

Springboard Data Science Capstone Project – Analysis on Traditional and Non-Traditional Burial in U.K. from 2016-2021

Yvonne Quiachon

September 10, 2021

Table of Contents

Introduction.....	3
Data Acquisition and Cleaning.....	3
Sources.....	3
Exploratory Data Analysis	3
Cremation vs. Burial	4
Crematoriums in the UK.....	5
Mapping and Analysis using Tableau: Crematoriums.....	7
Cemeteries in UK: Burial data.....	9
Mapping and Analysis using Tableau: Burials	10
Feature Engineering	12
Random Forest Regressor	12
Linear Regression	13
Skewness: Crematoriums.....	14
Fixing skewness	15
Random Forest Regressor: Crematoriums	16
Feature Scaling for Modeling: Crematoriums	17
Plotting attributes comparing normalization techniques.....	17
Z-score	17
Modeling	23
Drawing Insights from Analysis	24
Assumptions and Limitations	25

Introduction

Is there a growing demand for non-traditional burial in the UK?

Born out of the curiosity for numbers of COVID-19 related-deaths and concern for burial space in the UK, our goal here is to analyze trends in end-of-life disposition, predict prices for traditional and non-traditional burial, and ultimately provide insight for most reasonable solution to the problem of declining burial space.

According to the combined sources of [Wikipedia](#), [The New York Times](#), [JHU CSSE COVID-19 Data](#), and [Our World in Data](#), there have been 219 million COVID-19 cases and approximately 4.55 million deaths globally. In the UK, there have been 6.86 million cases and 133,178 deaths.

Burials as well as cemetery upkeep are very costly, both fiscally and environmentally. According to Thinkwillow.com, average funeral costs vary depending on the region in the UK. In South West of England funeral costs are around £4,500 while in London it's roughly £6,000. Cremations on average cost around £3,000-£4,000. Traditional burial, which includes basic coffin, hearse, collection and care of loved one, and a funeral director comes at about £5,000. The funeral director alone covers 50-60% of total cost of a funeral - so approximately £2,500. Other amenities and services have additional costs. It's important to remember these figures as I go on.

Take into account the environmental costs as well. An average cremation uses 92 cubic metres of natural gas - about 400 kg of CO₂, equivalent to a 500 miles car trip. Cremations also release toxic gases like dioxins, PCDFs, and mercury. Most coffins used in the UK for cremations are made from chipboard or MDF that release harmful nitric oxide and nitrogen dioxide when burned. Materials used for burials cannot be reused and most items won't decompose for decades. Lastly, it is unlikely or at least difficult to convert cemeteries to woodland or pasture and cities are running out of space for cemeteries (a quarter of England's cemeteries are expected to be full by 2023).

Data Acquisition and Cleaning

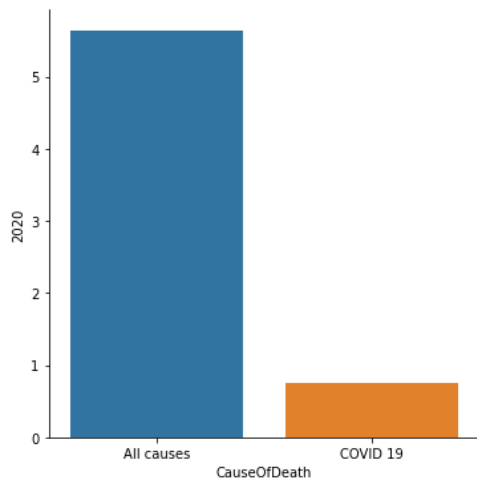
Sources

1. Office for National Statistics
2. Calderdale Metropolitan Borough Council
3. Wikipedia
4. Google maps
5. BBC News
6. Data.gov.uk
7. Individual borough's Ibsites (for burial fees)
8. Cemetery Ibsites

Exploratory Data Analysis

One of main goals of this capstone is to investigate how the covid-19 pandemic has affected funeral services and end-of-life body disposition in the UK, so we are looking at data from 2019 to 2021 (present) specifically. Data for full body burials or cemetery burials are from 2016-2021. The first data presented

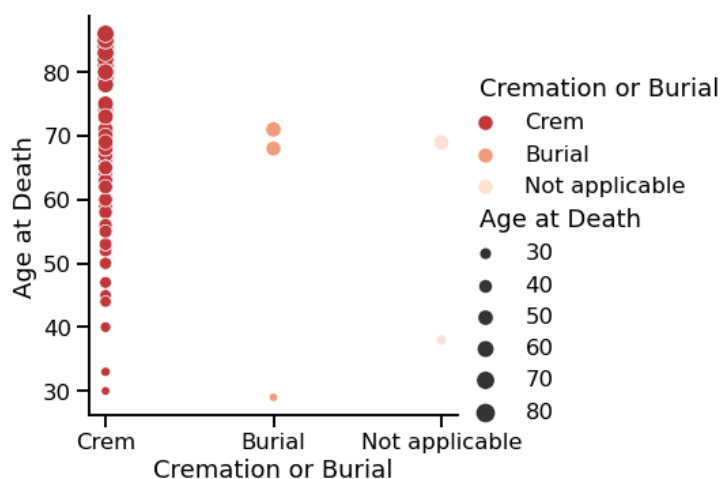
here is from 2020 from the Office for National Statistics, which includes number of deaths from districts across England, and cause of death recorded as either from Covid-19 or other causes simplified into a bar graph:



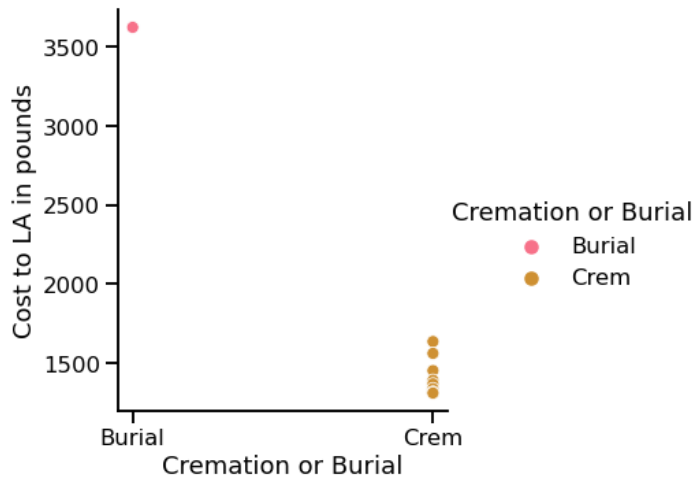
This dataset has 419,309 records, so it does not represent the entire population of UK. In 2020, according to BBC News and other sources, there are around 5.3 million recorded cases of Covid-19, with a total of 128,642 deaths.

Cremation vs. Burial

This data was published by Calderdale Metropolitan Borough Council, which contains public health funeral records in west Yorkshire, England. Public health funerals, according to the publisher and data.gov.uk, are arranged by local authorities for those who have died and have no known relatives to arrange or pay for their funeral or have relatives who do not wish to pay or these people are unable to arrange their funeral.



This figure here clearly shows that cremation was more popular than burial as a method of end-of-life disposition. It also shows the majority of the dead are aged 50 and above.

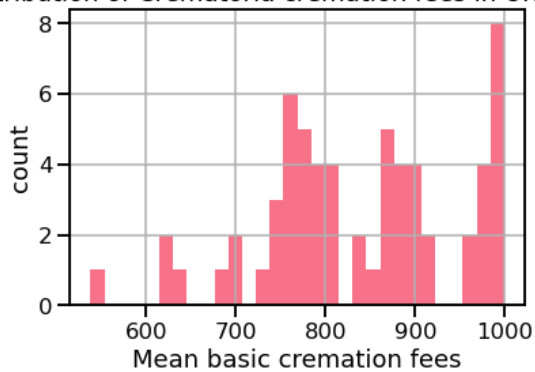


With the same data, I take a look at the costs local authorities/council spent to arrange the funerals. Here, burials exceeds cremations, with burials being $>£3500$ and cremations being $\leq £1700$.

Crematoriums in the UK

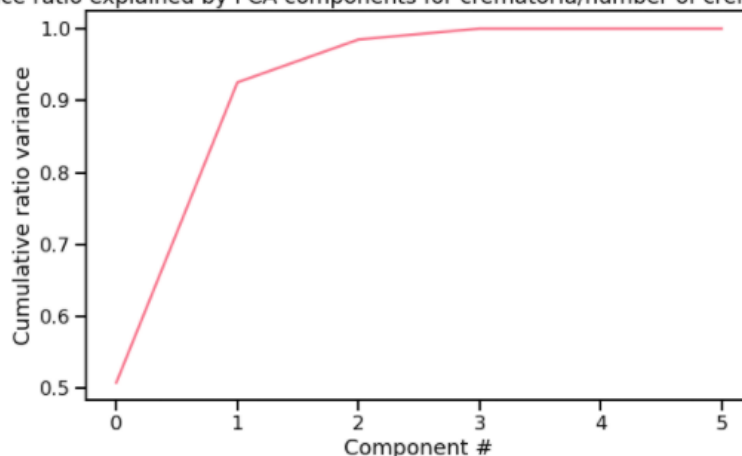
This is the dataset for cremations in the British Islands by Calderdale Metropolitan Borough Council.

