

W271 Assignment 7

Question-1: AIC and BIC and “Stringency”

(4 points) Part-1

In the async lecture, Jeffrey says “BIC is in general more stringent than AIC or AICc”. Let’s illustrate that and reason about it.

1. Produce a dataset, `d`, that includes 100 observations of pure white-noise.
 - The outcome variable should be a variable `y` that has 100 draws from `rnorm`, with `mean=0` and `sd=1`.
 - The input variables should be variables `x1 ... x10` that are also 100 draws from `rnorm` each with `mean=0` and `sd=1`.
 - There are fancy ways to write this code; the goal for this isn’t to place a clever coding task in front of you, so feel free to use copy-paste to create the data object in any way that you can.
2. After producing data, fit 11 models against that data, stored as `model0` through `model10`. (The number appended to `model` corresponds to the number of parameters that you have used in your estimation).
3. After estimating your models, create a new dataset, `results_data`, that contains the number of parameters that you have used in an estimation, and the AIC and BIC values that you calculated for that number of parameters.
 1. Note – this is another place where the way that you create the data, and the way that the data is the most useful to use are incongruent.
 2. When we created the data, we created a dataset that has a column called `parameters`, a column called `aic` and a column called `bic`.
 3. However, it is much more useful to have “tidy” data that has these values stacked. If you find yourself creating the dataset in the “wide” form that we have described above, you can use the `dplyr::pivot_longer` function to pivot your data into a tidy format. Specifically, we used this call `pivot_longer(cols = c('aic', 'bic'))` with our input data structure.
4. Finally, produce a plot that shows the AIC and BIC values on the y-axis and the number of estimated parameters on the x-axis. In the subtitle to your plot, note whether a relatively higher or lower AIC or BIC means that a model is performing better or worse (i.e. either “Higher values are better” or “Lower values are better”). What do you notice about these plots, and what does this tell you about the “stringency” of AIC vs. BIC?

```
d <- cbind.data.frame(rnorm(100), rnorm(100),
                      rnorm(100), rnorm(100),
                      rnorm(100), rnorm(100),
                      rnorm(100), rnorm(100),
                      rnorm(100), rnorm(100), rnorm(100))
names(d) <- c("y", "x1", "x2", "x3", "x4", "x5", "x6", "x7", "x8", "x9", "x10")
```

```
model0 <- lm(formula = y ~ rep(mean(d$y), 100), data = d)
model1 <- lm(formula = y ~ x1, data = d)
model2 <- lm(formula = y ~ x1 + x2, data = d)
model3 <- lm(formula = y ~ x1 + x2 + x3, data = d)
model4 <- lm(formula = y ~ x1 + x2 + x3 + x4, data = d)
model5 <- lm(formula = y ~ x1 + x2 + x3 + x4 + x5, data = d)
model6 <- lm(formula = y ~ x1 + x2 + x3 + x4 + x5 + x6, data = d)
```

```

model7 <- lm(formula = y ~ x1 + x2 +x3+ x4 + x5 +x6 + x7, data = d)
model8 <- lm(formula = y ~ x1 + x2 +x3+ x4 + x5 +x6 + x7 + x8, data = d)
model9 <- lm(formula = y ~ x1 + x2 +x3+ x4 + x5 +x6 + x7 + x8 + x9, data = d)
model10<-lm(formula = y ~ x1 + x2 +x3+ x4 + x5 +x6 + x7 + x8 + x9 + x10, data=d)

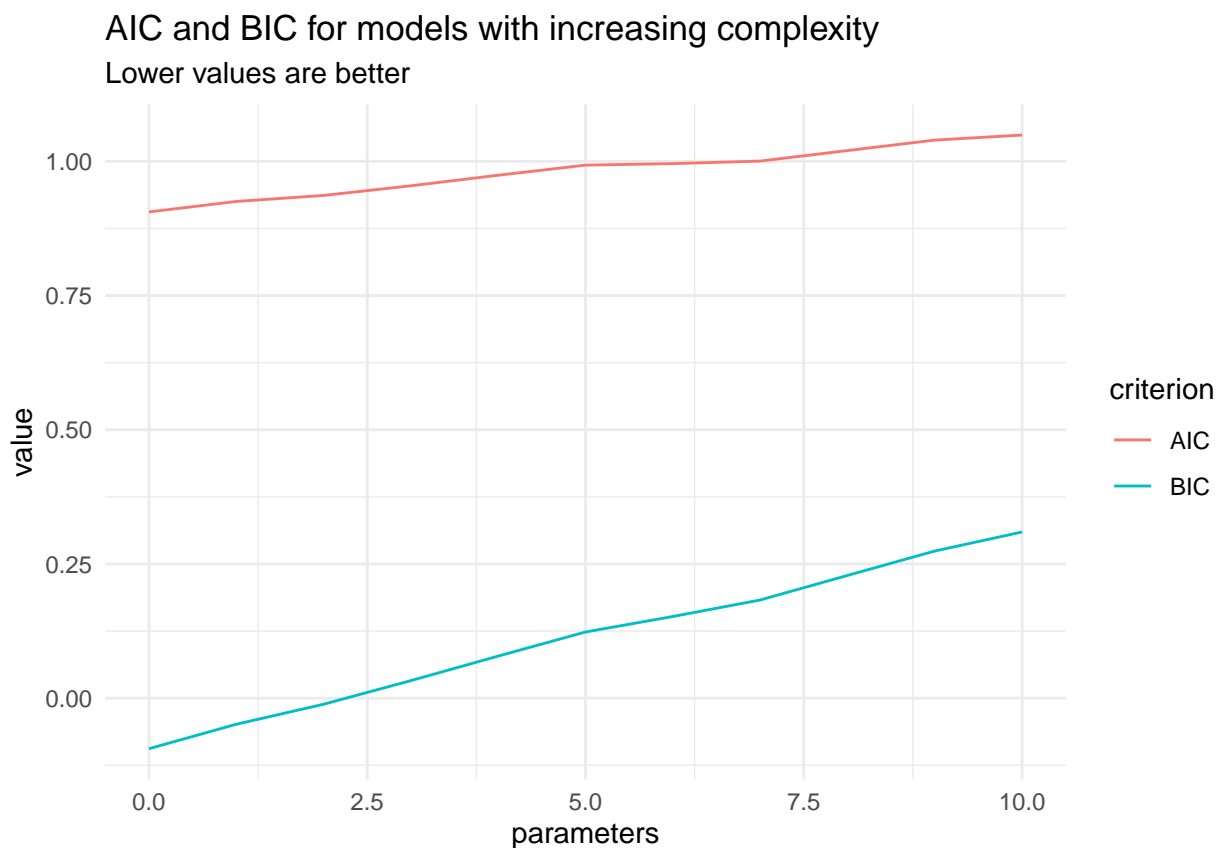
models <- list(model0, model1, model2, model3, model4, model5, model6, model7,
               model8, model9, model10)
results_data <- data.frame(matrix(nrow = 0, ncol = 3))
colnames(results_data) <- c("parameters", "criterion", "value")
counter = 0
for (model in models){
  AIC <- 1 + (2 * counter / 100) + log(sum(model$residuals * model$residuals)/100)
  results_data[nrow(results_data) + 1,] <- list(counter, "AIC", AIC)
  BIC <- (counter *log(100) / 100) + log(sum(model$residuals * model$residuals)/100)
  results_data[nrow(results_data) + 1,] <- list(counter, "BIC", BIC)
  counter = counter+1
}

