

Vaccination in the UK

Contents

1	Step 1 Find dataset	2
2	Step 2 Ask something	5
3	Step 3 Look at the datasets	6
3.1	Region	6
4	Step 4 Machine learning	13
4.1	Step 0: Look at and Modify the dataset	13
4.2	Step 1: Explore the dataset	16
4.3	Step 2: Split sets, train a Machine Learning Model and Evaluate performance	17
4.4	Step 3: Plot results	20
4.5	Step 4: Improve models by changing the dataset	22

1 Step 1 Find dataset

- Create a list of metrics for each dataset
- Look at the metrics

Lower Tier Local Authority (LTLA)

```
## [1] "New people receiving 2nd dose"
## [2] "New people vaccinated with a booster dose by publish date"
## [3] "New people vaccinated complete by publish date"
## [4] "New people fully vaccinated by vaccination date"
## [5] "New people vaccinated 1st dose by publish date"
## [6] "New people vaccinated with a first dose by vaccination date"
## [7] "New people vaccinated 2nd dose by publish date"
## [8] "New people vaccinated with a second dose by vaccination date"
## [9] "New people vaccinated with a third dose by publish date"
## [10] "New people vaccinated with a booster dose plus new people vaccinated with a third dose by publish date"
## [11] "New people vaccinated with a booster or third dose by vaccination date"
## [12] "New vaccines given by publish date"
```

Nation

```
## [1] "New people receiving 1st dose"
## [2] "New people receiving 2nd dose"
## [3] "New people vaccinated with a booster dose by publish date"
## [4] "New people vaccinated complete by publish date"
## [5] "New people fully vaccinated by vaccination date"
## [6] "New people vaccinated 1st dose by publish date"
## [7] "New people vaccinated with a first dose by vaccination date"
## [8] "New people vaccinated 2nd dose by publish date"
## [9] "New people vaccinated with a second dose by vaccination date"
## [10] "New people vaccinated with a third dose by publish date"
## [11] "New people vaccinated with a booster dose plus new people vaccinated with a third dose by publish date"
## [12] "New people vaccinated with a booster or third dose by vaccination date"
## [13] "New vaccines given by publish date"
```

So, as we can see, **some metrics are common**. I suggest finding out which metrics are the same for all datasets.

	ltla	msoa	nation	nhsRegion	nhsTrust	overview	region	utla
New people receiving 2nd dose	1	0	1	0	0	0	1	0
New people vaccinated with a booster dose by publish date	1	0	1	0	0	0	1	0
New people vaccinated complete by publish date	1	0	1	0	0	0	1	0
New people fully vaccinated by vaccination date	1	0	1	0	0	0	1	0
New people vaccinated 1st dose by publish date	1	0	1	0	0	0	1	0
New people vaccinated with a first dose by vaccination date	1	0	1	0	0	0	1	0
New people vaccinated 2nd dose by publish date	1	0	1	0	0	0	1	0
New people vaccinated with a second dose by vaccination date	1	0	1	0	0	0	1	0
New people vaccinated with a third dose by publish date	1	0	1	0	0	0	1	0
New people vaccinated with a booster dose plus new people vaccinated with a third dose by publish date	1	0	1	0	0	0	1	0
New people vaccinated with a booster or third dose by vaccination date	1	0	1	0	0	0	1	0
New vaccines given by publish date	1	0	1	0	0	0	1	0
New people receiving 1st dose	0	0	1	0	0	0	1	0

- Add new metrics in a common list
- Build zero-matrix, which dimension is the count of metrics x the count of area types
- Show links

Look at the result

First of all, I am interested in data about the **first jab**. **So, I need to look at the datasets:**

- Build zero-matrix, which dimension is the count of metrics x the count of area types

```
## [1] "Lower Tier Local Authority (LTLA)"  
## [1] "Nation"  
## [1] "Region"
```

2 Step 2 Ask something

I live in Bristol. What do I know about Bristol?

- This city is a part of the UK, England, and South West.
- There are two universities.
- The city rests every summer when students come back to their homes and works hardly elsewhen.

So, I am interested in data about the UK, England, South West, and Bristol.

Question 0: Are there dependencies between academic year events and vaccination waves? Does the vaccination depend on holidays?

I was vaccinated by

- the first dose on 8 August 2021,
- the second dose on 3 October 2021,
- the booster dose on 8 January 2022.

Question 1: How many people got their jabs with me?

I got the first and the second jabs on Sunday. There were fewer people in the vaccination centre. When I got the third jab on Saturday, there was a big queue.

Question 2: When do people prefer to get a jab: weekdays or weekends/Saturdays or Sundays?

Question 3: Is there something illogical in data?

3 Step 3 Look at the datasets

3.1 Region

As we can see on the website, Region metrics are available for regions of England. I am interested in the South West and metrics that start with “New”:

```
## [1] "areaCode"  
## [2] "areaName"  
## [3] "areaType"  
## [4] "date"  
## [5] "newPeopleVaccinatedFirstDoseByVaccinationDate"  
## [6] "newPeopleVaccinatedSecondDoseByVaccinationDate"  
## [7] "newPeopleVaccinatedThirdInjectionByVaccinationDate"
```

We have additional columns. Let’s look at them.

areaCode

```
## [1] "E12000009"
```

areaName

```
## [1] "South West"
```

areaType

```
## [1] "region"
```

So, we do not need to look at them in the future because these columns are used for filtering that we have already done on the website.

Let’s prepare data for the plotting.

