

# Predicting Cherry Blossom Bloom Dates

Joshua Write and Taehoon Song

2/26/2022

## Introduction

In this document, we demonstrate how we went about predicting the peak bloom dates in the coming decade for all four locations required by the competition. We use different modeling approaches due to different available data.

```
# Loading in packages and data
library(tidyverse)
library(forecast)
library(lubridate)
library(astsa)
library(here)
library(rnoaa)

cherry <- read_csv(here('peak-bloom-prediction', 'data', 'washingtondc.csv')) %>%
  bind_rows(read_csv(here('peak-bloom-prediction', 'data', 'liestal.csv'))) %>%
  bind_rows(read_csv(here('peak-bloom-prediction', 'data', 'kyoto.csv'))) %>%
  select(location, year, bloom_date, bloom_doy)
```

## Washington DC

In the field of phenology, plant growth is often modeled by using growing degree-day (GDD) as a predictor. Since the [National Park Service \(NPS\)](#) publishes dates of various growth stages for 2004 through 2021, we use this to build a growth model that helps us estimate the threshold GDD for which we can expect the plant to bloom. Details can be found in [Phenological Models of Flower Bud Stages and Fruit Growth of ‘Montmorency’ Sour Cherry Based on Growing Degree-day Accumulation](#) by C. Zavalloni, J. Adresen, and J. Flore. The authors use 4 degrees Celsius as the base for calculating GDD and fit a logistic function to the phenological stages from side green as a function of accumulated GDD.

$$y = \frac{k}{1 + \left[ \frac{k-n_0}{n_0} \cdot \exp(-r \cdot \text{GDD}) \right]}$$

where  $y$  is the phenological stage and  $k$  and  $r$  are empirical factors related to growth rate. Now, with just enough knowledge to be dangerous, we attempt to use this model to fit the growth stages published by the NPS as a function of cumulative GDD. Note that the authors use  $n_0 = 2$  because they were interested in phenological stages after side green (Stage 2). To be consistent with the paper’s authors, we label the phenological stages published by the NPS in a similar fashion.

## Obtaining temperature and phenological data for DC

Since the species in the paper is different from what we’re interested in, we must first identify the typical day of year when the cherry blossom trees in DC reach the green buds stage (Stage 3, as defined by the authors).

```
# Import historical weather data from DayMet
dc_tmin <- read_csv(here('peak-bloom-prediction', 'data', 'statistics_tmin.csv'))
dc_tmax <- read_csv(here('peak-bloom-prediction', 'data', 'statistics_tmax.csv'))
```

```

# DayMet data stops at 2020. Use rnoaa for 2021.
# Original code from demo analysis
# https://competition.statistics.gmu.edu/wp-content/uploads/2022/02/demo_analysis.html
getTemperature <- function (stationid,date_min,date_max) {
  tbls <- ghcn_d_search(stationid = stationid, var = c("tmax","tmin"),
                        date_min = date_min, date_max = date_max)
  tmax <- tbls[[1]] %>% select(id,tmax,date)
  tmin <- tbls[[2]] %>% select(tmin,date)
  temp <- inner_join(tmax,tmin,by="date") %>%
    group_by(date) %>%
    # Calculate average and convert unit from tenth degree Celsius to regular Celsius
    summarize(id,tavg=0.5*sum(tmax,tmin,na.rm=TRUE)/10) %>%
    ungroup() %>%
    mutate(year = year(date))
  return(temp)
}

dc_noaa <- getTemperature("USC00186350","2021-01-01","2022-01-31") %>%
  select(-id)

# Temperature data
dc_temp <- left_join(dc_tmin,dc_tmax, by = 'dt') %>%
  select(dt,value_mean.x,value_mean.y) %>%
  rename(tmin=value_mean.x,tmax=value_mean.y) %>%
  summarize(date=dt,tavg=0.5*(tmin+tmax), year=year(date)) %>%
  bind_rows(dc_noaa)

# Phenological stage data
dc_bloomdata <- cherry %>% filter(location=='washingtondc') %>%
  summarize(date=bloom_date,bloom_doy) %>%
  summarize(year=year(date),FullBloom=bloom_doy)

phenostage <- read_csv(here('peak-bloom-prediction','data',
                          'Phenological_Stage.csv')) %>%
  inner_join(dc_bloomdata,by='year')

summary(phenostage$GreenBuds)

```

```

##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##  50.00   59.00   64.00   63.94   69.50   77.00

```

This result suggests that Stage 3 occurs as early as 50 days from the start of the year. Thus, we shift the dates by 49 days to start calculating GDD to ensure that we start calculating GDD as soon as Stage 3 occurs. We also set  $n_0 = 3$  because we started from Stage 3 instead of 2.