Distribution of Crematoria cremation fees in UK in 2021

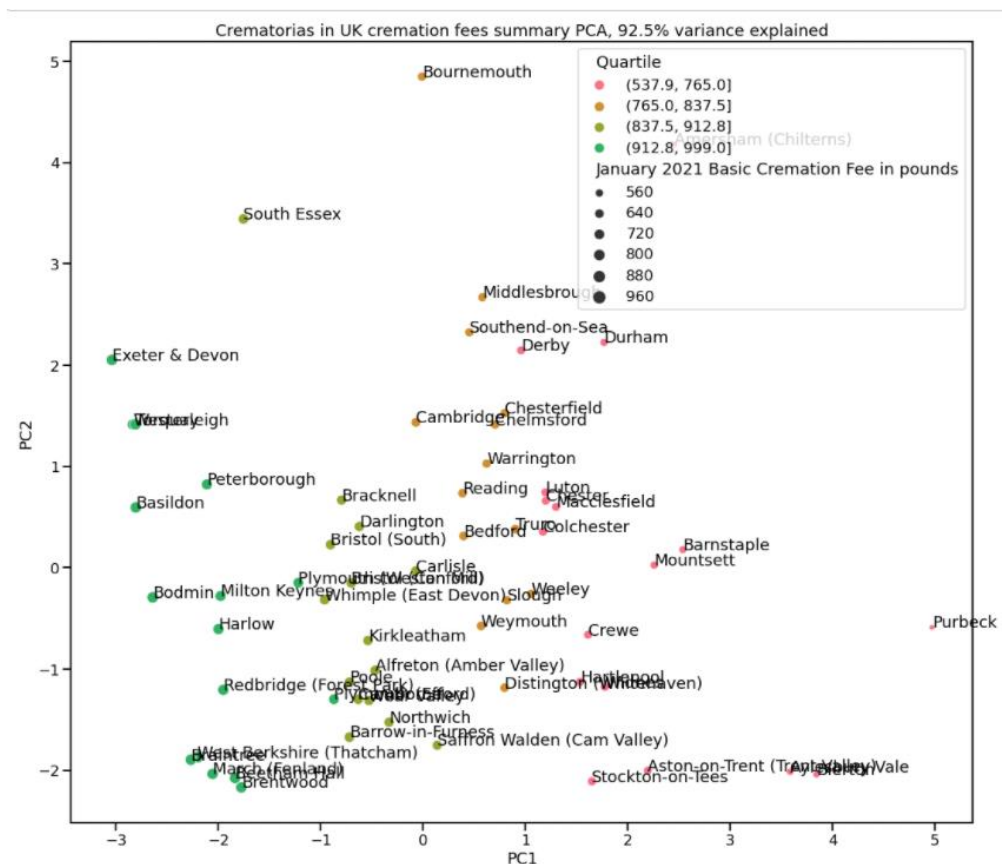


Taking the values of basic cremation fees from January 2021 and grouping by Crematoria, I get this figure for the distribution. Most of the cremation fees are $>£750$.

Cumulative variance ratio explained by PCA components for crematoria/number of cremations summary statistics



Principal component analysis was performed on January 2021 basic cremation fee in pounds. The figure above plots the cumulative variance ratio of number of components. Principal components are new variables that are constructed as linear combinations or mixtures of the initial variables; these combinations are done in a way that the new variables (i.e., principal components) are uncorrelated and most the information within the initial variables is squeezed or compressed into the first components. So, PC1 explains almost 95% of variance, and at PC2 it explains almost 100% total variance.

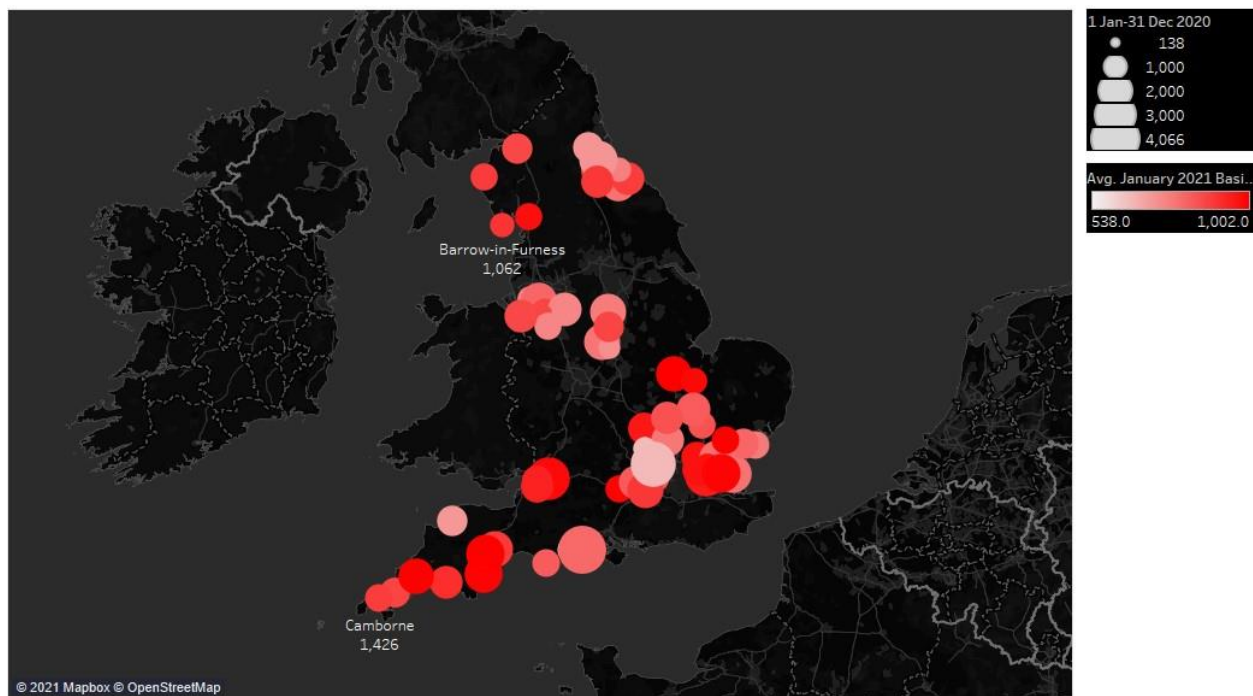


Here is a scatterplot of the summary PCA for cremation fees. They do not appear randomly but instead are clustered based on color; green points representing upper quartile prices are all to the left (negative side), olive green and brown in the middle, and pink to the right (positive side). Geographically, cities/boroughs in pink are smaller towns and farther away from London or in the rural parts near London. On the other hand, the green points are cities/boroughs closer to London; for instance, Redbridge is a London borough in East London, England. If I compare Exeter & Devon (left most) and Purbeck (rightmost) I see that they are both located southwest of England. This is a good visualization of each borough's cremation fees.

Mapping and Analysis using Tableau: Crematoriums

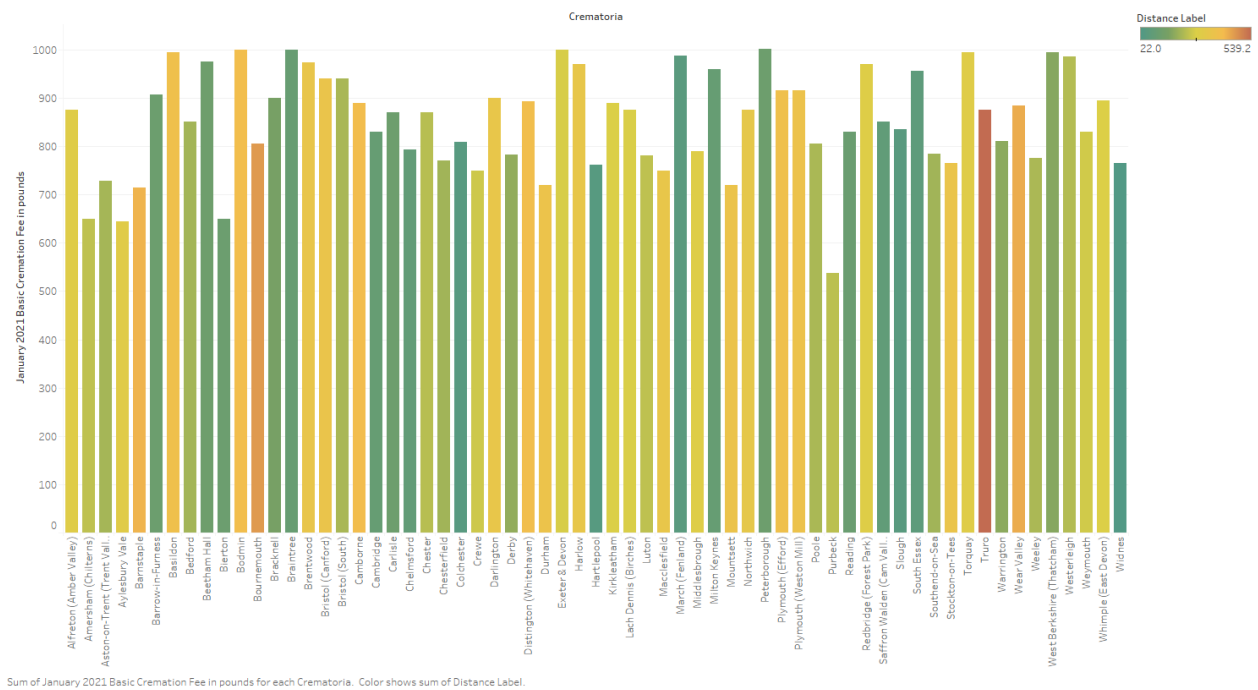
Using Tableau, I mapped the crematoriums. The figure below shows relationship between number of cremations in 2020 and Average Jan 2021 basic cremation fees. The size of the circle shows number of cremations, with larger being more. The color indicates range of cremation fees, with white being less expensive and red more expensive. It is consistent with the scatterplot PCA summary where crematoriums located in southwest England being more expensive and the north being less expensive.

