

1. Data preparation

In [1]:	<pre>import pandas as pd import seaborn as sns import matplotlib.pyplot as plt import random import time import joblib from sklearn.preprocessing import MinMaxScaler from sklearn.ensemble import RandomForestRegressor, GradientBoostingRegressor from sklearn.metrics import make_scorer from sklearn.model_selection import LeaveOneGroupOut, cross_val_score, GridSearchCV # Set random seed for reproducibility random.seed(123)</pre>																																																																																																																								
In [2]:	<pre>#Read csv and excel files cost = pd.read_csv('cost.csv') suppliers = pd.read_csv('suppliers.csv') tasks = pd.read_excel('tasks.xlsx')</pre>																																																																																																																								
Looking at the dataframes and checking their information (number of columns, rows, data types)																																																																																																																									
In [3]:	<pre>suppliers.head()</pre>																																																																																																																								
Out[3]:	<table><thead><tr><th></th><th>Supplier ID</th><th>SF1</th><th>SF2</th><th>SF3</th><th>SF4</th><th>SF5</th><th>SF6</th><th>SF7</th><th>SF8</th><th>SF9</th><th>SF10</th><th>SF11</th><th>SF12</th><th>SF13</th><th>SF14</th><th>SF15</th><th>SF16</th><th>SF17</th><th>SF18</th></tr></thead><tbody><tr><td>0</td><td>51</td><td>100</td><td>1000</td><td>1000</td><td>50</td><td>20</td><td>10</td><td>2</td><td>80</td><td>2000</td><td>100</td><td>1000</td><td>5</td><td>1000</td><td>1000</td><td>500</td><td>5000</td><td>100</td><td>96</td></tr><tr><td>1</td><td>52</td><td>100</td><td>1000</td><td>1000</td><td>50</td><td>20</td><td>10</td><td>0</td><td>80</td><td>2000</td><td>100</td><td>1000</td><td>5</td><td>1000</td><td>1000</td><td>500</td><td>5000</td><td>100</td><td>96</td></tr><tr><td>2</td><td>53</td><td>100</td><td>1000</td><td>1000</td><td>50</td><td>20</td><td>10</td><td>2</td><td>80</td><td>2000</td><td>100</td><td>1000</td><td>5</td><td>1000</td><td>1000</td><td>500</td><td>5000</td><td>0</td><td>96</td></tr><tr><td>3</td><td>54</td><td>100</td><td>1000</td><td>1000</td><td>50</td><td>20</td><td>10</td><td>2</td><td>80</td><td>2000</td><td>100</td><td>1000</td><td>5</td><td>1000</td><td>1000</td><td>500</td><td>5000</td><td>5000</td><td>96</td></tr><tr><td>4</td><td>55</td><td>10</td><td>1000</td><td>100</td><td>500</td><td>200</td><td>10</td><td>2</td><td>80</td><td>200</td><td>100</td><td>2000</td><td>8</td><td>2000</td><td>100</td><td>2000</td><td>5000</td><td>15000</td><td>90</td></tr></tbody></table>		Supplier ID	SF1	SF2	SF3	SF4	SF5	SF6	SF7	SF8	SF9	SF10	SF11	SF12	SF13	SF14	SF15	SF16	SF17	SF18	0	51	100	1000	1000	50	20	10	2	80	2000	100	1000	5	1000	1000	500	5000	100	96	1	52	100	1000	1000	50	20	10	0	80	2000	100	1000	5	1000	1000	500	5000	100	96	2	53	100	1000	1000	50	20	10	2	80	2000	100	1000	5	1000	1000	500	5000	0	96	3	54	100	1000	1000	50	20	10	2	80	2000	100	1000	5	1000	1000	500	5000	5000	96	4	55	10	1000	100	500	200	10	2	80	200	100	2000	8	2000	100	2000	5000	15000	90
	Supplier ID	SF1	SF2	SF3	SF4	SF5	SF6	SF7	SF8	SF9	SF10	SF11	SF12	SF13	SF14	SF15	SF16	SF17	SF18																																																																																																						
0	51	100	1000	1000	50	20	10	2	80	2000	100	1000	5	1000	1000	500	5000	100	96																																																																																																						
1	52	100	1000	1000	50	20	10	0	80	2000	100	1000	5	1000	1000	500	5000	100	96																																																																																																						
2	53	100	1000	1000	50	20	10	2	80	2000	100	1000	5	1000	1000	500	5000	0	96																																																																																																						
3	54	100	1000	1000	50	20	10	2	80	2000	100	1000	5	1000	1000	500	5000	5000	96																																																																																																						
4	55	10	1000	100	500	200	10	2	80	200	100	2000	8	2000	100	2000	5000	15000	90																																																																																																						
In [4]:	<pre>suppliers.info()</pre>																																																																																																																								
Out[4]:	<pre><class 'pandas.core.frame.DataFrame'> RangeIndex: 64 entries, 0 to 63 Data columns (total 19 columns): # Column Non-Null Count Dtype --- --- 0 Supplier ID 64 non-null object 1 SF1 64 non-null int64 2 SF2 64 non-null int64 3 SF3 64 non-null int64 4 SF4 64 non-null int64 5 SF5 64 non-null int64 6 SF6 64 non-null int64 7 SF7 64 non-null int64 8 SF8 64 non-null int64 9 SF9 64 non-null int64 10 SF10 64 non-null int64 11 SF11 64 non-null int64 12 SF12 64 non-null int64 13 SF13 64 non-null int64 14 SF14 64 non-null int64 15 SF15 64 non-null int64 16 SF16 64 non-null int64 17 SF17 64 non-null int64 18 SF18 64 non-null int64 dtypes: int64(18), object(1) memory usage: 9.6+ KB</pre>																																																																																																																								
In [5]:	<pre>tasks.head()</pre>																																																																																																																								
Out[5]:	<table><thead><tr><th></th><th>Task ID</th><th>TF1</th><th>TF2</th><th>TF3</th><th>TF4</th><th>TF5</th><th>TF6</th><th>TF7</th><th>TF8</th><th>TF9</th><th>...</th><th>TF107</th><th>TF108</th><th>TF109</th><th>TF110</th><th>TF111</th><th>TF112</th><th>TF113</th></tr></thead><tbody><tr><td>2019</td><td>0</td><td>05</td><td>706</td><td>2539</td><td>428536374</td><td>367</td><td>0.144545</td><td>35342375</td><td>0.08</td><td>829</td><td>0.326506</td><td>...</td><td>125</td><td>0</td><td>219407.09</td><td>903728.98</td><td>10174793</td><td>0</td></tr><tr><td>2019</td><td>1</td><td>09</td><td>697</td><td>2199</td><td>389831692</td><td>431</td><td>0.195998</td><td>20091114</td><td>0.05</td><td>460</td><td>0.209186</td><td>...</td><td>149</td><td>0</td><td>467568.10</td><td>1868010.80</td><td>19221510</td><td>0</td></tr><tr><td>2019</td><td>2</td><td>11</td><td>262</td><td>4156</td><td>500027098</td><td>1510</td><td>0.363330</td><td>89708355</td><td>0.18</td><td>1010</td><td>0.243022</td><td>...</td><td>394</td><td>0</td><td>228358.20</td><td>1358750.92</td><td>38657530</td><td>0</td></tr><tr><td>2020</td><td>3</td><td>01</td><td>469</td><td>4346</td><td>547810586</td><td>1376</td><td>0.316613</td><td>90478530</td><td>0.17</td><td>1097</td><td>0.252416</td><td>...</td><td>394</td><td>0</td><td>258487.96</td><td>1781016.93</td><td>39064840</td><td>0</td></tr><tr><td>2020</td><td>4</td><td>01</td><td>555</td><td>3934</td><td>521676289</td><td>1039</td><td>0.264108</td><td>69762831</td><td>0.13</td><td>943</td><td>0.239705</td><td>...</td><td>394</td><td>0</td><td>285558.89</td><td>1869912.30</td><td>39064840</td><td>0</td></tr></tbody></table> <p>5 rows x 17 columns</p>		Task ID	TF1	TF2	TF3	TF4	TF5	TF6	TF7	TF8	TF9	...	TF107	TF108	TF109	TF110	TF111	TF112	TF113	2019	0	05	706	2539	428536374	367	0.144545	35342375	0.08	829	0.326506	...	125	0	219407.09	903728.98	10174793	0	2019	1	09	697	2199	389831692	431	0.195998	20091114	0.05	460	0.209186	...	149	0	467568.10	1868010.80	19221510	0	2019	2	11	262	4156	500027098	1510	0.363330	89708355	0.18	1010	0.243022	...	394	0	228358.20	1358750.92	38657530	0	2020	3	01	469	4346	547810586	1376	0.316613	90478530	0.17	1097	0.252416	...	394	0	258487.96	1781016.93	39064840	0	2020	4	01	555	3934	521676289	1039	0.264108	69762831	0.13	943	0.239705	...	394	0	285558.89	1869912.30	39064840	0						