## Building the growth model

As mentioned above, we calculate the GDD using a base of 4 degrees Celsius and starting from 49 days from the beginning of the year. Then, we fit the known phenological stages for the cherry blossoms in DC to the calculated GDDs.

```

thresh <- 4
GDD_final <- dc_temp %>%
  mutate(NewYear=year(date %m-% period("49 day"))) %>%
  mutate(DD = ifelse(tavg - thresh > 0, tavg - thresh, 0)) %>%
  group_by(NewYear) %>%
  mutate(GDD = cumsum(DD),doy = yday(date)) %>%
  ungroup()

```

```

# obtain GDD for each phenological stage for DC
stages <- tibble(stage=names(phenostage)[2:7],stagenum=seq(3,8))
pheno_GDD <- phenostage %>%
  select(-GreenBuds) %>% #This stage is the starting point for our model.
  pivot_longer(cols=FloretsVisible:FullBloom,names_to = "stage",values_to="doy") %>%
  left_join(GDD_final,by=c("year","doy")) %>%
  summarize(year,stage,doy,GDD) %>%
  left_join(stages,by='stage')

# Fit growth model
n0 <- 3 #Starting phenological stage for our model
growth_model <- nls(stagenum ~ k/(1+(k-n0)/n0*exp(-r*GDD)),
  data=pheno_GDD,start=list(k=1,r=0.01))
summary(growth_model)

```

```

##
## Formula: stagenum ~ k/(1 + (k - n0)/n0 * exp(-r * GDD))
##
## Parameters:
##   Estimate Std. Error t value Pr(>|t|)
## k 10.223870   1.362473   7.504 4.72e-11 ***
## r  0.011973   0.001598   7.495 4.93e-11 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.9374 on 88 degrees of freedom
##
## Number of iterations to convergence: 8
## Achieved convergence tolerance: 7.797e-06

```

The model results show that both  $k$  and  $r$  are significant. However, inference is not of interest to us. We only care about prediction performance.

## Checking performance on historical data

Now that we have the growth model, we can test its prediction performance out on historical data that we did not use.

```

testdata <- GDD_final %>% filter(year<=2003) %>% #Used data from 2004 onward
  group_by(year) %>%
  filter(doy>=50 & doy <= 200) %>%
  ungroup() %>%
  select(date,GDD)

preds <- tibble(date=testdata$date, GDD=testdata$GDD,
  pred=predict(growth_model,newdata=testdata)) %>%
  mutate(year=year(date))

# Predicted date is the earliest date when predicted stage rounds to Stage 8
# Compare predicted date to actual historical dates
comp <- preds %>% group_by(year) %>%
  filter(round(pred) >= 8) %>%
  filter(row_number()==1) %>%
  ungroup() %>%
  mutate(pred.doy=yday(date)) %>%
  inner_join(dc_bloomdata,by='year') %>%
  summarize(year,GDD,pred.doy,act.doy=FullBloom,err_GDD=abs(pred.doy-act.doy)) %>%
  # compare to using simply using a lag predictor
  mutate(lag_bloom = lag(act.doy),err_lag=abs(act.doy - lag_bloom))