```

```

results_data %>% ggplot() + aes(x = parameters, y = value, colour = criterion) +
  geom_line() + labs(title = "AIC and BIC for models with increasing complexity", subtitle = "Lower values are better")

```



What do you note? Fill this in! Adding parameters beyond just constant inflates both criteria. BIC appears to be a more stringent criterion because it reacts much stronger to adding meaningless parameters. While AIC only increased by 10% upon adding 10 meaningless parameters, BIC went up almost by factor of 5.

(2 points) Part-2

Now, suppose that you had data that, *in the population model* actually held a relationship between the input features and the outcome feature. Specifically, suppose that for every unit increase in x_1 there was a 0.1 increase in the outcome, for every unit increase in x_2 there was a 0.2 increase in the outcome, ..., for every unit increase in x_{10} there was a 1.0 unit increase in the outcome. Suppose that if all $x_1 \dots x_{10}$ were zero, that the outcome would have an expectation of zero, but with white-noise around it with $\mu = 0$ and $\sigma = 1$.

- Modify the code that you wrote above to create data according to this schedule.
- Estimate 11 models as before.
- Produce a new dataset `results_data` that contains the AIC and BIC values from each of these models.
- Produce the same plot as you did before with the white noise series. Comment on what, if anything is similar or different between this plot, and the plot you created before.

```
coeffs <- seq(from = 0, to = 1, by = 0.1)
d_trend <- as_tibble( data.frame(matrix(nrow=100,ncol=0)) )
counter = 1
for (k in seq(1:10)){
  d_trend[as.character(k)] <- rnorm(100)
}

d_trend["y"] <- rnorm(100)
for (k in seq(1:10)){
  d_trend["y"] <- d_trend["y"] + 0.1*k*d_trend[as.character(k)]
}

colnames(d_trend) <- c("x1", "x2", "x3", "x4", "x5", "x6", "x7", "x8", "x9", "x10", "y")
```