	Task ID	TF1	TF2	TF3	TF4	TF5	TF6	TF7	TF8	TF9	...	TF107	TF108	TF109	TF110	TF111	TF112	TF113																																																																																																							
2019	0	05	706	2539	428536374	367	0.144545	35342375	0.08	829	0.326506	...	125	0	219407.09	903728.98	10174793	0																																																																																																							
2019	1	09	697	2199	389831692	431	0.195998	20091114	0.05	460	0.209186	...	149	0	467568.10	1868010.80	19221510	0																																																																																																							
2019	2	11	262	4156	500027098	1510	0.363330	89708355	0.18	1010	0.243022	...	394	0	228358.20	1358750.92	38657530	0																																																																																																							
2020	3	01	469	4346	547810586	1376	0.316613	90478530	0.17	1097	0.252416	...	394	0	258487.96	1781016.93	39064840	0																																																																																																							
2020	4	01	555	3934	521676289	1039	0.264108	69762831	0.13	943	0.239705	...	394	0	285558.89	1869912.30	39064840	0																																																																																																							
In [6]:	<pre>tasks.info(verbose=True)</pre>																																																																																																																								
Out[6]:	<pre><class 'pandas.core.frame.DataFrame'> RangeIndex: 130 entries, 0 to 129 Data columns (total 17 columns): # Column Dtype --- --- 0 Task ID object 1 TF1 int64 2 TF2 int64 3 TF3 int64 4 TF4 int64 5 TF5 float64 6 TF6 int64 7 TF7 float64 8 TF8 int64 9 TF9 float64 10 TF10 int64 11 TF11 int64 12 TF12 int64 13 TF13 float64 14 TF14 int64 15 TF15 float64 16 TF16 float64 17 TF17 float64 18 TF18 float64 19 TF19 int64 20 TF20 float64 21 TF21 float64 22 TF22 float64 23 TF23 int64 24 TF24 float64 25 TF25 float64 26 TF26 int64 27 TF27 int64 28 TF28 float64 29 TF29 float64 30 TF30 int64 31 TF31 int64 32 TF32 float64 33 TF33 float64 34 TF34 int64 35 TF35 int64 36 TF36 int64 37 TF37 int64 38 TF38 float64 39 TF39 int64 40 TF40 int64 41 TF41 int64 42 TF42 float64 43 TF43 float64 44 TF44 int64 45 TF45 int64 46 TF46 float64 47 TF47 float64 48 TF48 int64 49 TF49 int64 50 TF50 float64 51 TF51 float64 52 TF52 int64 53 TF53 float64 54 TF54 int64 55 TF55 float64 56 TF56 int64 57 TF57 float64 58 TF58 int64 59 TF59 float64 60 TF60 int64 61 TF61 float64 62 TF62 int64 63 TF63 float64 64 TF64 int64 65 TF65 int64 66 TF66 float64 67 TF67 float64 68 TF68 float64 69 TF69 int64 70 TF70 int64 71 TF71 int64 72 TF72 float64 73 TF73 float64 74 TF74 int64 75 TF75 int64 76 TF76 float64 77 TF77 float64 78 TF78 int64 79 TF79 int64 80 TF80 int64 81 TF81 float64 82 TF82 float64 83 TF83 int64 84 TF84 int64 85 TF85 float64 86 TF86 float64 87 TF87 int64 88 TF88 int64 89 TF89 float64 90 TF90 float64 91 TF91 int64 92 TF92 int64 93 TF93 float64 94 TF94 float64 95 TF95 int64 96 TF96 int64 97 TF97 float64 98 TF98 float64 99 TF99 int64 100 TF100 int64 101 TF101 float64 102 TF102 float64 103 TF103 int64 104 TF104 int64 105 TF105 float64 106 TF106 float64 107 TF107 int64 108 TF108 int64 109 TF109 float64 110 TF110 float64 111 TF111 int64 112 TF112 int64 113 TF113 int64 114 TF114 float64 115 TF115 int64 116 TF116 float64 dtypes: float64(18), int64(6), object(1) memory usage: 119.0+ KB</pre>																																																																																																																								
In [7]:	<pre>cost.head()</pre>																																																																																																																								
Out[7]:	<table><thead><tr><th></th><th>Task ID</th><th>Supplier ID</th><th>Cost</th></tr></thead><tbody><tr><td>0</td><td>30/05/2019</td><td>S1</td><td>0.478219</td></tr><tr><td>1</td><td>30/05/2019</td><td>S2</td><td>0.444543</td></tr><tr><td>2</td><td>30/05/2019</td><td>S3</td><td>0.521679</td></tr><tr><td>3</td><td>30/05/2019</td><td>S4</td><td>0.307331</td></tr><tr><td>4</td><td>30/05/2019</td><td>S5</td><td>0.357689</td></tr></tbody></table>		Task ID	Supplier ID	Cost	0	30/05/2019	S1	0.478219	1	30/05/2019	S2	0.444543	2	30/05/2019	S3	0.521679	3	30/05/2019	S4	0.307331	4	30/05/2019	S5	0.357689																																																																																																
	Task ID	Supplier ID	Cost																																																																																																																						
0	30/05/2019	S1	0.478219																																																																																																																						
1	30/05/2019	S2	0.444543																																																																																																																						
2	30/05/2019	S3	0.521679																																																																																																																						
3	30/05/2019	S4	0.307331																																																																																																																						
4	30/05/2019	S5	0.357689																																																																																																																						
In [8]:	<pre>cost.info()</pre>																																																																																																																								
Out[8]:	<pre><class 'pandas.core.frame.DataFrame'> RangeIndex: 7680 entries, 0 to 7679 Data columns (total 3 columns): # Column Non-Null Count Dtype --- --- 0 Task ID 7680 non-null object 1 Supplier ID 7680 non-null object 2 Cost 7680 non-null float64 dtypes: float64(1), object(2) memory usage: 180.1+ KB</pre>																																																																																																																								
In [9]:	<pre># Convert 'Task ID' to datetime format tasks['Task ID'] = pd.to_datetime(tasks['Task ID'], infer_datetime_format=True) cost['Task ID'] = pd.to_datetime(cost['Task ID'], infer_datetime_format=True)</pre>																																																																																																																								