```

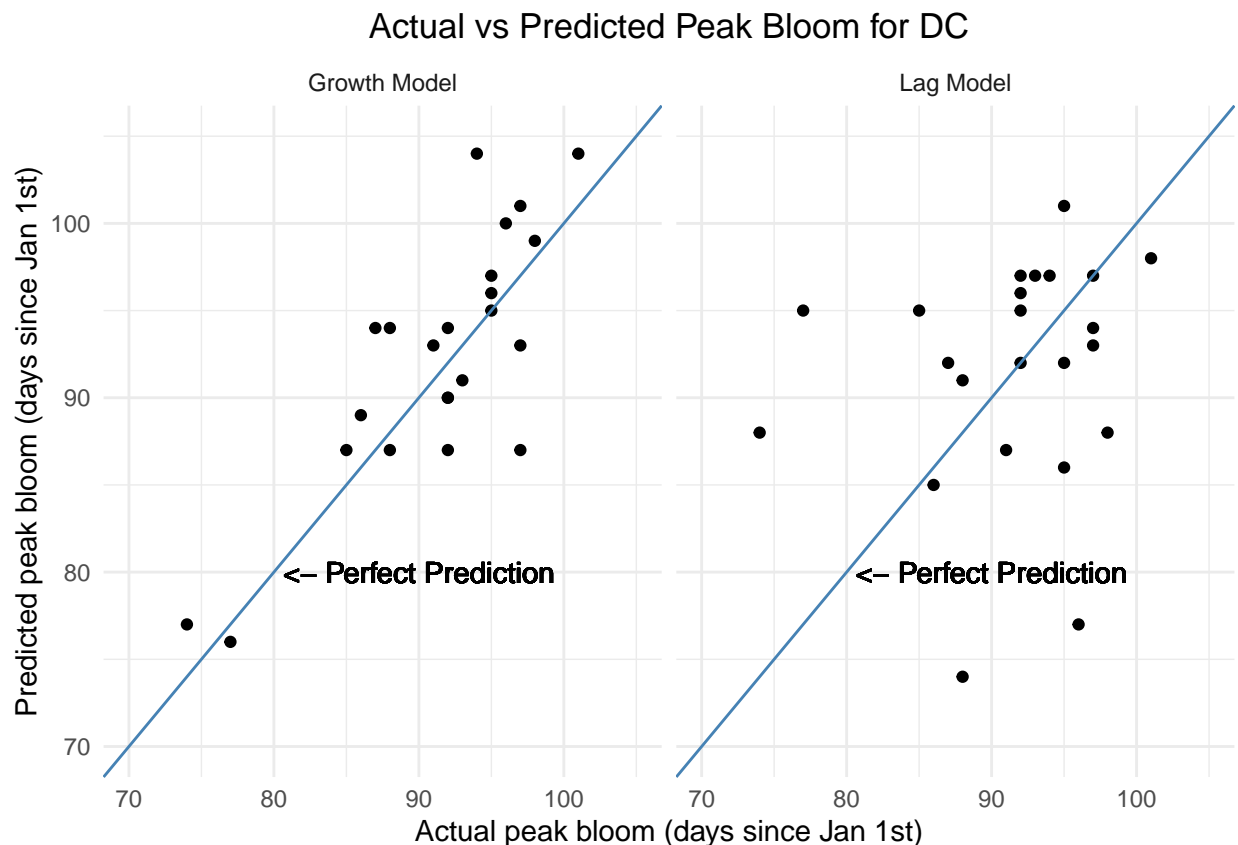
```
data.frame(MAE_gdd=mean(comp$err_GDD[which(!is.na(comp$err_lag))]),
           MAE_lag=mean(comp$err_lag,na.rm=TRUE))
```

```
##      MAE_gdd  MAE_lag
## 1  3.347826  6.304348
```

Our growth model has a much lower mean absolute error (MAE) compared to simply using the previous year's bloom date as a predictor. This is good because this is the bare minimum a useful model should be able to do. This is more apparent when we plot the two models against the actual peak bloom days since January 1st.

```
plt_data <- comp %>% rename(Actual=act.doy,GrowthModel=pred.doy,LagModel=lag_bloom) %>%
  pivot_longer(cols=c("GrowthModel","LagModel"),names_to="type",values_to="doy") %>%
  filter(!is.na(err_lag))
plt_data$type <- factor(plt_data$type,levels=c("GrowthModel","LagModel"),
                       labels=c("Growth Model","Lag Model"))

# Plot Actual vs Predicted for Growth Model
ggplot(plt_data,aes(x=Actual)) + geom_point(aes(y=doy)) +
  facet_grid(~type) +
  geom_abline(slope=1,col="steelblue") +
  geom_text(aes(x=90,y=80,label="<- Perfect Prediction")) +
  theme_minimal() +
  ggtitle("Actual vs Predicted Peak Bloom for DC") +
  xlab("Actual peak bloom (days since Jan 1st)") +
  ylab("Predicted peak bloom (days since Jan 1st)") +
  xlim(c(70,105)) + ylim(c(70,105)) +
  theme(plot.title=element_text(hjust=0.5))
```