- Rename columns and columns
- Create long table

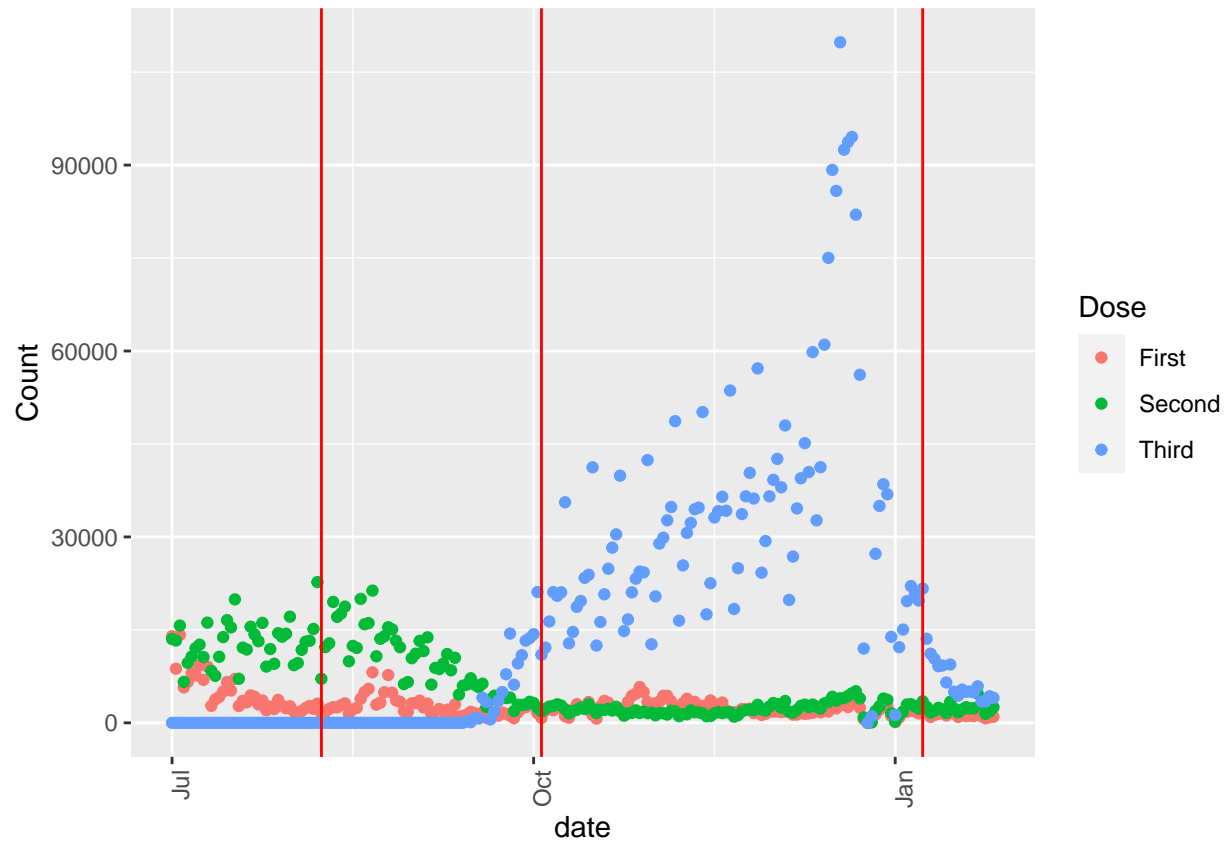
Let’s plot something.

areaCode	areaName	areaType	date	First	Second	Third	MonthYear
E12000009	South West	region	2022-01-26	986	2520	4034	1.2022
E12000009	South West	region	2022-01-25	899	1845	4283	1.2022
E12000009	South West	region	2022-01-24	723	1445	3441	1.2022
E12000009	South West	region	2022-01-23	1035	3007	3439	1.2022
E12000009	South West	region	2022-01-22	1822	4709	5896	1.2022
E12000009	South West	region	2022-01-21	1085	2362	4944	1.2022

date	MonthYear	Dose	Count
2022-01-26	1.2022	First	986
2022-01-25	1.2022	First	899
2022-01-24	1.2022	First	723
2022-01-23	1.2022	First	1035
2022-01-22	1.2022	First	1822
2022-01-21	1.2022	First	1085

3.1.1 Question 0

Are there dependencies between academic year events and vaccination waves? Does the vaccination depend on holidays?

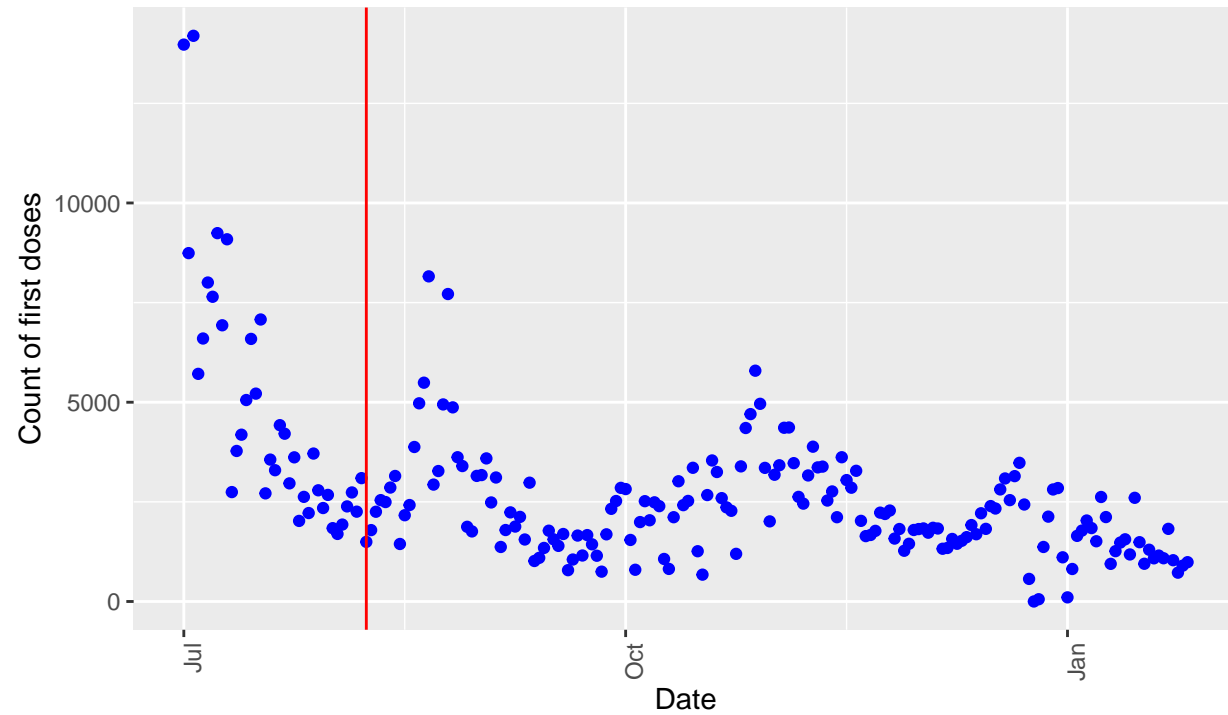


The result is not beautiful because of the active growth of the third jabs count at the end of 2021.

Let's plot them separately.