1.1 Data Verification

Missing Values

Calculate the number of null values in each dataset

In [10]:	<pre>print('Number of null values in cost:', cost.isnull().any().sum()) print('Number of null values in suppliers:', suppliers.isnull().any().sum()) print('Number of null values in tasks:', tasks.isnull().any().sum()) Number of null values in cost: 0 Number of null values in suppliers: 0 Number of null values in tasks: 0</pre>
Out[10]:	

Number of tasks, suppliers, features and cost values in all data sets

Count number of tasks, suppliers, features and cost values in all data sets

In [11]:	<pre>print('There are {} unique tasks in the Tasks dataset'.format(len(tasks['Task ID'].unique()))) print('There are {} unique tasks in the Costs dataset'.format(len(cost['Task ID'].unique()))) print('There are {} unique suppliers in the Suppliers dataset'.format(len(suppliers['Supplier ID'].unique()))) print('There are {} unique suppliers in the Costs dataset'.format(len(cost['Supplier ID'].unique()))) print('Each task is specified by {} features'.format(len(tasks.columns)-1)) print('Each supplier is specified by {} features'.format(len(suppliers.columns)-1)) print('There is information about {} unique suppliers in the Costs dataset'.format(len(suppliers['Supplier ID'].unique()))) There are 130 unique tasks in Tasks dataset There are 120 unique tasks in Costs dataset There are 64 unique suppliers in Suppliers dataset There are 64 unique suppliers in Costs dataset Each task is specified by 17 features Each supplier is specified by 18 features There is information about 64 unique suppliers associated with different tasks and suppliers</pre>
Out[11]:	

ID match: Task IDs

Making sure the Task IDs in Costs and Tasks dataset match.

Cost dataset is the most important for our analysis, since it contains the true value of the predicted variable.

That is why we should check whether other datasets contain features for all tasks presented in the cost dataset

In [12]:	<pre>cost_unique = pd.DataFrame(cost['Task ID'].unique()) tasks_unique = tasks['Task ID'].unique()</pre>
In [13]:	<pre>print('Number of unique Task ID in the Cost dataset:', len(cost_unique)) print('Number of unique Task ID in the Tasks dataset:', len(tasks_unique)) print('Number of unique Task ID in the Costs dataset that also appear in the Tasks dataset:', len(cost_unique.intersection(tasks_unique))) Number of unique Task ID in the Cost dataset: 120 Number of unique Task ID in the Tasks dataset: 130 Number of unique Task ID in the Costs dataset that also appear in the Tasks dataset: 120 Only 120 tasks from Cost dataset are specified in the Task dataset. Consequently, we should remove from both datasets tasks that are not matched</pre>
In [14]:	<pre>tasks_before = tasks.shape[0] print('Number of tasks before removal:', tasks_before) # Eliminating tasks with no costs specified from tasks dataset tasks = tasks[tasks['Task ID'].isin(cost['Task ID'].unique())] print('Number of tasks after removal:', tasks_after) print('Number of tasks removed from the Tasks dataset:', tasks_before - tasks_after) Number of tasks before removal: 130 Number of tasks after removal: 120 10 tasks were removed from the Tasks dataset</pre>
In [15]:	<pre>cost_tasks_before = cost.shape[0] print('Number of tasks before removal:', cost_tasks_before) # Eliminating tasks with no costs specified from cost dataset cost = cost[cost['Task ID'].isin(tasks['Task ID'].unique())] cost_tasks_after = cost.shape[0] print('Number of tasks after removal:', cost_tasks_after) print('Number of tasks removed from the Cost dataset:', cost_tasks_before - cost_tasks_after) Number of tasks before removal: 7680 Number of tasks after removal: 7680 0 tasks were removed from the Cost dataset</pre>
Out[15]:	

