

Attrition Analysis and Insights for Faculty of Engineering

Author: Omkar Measurekar

Description: This notebook addresses the HR Business Partner's request to analyze staff attrition within the Faculty of Engineering. It includes data cleaning, transformation, and explorations.

```
In [40]: 1 import pandas as pd
          2 import numpy as np
          3 pd.set_option('display.max_columns', None)
```

```
In [41]: 1 df= pd.read_csv(r"./test_dataset.csv")
```

```
In [42]: 1 df
```

Out[42]:

	FirstName	LastName	StartDate	ExitDate	PositionTitle	Supervisor	ADEmail
0	Aaden	Mercer	26-Jul-23	NaN	Senior Research Fellow	Victoria Jacobs	aaden.mercer@bilearner.com
1	Aaliyah	Watts	9-May-20	NaN	Senior Lecturer	Jared Whitehead	aaliyah.watts@bilearner.com
2	Aarav	Espinoza	11-Nov-19	NaN	Officer	David Ali	aarav.espinoza@bilearner.com
3	Aaron	Tapia	20-Jan-23	NaN	Officer	Michael Reed	aaron.tapia@bilearner.com
4	Aaron	Weber	6-Apr-20	24-Jul-23	Lecturer	Melanie Garcia	aaron.weber@bilearner.com
...
2995	Zoey	Page	6-May-19	9-Apr-22	Professor	William Newman	zoey.page@bilearner.com
2996	Zoey	Spence	6-May-20	NaN	Senior Lecturer	Ian Price	zoey.spence@bilearner.com
2997	Zoie	Logan	19-Apr-19	NaN	Lecturer	Heidi Terry	zoie.logan@bilearner.com
2998	Zoie	Mercado	1-Dec-20	NaN	Lecturer	David Young	zoie.mercado@bilearner.com
2999	Zoie	Rowland	8-Feb-22	22-Jul-22	Senior Lecturer	Mrs. Rachel Henry	zoie.rowland@bilearner.com

3000 rows × 21 columns

Time-Based Features

```
In [43]: 1 # Convert string dates to datetime
2 df['StartDate'] = pd.to_datetime(df['StartDate'], format='%d-%b-%y', errors='co
3 df['ExitDate'] = pd.to_datetime(df['ExitDate'], format='%d-%b-%y', errors='coer
4 df['Survey Date'] = pd.to_datetime(df['Survey Date'], format='%d/%m/%Y', errors
5
6 # For rows where the above format fails (e.g., 25/9/1983), try alternative form
7 df['StartDate'] = df['StartDate'].fillna(pd.to_datetime(df['StartDate'], format
8 df['ExitDate'] = df['ExitDate'].fillna(pd.to_datetime(df['ExitDate'], format='%
9 df['Survey Date'] = df['Survey Date'].fillna(pd.to_datetime(df['Survey Date'],
```

```
In [44]: 1 print(df[['StartDate', 'ExitDate']].dtypes)
```

```
StartDate    datetime64[ns]
ExitDate     datetime64[ns]
dtype: object
```

```
In [45]: 1 print("Failed StartDate conversions:", df['StartDate'].isna().sum())
2 print("Failed ExitDate conversions:", df['ExitDate'].isna().sum())
```

```
Failed StartDate conversions: 0
Failed ExitDate conversions: 1467
```

```
In [46]: 1 # For current employees (no ExitDate), use today's date
2 current_date = pd.to_datetime('2025-05-17') # Assuming analysis date
3 df['EmploymentDuration_months'] = np.where(
4     df['ExitDate'].isna(),
5     (current_date - df['StartDate']).dt.days / 30, # Months for current employ
6     (df['ExitDate'] - df['StartDate']).dt.days / 30 # Months for former employ
7 )
```

```
In [47]: 1 df['TerminationYear'] = df['ExitDate'].dt.year
2 df['TerminationYear'] = df['TerminationYear'].fillna('Active')
3 df['TerminationMonth'] = df['ExitDate'].dt.month
4 df['TerminationYearMonth'] = df['ExitDate'].dt.to_period('M')
```

```
In [48]: 1 # Create tenure buckets
2 df['TenureGroup'] = pd.cut(df['EmploymentDuration_months'],
3                             bins=[0,12,36,60,120,np.inf],
4                             labels=['<1yr', '1-3yrs', '3-5yrs', '5-10yrs', '10+yrs'])
5 df['TenureGroup'] = df['TenureGroup'].cat.add_categories('0')
6 df['TenureGroup'] = df['TenureGroup'].fillna('0')
7
```

```
In [49]: 1 print(df[['StartDate', 'ExitDate', 'EmploymentDuration_months', 'TenureGroup']])
```