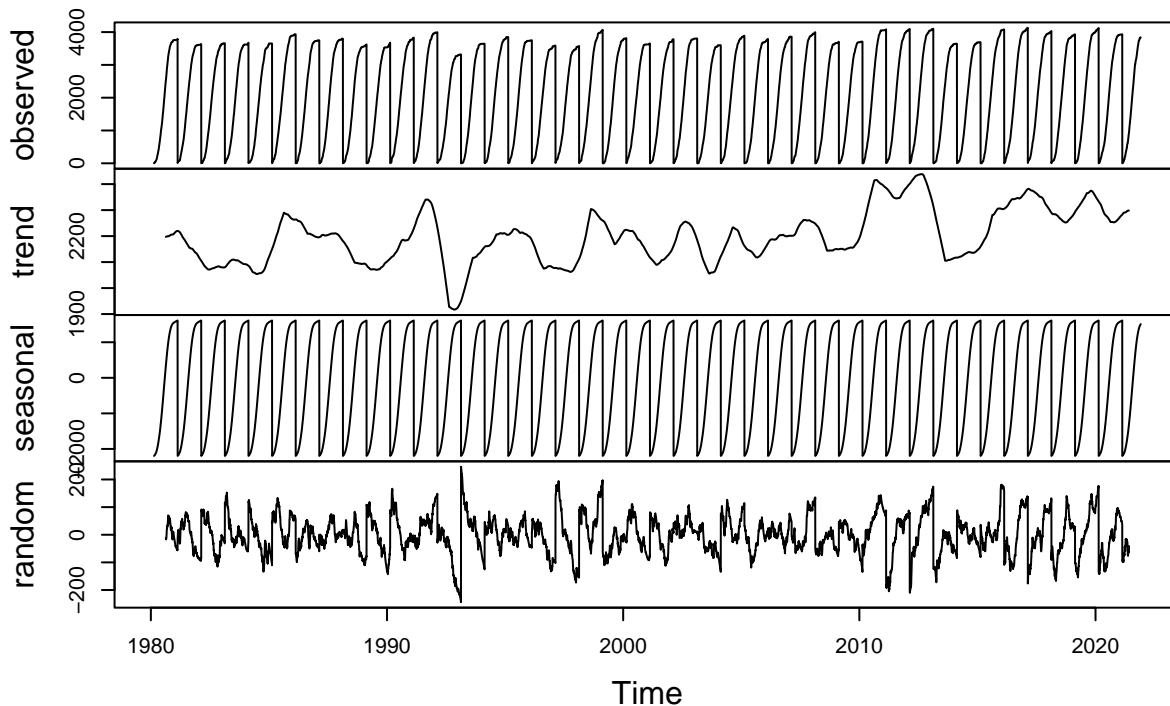


## Predicting future GDD and future peak bloom dates

For us to predict future peak bloom dates for DC, we must first predict future GDDs. We could predict future temperature and calculate GDD from that; however, that would be unnecessary because we always shift the dates by 49 days and use the base of 4 degrees Celsius to calculate GDD. We can directly predict GDD to avoid that step.

```
# Create time series
ts_gdd <- ts(GDD_final$GDD[50:nrow(GDD_final)],frequency = 365,
            #Only start counting GDD from day 50 & 2021 DC data only goes to November
            start=c(1980,50),end=c(2021,334))
plot(decompose(ts_gdd))
```

### Decomposition of additive time series



```
if (file.exists('pred_dc.rds')){
  pred_dc <- readRDS(here("pred_dc.rds"))
} else{
  # Fit time series to GDD
  ts_fit <- auto.arima(ts_gdd, ic = 'aicc')
  ar <- ts_fit$arma[1:2]
  ma <- ts_fit$arma[3:4]
  dd <- ts_fit$arma[5:7]

  # Use result from fitted time series to predict GDD until May 2031
  preddays <- difftime("2031-05-31", "2021-12-01", units="days")
  pred_dc <- sarima.for(ts_gdd, as.numeric(preddays), ar[1], dd[2], ar[2], ma[1], dd[3], ma[2], dd[1])
  saveRDS(pred_dc, here("pred_dc.rds"))
}
```

We can now predict the peak bloom day for DC using the predicted GDD as an input for the growth model.

```
# Forecasted GDD and growth stage
newdata <- tibble(date=seq.Date(as.Date("2021-12-01"),length.out=length(pred_dc$pred),by="day"),
                  GDD=pred_dc$pred,doy=yday(date),year=year(date))
future_growthstage <- predict(growth_model,newdata=newdata)
forecasts_dc <- newdata %>% mutate(pred.gs=future_growthstage) %>%
```

```

group_by(year) %>%
filter(doy >= 50 & doy <= 200) %>%
filter(round(pred.gs) >= 8) %>%
filter(row_number()==1) %>%
select(year,doy)

forecasts_dc

```

```

## # A tibble: 10 x 2
## # Groups:   year [10]
##   year    doy
##   <dbl> <dbl>
## 1  2022     86
## 2  2023     86
## 3  2024     86
## 4  2025     85
## 5  2026     85
## 6  2027     85
## 7  2028     84
## 8  2029     83
## 9  2030     83
## 10 2031     83