ID match: Supplier IDs

In [16]:	<pre>cost_unique = pd.DataFrame(cost['Supplier ID'].unique()) suppliers_unique = suppliers['Supplier ID'].unique() print('Number of unique Supplier ID in the Costs dataset:', len(cost_unique)) print('Number of unique Supplier ID in the Suppliers dataset:', len(suppliers_unique)) print('Number of unique Supplier ID in the Costs dataset that also appear in the Suppliers dataset:', len(cost_unique.intersection(suppliers_unique))) # All suppliers in the cost dataset also appear in the suppliers dataset, no adjustments required Number of unique Supplier ID in the Cost dataset: 64 Number of unique Supplier ID in the Suppliers dataset: 64 Number of unique Supplier ID in the Costs dataset that also appear in the Suppliers dataset: 64 Nothing was removed from the Costs and the Suppliers dataset</pre>
Out[16]:	

Saving modified DataFrames

In [17]:	<pre>cost_df = pd.DataFrame(cost) suppliers_df = pd.DataFrame(suppliers).set_index('Supplier ID') tasks_df = pd.DataFrame(tasks).set_index('Task ID')</pre>
Out[17]:	

1.2 Calculate the max, min, mean and var of each feature

1.2.1 Tasks Dataset

```
tasks_df = tasks_df.drop(cols_to_del, axis=1)

tf_after = tasks_df.shape[1]
print('Number of features after feature removal:', tf_after)
print('tf_before = tf_after, 'features were removed from the Tasks dataset')

tasks_df.head()
```

Number of features before feature removal: 116
Number of features after feature removal: 83
33 features were removed from the Tasks dataset

```
Out[20]:
```

	TF1	TF2	TF3	TF4	TF6	TF8	TF10	TF12	TF14	TF16	...	TF102	TF103	TF105	TF106	TF107	
Task ID																	
2019-05-30	706	2539	428536374	367	35342375	829	109388318	452	65396430	133010	...	983613171	52077023	6.15	25.10	125	2
2019-09-26	697	2199	389831692	431	20091114	460	72475671	287	47355481	140306	...	734279232	44103725	9.13	30.13	149	4
2019-11-29	262	4156	500027098	1510	89708355	1010	112404944	812	83864672	134716	...	100645909	50354895	5.06	23.75	394	2
2020-01-03	697	2199	389831692	431	20091114	460	72475671	287	47355481	140306	...	734279232	44103725	9.13	30.13	149	4
2020-01-07	436	4634	547810586	1376	90478330	1097	122674168	834	95043161	134570	...	912221784	49212552	5.25	26.07	394	2
2020-01-07	555	3934	521676289	1039	69762831	943	117268389	682	85212722	135606	...	866459129	48980163	5.80	27.35	394	2

5 rows x 17 columns

1.2.2 Suppliers Dataset

```
suppliers_df.head()
```

```
Out[21]:
```

	SF1	SF2	SF3	SF4	SF5	SF6	SF7	SF8	SF9	SF10	SF11	SF12	SF13	SF14	SF15	SF16	SF17	SF18
--	-----	-----	-----	-----	-----	-----	-----	-----	-----	------	------	------	------	------	------	------	------	------

Remove features with variance lower than 0.01

S1	100	1000	1000	500	20	10	2	80	2000	100	1000	5	1000	1000	500	5000	100	96
S2	100	1000	1000	500	20	10	0	80	2000	100	1000	5	1000	1000	500	5000	100	96
S3	100	1000	1000	500	20	10	2	80	2000	100	1000	5	1000	1000	500	5000	0	96
S4	100	1000	1000	500	20	10	2	80	2000	100	1000	5	1000	1000	500	5000	5000	96
S5	10	1000	100	500	200	10	2	80	2000	100	2000	8	2000	100	2000	5000	15000	90

```
In [22]: # Calculating the variance of each task feature
var_suppliers = pd.DataFrame(suppliers_df.var()).T.rename(index={0:'var'})

# Getting summary statistics of each task feature
suppliers_desc = suppliers_df.describe()

# Concatenating both result for an easier read
suppliers_summary = pd.concat([suppliers_desc, var_suppliers], axis=0)
suppliers_summary

Out[22]:
```

SF1	SF2	SF3	SF4	SF5	SF6	SF7	SF8	SF9	SI
-----	-----	-----	-----	-----	-----	-----	-----	-----	----

1.2.2 Suppliers Dataset

	std	441.163235	757.371185	757.371185	220.581678	88.232647	44.116324	0.835117	213.615141	1.514742e+03	209.359
	min	10.000000	100.000000	100.000000	5.000000	2.000000	1.000000	0.000000	8.000000	2.000000e+02	10.000
	25%	100.000000	100.000000	100.000000	5.000000	2.000000	1.000000	1.000000	8.000000	2.000000e+02	10.000
	50%	100.000000	1000.000000	1000.000000	50.000000	20.000000	10.000000	2.000000	80.000000	2.000000e+03	100.000
	75%	100.000000	2000.000000	2000.000000	500.000000	200.000000	100.000000	3.000000	500.000000	4.000000e+03	500.000
	max	1000.000000	10000.000000	10000.000000	5000.000000	2000.000000	1000.000000	3.000000	5000.000000	4.000000e+03	500.000
	var	194625.000000	573611.111111	573611.111111	48656.250000	7785.000000	1946.250000	0.697421	45631.482571	2.294444e+06	43831.3495

```
In [23]: print('Number of features that have variance less than 0.01:', (suppliers_summary.loc['var',:] < 0.01).sum())
```

Number of features that have variance less than 0.01: 0

There are no features that have variance less than 0.01 in the Suppliers dataset, so no features were removed

1.3 Feature Scaling

```
In [24]: # Feature scaling using MinMaxScaler
scaler = MinMaxScaler(feature_range=(-1,1))
tasks_df[1:] = scaler.fit_transform(tasks_df)
suppliers_df[1:] = scaler.fit_transform(suppliers_df)
```

Tasks and Suppliers dataframe after scaling

```
In [25]: tasks_df.head()
```

```
Out[25]:
```