	StartDate	ExitDate	EmploymentDuration_months	TenureGroup
0	2023-07-26	NaT	22.033333	1-3yrs
1	2020-05-09	NaT	61.133333	5-10yrs
2	2019-11-11	NaT	67.133333	5-10yrs
3	2023-01-20	NaT	28.266667	1-3yrs
4	2020-04-06	2023-07-24	40.133333	3-5yrs

Termination Analysis Features

```
In [50]: 1 df['EmploymentStatus'] = np.where(
2         df['ExitDate'].isna(),
3         'Active',
4         'Terminated'
5     )
```

```
In [51]: 1 df['StatusValidated'] = np.where(
2         (df['EmploymentStatus'] == 'Active') & df['ExitDate'].isna(),
3         'Active',
4         np.where(
5             (df['EmploymentStatus'] != 'Active') & df['ExitDate'].notna(),
6             'Terminated',
7             'Status Conflict' # Flag inconsistencies
8         )
9     )
```

```
In [52]: 1 print(df['StatusValidated'].value_counts())
```

```
Terminated    1533
Active         1467
Name: StatusValidated, dtype: int64
```

```
In [53]: 1 df['TerminationCategory'] = np.where(
2         df['ExitDate'].isna(), # No exit date = active staff
3         'Active',
4         np.where(
5             df['TerminationType'].isin(['Voluntary', 'Resignation']),
6             'Voluntary Departure',
7             np.where(
8                 df['TerminationType'] == 'Retirement',
9                 'Retirement',
10                np.where(
11                    df['TerminationType'] == 'Involuntary',
12                    'Involuntary Departure',
13                    'Unclassified Departure' # For terminated with unknown type
14                )
15            )
16        )
17 )
18
19 df['VoluntaryTermination'] = np.where(
20     df['TerminationCategory'] == 'Voluntary Departure', 1, 0)
21 df['InvoluntaryTermination'] = np.where(
22     df['TerminationCategory'] == 'Involuntary Departure', 1, 0)
23 df['ActiveStaff'] = np.where(
24     df['TerminationCategory'] == 'Active', 1, 0)
25 df['RetirementCategory'] = np.where(
26     df['TerminationCategory'] == 'Retirement', 1, 0)
27 df['TerminatedStaff'] = df['ExitDate'].notna().astype(int)
28
29 # Employees who stayed >2 years
30 df['RetainedOver2Years'] = np.where(
31     (df['EmploymentDuration_months'] > 24) & (df['TerminatedStaff'] == 0), 1, 0
```

```
In [54]: 1 print(df['TerminationCategory'].value_counts())
```

```
Active          1467
Voluntary Departure  768
Involuntary Departure  388
Retirement      377
Name: TerminationCategory, dtype: int64
```

```
In [55]: 1 #df.to_csv("./cleaned_test_csv.csv")
```

```
In [56]: 1 monthly_trends = (  
2     df.groupby('TerminationYearMonth')  
3         .agg(Terminations=('Employee ID', 'count'))  
4         .reset_index()  
5 )
```

In [57]:

1	monthly_trends
---	----------------

Out[57]:

	TerminationYearMonth	Terminations
0	2018-11	2
1	2018-12	2
2	2019-01	4
3	2019-02	3
4	2019-03	5
5	2019-04	6
6	2019-05	2
7	2019-06	2
8	2019-07	6
9	2019-08	7
10	2019-09	2
11	2019-10	12
12	2019-11	4
13	2019-12	9
14	2020-01	14
15	2020-02	5
16	2020-03	11
17	2020-04	7
18	2020-05	12
19	2020-06	9
20	2020-07	9
21	2020-08	12
22	2020-09	12
23	2020-10	14
24	2020-11	8
25	2020-12	20
26	2021-01	13
27	2021-02	12
28	2021-03	19
29	2021-04	15
30	2021-05	22
31	2021-06	21
32	2021-07	23
33	2021-08	20
34	2021-09	31
35	2021-10	37
36	2021-11	31
37	2021-12	34
38	2022-01	25
39	2022-02	21
40	2022-03	31

	TerminationYearMonth	Terminations
41	2022-04	27
42	2022-05	31
43	2022-06	37
44	2022-07	44
45	2022-08	49
46	2022-09	49
47	2022-10	48
48	2022-11	52
49	2022-12	46
50	2023-01	61
51	2023-02	62
52	2023-03	67
53	2023-04	62
54	2023-05	101
55	2023-06	90
56	2023-07	127
57	2023-08	26

```
In [58]: 1 # Calculate termination rates by faculty
2 faculty_terminations = df.groupby('Faculty').agg(
3     TotalEmployees=('Employee ID', 'count'),
4     Terminations=('TerminatedStaff', 'sum'),
5     Voluntary=('VoluntaryTermination', 'sum'),
6     Involuntary=('InvoluntaryTermination', 'sum'),
7     Retirement=('RetirementCategory', 'sum'),
8     ActiveStaff=('ActiveStaff', 'sum'),
9     Retention=('RetainedOver2Years', 'mean')
10 ).reset_index()
11
12 faculty_terminations['AttritionRate'] = (
13     faculty_terminations['Terminations'] / faculty_terminations['TotalEmployees']
14 )
15 faculty_terminations['ActiveRate'] = (
16     faculty_terminations['ActiveStaff'] / faculty_terminations['TotalEmployees']
17 )
18 faculty_terminations['VoluntaryRate'] = (
19     faculty_terminations['Voluntary'] / faculty_terminations['Terminations']*10
20 )
21 faculty_terminations['InvoluntaryRate'] = (
22     faculty_terminations['Involuntary'] / faculty_terminations['Terminations']*
23 )
24 faculty_terminations['RetirementRate'] = (
25     faculty_terminations['Retirement'] / faculty_terminations['Terminations']*1
26 )
27 faculty_terminations['Retention_rate']=(faculty_terminations['Retention']*100)
```

In [59]: 1 faculty_terminations

Out[59]:

	Faculty	TotalEmployees	Terminations	Voluntary	Involuntary	Retirement	ActiveStaff	Retention
0	Faculty of Arts	591	306	142	85	79	285	0.458545
1	Faculty of Engineering	904	467	247	104	116	437	0.466814
2	Faculty of Medicine	1505	760	379	199	182	745	0.475083

In [60]: 1 #faculty_terminations.to_csv("./faculty_termibation_metrics.csv")

Performance & Engagement Metrics

In [61]: 1 survey_metrics = df.melt(id_vars=['Employee ID', 'Faculty'],
2 value_vars=['Engagement Score', 'Satisfaction Score', 'Wo
3 var_name='Metric',
4 value_name='Score')

In [62]: 1 survey_metrics

Out[62]:

	Employee ID	Faculty	Metric	Score
0	1808	Faculty of Medicine	Engagement Score	5
1	3783	Faculty of Medicine	Engagement Score	3
2	2956	Faculty of Arts	Engagement Score	5
3	1538	Faculty of Arts	Engagement Score	2
4	3541	Faculty of Engineering	Engagement Score	1
...
8995	1817	Faculty of Engineering	Work-Life Balance Score	5
8996	1492	Faculty of Medicine	Work-Life Balance Score	2
8997	3388	Faculty of Engineering	Work-Life Balance Score	3
8998	1489	Faculty of Medicine	Work-Life Balance Score	2
8999	3772	Faculty of Medicine	Work-Life Balance Score	1

9000 rows × 4 columns

In [63]: 1 df[['Engagement Score', 'Satisfaction Score']].corr()

Out[63]:

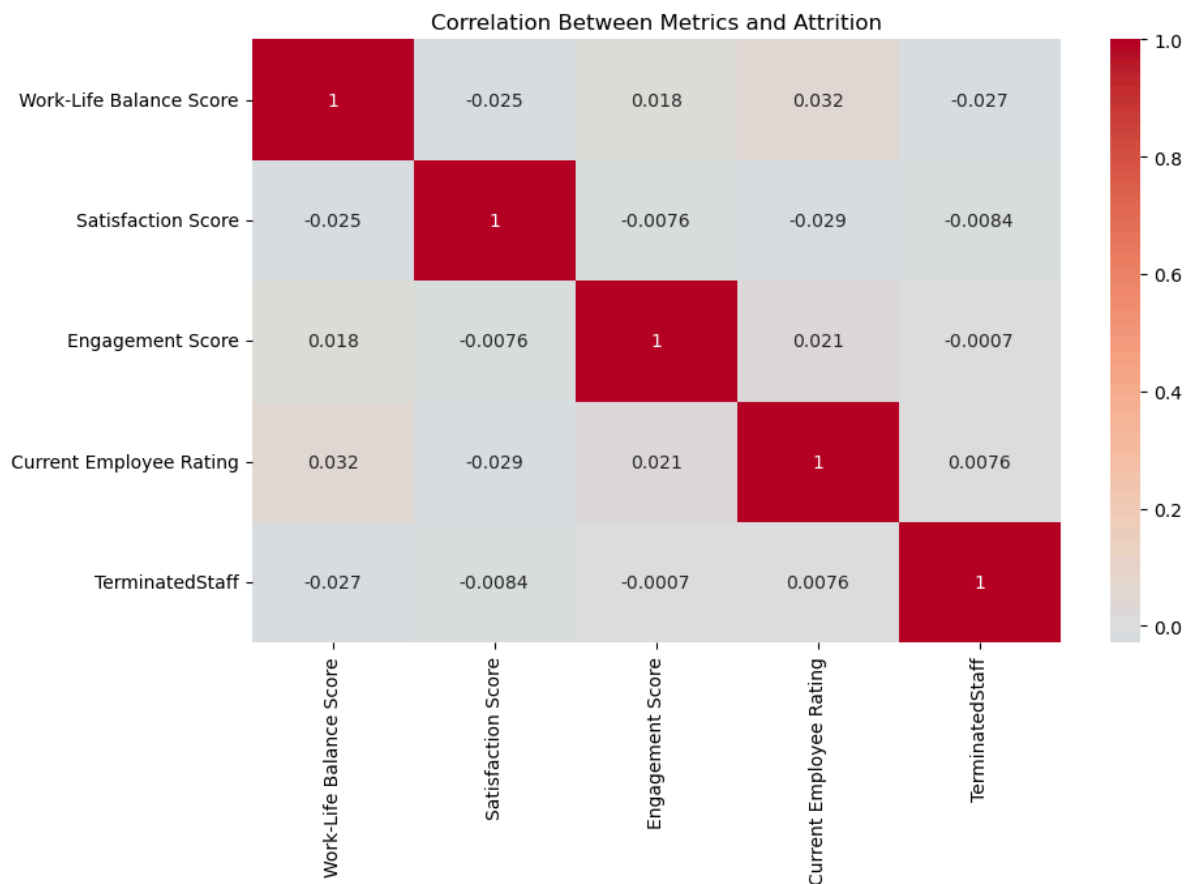
	Engagement Score	Satisfaction Score
Engagement Score	1.000000	-0.007598
Satisfaction Score	-0.007598	1.000000


```
In [64]: 1 df.groupby('Faculty')[['Engagement Score', 'Satisfaction Score']].corr().unstack()
2
```

```
Out[64]: Faculty
Faculty of Arts      -0.005539
Faculty of Engineering -0.013304
Faculty of Medicine   -0.004322
Name: (Engagement Score, Satisfaction Score), dtype: float64
```

```
In [65]: 1 # Calculate correlation matrix
2 corr_matrix = df[['Work-Life Balance Score', 'Satisfaction Score',
3                  'Engagement Score', 'Current Employee Rating',
4                  'TerminatedStaff']].corr()
5
6 # Visualize with heatmap
7 import matplotlib.pyplot as plt
8
9 import seaborn as sns
10 plt.figure(figsize=(10,6))
11 sns.heatmap(corr_matrix, annot=True, cmap='coolwarm', center=0)
12 plt.title("Correlation Between Metrics and Attrition")
```

```
Out[65]: Text(0.5, 1.0, 'Correlation Between Metrics and Attrition')
```