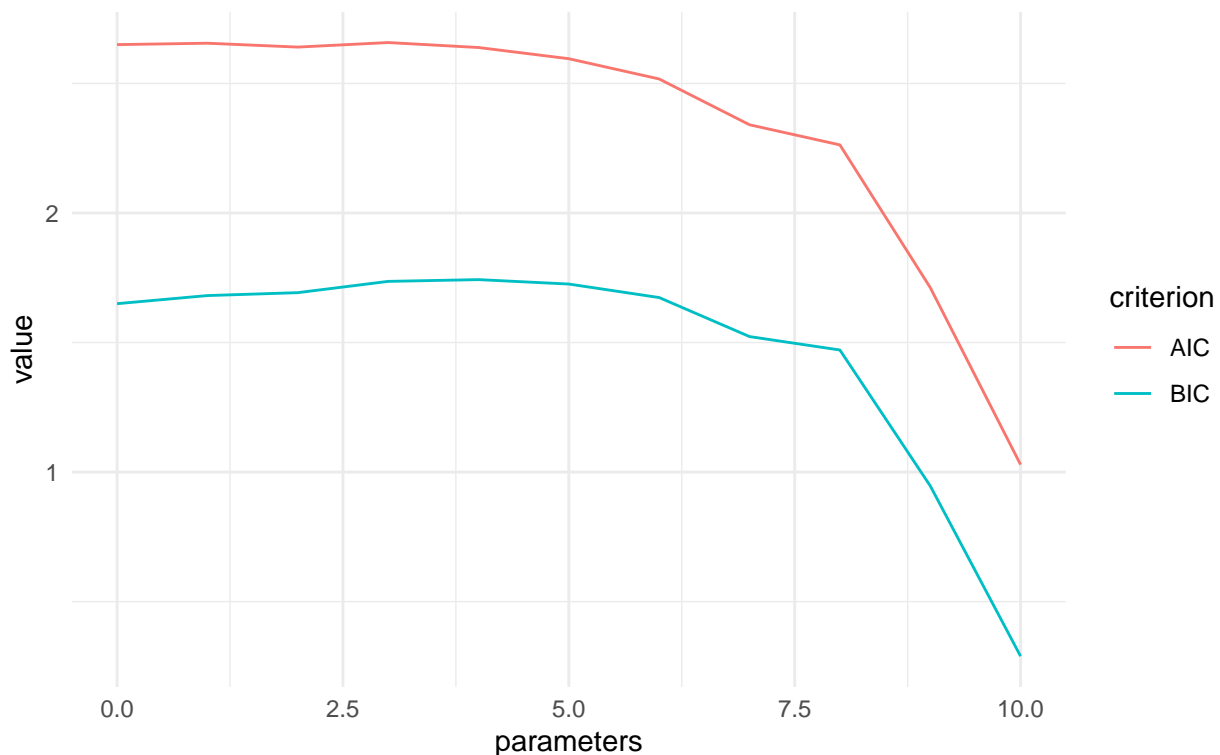
```
model0 <- lm(formula = y ~ rep(mean(d_trend$y), 100), data = d_trend)
model1 <- lm(formula = y ~ x1, data = d_trend)
model2 <- lm(formula = y ~ x1 + x2, data = d_trend)
model3 <- lm(formula = y ~ x1 + x2 +x3, data = d_trend)
model4 <- lm(formula = y ~ x1 + x2 +x3 + x4, data = d_trend)
model5 <- lm(formula = y ~ x1 + x2 +x3 + x4+ x5, data = d_trend)
model6 <- lm(formula = y ~ x1 + x2 +x3 + x4 + x5 +x6, data = d_trend)
model7 <- lm(formula = y ~ x1 + x2 +x3+ x4 + x5 +x6 + x7, data = d_trend)
model8 <- lm(formula = y ~ x1 + x2 +x3+ x4 + x5 +x6 + x7 + x8, data = d_trend)
model9 <- lm(formula = y ~ x1 + x2 +x3+ x4 + x5 +x6 + x7 + x8 + x9, data = d_trend)
model10<-lm(formula = y ~ x1 + x2 +x3+ x4 + x5 +x6 + x7 + x8 + x9 + x10, data=d_trend)
```

```
models <- list(model0, model1, model2, model3, model4, model5, model6, model7,
              model8, model9, model10)
results_data_trend <- data.frame(matrix(nrow = 0, ncol = 3))
colnames(results_data_trend) <- c("parameters", "criterion", "value")
counter = 0
for (model in models){
  AIC <- 1 + (2 * counter / 100) + log(sum(model$residuals * model$residuals)/100)
  results_data_trend[nrow(results_data_trend) + 1,] <- list(counter, "AIC", AIC)
  BIC <- (counter *log(100) / 100) + log(sum(model$residuals * model$residuals)/100)
  results_data_trend[nrow(results_data_trend) + 1,] <- list(counter, "BIC", BIC)
  counter = counter+1
}
```

```
results_data_trend %>% ggplot() + aes(x = parameters, y = value, colour = criterion) +
  geom_line() + labs(title = "AIC and BIC for models with increasing complexity", subtitle = "Lower val
```

AIC and BIC for models with increasing complexity

Lower values are better



What do you notice about this plot and the model performance overall? Adding parameters that affect the outcome strongly influences the value for both criteria, while adding weakly affecting parameters does not change AIC and BIC much. BIC appears to be a more stringent criterion because it takes even more meaningful parameters to make a significant change in BIC.

Question-2: Weather in NYC

Our goals with this question are to:

- (If necessary) Clean up code that we've written before to re-use. This task of writing code, and then coming back and using it later is often overlooked in the MIDS program. Here's a chance to practice!
- Estimate several different polynomial regressions against a time series and evaluate at what point we have produced a model with "enough complexity" that the model evaluation scores cease to tell us that additional model parameters are improving the model fit.