	TF1	TF2	TF3	TF4	TF6	TF8	TF10	TF12	TF14	...	TF102	TF103	
Task ID													
2019-05-30	1.000000	-0.583542	0.022445	-1.000000	-0.821200	-0.454300	0.059902	-0.697013	-0.090516	-0.559193	-0.712664	-0.401560	-0.9
2019-06-01	0.974504	-0.668382	-0.088391	-0.957004	-1.000000	-0.786583	-0.333066	-0.931721	-0.394759	0.017794	-0.903937	-0.553691	-0.8


```
Task ID 0
Supplier 0
Cost 0
dtype: int64

[35]: # Removing the flagged supplier from Suppliers dataframe
final_supplier = suppliers_df[~(suppliers_df.index == sup_to_del)]
final_supplier

Out[35]:
Supplier ID
0 -0.818182 -0.052632 -0.052632 -0.818182 -0.818182 -0.818182 0.333333 -0.707317 -0.052632 -0.632653 -0.052632 0.142857 -0.052632
1 -0.818182 -0.052632 -0.052632 -0.818182 -0.818182 -0.818182 0.333333 -0.707317 -0.052632 -0.632653 -0.052632 0.142857 -0.052632
3 -0.818182 -0.052632 -0.052632 -0.818182 -0.818182 -0.818182 0.333333 -0.707317 -0.052632 -0.632653 -0.052632 0.142857 -0.052632
4 -0.818182 -0.052632 -0.052632 -0.818182 -0.818182 -0.818182 0.333333 -0.707317 -0.052632 -0.632653 -0.052632 0.142857 -0.052632
5 -1.000000 -0.052632 -1.000000 1.000000 1.000000 -0.818182 0.333333 -0.707317 -1.000000 -0.632653 1.000000 1.000000 1.000000
6 -1.000000 -1.000000 1.000000 -1.000000 1.000000 1.000000 -0.333333 -1.000000 1.000000 -1.000000 1.000000 -1.000000 -1.000000
...
560 -1.000000 -1.000000 -1.000000 -1.000000 -1.000000 -0.818182 0.333333 -1.000000 -1.000000 -0.632653 -1.000000 1.000000 -0.052632
561 -0.818182 1.000000 -0.052632 -1.000000 -1.000000 -1.000000 -0.818182 0.333333 1.000000 1.000000 1.000000 -0.052632 -1.000000
562 1.000000 -0.052632 -0.052632 -1.000000 -1.000000 -1.000000 -0.818182 0.333333 1.000000 1.000000 -1.000000 1.000000 1.000000
563 -0.818182 -0.052632 -0.052632 -0.818182 1.000000 -1.000000 0.333333 1.000000 1.000000 -1.000000 1.000000 -1.000000 -1.000000
564 -0.818182 1.000000 -0.052632 1.000000 -1.000000 -1.000000 1.000000 -1.000000 1.000000 1.000000 -0.052632 0.142857 1.000000

63 rows x 18 columns

In [36]: # Checking that the flagged supplier ID was deleted
final_supplier[final_supplier.index==sup_to_del].count()

Out[36]:
SFI 0
SF2 0
SF3 0
SF4 0
SFE 0
SFF 0
SFG 0
SFI0 0
SFI1 0
SFI2 0
SFI3 0
SFI4 0
SFI5 0
SFI6 0
SFI7 0
SFI8 0
SFI9 0
SFI10 0
SFI11 0
SFI12 0
dtype: int64

In [37]: final_cost.to_csv(r'cost_final.csv')
final_supplier.to_csv(r'supplier_final.csv')
```

Exploratory Data Analysis

2.1 Box Plot for Task Dataset

```
In [38]: # Set the figure size
plt.rcParams['figure.figsize'] = [30,20]
plt.rcParams['figure.autolayout'] = True

# Set plot style
sns.set_style('whitegrid')
sns.set_context('poster')

# Load the cleaned Tasks dataframe
tasks_cleaned = pd.read_csv(r'tasks_cleaned.csv').set_index('Task ID')

# Make a boxplot of the Tasks dataframe
sns.boxplot(data=tasks_cleaned, palette='deep')
plt.title('Tasks Data')
plt.xlabel('Task features')
plt.ylabel('Normalized feature value')

# Save plot
plt.savefig(r'boxplot1.png')

# Display the plot
plt.show()

Tasks Data
Normalized feature value
1.00
0.75
0.50
0.25
-0.25
-0.50
-0.75
-1.00
Task feature
```