```

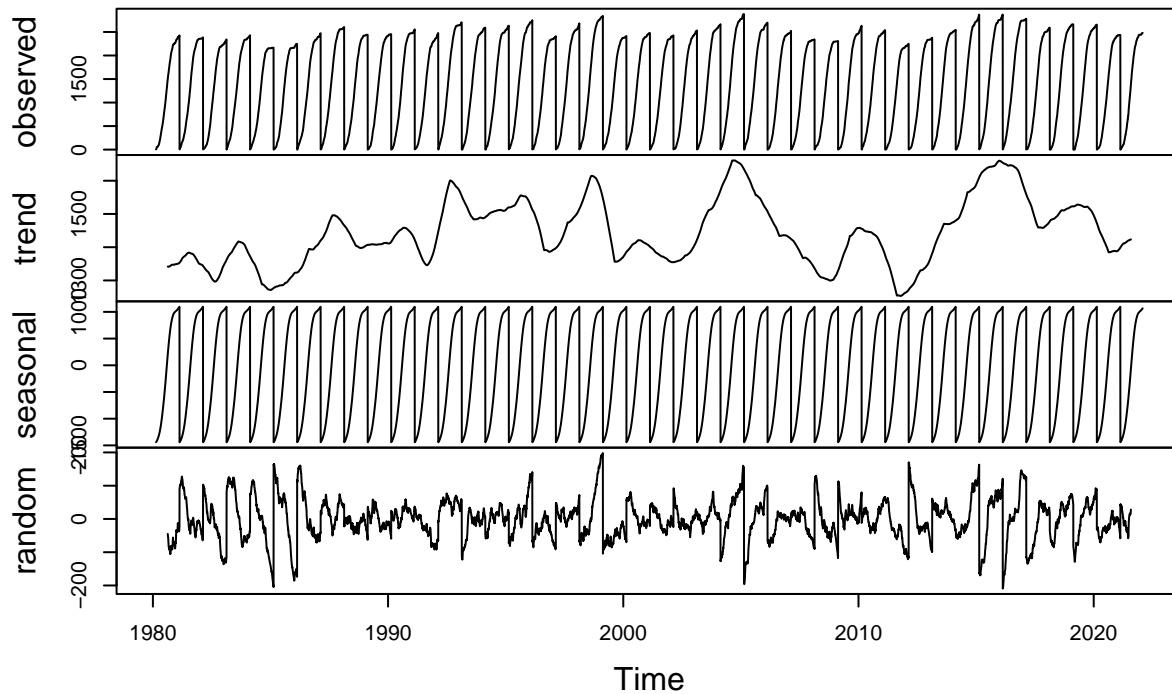
## Vancouver

Vancouver doesn't provide historical peak bloom dates. However, we can use the growth model from DC to predict future peak bloom dates because the tree species are identical (or close to identical) and the peak bloom definition is the same. We simply need to calculate historical GDD using the same definitions and predict future GDD for Vancouver to be able to predict the peak bloom dates.

### Obtaining historical weather data

Similar to DC, we obtained the historical weather data for Vancouver via DayMet. However, this dataset only goes to December 2019, so we use `rnoaa` to obtain the remaining data up to Jan 2022.

## Decomposition of additive time series



```
## # A tibble: 10 x 2
## # Groups:   year [10]
##   year   doy
##   <dbl> <dbl>
## 1  2022   109
## 2  2023   108
## 3  2024   108
## 4  2025   107
## 5  2026   106
## 6  2027   106
## 7  2028   106
## 8  2029   105
## 9  2030   104
## 10 2031   104
```

## Liestal

The cherry tree in Liestal is a different species from the ones in DC and Vancouver. Ideally, we would want to fit a new growth model using historical phenological stages. Unfortunately, such a dataset is not readily available. So we use the same growth model from DC to calculate the mean growth stage for historical peak bloom dates. We use the average growth stage for the past 5 years as the cutoff for our prediction.

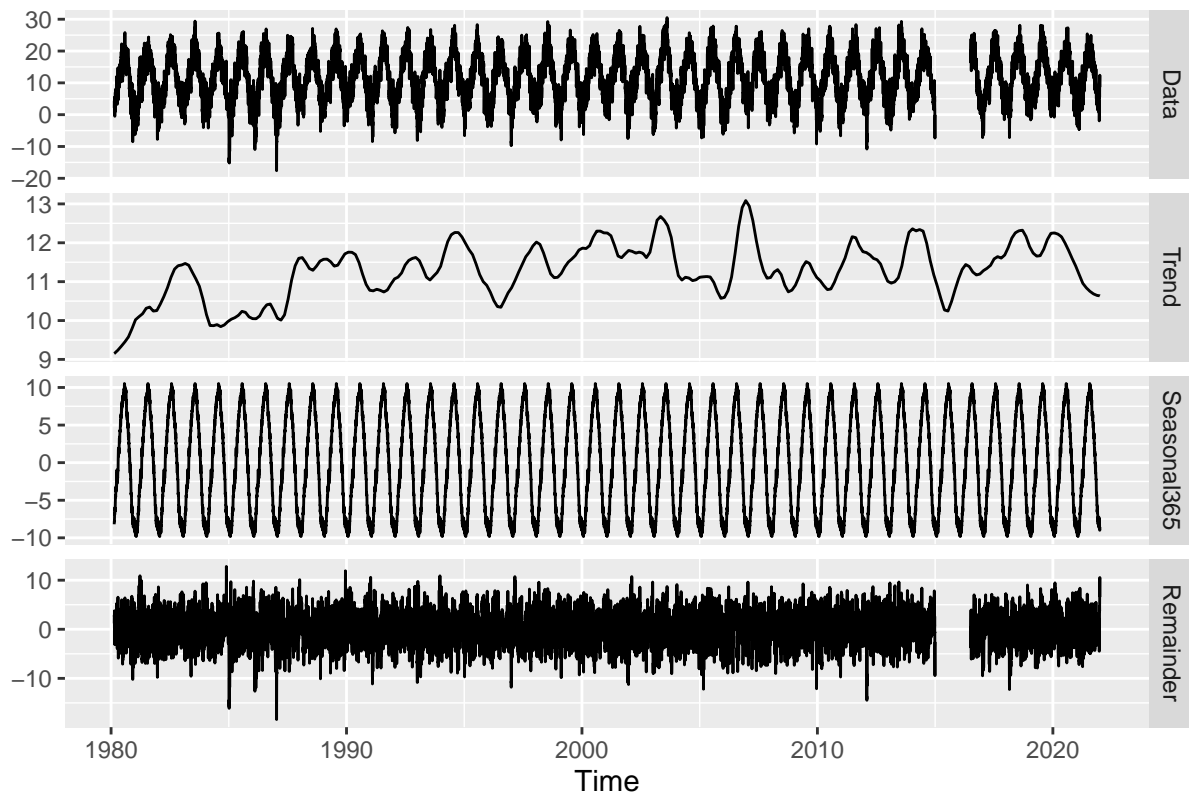
```
# Check historical GDD and corresponding predicted growth stages for actual bloom dates in Liestal
testdata <- cherry %>%
  left_join(GDD_final %>% select(date,GDD), by=c('bloom_date'='date')) %>%
  filter(year >= 2016 & location=='liestal') %>%
  select(GDD)

summary(predict(growth_model,newdata=testdata))
```

```
##   Min. 1st Qu.  Median    Mean 3rd Qu.    Max.   NA's
##  5.995  6.532   6.773   6.799   7.170   7.524     1
```

```
li_thresh <- mean(predict(growth_model,newdata=testdata),na.rm=TRUE)
```