(1 point) Part-1: Load the Weather Data

Load the weather data in the same way as you did in the previous assignment, recalling that there was some weird duplication of data for one of the days. Then, create an object, `weather_weekly` that aggregates the data to have two variables `average_temperature` and `average_dewpoint` at the year-week level, for each airport. After your aggregation is complete, you should have a `tsibble` that has the following shape:

```
A tsibble: 159 x 4 [1W]
# Key:      origin [3]
  origin week_index average_temperature average_dewpoint
  <chr>    <week>          <dbl>             <dbl>
1 EWR     2013 W01         34.3              19.4
```

2	EWR	2013 W02	42.7	33.3
3	EWR	2013 W03	39.6	26.5

```

# Load the data
weather <- nycflights13::weather

# Make time index
weather$time_index <- make_datetime(year = weather$year,
                                   month = weather$month,
                                   day = weather$day,
                                   hour = weather$hour)

# Remove duplicates by taking a median among them, omit empty rows
weather <- weather %>%
  group_by(origin, time_index) %>%
  summarise(across(5:6, ~median(.x))) %>% na.omit()

# Coerce into tsibble
weather_ts <- as_tsibble(weather, key = origin, index = time_index)

# Create weekly aggregation
weather_weekly <- weather_ts %>%
  index_by(week = ~week(.)) %>% summarise(across(2:3, ~mean(.x)))

weather_weekly %>% head()

## # A tsibble: 6 x 4 [1]
## # Key:      origin [1]
##   origin week temp dewp
##   <chr>  <dbl> <dbl> <dbl>
## 1 EWR      1  35.1  20.0
## 2 EWR      2  44.2  35.9
## 3 EWR      3  36.5  21.9
## 4 EWR      4  21.7   5.9
## 5 EWR      5  35.4  23.4
## 6 EWR      6  31.5  21.4

```

(2 points) Part-2: Fit Polynomial Regression Models

For each of the `average_temperature` and `average_dewpoint` create ten models that include polynomials of increasing order.

- One issue that you're likely to come across is dealing with how to make the time index that you're using in your `tsibble` work with either `poly` or some other function to produce the polynomial terms; this arises because although the time index is ordered, it isn't really a "numeric" feature so when you call for something like, `poly(week_index, degree=2)` you will be met with an error.
- Cast the index to a numeric variable, where the first week is indexed to be 0. Recall that Jeffrey notes that this form of translation only changes the way that the intercept is interpreted; we will note that because the `as.numeric(week_index)` creates input variables that are in the vicinity, it also changes the magnitude of the higher-order polynomial terms that are estimated, though it does not change the regression diagnostics and model scoring to transform (or not) these time index variables.

Additionally, you might recall that in 203, we actually recommended you away from using the `poly` function. That was a recommendation based on students' knowledge at the time, when we were considering fitting log and square root transformations of data. At this point, you can handle the additional complexity and can take the recommendation that `poly` is nice for working with polynomial translations of time.

```

weather_weekly$week <- weather_weekly$week - 1

degrees <- seq(1:10)
models_temp <- list()
models_dew <- list()

for (degree in degrees){
  models_temp[[degree]] <- lm(formula = temp ~ poly(week, degree=degree), data = weather_weekly)
  models_dew[[degree]] <- lm(formula = dewp ~ poly(week, degree=degree), data = weather_weekly)
}

```

(2 points) Part-3: Evaluate the model fits best for each outcomes

For each of the outcomes – `average_temperature` at the weekly level, and `average_dewpoint` at the weekly level – make an assessment based on either AIC or BIC for why one polynomial degree produces the best fitting model. In doing so, describe why you have chosen to use either AIC or BIC, what the particular scoring of this metric is doing (i.e. write the formula, and explain to your reader what is happening in that formula). Especially compelling in producing your argument for why you prefer a particular model form is to create a plot of the polynomial degree on the x-axis and the metric score on the y-axis.

```

temp_IC <- data.frame(matrix(nrow = 0, ncol = 3))
colnames(temp_IC) <- c("degrees", "criterion", "value")
counter = 1
for (model in models_temp){
  AIC <- 1 + (2 * counter / 100) + log(sum(model$residuals * model$residuals)/100)
  temp_IC[nrow(temp_IC) + 1,] <- list(counter, "AIC", AIC)
  BIC <- (counter * log(100) / 100) + log(sum(model$residuals * model$residuals)/100)
  temp_IC[nrow(temp_IC) + 1,] <- list(counter, "BIC", BIC)
  counter = counter+1
}

```

```

temp_plot <- temp_IC %>% ggplot() +
  aes(x = degrees, y = value, colour = criterion) +
  geom_point() + geom_line() + labs(title = "Temp models") +
  theme(legend.position = "none") +
  scale_x_continuous(breaks=c(1, 2, 3, 4, 5, 6, 7, 8, 9, 10))