2.2 Distribution of Errors

```
In [39]: # Define RMSE function
def rmse(error):
    return np.sqrt(np.mean(error**2)/len(error))

In [40]: # Load the cleaned Cost dataframe
cost = pd.read_csv(r'cost_final.csv').drop('Unnamed: 0', axis=1)

# Convert the Cost to a wide format
cost_error = cost.pivot(index='Task ID', values='Cost', columns='Supplier ID')

# Get the minimum cost for each task
cost_error['min_costs'] = cost_error.min(axis=1)

# Calculating error distribution for all suppliers
for i in cost_error.columns[1:]:
    cost_error[i] = cost_error['min_costs'] - cost_error[i]

In [41]: cost_error.head()

Out[41]:
Supplier ID S1 S10 S11 S12 S13 S14 S15 S16 S17 S18 ... S6 S60
Task ID
2019-05-30 -0.187252 -0.040986 -0.043396 -0.022542 -0.018435 -0.029696 -0.016725 -0.135365 -0.058547 -0.035238 ... -0.061014 -0.170021
2019-09-26 -0.124033 -0.043402 -0.048565 -0.041030 -0.026667 -0.005421 -0.040261 -0.037321 -0.031454 -0.023414 ... -0.020424 -0.048707
11-29 -0.153309 -0.094888 -0.024619 -0.018330 -0.020813 -0.069651 -0.028824 -0.043895 -0.074065 -0.074915 ... -0.020996 -0.018257
2020-01-03 -0.163084 -0.075884 -0.024619 -0.035677 -0.028765 -0.046851 -0.032314 -0.074921 -0.056704 -0.067153 ... -0.021011 -0.023527
2020-01-07 -0.149066 -0.031547 -0.023474 -0.031520 -0.018434 -0.031524 -0.009770 -0.015512 -0.054461 -0.017052 ... -0.023803 -0.028648

5 rows x 64 columns

In [42]: rmse_df = pd.DataFrame(
    [rmse(cost_error.iloc[:,i-1])
     for i in cost_error.columns[1:]]
    , columns=[f'rmse_{cost_error.columns[i+1]}'], round(2)
    , rmse_df

Out[42]:
Supplier ID S1 S10 S11 S12 S13 S14 S15 S16 S17 S18 ... S6 S60
Task ID
0 0.08 0.04 0.03 0.05 0.04 0.04 0.03 0.06 0.06 0.05 ... 0.06 0.04
1 0.08 0.04 0.03 0.05 0.04 0.04 0.03 0.06 0.06 0.05 ... 0.06 0.04
3 0.08 0.04 0.03 0.05 0.04 0.04 0.03 0.06 0.06 0.05 ... 0.06 0.04
4 0.08 0.04 0.03 0.05 0.04 0.04 0.03 0.06 0.06 0.05 ... 0.06 0.04
5 0.08 0.04 0.03 0.05 0.04 0.04 0.03 0.06 0.06 0.05 ... 0.06 0.04
...
560 0.08 0.04 0.03 0.05 0.04 0.04 0.03 0.06 0.06 0.05 ... 0.06 0.04
561 0.08 0.04 0.03 0.05 0.04 0.04 0.03 0.06 0.06 0.05 ... 0.06 0.04
562 0.08 0.04 0.03 0.05 0.04 0.04 0.03 0.06 0.06 0.05 ... 0.06 0.04
563 0.08 0.04 0.03 0.05 0.04 0.04 0.03 0.06 0.06 0.05 ... 0.06 0.04
564 0.08 0.04 0.03 0.05 0.04 0.04 0.03 0.06 0.06 0.05 ... 0.06 0.04
565 0.08 0.04 0.03 0.05 0.04 0.04 0.03 0.06 0.06 0.05 ... 0.06 0.04
566 0.08 0.04 0.03 0.05 0.04 0.04 0.03 0.06 0.06 0.05 ... 0.06 0.04
567 0.08 0.04 0.03 0.05 0.04 0.04 0.03 0.06 0.06 0.05 ... 0.06 0.04
568 0.08 0.04 0.03 0.05 0.04 0.04 0.03 0.06 0.06 0.05 ... 0.06 0.04
569 0.08 0.04 0.03 0.05 0.04 0.04 0.03 0.06 0.06 0.05 ... 0.06 0.04
570 0.08 0.04 0.03 0.05 0.04 0.04 0.03 0.06 0.06 0.05 ... 0.06 0.04
571 0.08 0.04 0.03 0.05 0.04 0.04 0.03 0.06 0.06 0.05 ... 0.06 0.04
572 0.08 0.04 0.03 0.05 0.04 0.04 0.03 0.06 0.06 0.05 ... 0.06 0.04
573 0.08 0.04 0.03 0.05 0.04 0.04 0.03 0.06 0.06 0.05 ... 0.06 0.04
574 0.08 0.04 0.03 0.05 0.04 0.04 0.03 0.06 0.06 0.05 ... 0.06 0.04
575 0.08 0.04 0.03 0.05 0.04 0.04 0.03 0.06 0.06 0.05 ... 0.06 0.04
576 0.08 0.04 0.03 0.05 0.04 0.04 0.03 0.06 0.06 0.05 ... 0.06 0.04
577 0.08 0.04 0.03 0.05 0.04 0.04 0.03 0.06 0.06 0.05 ... 0.06 0.04
578 0.08 0.04 0.03 0.05 0.04 0.04 0.03 0.06 0.06 0.05 ... 0.06 0.04
579 0.08 0.04 0.03 0.05 0.04 0.04 0.03 0.06 0.06 0.05 ... 0.06 0.04
580 0.08 0.04 0.03 0.05 0.04 0.04 0.03 0.06 0.06 0.05 ... 0.06 0.04
581 0.08 0.04 0.03 0.05 0.04 0.04 0.03 0.06 0.06 0.05 ... 0.06 0.04
582 0.08 0.04 0.03 0.05 0.04 0.04 0.03 0.06 0.06 0.05 ... 0.06 0.04
583 0.08 0.04 0.03 0.05 0.04 0.04 0.03 0.06 0.06 0.05 ... 0.06 0.04
584 0.08 0.04 0.03 0.05 0.04 0.04 0.03 0.06 0.06 0.05 ... 0.06 0.04
585 0.08 0.04 0.03 0.05 0.04 0.04 0.03 0.06 0.06 0.05 ... 0.06 0.04
586 0.08 0.04 0.03 0.05 0.04 0.04 0.03 0.06 0.06 0.05 ... 0.06 0.04
587 0.08 0.04 0.03 0.05 0.04 0.04 0.03 0.06 0.06 0.05 ... 0.06 0.04
588 0.08 0.04 0.03 0.05 0.04 0.04 0.03 0.06 0.06 0.05 ... 0.06 0.04
589 0.08 0.04 0.03 0.05 0.04 0.04 0.03 0.06 0.06 0.05 ... 0.06 0.04
590 0.08 0.04 0.03 0.05 0.04 0.04 0.03 0.06 0.06 0.05 ... 0.06 0.04
591 0.08 0.04 0.03 0.05 0.04 0.04 0.03 0.06 0.06 0.05 ... 0.06 0.04
592 0.08 0.04 0.03 0.05 0.04 0.04 0.03 0.06 0.06 0.05 ... 0.06 0.04
593 0.08 0.04 0.03 0.05 0.04 0.04 0.03 0.06 0.06 0.05 ... 0.06 0.04
594 0.08 0.04 0.03 0.05 0.04 0.04 0.03 0.06 0.06 0.05 ... 0.06 0.04
595 0.08 0.04 0.03 0.05 0.04 0.04 0.03 0.06 0.06 0.05 ... 0.06 0.04
596 0.08 0.04 0.03 0.05 0.04 0.04 0.03 0.06 0.06 0.05 ... 0.06 0.04
597 0.08 0.04 0.03 0.05 0.04 0.04 0.03 0.06 0.06 0.05 ... 0.06 0.04
598 0.08 0.04 0.03 0.05 0.04 0.04 0.03 0.06 0.06 0.05 ... 0.06 0.04
599 0.08 0.04 0.03 0.05 0.04 0.04 0.03 0.06 0.06 0.05 ... 0.06 0.04
600 0.08 0.04 0.03 0.05 0.04 0.04 0.03 0.06 0.06 0.05 ... 0.06 0.04

1 rows x 63 columns

In [43]: cost_error

Out[43]:
Supplier ID S1 S10 S11 S12 S13 S14 S15 S16 S17 S18 ... S6 S60
Task ID
2019-05-30 -0.187252 -0.040986 -0.043396 -0.022542 -0.018435 -0.029696 -0.016725 -0.135365 -0.058547 -0.035238 ... -0.061014 -0.170021
2019-09-26 -0.124033 -0.043402 -0.048565 -0.041030 -0.026667 -0.005421 -0.040261 -0.037321 -0.031454 -0.023414 ... -0.020424 -0.048707
11-29 -0.153309 -0.094888 -0.024619 -0.018330 -0.020813 -0.069651 -0.028824 -0.043895 -0.074065 -0.074915 ... -0.020996 -0.018257
2020-01-03 -0.163084 -0.075884 -0.024619 -0.035677 -0.028765 -0.046851 -0.032314 -0.074921 -0.056704 -0.067153 ... -0.021011 -0.023527
2020-01-07 -0.149066 -0.031547 -0.023474 -0.031520 -0.018434 -0.031524 -0.009770 -0.015512 -0.054461 -0.017052 ... -0.023803 -0.028648

7560 rows x 64 columns

In [44]: # Set the figure size
plt.rcParams['figure.figsize'] = [50,30]
plt.rcParams['figure.autolayout'] = True

# Set plot style
sns.set_style('whitegrid')
sns.set_context('poster')

# Plot the dataframe
ax = sns.boxplot(data=cost_error.iloc[:,1:], palette='deep')
plt.title('Distribution of Errors per Supplier')
plt.xlabel('Suppliers')
plt.ylabel('Error')

# Annotate the dataframe with RMSE values
for i, x in enumerate(cost_error.columns):
    ax.text(ax.get_xlim()[0]+i, ax.get_ylim()[0], f'{rmse_df[x]}')

# Save plot
plt.savefig(r'boxplot22.png')

# Display the plot
plt.show()

Distribution of Errors per Supplier
Error
-0.40
-0.30
-0.20
-0.10
0.00
0.10
0.20
0.30
0.40
0.50
0.60
0.70
0.80
0.90
1.00
Suppliers
```