We can approximate the peak bloom for the cherry trees in Liestal (as defined by MeteoSwiss) using the growth model for the trees in DC if we use 6.7988151 as a threshold for the predicted growth stage. In other words, we estimate the peak bloom date by finding the earliest date for which the predicted growth stage is greater than or equal to 6.7988151.



```
## # A tibble: 10 x 2
## # Groups:   year [10]
##   year   doy
##   <dbl> <dbl>
## 1 2022    88
## 2 2023    88
## 3 2024    88
## 4 2025    87
## 5 2026    87
## 6 2027    87
## 7 2028    87
## 8 2029    86
## 9 2030    86
## 10 2031    86
```

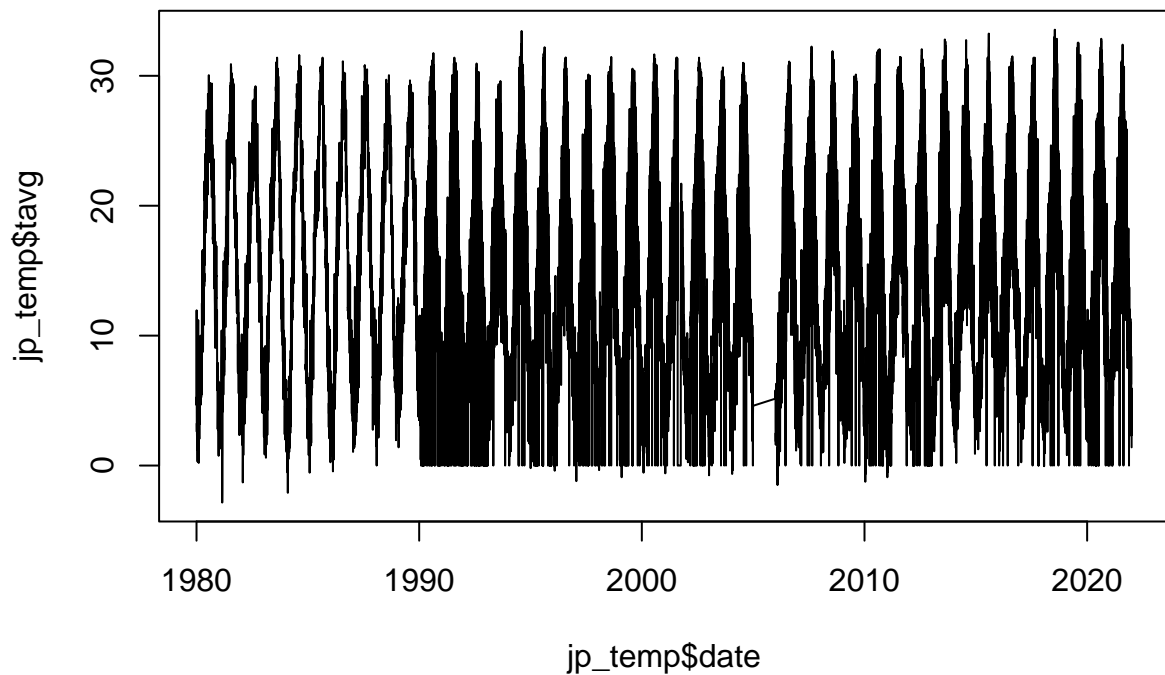
## Kyoto

Similar to predicting peak bloom dates for Liestal, the peak bloom dates for Kyoto (as defined by a local newspaper in Arashiyama) need to be approximated using the growth model for the DC cherry trees.

```
jp_temp <- temps %>% filter(city == 'Kyoto') %>%
  mutate(doy = yday(date)) %>%
  filter(doy != 366) #Leap year mess

# Historical temperature plot for Kyoto
plot(jp_temp$date, jp_temp$tavg, type="l")
```