```

```

dew_IC <- data.frame(matrix(nrow = 0, ncol = 3))
colnames(dew_IC) <- c("degrees", "criterion", "value")
counter = 1
for (model in models_dew){
  AIC <- 1 + (2 * counter / 100) + log(sum(model$residuals * model$residuals)/100)
  dew_IC[nrow(dew_IC) + 1,] <- list(counter, "AIC", AIC)
  BIC <- (counter * log(100) / 100) + log(sum(model$residuals * model$residuals)/100)
  dew_IC[nrow(dew_IC) + 1,] <- list(counter, "BIC", BIC)
  counter = counter+1
}

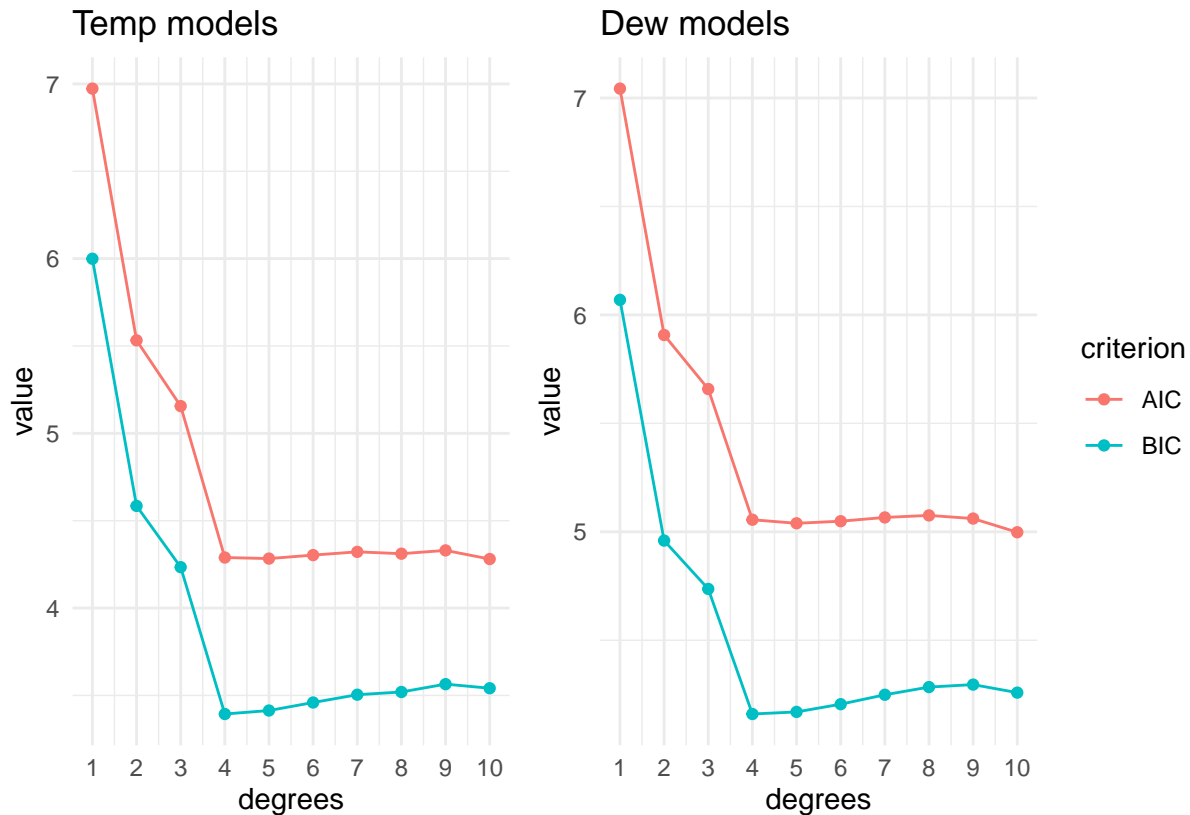
```

```

dew_plot <- dew_IC %>% ggplot() +
  aes(x = degrees, y = value, colour = criterion) +
  geom_point() + geom_line() + labs(title = "Dew models") +
  scale_x_continuous(breaks=c(1, 2, 3, 4, 5, 6, 7, 8, 9, 10))

```

```
temp_plot | dew_plot
```



Looking at these two BIC scoring criteria there seems to be a clear **lack** of improvement beyond a polynomial order of four. *Perhaps* moving from four to five would still increase the model's performance, but it is small compared to the polynomials 2-4. For us, if we were fitting this model, we would be likely to stop at `poly(, degree = 4)`. `> BIC` appears to be a preferred method of assesment as it shows a clear minimum *in the graph above, while AIC hardly distinguhes between polynomes above 4'th degree. Using BIC as a criterion, we can see that polynome of 4'th degree has the minimum IC value and thus is the optimal choice.