Vaccination in South West

The first dose

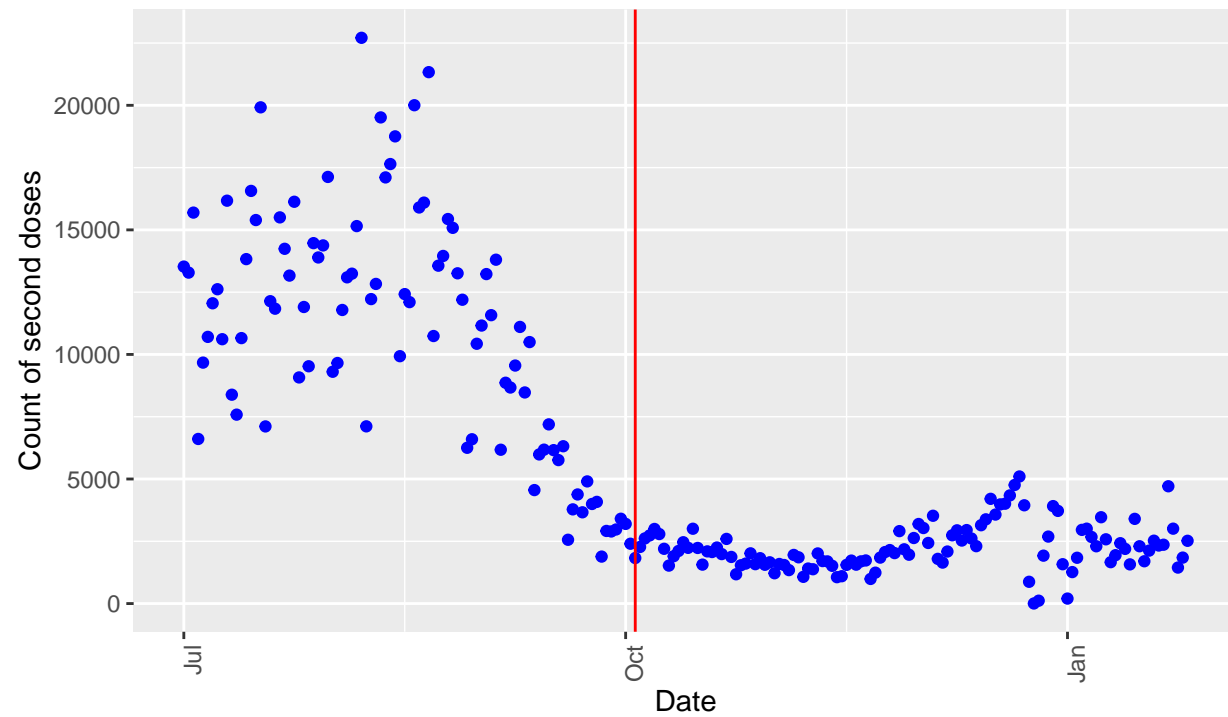


More information <https://coronavirus.data.gov.uk/details/about-data>

It is so interesting why the graph is wavy.