Sheet 1



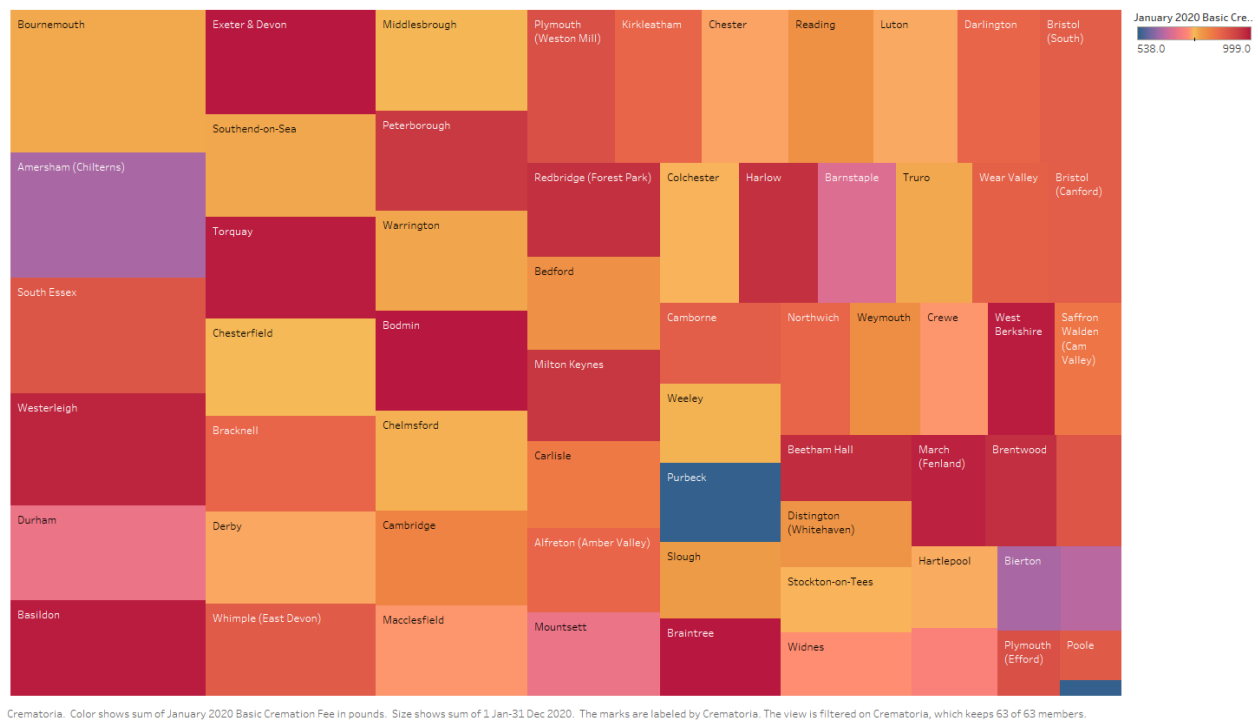
Map based on average of Longitude and average of Latitude. Color shows average of January 2021 Basic Cremation Fee in pounds. Size shows sum of 1 Jan-31 Dec 2020. The marks are labeled by Crematoria and sum of 1 Jan-31 Dec 2020. The view is filtered on Crematoria, which keeps 63 of 63 members.

Sheet 2



A bar chart of Crematoria, January 2021 basic cremation fees and relative distances between crematoria is shown here. Truro, Bournemouth, and Wear Valley stand out. Again, these are in south England.

Sheet 1



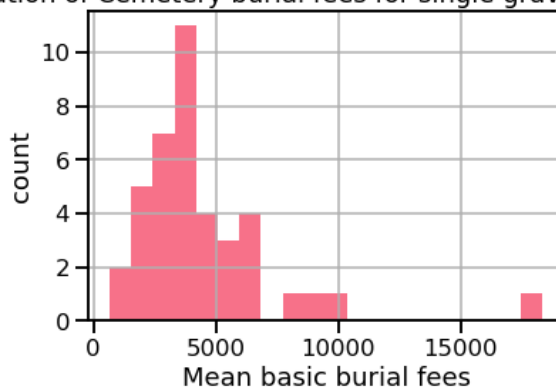
Using Tableau's treemap, I can identify which Crematoria in which district in the UK had the highest basic (average) cremation fees in January 2021 and the most cremations in 2020. The larger the square, the more cremations, and the color bar at the top shows that blue or purple indicates less expensive basic

cremation fees and red is more expensive. According to this, Bournemouth had the most cremations in 2020 and Lach Dennis (Birches Remembrance Park & Crematorium) had the least. Based on the scale at the top right, Exter & Devon, Torquay, Bodmin, West Berkshire, March (Fenland), Westerleigh, Basildon, and Braintree have the highest cremation fees in January 2021. Vice versa, Purbeck, Lach Dennis, Amersham (Chilterns), Bierton, and Aylesbury Vale had the lowest cremation fees. Again, these are consistent with the PCA summary.

Cemeteries in UK: Burial data

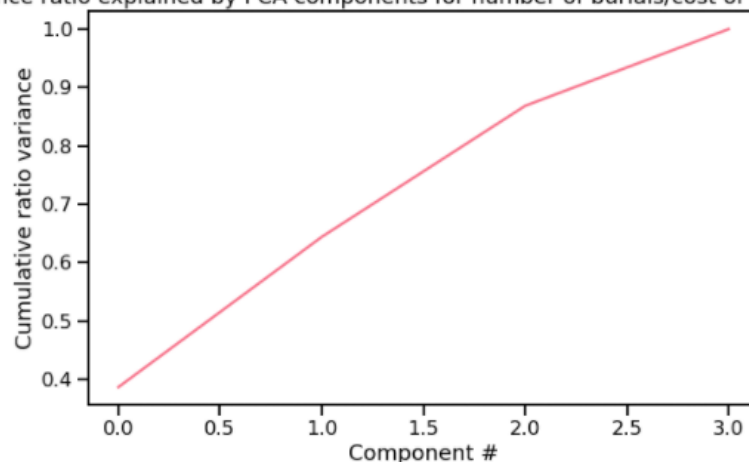
There were no available free data for number of graves for cemeteries, burial sites, and/or graveyards in the UK. Here, instead, I performed secondary data collection from sources such as Wikipedia, Google maps, data.gov.uk, UK cemetery websites, and UK boroughs council's websites. I did in order to: 1) produce a list of burial sites 2) obtain number of graves/bodies buried for each burial site 3) get the coordinates for each burial site 4) get the burial fees for each burial site. Information obtained is from 2016-2021, and only cemeteries that are still open to burying more bodies (not closed) was included. In addition, these cemeteries are in the south of England only.