2.3 Heatmap

```
In [45]: final_cost

Out[45]:
Task ID Supplier ID Cost
0 2019-05-30 S1 0.478219
3 2019-05-30 S4 0.307331
4 2019-05-30 S5 0.357689
5 2019-05-30 S6 0.351982
...
7675 2021-12-22 S60 0.410605
7676 2021-12-22 S61 0.410376
7677 2021-12-22 S62 0.407894
7678 2021-12-22 S63 0.420536
7679 2021-12-22 S64 0.423008

7560 rows x 3 columns

In [46]: # Convert the Cost dataframe into a wide format
final_cost_matrix = final_cost.pivot(index='Task ID', values='Cost', columns='Supplier ID')
final_cost_matrix

Out[46]:
Supplier ID S1 S10 S11 S12 S13 S14 S15 S16 S17 S18 ... S59 S60
Task ID
0 0.478219 0.31953 0.334363 0.313509 0.309403 0.320663 0.307692 0.426332 0.349514 0.326205 ... 0.371470 0.315192 0.460988 0
2019-05-30 0.478219 0.31953 0.334363 0.313509 0.309403 0.320663 0.307692 0.426332 0.349514 0.326205 ... 0.371470 0.315192 0.460988 0
2019-09-26 0.407784 0.327153 0.332316 0.324781 0.310419 0.289172 0.324012 0.321072 0.315205 0.307165 ... 0.298476 0.304175 0.332458 0
11-29 0.417884 0.379463 0.309194 0.302905 0.305388 0.354226 0.313399 0.328469 0.358640 0.359490 ... 0.369166 0.305571 0.302832 0
2020-01-03 0.473749 0.347978 0.335874 0.346072 0.339160 0.372445 0.342708 0.385315 0.367099 0.377548 ... 0.401964 0.331406 0.333921 0
2020-01-07 0.466028 0.333509 0.340306 0.330482 0.335396 0.348486 0.326732 0.332474 0.371423 0.334014 ... 0.370986 0.340765 0.345610 0
...
7675 0.410605 0.410376 0.407894 0.420536 0.423008 0.423008 0.423008 0.423008 0.423008 0.423008 ... 0.423008 0.423008 0.423008 0
7676 0.410376 0.410376 0.410376 0.410376 0.410376 0.410376 0.410376 0.410376 0.410376 0.410376 ... 0.410376 0.410376 0.410376 0
7677 0.407894 0.407894 0.407894 0.407894 0.407894 0.407894 0.407894 0.407894 0.407894 0.407894 ... 0.407894 0.407894 0.407894 0
7678 0.420536 0.420536 0.420536 0.420536 0.420536 0.420536 0.420536 0.420536 0.420536 0.420536 ... 0.420536 0.420536 0.420536 0
7679 0.423008 0.423008 0.423008 0.423008 0.423008 0.423008 0.423008 0.423008 0.423008 0.423008 ... 0.423008 0.423008 0.423008 0
...
7559 0.423008 0.423008 0.423008 0.423008 0.423008 0.423008 0.423008 0.423008 0.423008 0.423008 ... 0.423008 0.423008 0.423008 0
7560 0.423008 0.423008 0.423008 0.423008 0.423008 0.423008 0.423008 0.423008 0.423008 0.423008 ... 0.423008 0.423008 0.423008 0

120 rows x 63 columns

In [47]: # Set plot style
sns.set_context('notebook')

# Generate a custom diverging colormap
cm = sns.diverging_palette(220, 10, as_cmap=True)

# Draw the heatmap with the mask and correct aspect ratio
sns.heatmap(
    final_cost_matrix, # The data to plot
    square=True, # What colors to plot the heatmap as
    annot=True, # Force cells to be square
    linewidth=5, # Width of lines that divide cells

# Save plot
plt.savefig(r'heatmapcorr_new.png')

plt.show()

Heatmap of Cost matrix
Cost
-0.40
-0.30
-0.20
-0.10
0.00
0.10
0.20
0.30
0.40
0.50
0.60
0.70
0.80
0.90
1.00
Suppliers
```