Vaccination in England

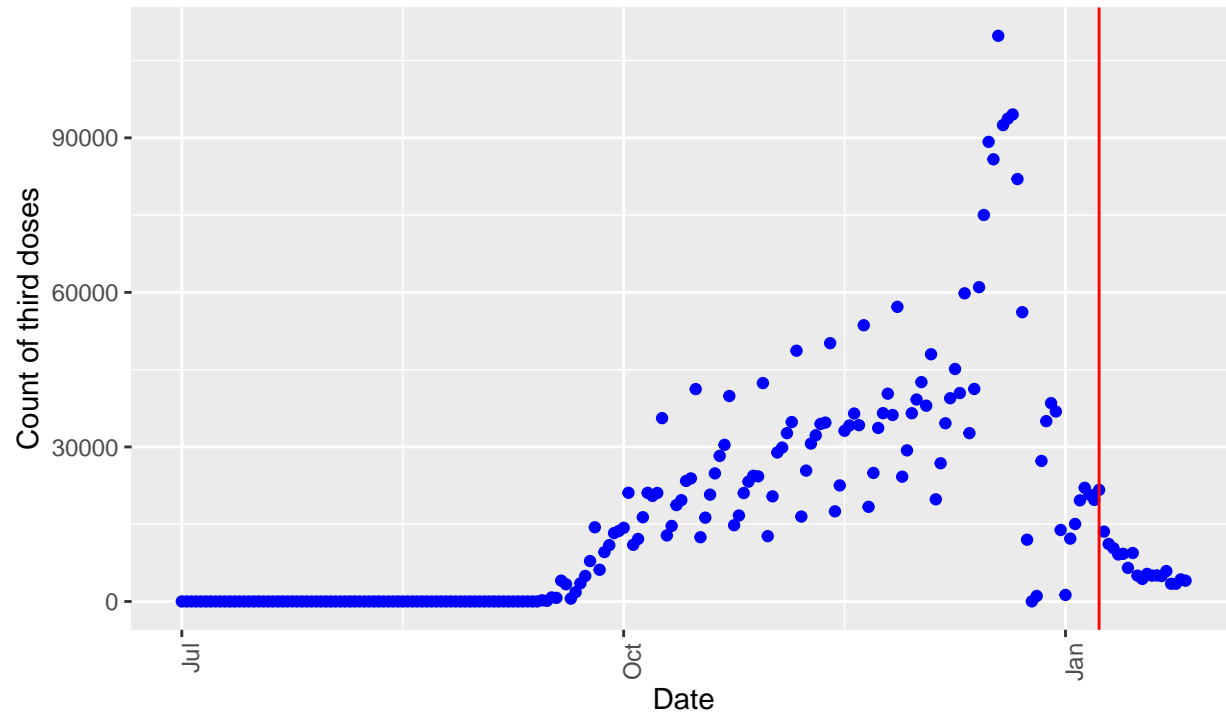
The second dose



More information <https://coronavirus.data.gov.uk/details/about-data>

Vaccination in England

The third dose

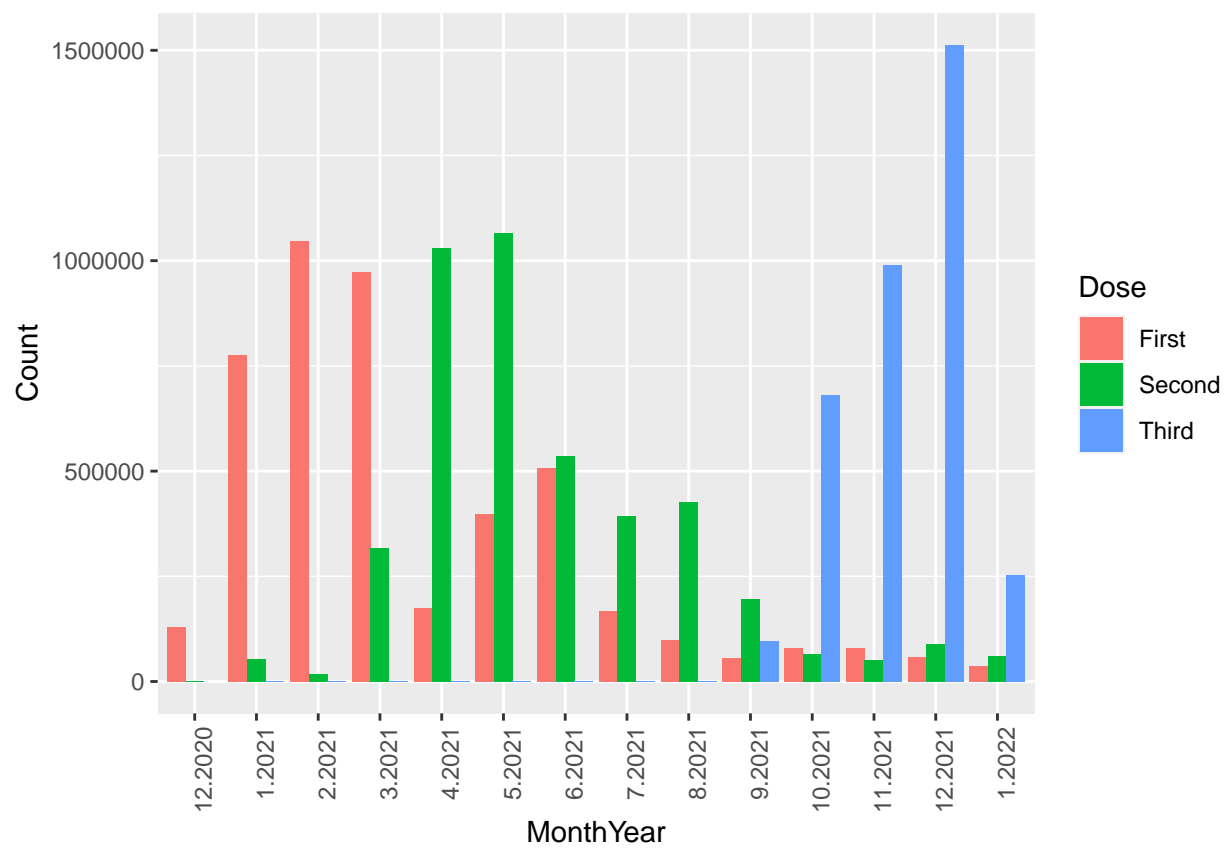


More information <https://coronavirus.data.gov.uk/details/about-data>

We can see when the active phase of vaccination by the third dose started.

Let's calculate the date.

```
## Warning: Removed 1 rows containing missing values (geom_col).
```



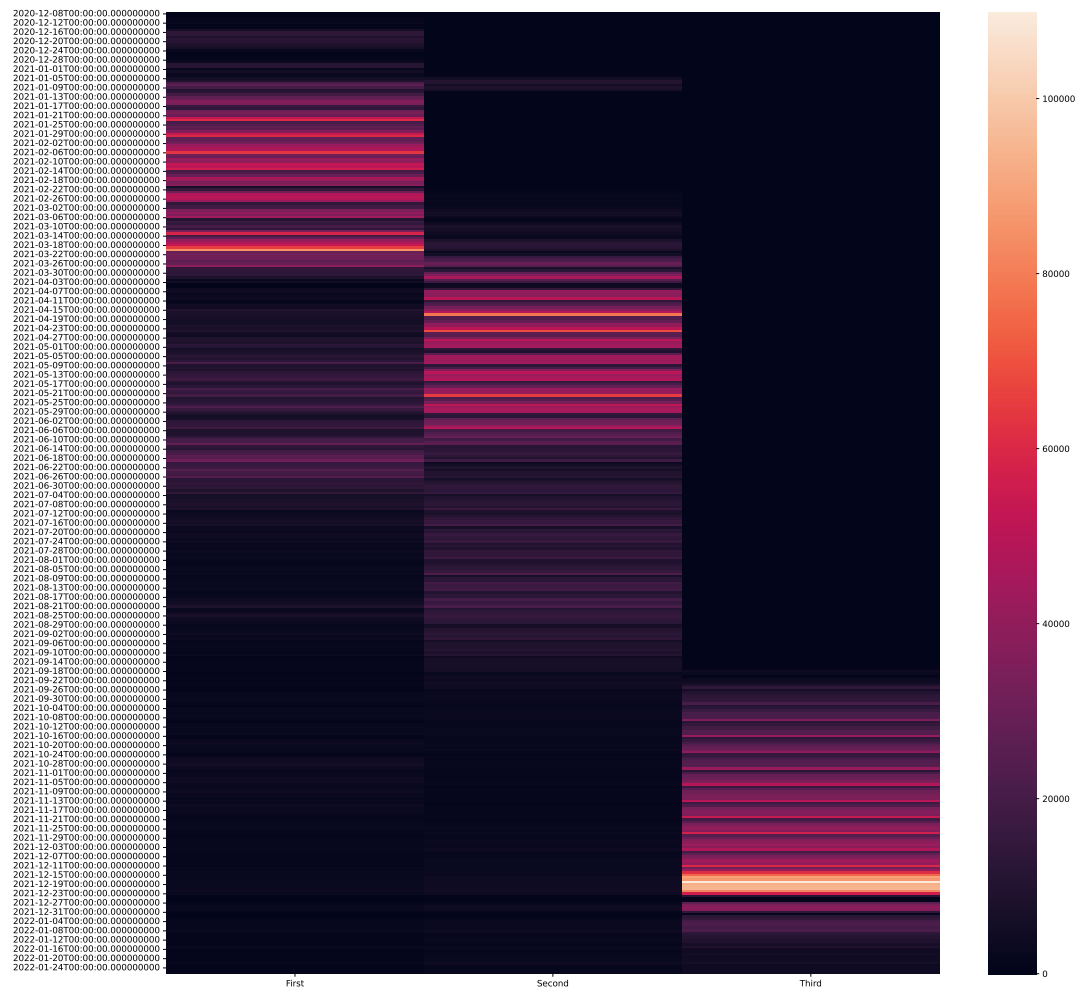
4 Step 4 Machine learning

4.1 Step 0: Look at and Modify the dataset

So, I am curious. Can I predict vaccination data?

I will work with the South West's vaccination data.

	First	Second	Third
2022-01-26	986	2520	4034
2022-01-25	899	1845	4283
2022-01-24	723	1445	3441
2022-01-23	1035	3007	3439
2022-01-22	1822	4709	5896



As we can see, there are waves. So, the count of jabs depends on dates.