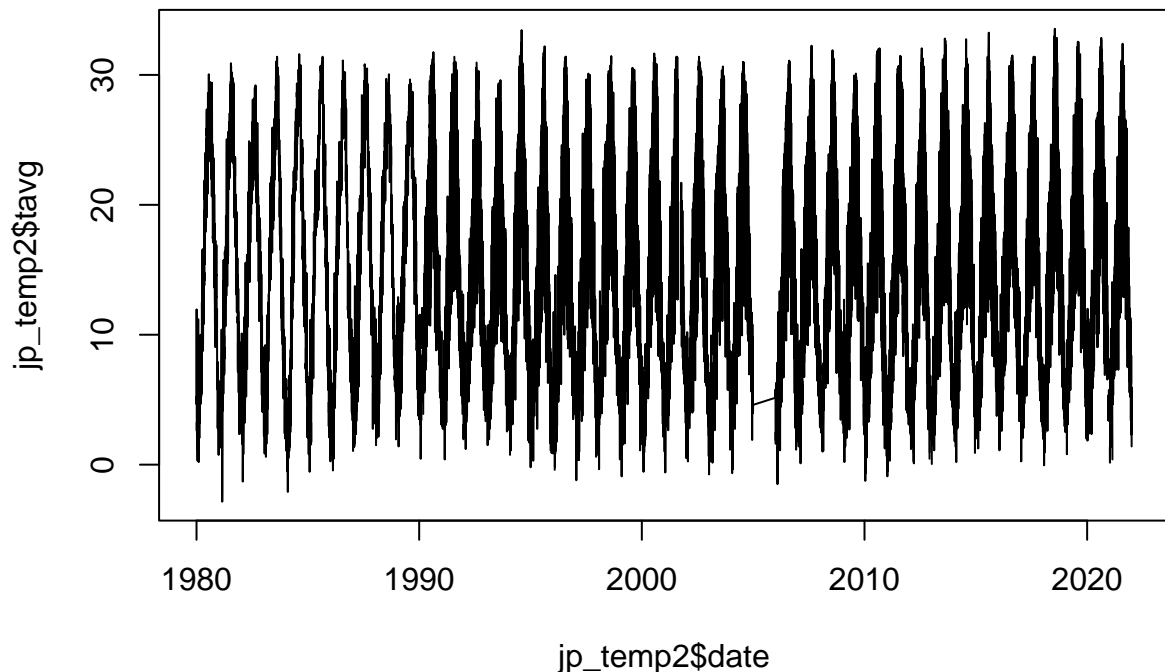




Notice that the data abruptly cuts off at 0 throughout the years. This is likely due to missing values. We remove 0s to NAs so that the algorithm knows that the value is actually missing. Kyoto is a very warm city and rarely dips below 0 degrees Celsius, so we shouldn't lose too much information with this substitution. We're missing a month of data for September 2009, but that shouldn't be too much of a problem for predicting 2022 and onward.

```
# Substitute 0s for NA
jp_temp2 <- jp_temp %>%
  mutate(year=year(date), month=month(date)) %>%
  group_by(year,month) %>%
  mutate(tavg=ifelse(tavg==0,NA,tavg))

plot(jp_temp2$date,jp_temp2$tavg,type="l")
```



With the substitution, the odd behavior is gone. We can proceed to calculate GDD for Kyoto.

```
# GDD calculations
GDD_final <- jp_temp2 %>%
  mutate(NewYear=year(date %m-% period("49 day"))) %>%
  mutate(DD = ifelse(tavg - thresh > 0, tavg - thresh, 0)) %>%
  group_by(NewYear) %>%
  mutate(GDD = cumsum(DD),doy = yday(date)) %>%
  ungroup()

# Check historical GDD and corresponding predicted growth stages for actual bloom dates in Liestal
testdata <- cherry %>%
  left_join(GDD_final %>% select(date,GDD), by=c('bloom_date'='date')) %>%
  filter(year >= 2016 & location=='kyoto') %>%
  select(GDD)

summary(predict(growth_model,newdata=testdata))
```