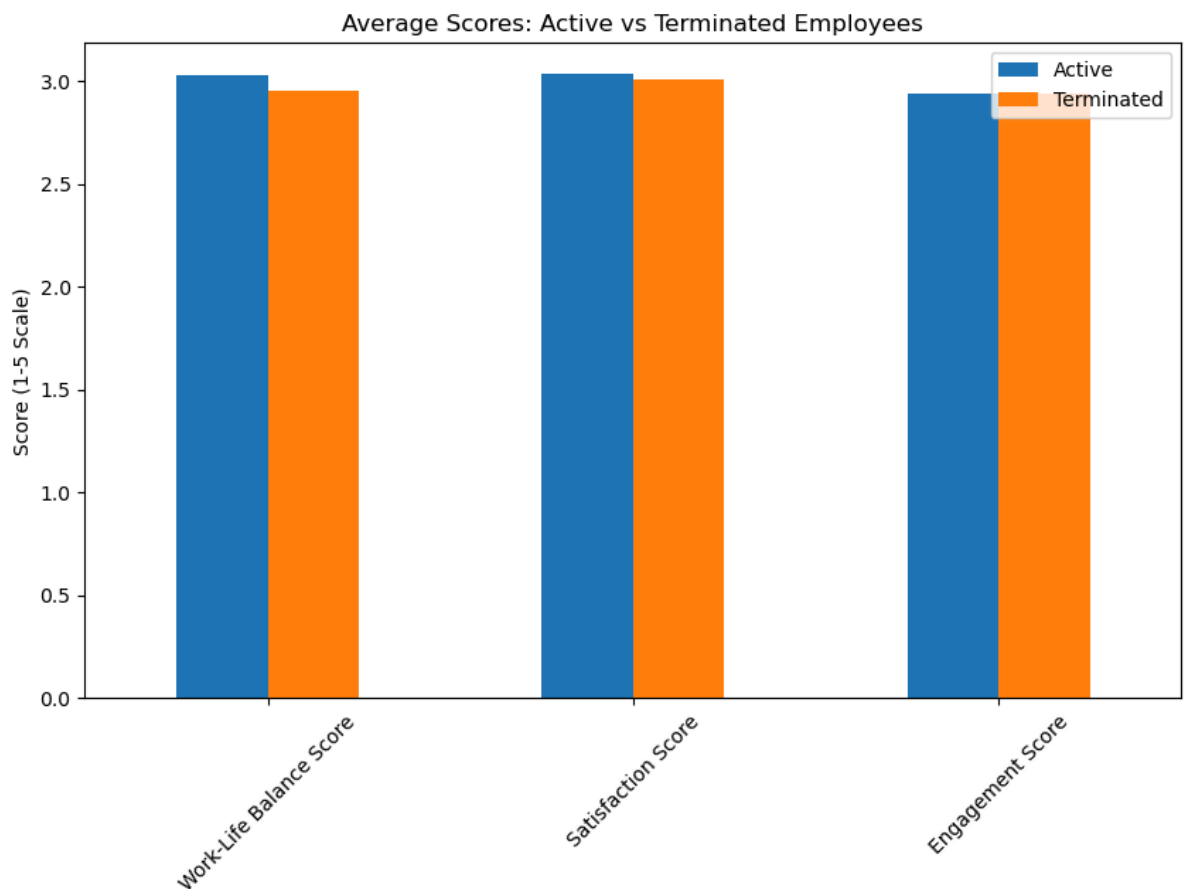
```
In [ ]: 1
```

```
In [66]: 1 # Compare scores between groups
2 terminated_scores = df[df['TerminatedStaff']==1][['Work-Life Balance Score',
3                                                    'Satisfaction Score',
4                                                    'Engagement Score']].mean()
5
6 active_scores = df[df['TerminatedStaff']==0][['Work-Life Balance Score',
7                                                  'Satisfaction Score',
8                                                  'Engagement Score']].mean()
9
10 comparison = pd.DataFrame({'Active': active_scores,
11                             'Terminated': terminated_scores,
12                             'Difference': active_scores - terminated_scores})
13 print(comparison)
```