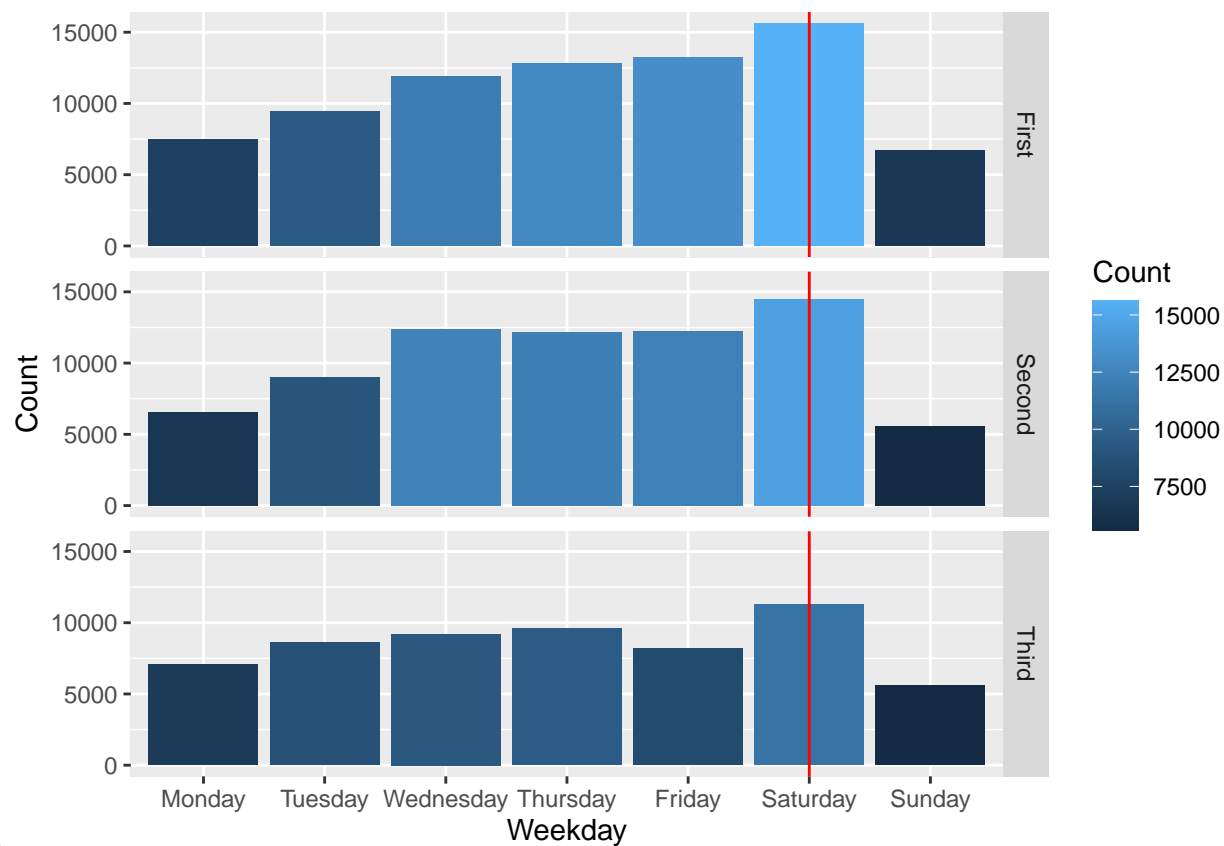
Let's get features: 1) Year 2) Month 3) Day etc.

	First	Second	Third	Year	Month	Day	DayOfYear	Weekday	Quarter	IsMonthStart	IsMonthEnd
2022-01-26	986	2520	4034	2022	1	26	26	2	1	FALSE	FALSE
2022-01-25	899	1845	4283	2022	1	25	25	1	1	FALSE	FALSE
2022-01-24	723	1445	3441	2022	1	24	24	0	1	FALSE	FALSE
2022-01-23	1035	3007	3439	2022	1	23	23	6	1	FALSE	FALSE
2022-01-22	1822	4709	5896	2022	1	22	22	5	1	FALSE	FALSE

4.2 Step 1: Explore the dataset

4.2.1 Weekdays

As you remember, I have a question. When do people prefer to get a job: weekdays or weekends/Saturdays or Sundays?



Let's answer.

So, most of South West's people prefer to get a job on Saturdays.

4.2.2 Missing values

Calculate a count of dates in the dataset.


```
## 415
```

Calculate a count of dates between maximum and minimum dates.

```
## 415
```

There are no missing dates.

4.3 Step 2: Split sets, train a Machine Learning Model and Evaluate performance

Define necessary variables

Prepare sets and train models using parameters.

Compare the score with the mean value of the column that we predicted.