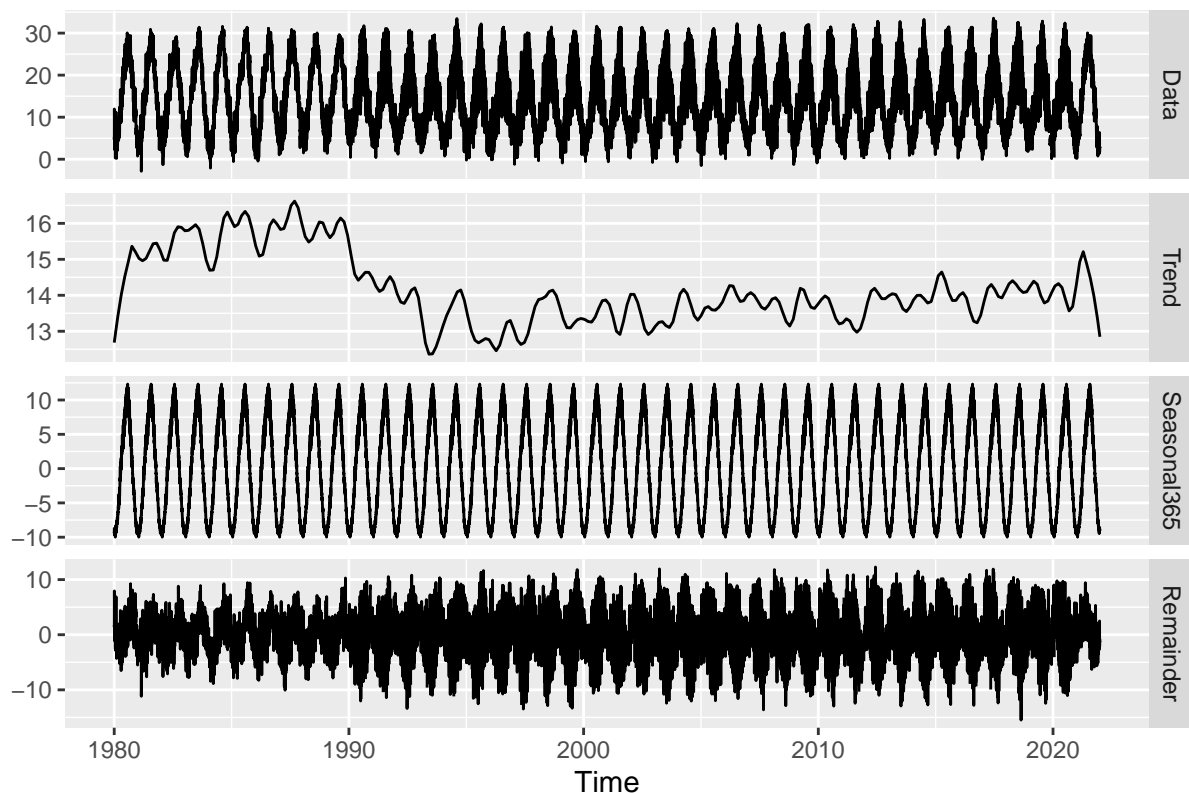
```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.     NA's
##      8.234   8.268   8.303   8.303   8.337   8.372         4
```

```
jp_thresh <- mean(predict(growth_model,newdata=testdata),na.rm=TRUE)
```

Based on the above result, we estimate the peak bloom date by finding the earliest date for which the predicted growth stage is greater than or equal to 8.3028025.

```
# Create time series
ts_temp <- ts(GDD_final$tavg[1:nrow(GDD_final)],frequency = 365,
             start=c(1980,1),end=c(2021,365))

ts_temp %>%
  mstl(s.window="periodic") %>%
  autoplot()
```



```
## # A tibble: 10 x 2
## # Groups:   year [10]
##   year    doy
##   <dbl> <dbl>
## 1 2022     79
## 2 2023     79
## 3 2024     79
## 4 2025     79
## 5 2026     79
## 6 2027     79
## 7 2028     79
## 8 2029     78
## 9 2030     78
## 10 2031     78
```

## Final predictions

The summary table below shows the final predictions for the four locations from 2022 to 2031.

```
## # A tibble: 10 x 5
##   year kyoto liestal washingtondc vancouver
##   <int> <dbl>   <dbl>         <dbl>     <dbl>
## 1 2022     79     88           86       109
## 2 2023     79     88           86       108
## 3 2024     79     88           86       108
## 4 2025     79     87           85       107
## 5 2026     79     87           85       106
## 6 2027     79     87           85       106
## 7 2028     79     87           84       106
## 8 2029     78     86           83       105
## 9 2030     78     86           83       104
## 10 2031     78     86           83       104
```

## Closing thoughts

If historical phenological stages were recorded for all locations in this competition, we would have been able to fit a growth model for each species and location. This would have allowed us to create more accurate models. Additionally, weather data was extremely limited or unreliable for locations outside the United States. If historical weather data was free and readily available to the public, we would have been able to create a better time series model for GDD. Lastly, using weather forecasts from meteorologists will likely improve the accuracy of the peak bloom predictions.