	Active	Terminated	Difference
Work-Life Balance Score	3.027267	2.952381	0.074886
Satisfaction Score	3.034083	3.010437	0.023646
Engagement Score	2.940695	2.938682	0.002013

```
In [67]: 1 comparison[['Active', 'Terminated']].plot(kind='bar', figsize=(10,6))
2 plt.title("Average Scores: Active vs Terminated Employees")
3 plt.ylabel("Score (1-5 Scale)")
4 plt.xticks(rotation=45)
```

```
Out[67]: (array([0, 1, 2]),
[Text(0, 0, 'Work-Life Balance Score'),
Text(1, 0, 'Satisfaction Score'),
Text(2, 0, 'Engagement Score')])
```



Classification of Survey Scores into High, Medium, and Low Groups

```
In [68]: 1 def group(x):
2         if x < 3:
3             return 'Low'
4         elif x < 4:
5             return 'Medium'
6         else:
7             return 'High'
8
9 df['Total_Survey_Score'] = (df['Engagement Score'] + df['Satisfaction Score'] +
10 df['PerformanceGroup'] = df['Current Employee Rating'].apply(group)
11 df['EngagementGroup'] = df['Engagement Score'].apply(group)
12 df['SatisfactionGroup'] = df['Satisfaction Score'].apply(group)
13 df['Work-Life_BalanceGroup'] = df['Work-Life Balance Score'].apply(group)
14 df['Total_Survey_Group'] = df['Total_Survey_Score'].apply(group)
```

```
In [69]: 1 df.to_csv("./cleaned_test_csv.csv")
```

```
In [70]: 1 grouped = df.groupby(['Faculty', 'PerformanceGroup', 'Total_Survey_Group']).agg(
2     Total=('Employee ID', 'count'),
3     Terminated=('TerminatedStaff', 'sum')
4 ).reset_index()
5
6 grouped['AttritionRate'] = grouped['Terminated'] / grouped['Total'] * 100
```

In [71]:

```
1 grouped
```

Out[71]:

	Faculty	PerformanceGroup	Total_Survey_Group	Total	Terminated	AttritionRate
0	Faculty of Arts	High	High	22	12	54.545455
1	Faculty of Arts	High	Low	55	32	58.181818
2	Faculty of Arts	High	Medium	52	28	53.846154
3	Faculty of Arts	Low	High	18	8	44.444444
4	Faculty of Arts	Low	Low	66	35	53.030303
5	Faculty of Arts	Low	Medium	70	36	51.428571
6	Faculty of Arts	Medium	High	43	20	46.511628
7	Faculty of Arts	Medium	Low	145	74	51.034483
8	Faculty of Arts	Medium	Medium	120	61	50.833333
9	Faculty of Engineering	High	High	38	22	57.894737
10	Faculty of Engineering	High	Low	97	55	56.701031
11	Faculty of Engineering	High	Medium	85	47	55.294118
12	Faculty of Engineering	Low	High	35	18	51.428571
13	Faculty of Engineering	Low	Low	105	45	42.857143
14	Faculty of Engineering	Low	Medium	91	54	59.340659
15	Faculty of Engineering	Medium	High	78	37	47.435897
16	Faculty of Engineering	Medium	Low	193	98	50.777202
17	Faculty of Engineering	Medium	Medium	182	91	50.000000
18	Faculty of Medicine	High	High	52	25	48.076923
19	Faculty of Medicine	High	Low	151	75	49.668874
20	Faculty of Medicine	High	Medium	137	61	44.525547
21	Faculty of Medicine	Low	High	48	25	52.083333
22	Faculty of Medicine	Low	Low	176	98	55.681818
23	Faculty of Medicine	Low	Medium	172	79	45.930233
24	Faculty of Medicine	Medium	High	123	67	54.471545
25	Faculty of Medicine	Medium	Low	330	184	55.757576
26	Faculty of Medicine	Medium	Medium	316	146	46.202532