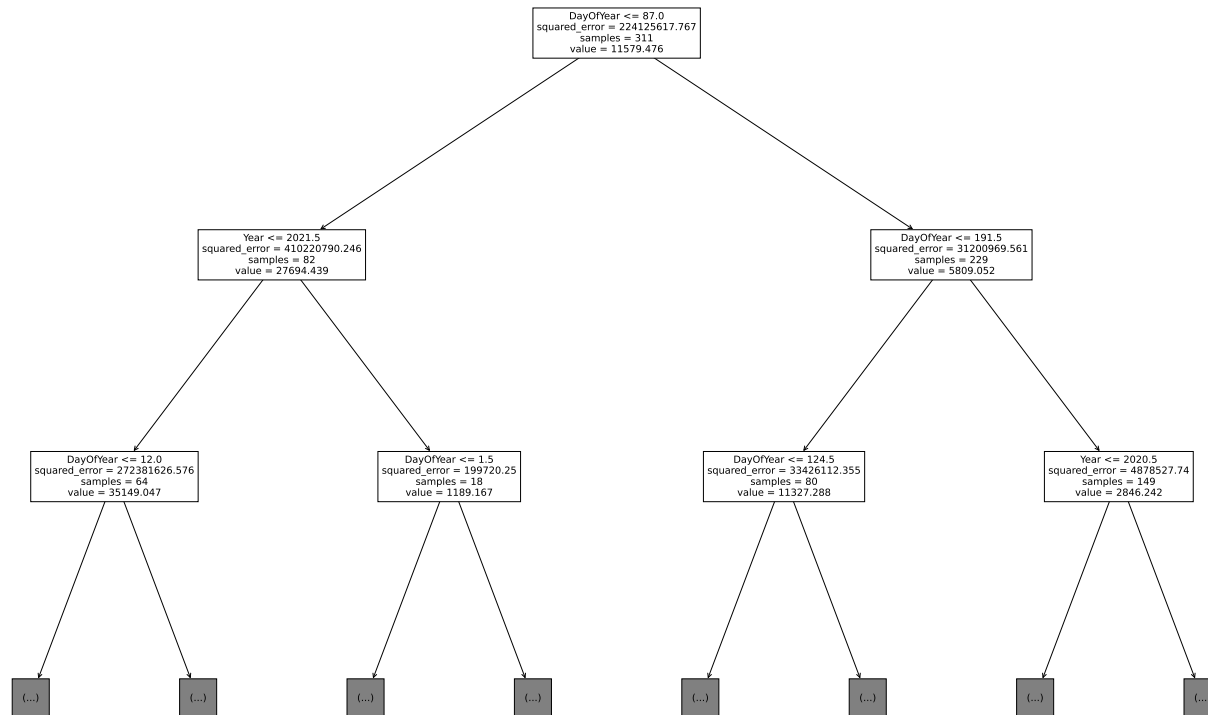
```
## 0.719657929335243
```

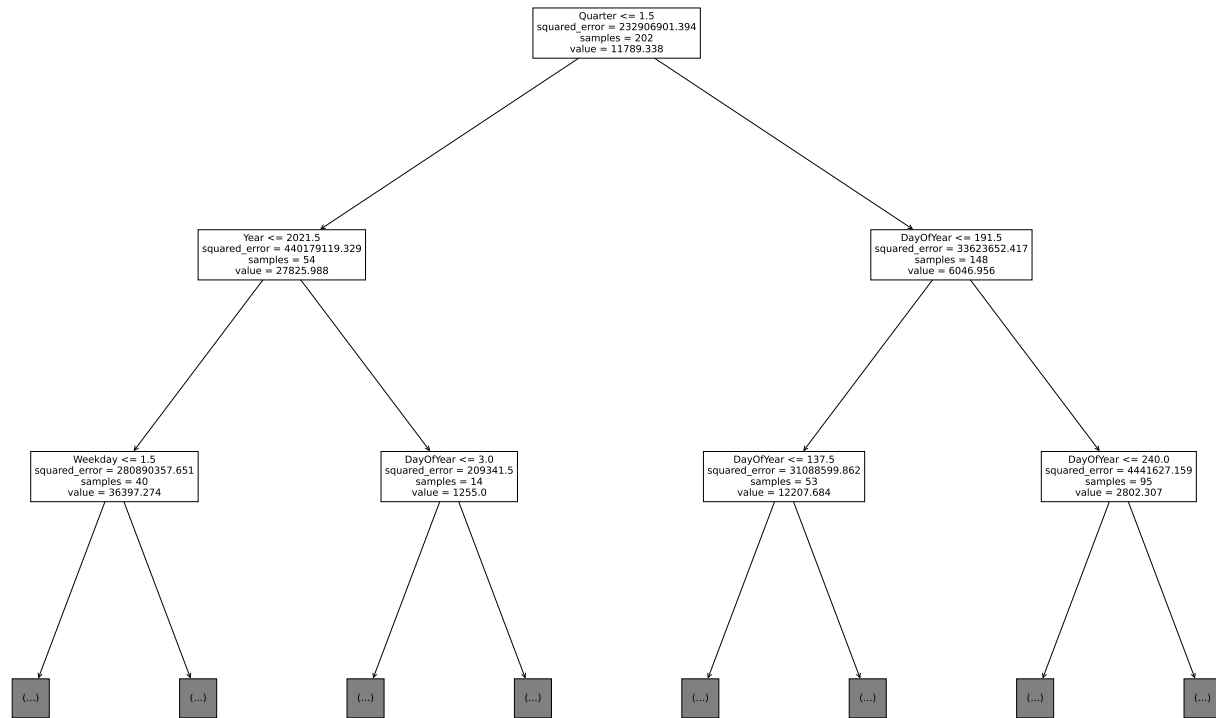
```
## 0.774580856609961
```

```
## 0.7445280765098062
```