Distribution of Cemetery burial fees for single grave (adult) in UK

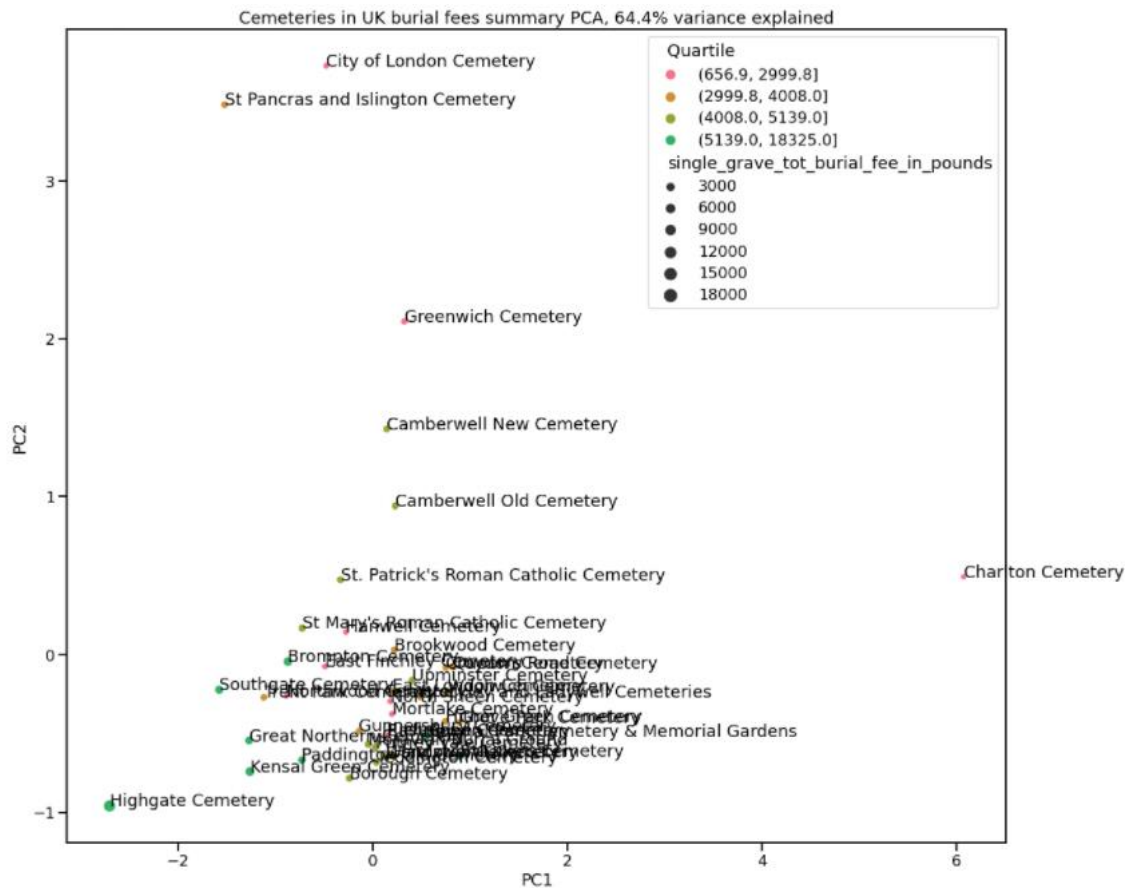


Using the same method for the crematorium data, I find the distribution of cemetery burial fees for a single, adult grave in the UK. Most lie <£5,000, around £2,000 to £4,000.

Cumulative variance ratio explained by PCA components for number of burials/cost of burials summary statistics



Another PCA was performed, this time for the same burials dataset I collected. Again, I was trying to preserve as much variance as possible. The above figure is the explained variance for cost of burials. Here, PC1 explains only 65% variance, PC2 explains almost 90% variance, and at PC3 100% variance.



Using seaborn, I do the same thing as for crematoriums data onto burials data, creating a scatterplot of PCA summary showing quartile ranges. Most of the cemeteries lie between 0 and 1, with about 4 outliers. Highgate Cemetery was in the negative side and represents the upper quartile range, while Charlton Cemetery was in the positive side and represents the lower quartile range. Geographically, Highgate Cemetery is in London. It seems, with Charlton Cemetery as an outlier, I may have chosen the wrong cemetery (there is a Charlton Cemetery in London and another in Dover, Kent).

Mapping and Analysis using Tableau: Burials

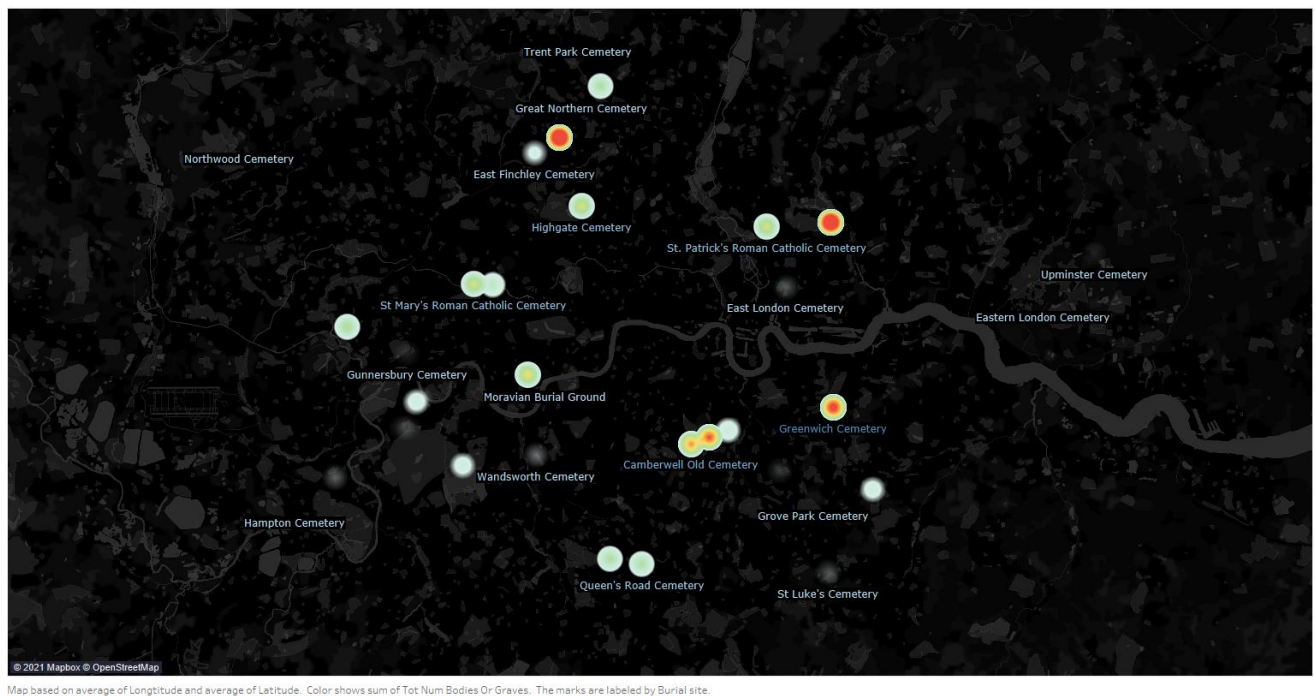
The figure below shows a relationship between number of burials/graves and single grave total burial fees (for adult). Highgate Cemetery shows to be the most expensive, although not the most in terms of bodies buried. It seems St. Pancras and Islington Cemetery has the most bodies buried and/or most graves. Charlton Cemetery is not shown in the map, which indicates I have the wrong coordinates for this cemetery. The correct Charlton Cemetery is in the Royal Borough of Greenwich, which is the data we used for burial fees (both Charlton Cemetery and Greenwich Cemetery are pink in the PCA summary).

Sheet 1



Below is a density map where an intense orange color signifies more burials.

Sheet 1



A density map is shown here, where an intense orange color signifies more burials. Again, City of London Cemetery and St. Pancras & Islington Cemetery have the most burials.

Feature Engineering

After cleaning data and exploring the data, I did some feature engineering on the cremation vs. burial data and the crematoriums data.

Random Forest Regressor

I one-hot-encoded the cremation vs burial data here and performed random forest regressor. I chose random forest because it requires few parameters to set, it is simple, and performs well. As shown in the code, the final data contains the costs local authorities/council spent on funerals, cremation or burial, and genders.

Convert categorical variable, Cremation or Burial, into indicator variables or dummy

```
burial_subset = df_final[['Age at Death', 'Gender', 'Cost to LA in pounds', 'Cremation or Burial']]
```

```
burial_dummies = pd.get_dummies(burial_subset['Cremation or Burial'], prefix='Cremation or Burial')  
ds = pd.concat([burial_subset, burial_dummies], axis = 1)
```

ds

	Age at Death	Gender	Cost to LA in pounds	Cremation or Burial	Cremation or Burial_Burial	Cremation or Burial_Crem
0	53	Male	214.00	Crem	0	1
1	40	Male	1002.00	Crem	0	1
2	63	Male	150.00	Crem	0	1
3	70	Male	187.00	Crem	0	1
4	62	Male	68.00	Crem	0	1
5	81	Male	1392.00	Crem	0	1
6	44	Female	1005.00	Crem	0	1
7	79	Female	1060.00	Crem	0	1
8	60	Male	1103.00	Crem	0	1

```
from sklearn.preprocessing import OneHotEncoder
```

```
one_hot = OneHotEncoder()  
encoded = one_hot.fit_transform(ds[['Cremation or Burial']])
```

```
burial_subset[one_hot.categories_[0]] = encoded.toarray()
```

burial_subset

	Age at Death	Gender	Cost to LA in pounds	Cremation or Burial	Burial	Crem
0	53	Male	214.00	Crem	0.0	1.0
1	40	Male	1002.00	Crem	0.0	1.0
2	63	Male	150.00	Crem	0.0	1.0
3	70	Male	187.00	Crem	0.0	1.0
4	62	Male	68.00	Crem	0.0	1.0
5	81	Male	1392.00	Crem	0.0	1.0
6	44	Female	1005.00	Crem	0.0	1.0
7	79	Female	1060.00	Crem	0.0	1.0
8	60	Male	1103.00	Crem	0.0	1.0

```
# Repeat for Gender
gender_dummies = pd.get_dummies(ds['Gender'], prefix='Gender')
ds = pd.concat([ds, gender_dummies], axis = 1)
ds.head()
```