3. ML Model Fitting and Scoring (Random Forest Regressor)

3.1 Combine Data

```
In [48]: # Load the cleaned datasets
cost_df = pd.read_csv(r'cost_final.csv').drop('Unnamed: 0', axis=1)
tasks_df = pd.read_csv(r'tasks_cleaned.csv')
suppliers_df = pd.read_csv(r'supplier_final.csv')

Out[48]:
Task ID Supplier ID Cost
0 2019-05-30 S1 0.478219
3 2019-05-30 S4 0.307331
4 2019-05-30 S5 0.357689
5 2019-05-30 S6 0.351982
...
7555 2021-12-22 S60 0.410605
7556 2021-12-22 S61 0.410376
7557 2021-12-22 S62 0.407894
7558 2021-12-22 S63 0.420536
7559 2021-12-22 S64 0.423008

7560 rows x 3 columns

In [50]: tasks_df

Out[50]:
Task ID TF1 TF6 TF8 TF17 TF19 TF20 TF22 TF26 TF27 ... TF43 TF45 TF52
0 2019-05-30 1.000000 -0.821200 -0.454300 -0.350672 1.000000 0.961792 0.790235 0.547957 -0.599289 ... -0.142857 0.608042 -0.444986
1 2019-09-26 0.974504 -1.000000 -0.786583 -0.560807 0.662732 0.961792 0.386772 0.444828 -0.599289 ... -0.428571 0.214577 -0.380919
...
118 2021-12-21 -0.997167 -0.600934 -0.108510 0.458689 1.000000 -0.951252 0.172824 -0.751655 -0.411680 ... -0.918367 0.477116 -0.481894
119 2021-12-22 -0.997167 -0.418980 0.002251 0.445513 1.000000 -0.951252 0.094291 -0.590019 -0.411680 ... -0.877551 0.551804 -0.447772
120 2021-12-22 -0.997167 -0.387289 0.010356 0.439078 1.000000 -0.951252 0.090674 -0.535657 -0.414927 ... -0.877551 0.636270 -0.444986
...
7555 2021-12-22 0.410605 -0.997167 -0.253683 0.062584 0.379247 1.0 -0.951252 -0.103591 -0.541969 ... -0.052632 -0.632653 -0.052632
7556 2021-12-22 0.410376 -0.997167 -0.253683 0.062584 0.379247 1.0 -0.951252 -0.103591 -0.541969 ... -0.052632 -0.632653 -0.052632
7557 2021-12-22 0.407894 -0.997167 -0.253683 0.062584 0.379247 1.0 -0.951252 -0.103591 -0.541969 ... -0.052632 -0.632653 -0.052632
7558 2021-12-22 0.420536 -0.997167 -0.253683 0.062584 0.379247 1.0 -0.951252 -0.103591 -0.541969 ... -0.052632 -0.632653 -0.052632
7559 2021-12-22 0.423008 -0.997167 -0.253683 0.062584 0.379247 1.0 -0.951252 -0.103591 -0.541969 ... -0.052632 -0.632653 -0.052632

120 rows x 24 columns

In [51]: suppliers_df

Out[51]:
Supplier ID SFI SFI2 SFI3 SFI4 SFI5 SFI6 SFI7 SFI8 SFI9 SFI10 SFI11 SFI12
0 S1 -0.818182 -0.052632 -0.052632 -0.818182 -0.818182 -0.818182 0.333333 -0.707317 -0.052632 -0.632653 -0.052632 0.142857
1 S2 -0.818182 -0.052632 -0.052632 -0.818182 -0.818182 -0.818182 0.333333 -0.707317 -0.052632 -0.632653 -0.052632 0.142857
2 S3 -0.818182 -0.052632 -0.052632 -0.818182 -0.818182 -0.818182 0.333333 -0.707317 -0.052632 -0.632653 -0.052632 0.142857
3 S4 -1.000000 -0.052632 -1.000000 1.000000 1.000000 -0.818182 0.333333 -0.707317 -1.000000 -0.632653 1.000000 1.000000
4 S6 -1.000000 -1.000000 1.000000 -1.000000 1.000000 1.000000 -0.333333 -1.000000 1.000000 -1.000000 1.000000 -1.000000
...
58 S60 -1.000000 -1.000000 -1.000000 -1.000000 -1.000000 1.000000 -0.818182 0.333333 -1.000000 -1.000000 -0.632653 -1.000000
59 S61 -0.818182 1.000000 -0.052632 -1.000000 -1.000000 -1.000000 -0.818182 0.333333 1.000000 1.000000 -0.052632 -1.000000
60 S62 1.000000 -0.052632 -0.052632 -1.000000 -1.000000 -1.000000 -0.818182 0.333333 1.000000 1.000000 -1.000000 -1.000000
61 S63 -0.818182 -0.052632 -0.052632 -0.818182 1.000000 -1.000000 0.333333 1.000000 1.000000 -1.000000 1.000000 -1.000000
62 S64 -0.818182 1.000000 -0.052632 1.000000 -1.000000 -1.000000 1.000000 -1.000000 1.000000 1.000000 -0.052632 0.142857

63 rows x 19 columns

In [52]: # Convert the Task ID to datetime format in both dataframes
cost_df['Task ID'] = pd.to_datetime(cost_df['Task ID'], infer_datetime_format=True)
tasks_df['Task ID'] = pd.to_datetime(tasks_df['Task ID'], infer_datetime_format=True)

# Merge all dataframes
full_data = pd.merge(cost_df, tasks_df, on='Task ID', how='left')
full_data = pd.merge(full_data, suppliers_df, on='Supplier ID', how='left')
full_data.drop('Supplier ID', axis=1, inplace=True)
full_data

Out[52]:
Task ID Cost TF6 TF8 TF17 TF19 TF20 TF22 TF26 ... SFI9 SFI10 SFI11
0 2019-05-30 0.478219 1.000000 -0.821200 -0.454300 -0.350672 1.0 0.961792 0
```


[illegible]

[illegible]