 1000))
by_area
```

Residuals Produced by Community Area



In most areas, the model performs very well with relatively low errors. However, there are specific areas in which the predicted counts are far higher or lower than the actual counts.

4.3.4 Adding Categorical Data

On further exploratory analysis we found that adding the neighbourhood list did not significantly improve upon the model which took community area as a factor with no neighbourhood list provided. In this final section we fit a GAM without the neighbourhood list, with an additional factor indicating the type of crime reported.

```
# Construct training data set
spatial_fbi_train <- gam_train %>%
  get_count_data(time_period = "month", region = "community_area",
                 crime_type = "fbi_code")
# Construct test data set
spatial_fbi_test <- gam_test %>%
  get_count_data(time_period = "month", region = "community_area",
                 crime_type = "fbi_code")
# Split into data matrix and response vector
X_train <- spatial_fbi_train %>% select(-n)
y_train <- spatial_fbi_train$n
X_test <- spatial_fbi_test %>% select(-n)

# Initialise PoissonGAM
mgam_com_fbi <- PoissonGAM$new(time_period = "month", region = "community_area",
                              crime_type = "fbi_code", include_nb = FALSE)
# Find CV error for given metrics
```

```
mgam_fbi_results <- kfold_cv(mgam_com_fbi, X = X, y = y, error_funcs = metrics,
                             k = 5, n_reps = 5, parallel = TRUE, n_threads = 5)
names(mgam_fbi_results) <- c("RMSE", "R2")

tab <- xtable::xtable(as_tibble(mgam_fbi_results), digits = 4)
print(tab, type = "latex", comment = FALSE)
```

	RMSE	R2
1	16.8159	0.7177
2	16.8304	0.7167
3	16.8168	0.7166
4	16.8706	0.7156
5	16.8499	0.7165

```
# Fit full model
mgam_com_fbi$fit(X_train, y_train, n_threads = 7)

spatial_fbi_test$predicted <- mgam_com_fbi$predict(X = X_test)

# Compute evaluation metrics
rmse <- rmse_loss(spatial_fbi_test$n, spatial_fbi_test$predicted)
rsq <- r_squared(spatial_fbi_test$n, spatial_fbi_test$predicted)
tab <- xtable::xtable(data.frame("RMSE" = rmse, "R2" = rsq))
print(tab, type = "latex", comment = FALSE)
```

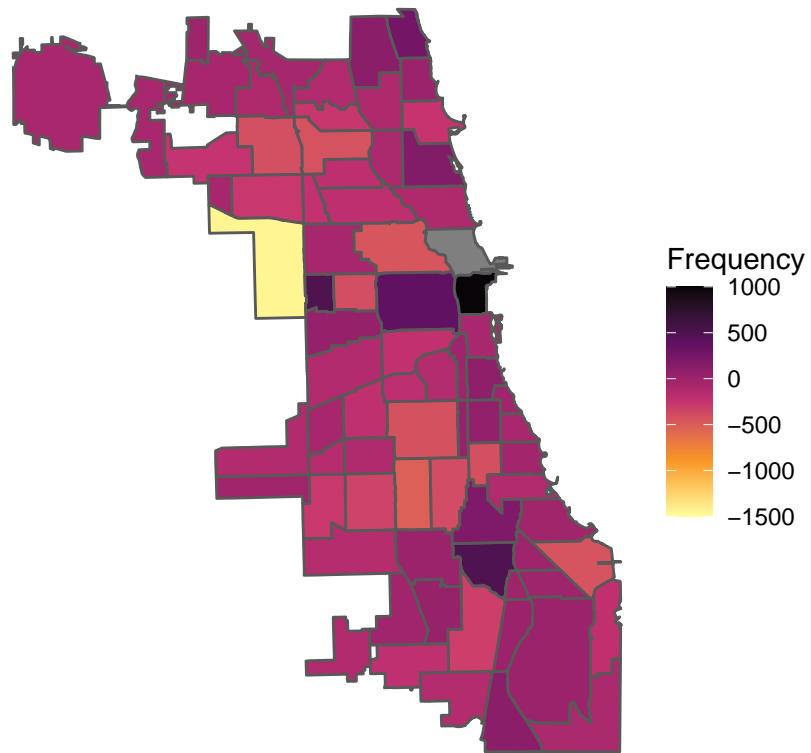
	RMSE	R2
1	17.63	0.70

The updated model has a lower RMSE value but also a lower R^2 value. Below we produce a map in the same way as the previous section.

```
com_area <- spatial_fbi_test %>% group_by(community_area) %>%
  summarise(n = sum(n), n_pred = sum(predicted))
com_area$residual <- (com_area$n - com_area$n_pred)
community_bounds$residual <- com_area$residual

# Plot the residuals for each Community Area
by_area <- ggplot(community_bounds) +
  geom_sf(aes(fill = residual)) +
  ggtitle("Residuals Produced by Community Area") +
  colorspace::scale_fill_continuous_sequential(5, palette = "Inferno") +
  theme_void() +
  guides(fill = guide_colorbar(title = "Frequency")) +
  scale_fill_gradientn(colors = colorspace::sequential_hcl(5, palette = "Inferno")[5:1],
                      breaks = seq(-1500, 1000, 500), limits = c(-1500, 1000))
by_area
```

Residuals Produced by Community Area



The maps are quite similar, which may indicate little improvement in the model from adding the categorical component.