	Age at Death	Gender	Cost to LA in pounds	Cremation or Burial	Cremation or Burial_Burial	Cremation or Burial_Crem	Gender_Female	Gender_Male
0	53	Male	214.0	Crem	0	1	0	1
1	40	Male	1002.0	Crem	0	1	0	1
2	63	Male	150.0	Crem	0	1	0	1
3	70	Male	187.0	Crem	0	1	0	1
4	62	Male	68.0	Crem	0	1	0	1

Get numerical values into new dataset

```
num_set = ds.select_dtypes(include=['int', 'float'])
num_set.head()
```

	Cost to LA in pounds
0	214.0
1	1002.0
2	150.0
3	187.0
4	68.0

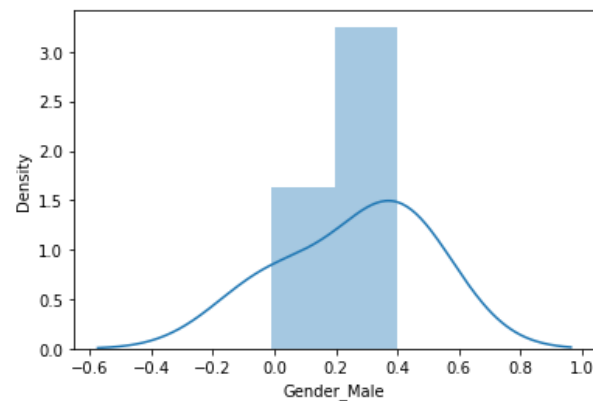
```
data = ds[['Cost to LA in pounds', 'Cremation or Burial_Burial', 'Cremation or Burial_Crem', 'Gender_Female', 'Gender_Male']]
```

```
# split into inputs and outputs
X, y = data.iloc[:, :-1], ds_.iloc[:, -1]
```

```
print(X.shape, y.shape)
```

```
(10, 4) (10,)
```

Below is a result for density plot of the random forest predicted probabilities, with gender as y variable and cost to LA in pounds and cremation or burial as the set of predictor X variables.



Linear Regression

Using Standard Scaler on numeric features Cost to LA in pounds and Age at Death, I did a Linear Regression.

```
num_subset = df_final[['Cremation or Burial', 'Age at Death', 'Cost to LA in pounds']]
num_subset.head()
```

	Cremation or Burial	Age at Death	Cost to LA in pounds
0	Crem	53	214.0
1	Crem	40	1002.0
2	Crem	63	150.0
3	Crem	70	187.0
4	Crem	62	68.0

```
# Repeat splitting into train and test set
x = num_subset.iloc[:, -1].values
y = num_subset.iloc[:, 1].values
```

y

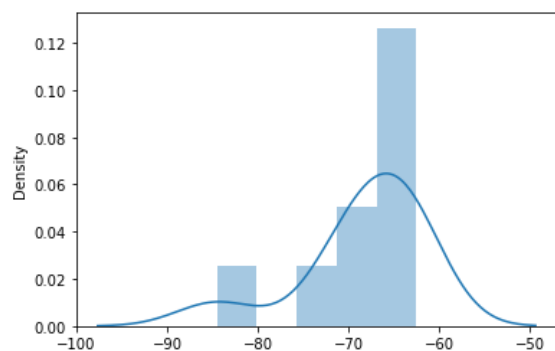
```
array([53, 40, 63, 70, 62, 81, 44, 79, 60, 84, 71, 67, 74, 52, 79, 52, 60,
       65, 70, 70, 83, 45, 64, 56, 85, 68, 59, 33, 56, 78, 63, 65, 65, 86,
       71, 60, 71, 68, 80, 80, 69, 58, 70], dtype=int64)
```

x

```
array([ 214. , 1002. ,  150. ,  187. ,   68. , 1392. , 1005. ,
       1060. , 1103. , 1216. , 1167. ,  195. , 1081. , 1197. ,
       1103. , 1174. , 1114. ,  396. ,  427. ,  346. ,    8. ,
       1173. , 1369. , 1247. ,  122.7, 3625. , 1313. , 1210. ,
       1227. , 1562. , 1311. , 1209. , 1234. , 1637. ,   55. ,
       1453. , 1271. , 1266. , 1142.75, 1084.25, 1317. , 1270. ,
       1339.  ])
```

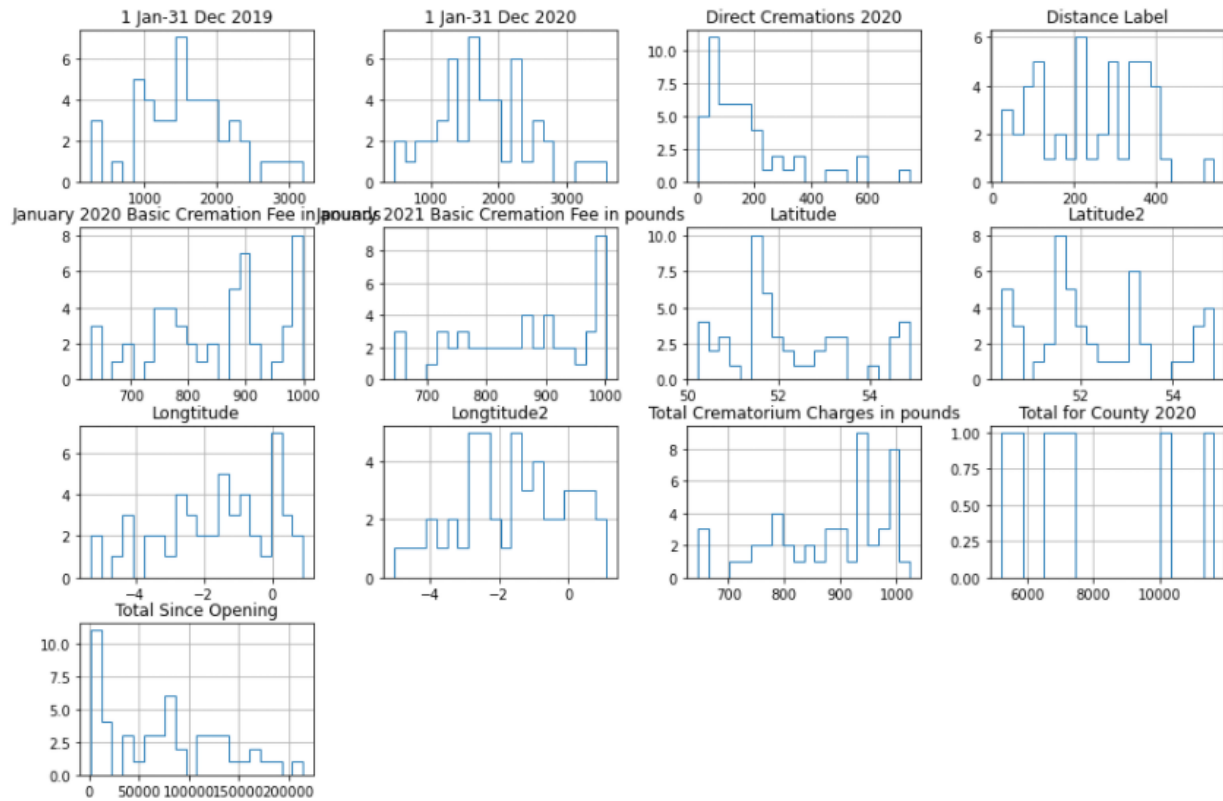
```
X_train, X_test, y_train, y_test = train_test_split(x, y, test_size = 0.2, random_state = 1)
```

Here is a density curve of linear regression probabilities, with age at death as y variable and age at death the x variable.



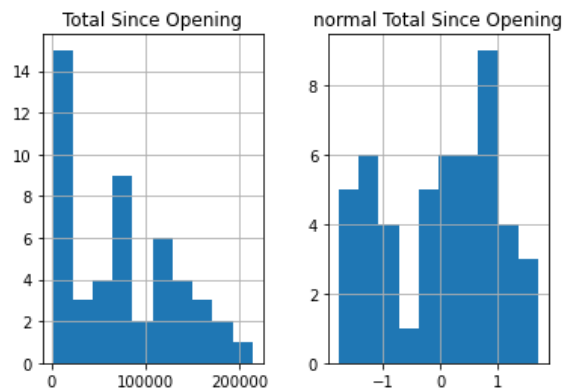
Skewness: Crematoriums

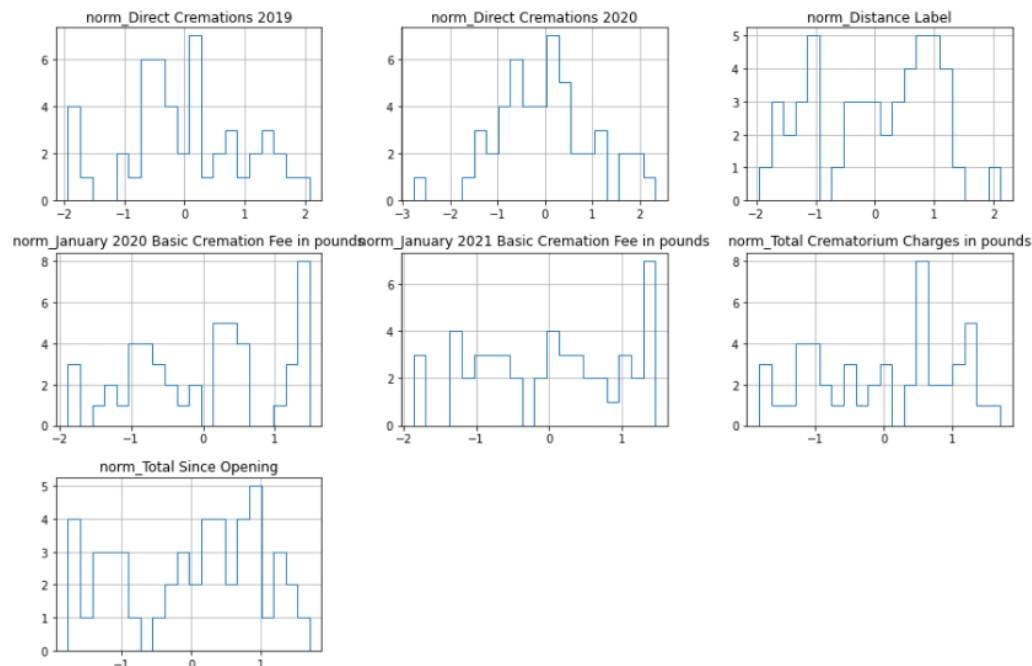
I wanted to visualize the skewness of the crematorium data, so I created a histogram for each feature.