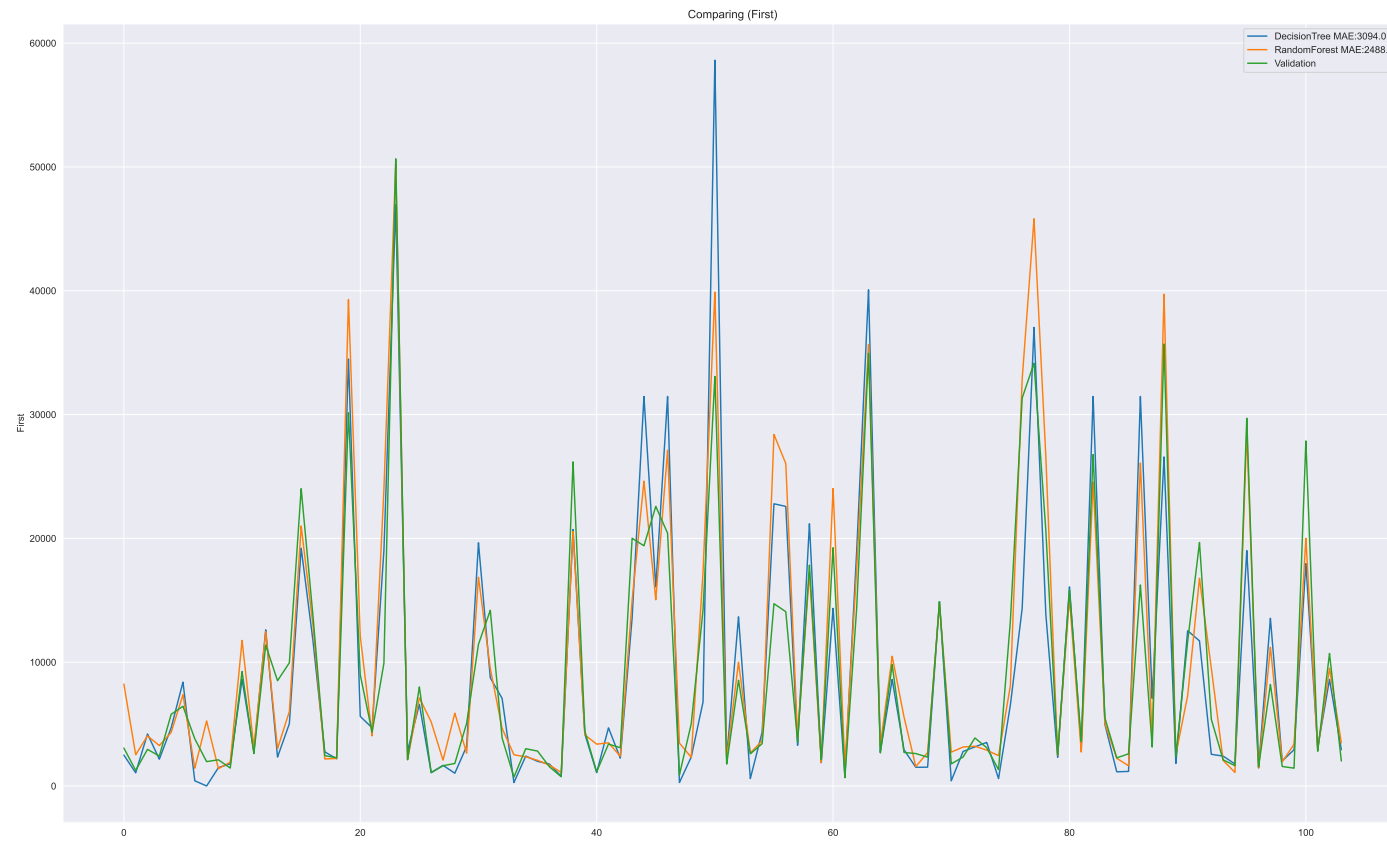
```
## 0.8144473412230163
```

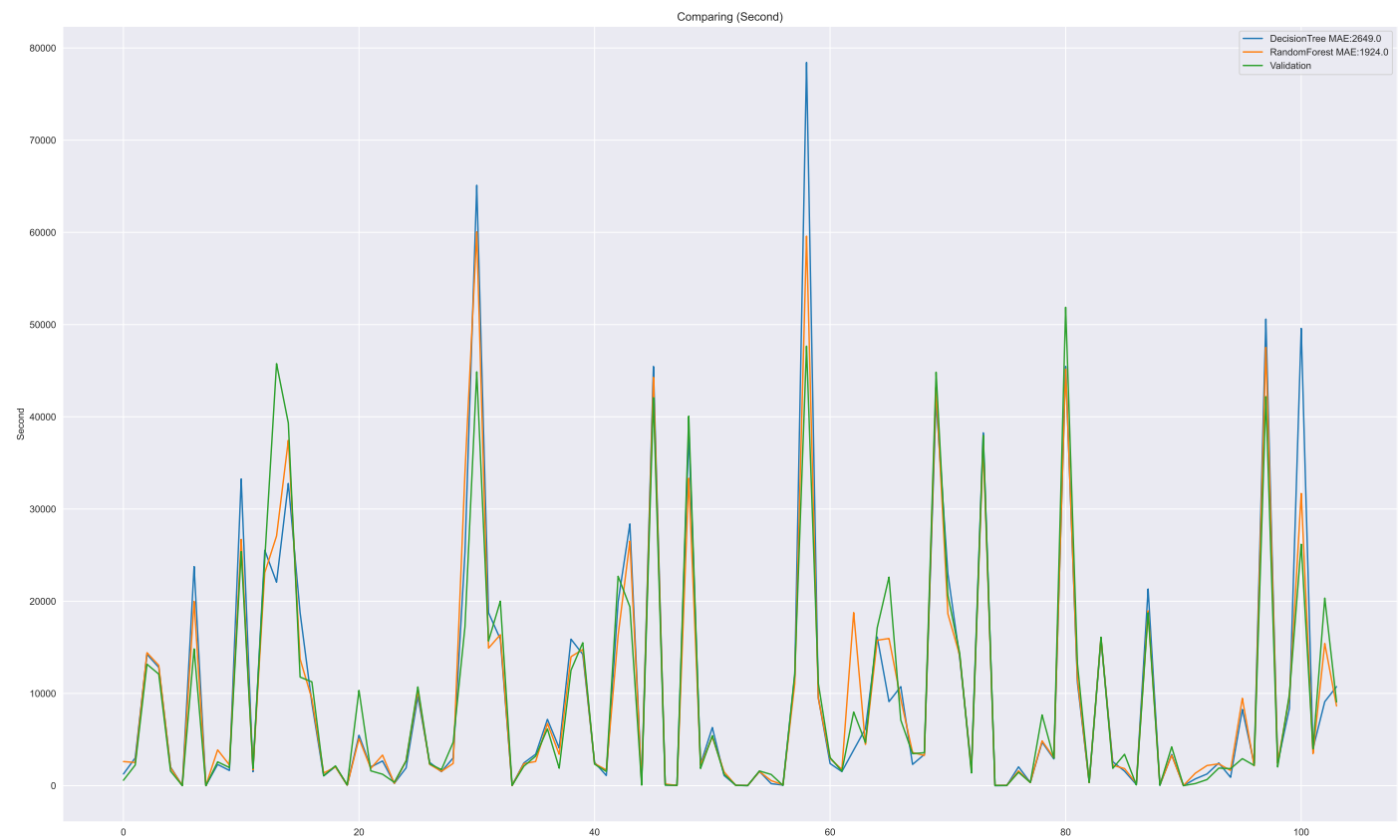
Look at the tree





4.4 Step 3: Plot results





4.5 Step 4: Improve models by changing the dataset

I am going to work with features.

Define necessary variables

Prepare sets and Train models

Compare the score with the mean value of the column that we predicted.