Question-3: Smooth Moves

In the async lecture, Jeffrey proposes four different smoothers that might be used:

1. **Moving Average:** These moving average smoothers can be either symmetric or, often preferred, backward smoothers. Please use a backward smoother, and make the choice about the number of periods based off of some evaluation of different choices. You might consult [this page] in *Forecasting Principles and Practice 3*.
2. **Regression Smoothers:** Please use the polynomial regression that you stated you most preferred from your BIC analysis to the last question.
3. (Optional) **Spline Smoothers:** There is a reading in the repository that provides some more background (it is a review from 2019) on using spline smoothers. The current implementation that we prefer in R is the `splines2` library. For your spline smoother, use the `splines2::naturalSpline` function. Once you have fitted this spline, you can use the `predict` method to produce values. A good starting place for this is [here]. We'll note that this is the most challenging of the smoothers to get running in this assignment, and so getting it running successfully is optional.
4. **Kernel Smoothers.:** Please use the `ksmooth` function that is available to you in the `stats` library. Because `stats` is always loaded in R, we have not referred to it using the `::` notation.

(6 points, with 2 optional) Part-1: Create Smoothers

With the weekly weather data that you used for the previous question, produce a smoothed variable for `average_temperature` and `average_dewpoint` using each of the four smoothers described in the `async`. Three smoothers are required of this question – (1) Moving Average; (2) Regression Smoothers; and, (3) Kernel Smoothers. The fourth, splines, is optional but if you produce a spline smoother that is working effectively, you can earn two bonus points. (Note that the homework maximum score is still 100%.)

When you are done with this task, you should have created eight new variables that are each a smoothed version of this series.

For each smoother that you produce:

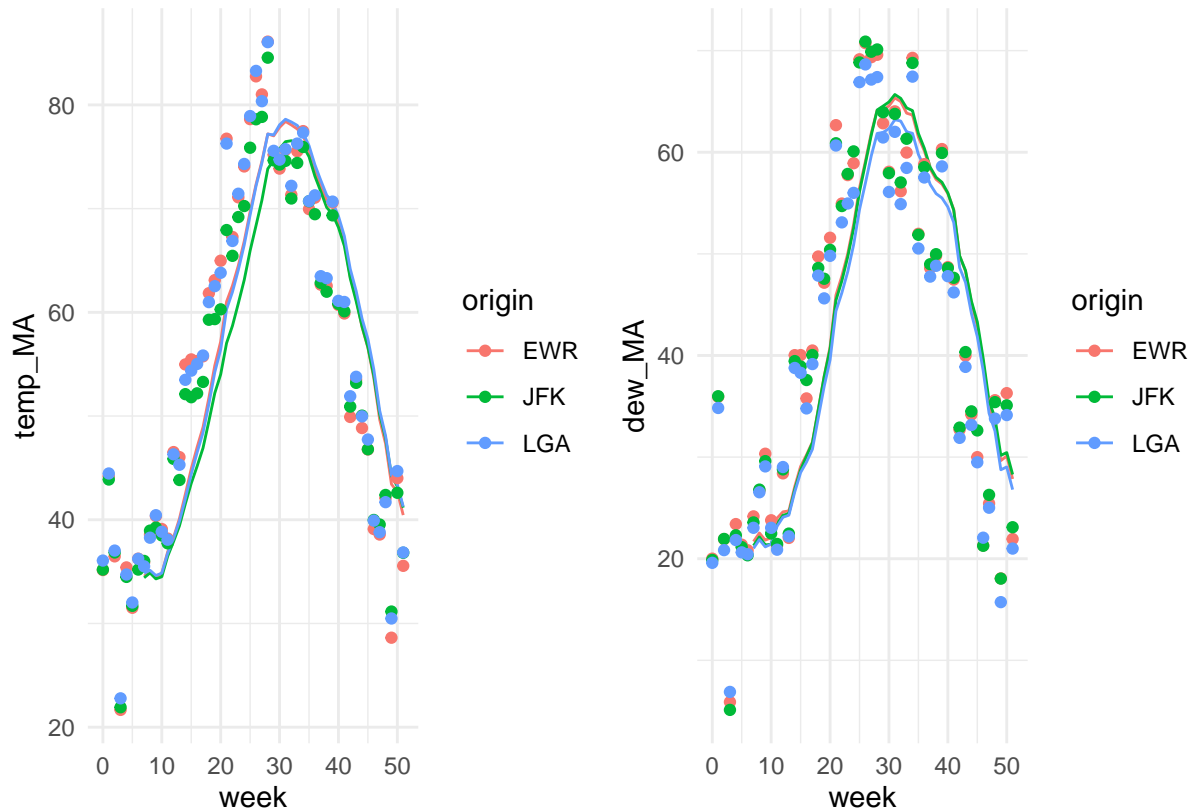
- Fit the smoother **within** each origin. That is, fit the smoother for JFK separately from LaGuardia and Newark.
- Attach the values that are produced by the smoother onto the `weekly_weather` dataframe.
- Produce a plot that shows the original data as `geom_point()`, and the smoother's predictions as `geom_line()`.
- Your goal is not to produce **any** smoother, but instead, for each class of smoother, the version that is do the best job that is possible by this smoother. That is, you are working through the hyperparameters to these algorithms to produce their most effective output.

```
weather_weekly <- weather_weekly %>%
  group_by(origin) %>%
  mutate(
    `temp_MA` = slider::slide_dbl(temp, mean,
                                   .before = 7, .after = 0, .complete = TRUE),
    `dew_MA` = slider::slide_dbl(dewp, mean,
                                   .before = 7, .after = 0, .complete = TRUE)
  )

temp_plot_av <- weather_weekly %>% ggplot() +
  geom_line(aes(x = week, y = temp_MA , colour = origin)) +
  geom_point(aes(x = week, y = temp , colour = origin)) +
  theme(legend.position = "none") +
  theme_minimal()

dew_plot_av <- weather_weekly %>% ggplot() +
  geom_line(aes(x = week, y = dew_MA , colour = origin)) +
  geom_point(aes(x = week, y = dewp , colour = origin)) +
  theme_minimal()

temp_plot_av | dew_plot_av
```