The features that are left-skewed are: direct cremations 2020, distance label, and total since opening. Features that are right-skewed are: January 2020 basic cremation fee in pounds, January 2021 basic cremation fee in pounds and total crematorium charges in pounds. The following histograms is where I fix skewness for these features.

Fixing skewness





Random Forest Regressor: Crematoriums

As I did before with the public health burials dataset (cremation vs. burial), I performed random forest regression this time for crematoriums data. I first one-hot encoded the data, then split the data into training and test sets, with the following normalized features: total since opening, January 2020 basic cremation fees, January 2021 basic cremation fees, and total crematorium charges.

```
# create dummy variables for Regional indicator
dummy_df = pd.get_dummies(xl3_noskew['Crematoria'])
dummy_df
```

	Alfreton (Amber Valley)	Amersham (Chilterns)	Aston- on- Trent (Trent Valley)	Aylesbury Vale	Barnstaple	Barrow- in- Furness	Basildon	Bedford	Blerton	Bodmin	...	South Essex	Southend- on-Sea
0	1	0	0	0	0	0	0	0	0	0	...	0	0
1	0	1	0	0	0	0	0	0	0	0	...	0	0
2	0	0	1	0	0	0	0	0	0	0	...	0	0
3	0	0	0	1	0	0	0	0	0	0	...	0	0
4	0	0	0	0	1	0	0	0	0	0	...	0	0

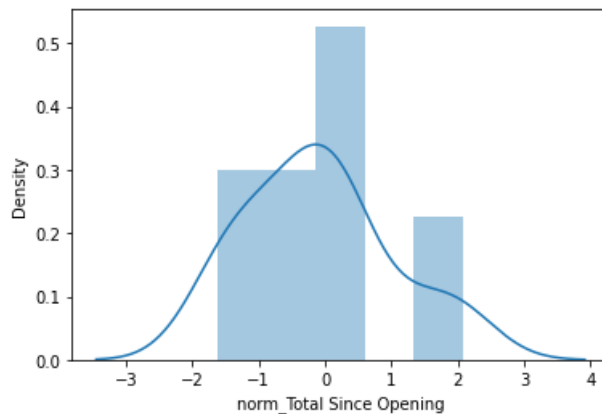

```
x13_final_train = x13_final[['norm_Total Crematorium Charges in pounds', 'norm_January 2020 Basic Cremation Fee in pounds', 'norm_January 2021 Basic Cremation Fee in pounds', 'norm_Total Since Opening']]
```

```
x13_final_train = pd.DataFrame(x13_final_train).fillna(0)
```

```
x13_final_train['norm_Total Since Opening'].unique()
```

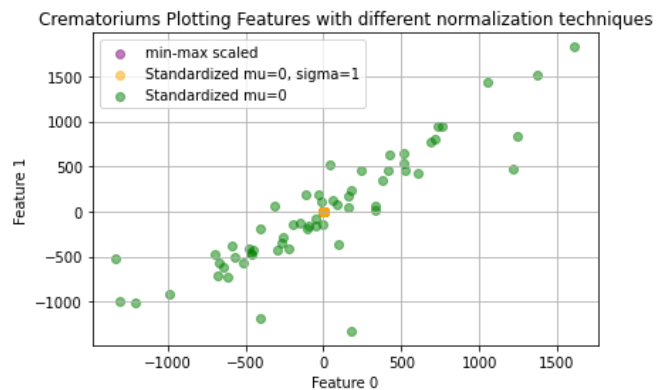
```
array([-1.253,  1.118, -1.697, -1.761,  0.13 , -0.126, -0.253,  0.408,
        -1.772, -0.332,  0.188, -1.117, -1.328,  0.976,  0.418, -1.062,
         1.319,  0.496,  0.794,  0.698,  0.116,  0.874,  1.526,  0.812,
         1.288,  0.381, -0.045, -1.209,  0.314, -1.199,  0.163, -1.296,
         0.839,  0.712,  0.302,  0.292,  0.994, -0.821, -1.578,  1.704,
         1.395, -1.713,  1.193,  0.965,  0.531, -0.949, -0.037, -0.981,
        -0.42 ,  0.    ])
```

Here is a normal density curve for random forest probabilities, with total since opening as the response variable y and total crematorium charges, January 2020 basic cremation fees and January 2021 basic cremation fees as set of predictor variables x.



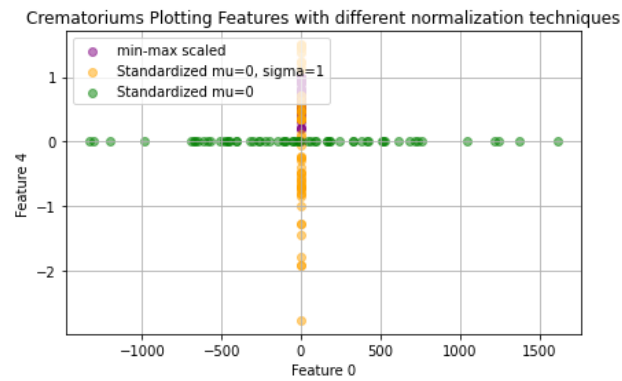
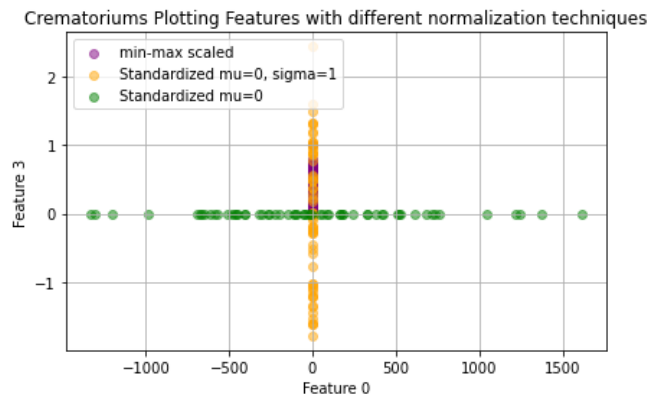
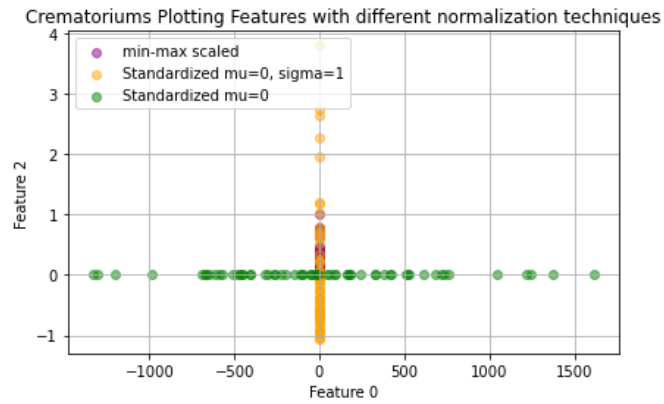
Feature Scaling for Modeling: Crematoriums