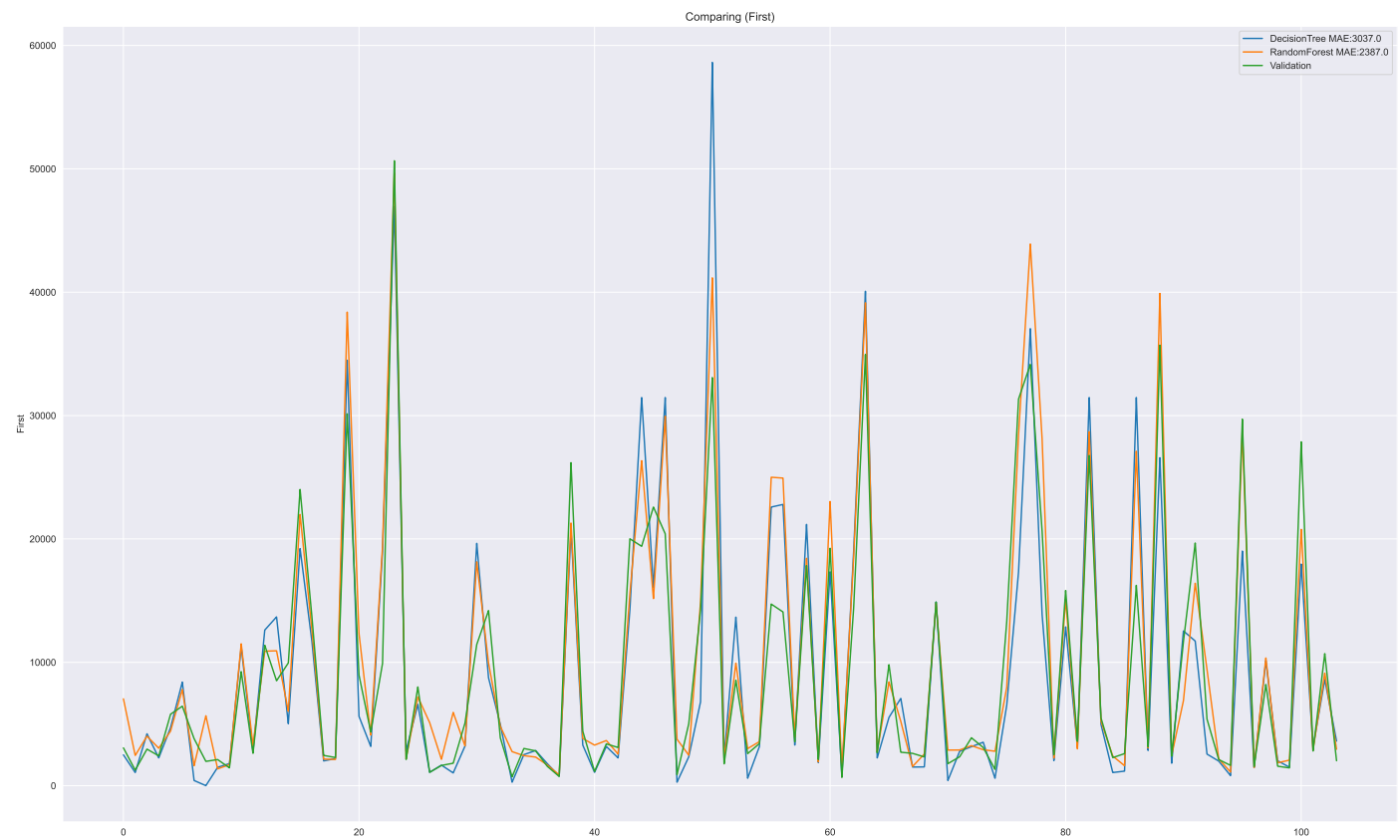
```
## 0.7248630326024768
```

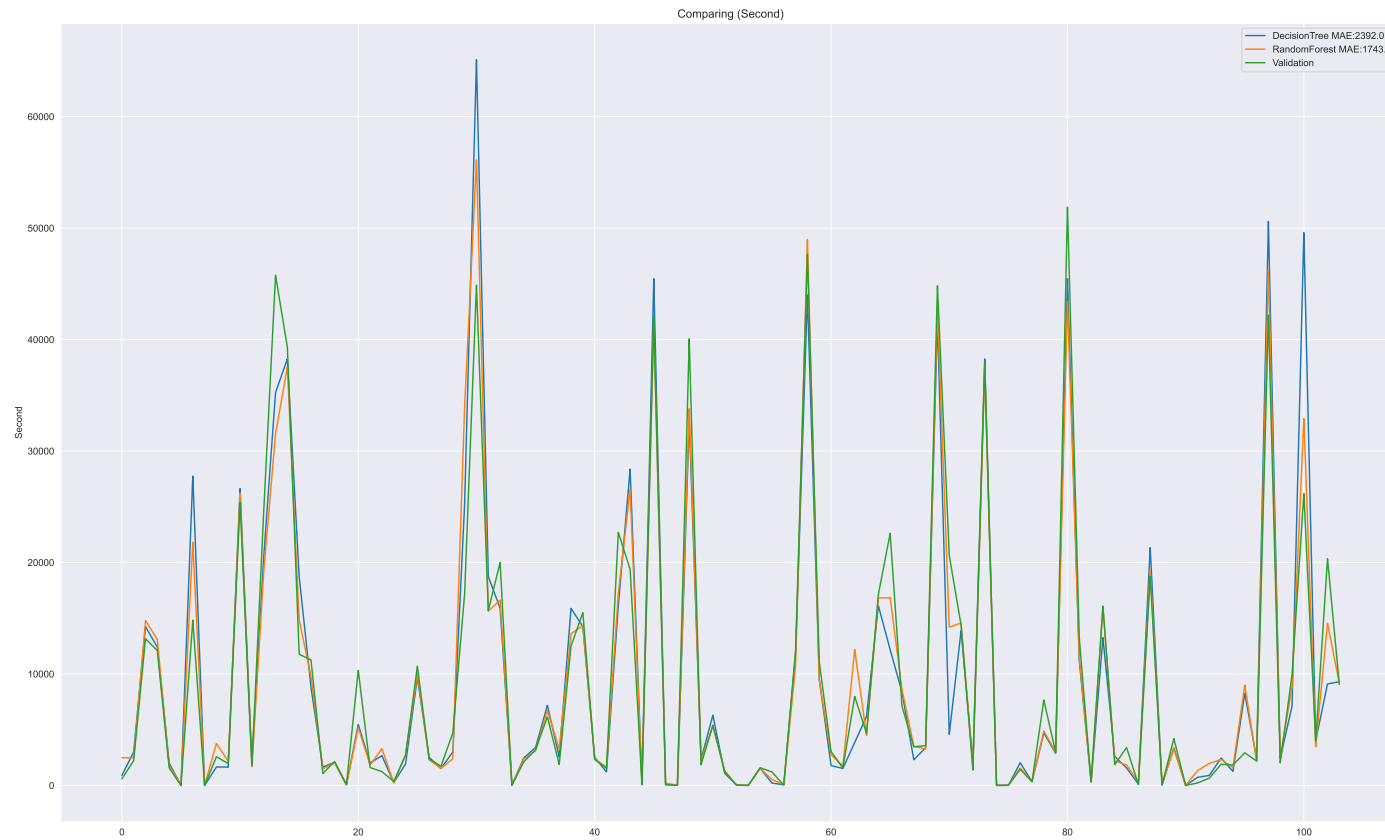
```
## 0.7837038702898657
```

```
## 0.7692636418171874
```

```
## 0.8318561653086721
```

Plot the result.





A combination of the following features give us the best result: * Weekday, * Year, * DayOfYear