In [72]:

```
1 grouped.to_csv('./survey_analysis.csv')
```

In [73]:

```
1 pivot = grouped.pivot_table(  
2     index=['PerformanceGroup', 'Total_Survey_Group'],  
3     values='AttritionRate',  
4     aggfunc='mean'  
5 ).reset_index()
```

In [74]:

1 pivot

Out[74]:

	PerformanceGroup	Total_Survey_Group	AttritionRate
0	High	High	53.505705
1	High	Low	54.850574
2	High	Medium	51.221940
3	Low	High	49.318783
4	Low	Low	50.523088
5	Low	Medium	52.233154
6	Medium	High	49.473023
7	Medium	Low	52.523087
8	Medium	Medium	49.011955

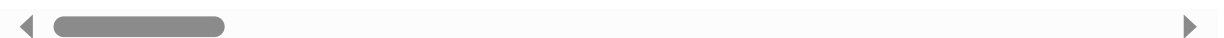
In [75]:

1 df

Out[75]:

	FirstName	LastName	StartDate	ExitDate	PositionTitle	Supervisor	ADEmail
0	Aaden	Mercer	2023-07-26	NaT	Senior Research Fellow	Victoria Jacobs	aaden.mercer@bilearner.com
1	Aaliyah	Watts	2020-05-09	NaT	Senior Lecturer	Jared Whitehead	aaliyah.watts@bilearner.com
2	Aarav	Espinoza	2019-11-11	NaT	Officer	David Ali	aarav.espinoza@bilearner.com
3	Aaron	Tapia	2023-01-20	NaT	Officer	Michael Reed	aaron.tapia@bilearner.com
4	Aaron	Weber	2020-04-06	2023-07-24	Lecturer	Melanie Garcia	aaron.weber@bilearner.com
...
2995	Zoey	Page	2019-05-06	2022-04-09	Professor	William Newman	zoey.page@bilearner.com
2996	Zoey	Spence	2020-05-06	NaT	Senior Lecturer	Ian Price	zoey.spence@bilearner.com
2997	Zoie	Logan	2019-04-19	NaT	Lecturer	Heidi Terry	zoie.logan@bilearner.com
2998	Zoie	Mercado	2020-12-01	NaT	Lecturer	David Young	zoie.mercado@bilearner.com
2999	Zoie	Rowland	2022-02-08	2022-07-22	Senior Lecturer	Mrs. Rachel Henry	zoie.rowland@bilearner.com

3000 rows × 41 columns



Staff-Type Analysis Transformations

```
In [76]: 1 df['ExitYear'] = pd.to_datetime(df['ExitDate']).dt.year
2 # Staff-type attrition trends
3 staff_attrition = df.groupby(['Faculty', 'StaffType']).agg(
4     TerminationCount=('TerminatedStaff', 'sum'),
5     VoluntaryTermination=('VoluntaryTermination', 'sum'),
6     InvoluntaryTermination=('InvoluntaryTermination', 'sum'),
7     RetirementCategory=('RetirementCategory', 'sum'),
8     ActiveStaff=('ActiveStaff', 'sum'),
9     TotalStaff=('Employee ID', 'count')
10 ).reset_index()
```

```
In [77]: 1 staff_attrition
```

Out[77]:

	Faculty	StaffType	TerminationCount	VoluntaryTermination	InvoluntaryTermination	RetirementC
0	Faculty of Arts	Academic	209	100	53	
1	Faculty of Arts	Professional	97	42	32	
2	Faculty of Engineering	Academic	326	170	79	
3	Faculty of Engineering	Professional	141	77	25	
4	Faculty of Medicine	Academic	489	235	131	
5	Faculty of Medicine	Professional	271	144	68	

```
In [78]: 1 df['Year'] = pd.to_datetime(df['ExitDate']).dt.year
2 terminated_df = df[df['TerminatedStaff'] == 1]
3
4 grouped = terminated_df.groupby(['Faculty', 'Year']).size().reset_index(name='T
5 grouped['CumulativeTerminations'] = grouped.groupby('Faculty')['Terminations'].
```