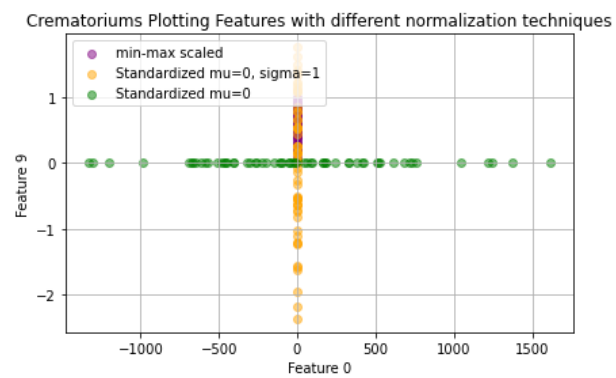
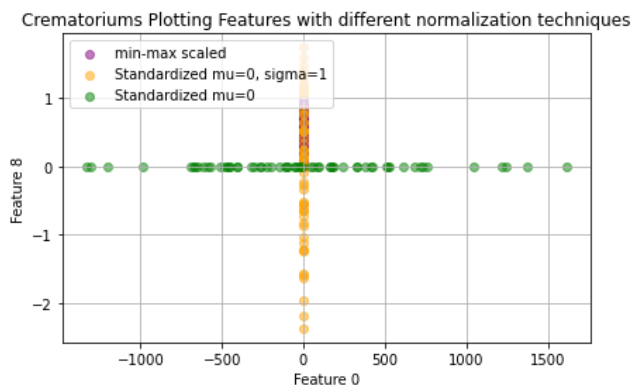
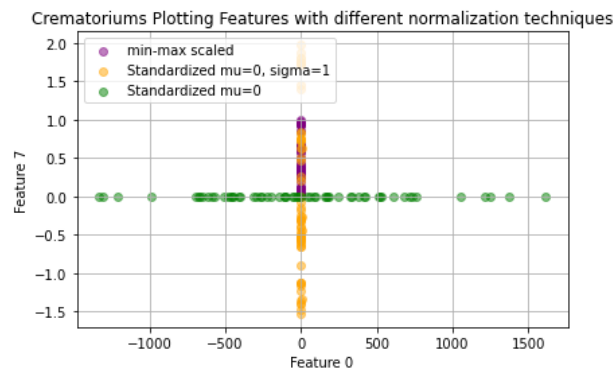
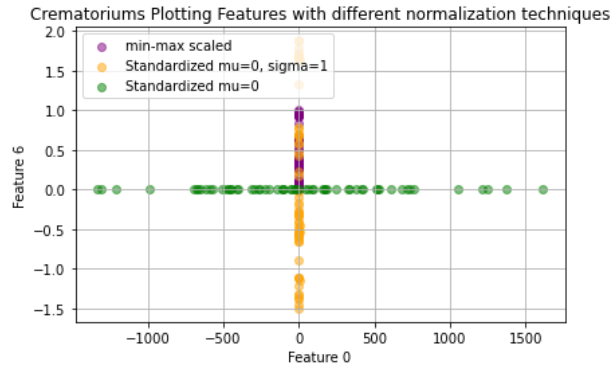
Plotting attributes comparing normalization techniques

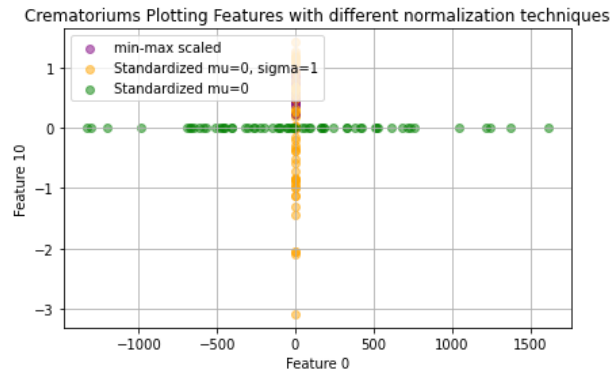


Z-score

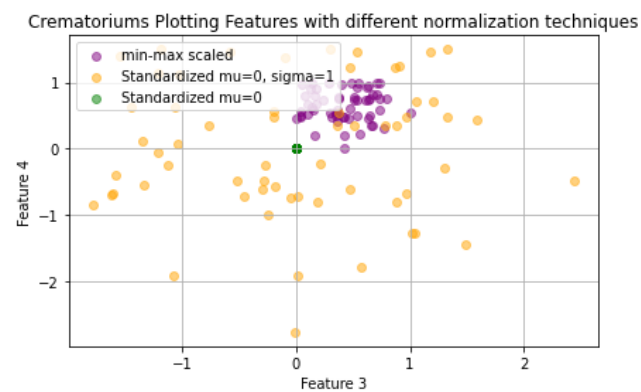
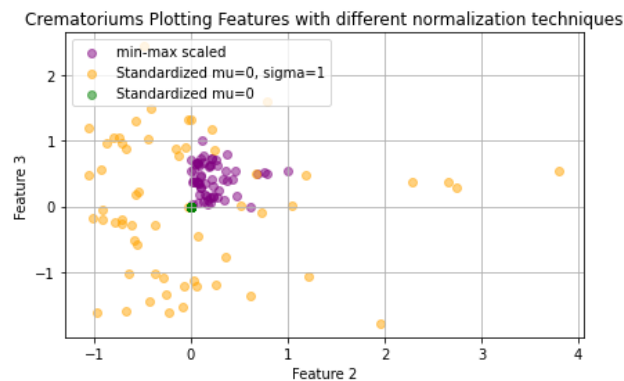
Z-score here is 0, which indicates the value is exactly the mean. This is features 1 Jan-31 Dec 2019 and 1 Jan-31 Dec 2020. The following plots shows feature 1 Jan-31 Dec 2019 and its relationship with the other features. For all, the z-score is 0.



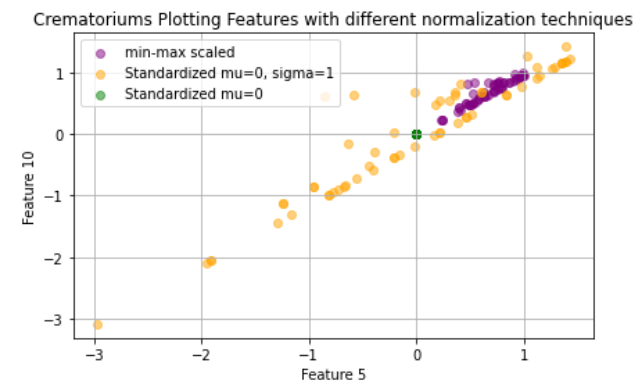
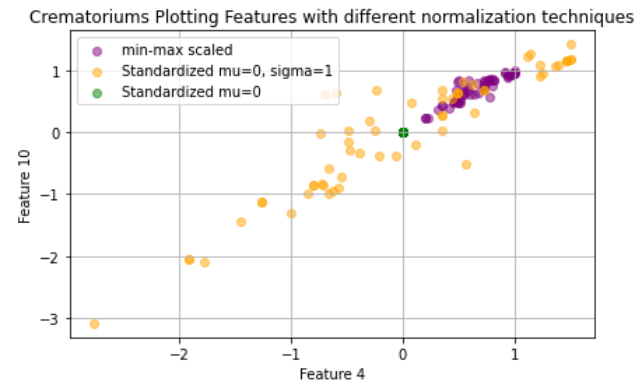
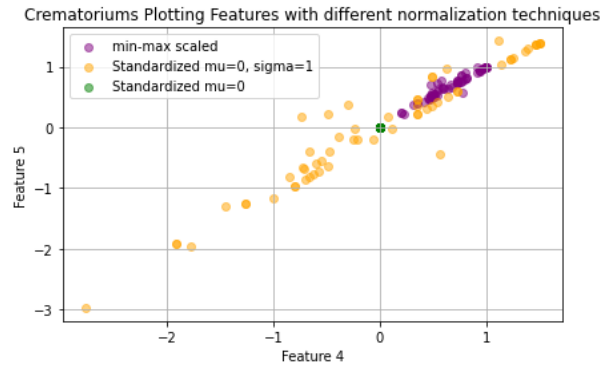




The results are similar for Feature 1 (1 Jan-31 Dec 2020) on the x axis and Features 2 to 10 on the y axis. The next scatterplots shows the mean subtracted (standardized) at 0 and the min-max range (normalized) 0 to 1.

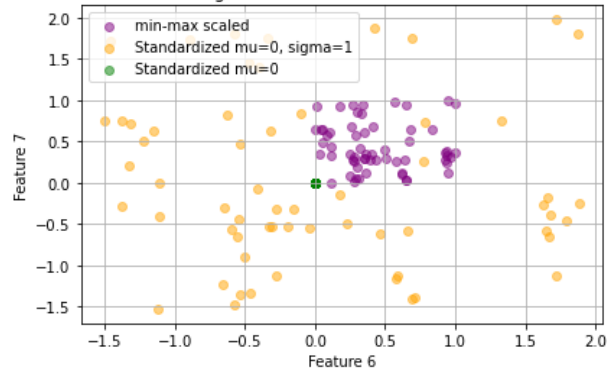


The following plots shows z-scores and min-max ranges in line with each other, forming a positive slope between ranges -1 to 1. This is where standardized features and normalized features fit (match). Feature 4 is January 2020 basic cremation fee and feature 5 is January 2021 basic cremation fee.

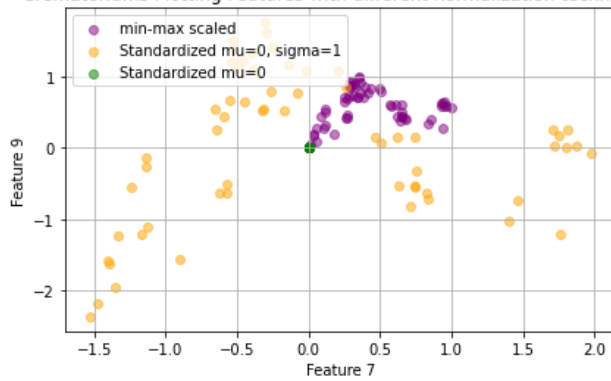


In the next scatterplots, normalized features are between 0 and 1, with z-scores having random patterns. The mean subtracted standardized features is at 0. Here, feature 6 is latitude and feature 7 is longitude. Feature 8 is total crematorium charges and feature 9 is total since opening.