> Window of 7 points was selected as it seems to strike a balance between removing noise and preserving features of the series.

```
temp_JFK <- weather_weekly %>%
  filter(origin == "JFK") %>%
  mutate(
    `temp_poly_4` = predict.lm(lm(formula = temp ~ poly(week, 4), data = .), .),
    `dew_poly_4` = predict.lm(lm(formula = dewp ~ poly(week, 4), data = .), .)
  )

temp_EWR <- weather_weekly %>%
  filter(origin == "EWR") %>%
  mutate(
    `temp_poly_4` = predict.lm(lm(formula = temp ~ poly(week, 4), data = .), .),
    `dew_poly_4` = predict.lm(lm(formula = dewp ~ poly(week, 4), data = .), .)
  )

temp_LGA <- weather_weekly %>%
  filter(origin == "LGA") %>%
  mutate(
    `temp_poly_4` = predict.lm(lm(formula = temp ~ poly(week, 4), data = .), .),
    `dew_poly_4` = predict.lm(lm(formula = dewp ~ poly(week, 4), data = .), .)
  )

weather_weekly <- rbind.data.frame(temp_JFK, temp_EWR, temp_LGA)

temp_plot_poly <- weather_weekly %>% ggplot() +
  geom_line(aes(x = week, y = temp_poly_4, colour = origin)) +
```

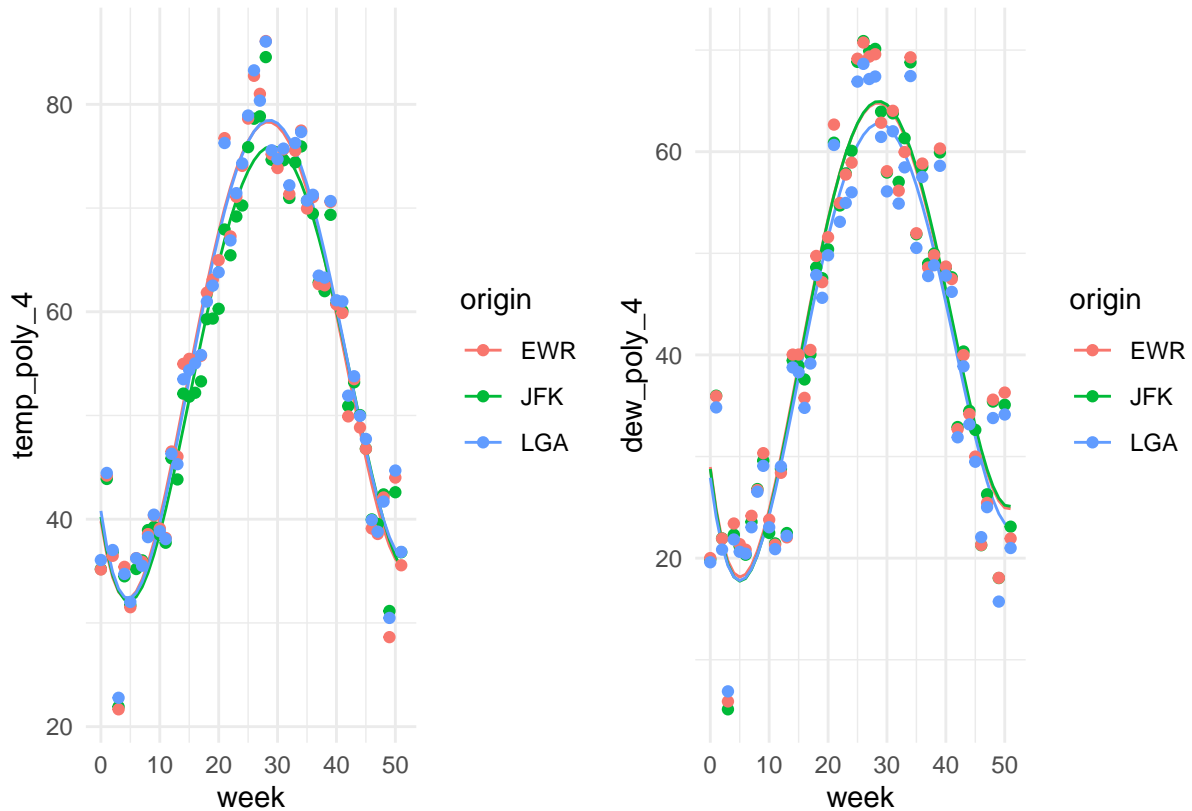
```

geom_point(aes(x = week, y = temp , colour = origin)) +
theme(legend.position = "none") +
theme_minimal()

dew_plot_poly <- weather_weekly %>% ggplot() +
  geom_line(aes(x = week, y = dew_poly_4 , colour = origin)) +
  geom_point(aes(x = week, y = dewp , colour = origin)) +
  theme_minimal()

temp_plot_poly |> dew_plot_poly