In [79]:

1 grouped

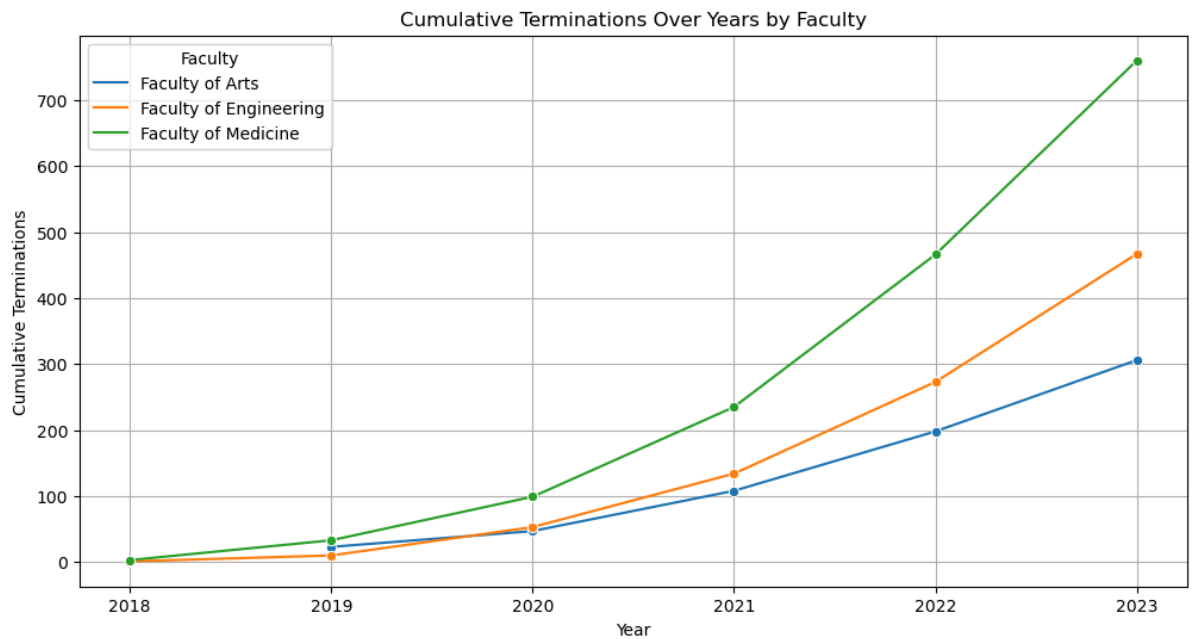
Out[79]:

	Faculty	Year	Terminations	CumulativeTerminations
0	Faculty of Arts	2019.0	23	23
1	Faculty of Arts	2020.0	24	47
2	Faculty of Arts	2021.0	61	108
3	Faculty of Arts	2022.0	90	198
4	Faculty of Arts	2023.0	108	306
5	Faculty of Engineering	2018.0	1	1
6	Faculty of Engineering	2019.0	9	10
7	Faculty of Engineering	2020.0	43	53
8	Faculty of Engineering	2021.0	81	134
9	Faculty of Engineering	2022.0	139	273
10	Faculty of Engineering	2023.0	194	467
11	Faculty of Medicine	2018.0	3	3
12	Faculty of Medicine	2019.0	30	33
13	Faculty of Medicine	2020.0	66	99
14	Faculty of Medicine	2021.0	136	235
15	Faculty of Medicine	2022.0	231	466
16	Faculty of Medicine	2023.0	294	760

```

In [80]: 1 import matplotlib.pyplot as plt
2 import seaborn as sns
3
4 plt.figure(figsize=(12, 6))
5 sns.lineplot(data=grouped, x='Year', y='CumulativeTerminations', hue='Faculty',
6
7 plt.title('Cumulative Terminations Over Years by Faculty')
8 plt.ylabel('Cumulative Terminations')
9 plt.xlabel('Year')
10 plt.legend(title='Faculty')
11 plt.grid(True)
12 plt.show()
13

```



In []: 1