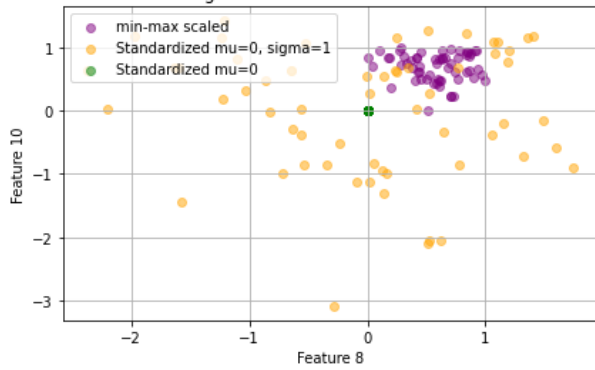
Crematoriums Plotting Features with different normalization techniques



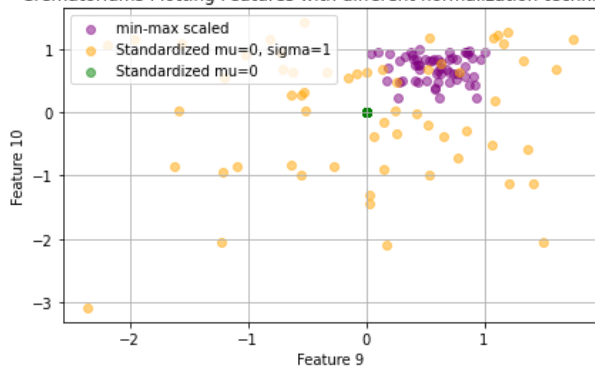
Crematoriums Plotting Features with different normalization techniques



Crematoriums Plotting Features with different normalization techniques



Crematoriums Plotting Features with different normalization techniques



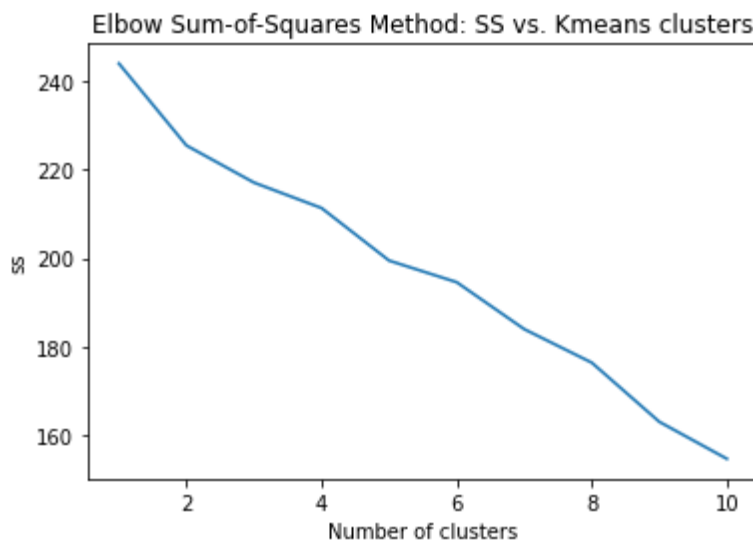
Modeling

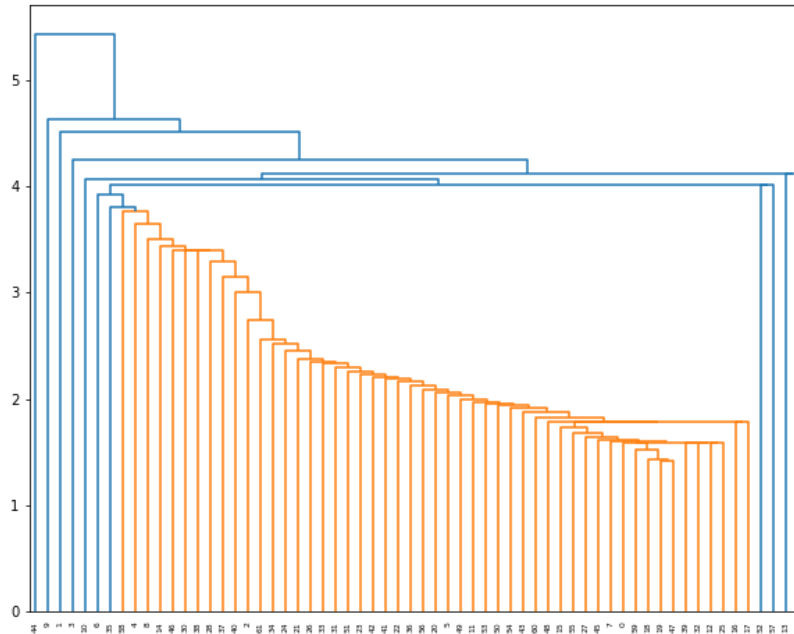
After feature engineering dataset for cremation vs. burial and crematoriums, I realized it would be very difficult to create a model with the limited data, as it only spans one year and a time-series analysis or synthetic control model would be infeasible. Although some of the crematorium vs burial data is from 2016 to 2020, it is only a cumulative amount and does not specify number of cremations or burials per year.

Instead, a point-biserial correlation was performed for cremation vs. burial dataset. Point-biserial correlation is the same as the Pearson correlation coefficient used in linear regression, with the difference being dichotomous data is being compared to continuous data instead of continuous data to continuous data. From this point on I assume that our dichotomous data is composed of items from two groups (group 0 and group 1) and the continuous data as “y”. Our results are:

PointbiserialResult(correlation = 0.291, p-value=0.06)

I ran a K-means and agglomerative/hierarchical clustering on crematorium dataset.





Drawing Insights from Analysis

With the data and data visualizations presented, we can draw insights on traditional and non-traditional burials in the UK. Our present data suggests that cremation is the leading method of end-of-life disposition. Multiple factors play into the method of choice, but here we can confidently say that cost is a key factor. Average cremation fees in January 2021 were between £750 to £1500, while average burial fees for a single grave for an adult from 2016 to 2021 were between £1000 to £5000. With that said, burial fees are approximately 20% more expensive than cremation fees. For volumes of bodies disposed, the highest number of cremations by a crematoria (in this data, it was Bournemouth Crematorium) in 2020 was 4,066. The cemetery with the most bodies buried, St. Pancras & Islington Cemetery, has 812,000 burials. Even though this cemetery is 190 acres, it will eventually reach its full capacity. Graves can be reused, but an exclusive right of burial or cremated remains plot is sold for at least 50 to 100 years, so cemeteries can get crowded whenever there is a surge of deaths within a population as the demand for plots increases while the supply is limited. On the other hand, there are various ways to preserve cremation ashes other than burial.

When looking at the mapping analysis, we can see that cremation costs are highest in southwest England. Peterborough had the highest basic cremation fees in 2021 (£1,002), and Purbeck had the lowest basic cremation fees (£538). These two crematoriums are located in South England. In the bar graph where distances between crematoriums were randomized, Truro had the longest distance because it was relative to Middlesbrough, where Truro is at the very south west of England and Middlesbrough in the North. It was not a very successful graph in that it did not show a clear relationship between distance and cremation fees. In the burials data, the area was concentrated in South England around London area. From 2016-2021, Highgate Cemetery had the highest burial costs (approximately £18,325) and Charlton Cemetery the lowest burial costs (about £657). The Magnificent Seven are the seven large private cemeteries in London, and these were established in the 19th century to alleviate overcrowding in existing

parish burial grounds. Besides Highgate Cemetery, these include Kensal Green Cemetery, West Norwood Cemetery, Brompton Cemetery, Abney Park Cemetery, Nunhead Cemetery and Tower Hamlets Cemetery. Abney Park Cemetery, Tower Hamlets Cemetery, and Nunhead Cemetery are closed cemeteries. Interestingly, these seven cemeteries form a circle around central London. So, Kensal Green and Brompton cemeteries, according to our mapping analysis, are in the mid range of burial costs, between £8,000-£10,000, which is still higher than the majority cemeteries. Hence, privately owned cemeteries are more expensive. My mapping analysis is congruent with the scatter plot PCA summaries for both crematoriums and burial datasets.

Assumptions and Limitations

There were several limitations that prevented my from creating a good model. For instance, the data for crematorium spans only one year, so it was impossible to create a time-series analysis. The burial data spans from 2016 to 2021, but the data was cumulative and there was no data present for number of burials or burial fees per year. My ultimate goal was to create a robust synthetic control model for a causal impact model to make predictions for burial and cremation rates in scenario where Covid-19 did not happen. This proved to be difficult to achieve, but a new perspective was discovered during this project that could be applied to future work or make improvements. This includes collection of more secondary data for cemeteries outside of London area, further research for burial and cremation records divided years spanning more than 8 to 10 years for building a good model, and ensuring that coordinates for locations are precise.