```



> 7'th degree polynome was selected based on the previous BIC analysis.

```

weather_weekly <- weather_weekly %>%
  group_by(origin) %>%
  mutate (
    `temp_kern` = ksmooth(week, temp, kernel = "normal", bandwidth = 7, n.points = length(week))$y,
    `dew_kern` = ksmooth(week, dewp, kernel = "normal", bandwidth = 7, n.points = length(week))$y
  )

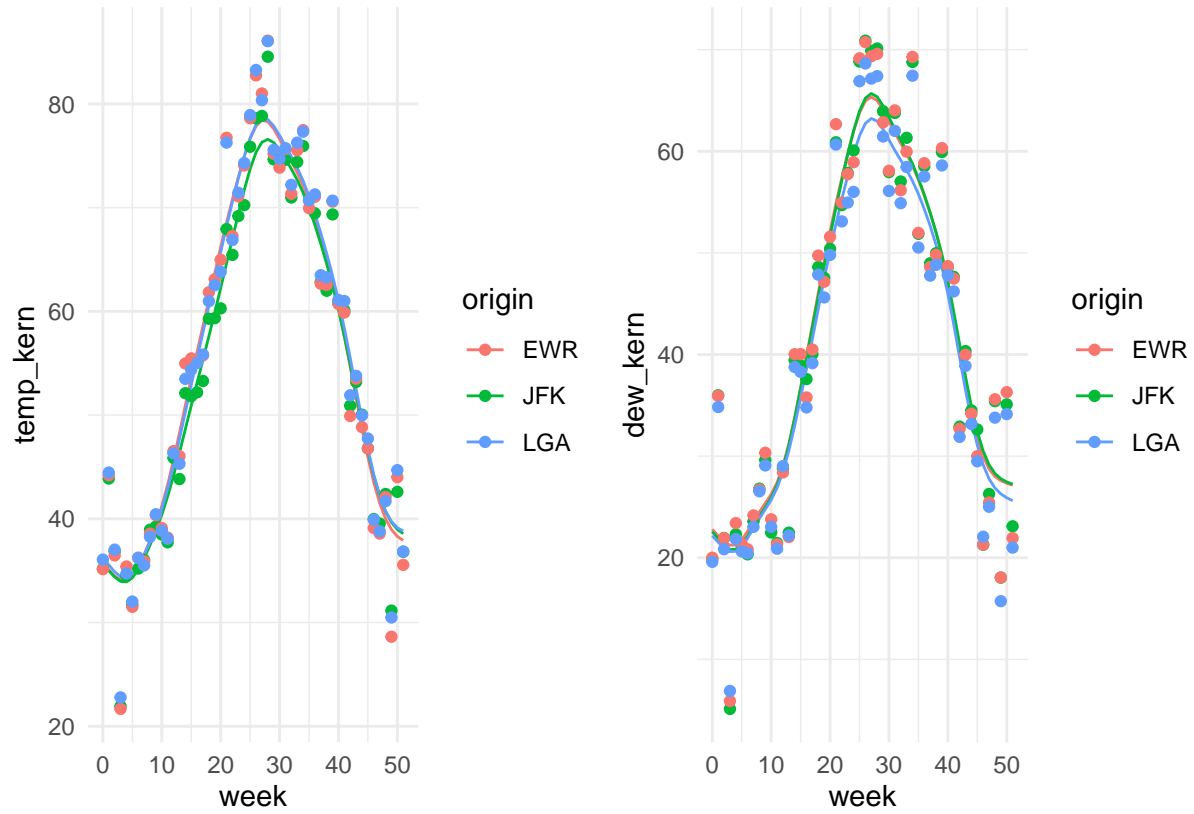
temp_plot_kern <- weather_weekly %>% ggplot() +
  geom_line(aes(x = week, y = temp_kern , colour = origin)) +
  geom_point(aes(x = week, y = temp , colour = origin)) +
  theme(legend.position = "none") +
  theme_minimal()

dew_plot_kern <- weather_weekly %>% ggplot() +

```

```
geom_line(aes(x = week, y = dew_kern , colour = origin)) +  
geom_point(aes(x = week, y = dewp , colour = origin)) +  
theme_minimal()
```

temp_plot_kern|dew_plot_kern



Bandwidth 7 and kernel “normal” were selected as they offer a balance between reduction of noise and preserving overall features of the series.