

# Problem Statement

The management at Expert Sales Enterprise, a fictitious marketing agency, consistently invests in a variety of training and development schemes aimed at enhancing the success rate of their sales employees in achieving their set targets. The management is determined to comprehend the efficacy of these training initiatives and their contributions to the performance and efficiency of the employees. They are convinced that insights derived from data analysis focused on the impact will empower them to make well-informed choices concerning the optimization of the training programs and the efficient allocation of resources.

In your role as an HR Analyst, you are tasked with evaluating the influence of the training and development programs on the performance and productivity of the employees. Your responsibility involves conducting a comprehensive analysis of the impact of these programs using both pre-training and post-training data.

In [1]:

```
1 #import needed Liberaries
2
3 import pandas as pd
4 import numpy as np
5 import matplotlib.pyplot as plt
6 %matplotlib inline
7 from scipy.stats import ttest_rel
8 import seaborn as sns
```

In [2]:

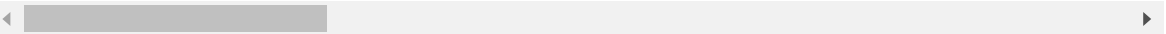
```
1 # Load the CSV data and create a dataframe
2 df = pd.read_csv(r'C:\Users\agozi\Desktop\DATA\HR_Data.CSV')
```

In [3]:

```
1 # print the data to see if it Loaded properly
2 df.head()
```

Out[3]:

	Emp_ID	Training_ID	Gender	Age	Edu_background	Experience_level	Pre_training_Sale
0	ENG00234	Training_2	Male	35	High school	senior	
1	ENG00073	Training_1	Female	27	High school	senior	
2	ENG00078	Training_3	Male	35	Masters	senior	
3	ENG00075	Training_3	Female	30	College_degree	senior	
4	ENG00127	Training_2	Male	23	College_degree	senior	



In [4]:

```
1 # use the describe command to check the distribution of the dataset
2 df.describe()
```

Out[4]:

	Age	Pre_training_Sales_revenue	Pre_training_data_year	Training_evaluation_scc
count	350.000000	3.500000e+02	350.0000	350.0000
mean	39.237143	2.032401e+07	2020.9800	52.5114
std	12.819813	1.141847e+07	0.7919	28.1334
min	18.000000	6.060100e+05	2020.0000	1.0000
25%	28.000000	1.063604e+07	2020.0000	28.0000
50%	39.000000	1.956613e+07	2021.0000	52.0000
75%	51.000000	3.052556e+07	2022.0000	77.7500
max	60.000000	3.993977e+07	2022.0000	100.0000

In [5]:

```
1 # have already cleaned the data in the SQL process but lets check for null values age
2 df.isnull().sum()
```

Out[5]:

Emp_ID	0
Training_ID	0
Gender	0
Age	0
Edu_background	0
Experience_level	0
Pre_training_Sales_revenue	0
Pre_training_data_year	0
Department	0
Role	0
Pre_training_target	0
Training_evaluation_score	0
Training_completion_status	0
Training_duration_in_days	0
Post_training_Sales_revenue	0
Post_training_data_year	0
Post_training_target	0
Employment_type	0

dtype: int64

## lets see how well each Traning Program performed

To do this, we will be comparing the sucess matrices to see which program performed better

In [6]:

```
1 # Replace 'Failed' with No and Achieved with yes
2 df['Pre_training_target'] = df['Pre_training_target'].replace('Failed', 'No')
3 df['Pre_training_target'] = df['Pre_training_target'].replace('Achieved', 'Yes')
4 df['Post_training_target'] = df['Post_training_target'].replace('Achieved', 'Yes')
5 df['Post_training_target'] = df['Post_training_target'].replace('Failed', 'No')
6 df.head(10)
```

Out[6]:

	Emp_ID	Training_ID	Gender	Age	Edu_background	Experience_level	Pre_training_Sale
0	ENG00234	Training_2	Male	35	High school	senior	
1	ENG00073	Training_1	Female	27	High school	senior	
2	ENG00078	Training_3	Male	35	Masters	senior	
3	ENG00075	Training_3	Female	30	College_degree	senior	
4	ENG00127	Training_2	Male	23	College_degree	senior	
5	ENG00027	Training_2	Female	30	Masters	senior	
6	ENG00074	Training_1	Male	28	Masters	Entry	
7	ENG00165	Training_3	Male	59	College_degree	mid-level	
8	ENG00205	Training_1	Female	51	College_degree	senior	
9	ENG00214	Training_1	Female	38	High school	mid-level	

In [7]:

```
1 # Lets check for unique values in the columns
2 unique_pre_training_targets = df['Pre_training_target'].unique()
3 unique_post_training_targets = df['Post_training_target'].unique()
4
5 print("Unique values in Pre_training_target:")
6 print(unique_pre_training_targets)
7
8 print("Unique values in Post_training_target:")
9 print(unique_post_training_targets)
10
```

Unique values in Pre\_training\_target:

['No' 'Yes']

Unique values in Post\_training\_target:

['Yes' 'No']

In [22]:

```
1 # Selecting specific columns for analysis
2 df1 = df[['Training_ID', 'Edu_background', 'Experience_level', 'Pre_training_Sales_revenue',
3           'Pre_training_target', 'Training_evaluation_score', 'Training_completion_status',
4           'Post_training_Sales_revenue', 'Post_training_data_year', 'Post_training_target']]
5
6 # Converting 'Pre_training_target' and 'Post_training_target' to binary using .loc indexer
7 df1.loc[:, 'Pre_training_target'] = df1['Pre_training_target'].map({'No': 0, 'Yes': 1})
8 df1.loc[:, 'Post_training_target'] = df1['Post_training_target'].map({'No': 0, 'Yes': 1})
9
10 # Display the updated DataFrame
11 #df1.head(10)
12
```

C:\Users\agozi\AppData\Local\Temp\ipykernel\_12776\729497627.py:7: SettingWithCopyWarning:  
A value is trying to be set on a copy of a slice from a DataFrame.  
Try using .loc[row\_indexer,col\_indexer] = value instead

See the caveats in the documentation: [https://pandas.pydata.org/pandas-docs/stable/user\\_guide/indexing.html#returning-a-view-versus-a-copy](https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy) ([https://pandas.pydata.org/pandas-docs/stable/user\\_guide/indexing.html#returning-a-view-versus-a-copy](https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy))

```
df1.loc[:, 'Pre_training_target'] = df1['Pre_training_target'].map({'No': 0, 'Yes': 1})
```

C:\Users\agozi\AppData\Local\Temp\ipykernel\_12776\729497627.py:8: SettingWithCopyWarning:  
A value is trying to be set on a copy of a slice from a DataFrame.  
Try using .loc[row\_indexer,col\_indexer] = value instead

See the caveats in the documentation: [https://pandas.pydata.org/pandas-docs/stable/user\\_guide/indexing.html#returning-a-view-versus-a-copy](https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy) ([https://pandas.pydata.org/pandas-docs/stable/user\\_guide/indexing.html#returning-a-view-versus-a-copy](https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy))

```
df1.loc[:, 'Post_training_target'] = df1['Post_training_target'].map({'No': 0, 'Yes': 1})
```

In [9]:

```
1 # Lets check for unique values in the columns
2 unique_pre_training_targets = df1['Pre_training_target'].unique()
3 unique_post_training_targets = df1['Post_training_target'].unique()
4
5 print("Unique values in Pre_training_target:")
6 print(unique_pre_training_targets)
7
8 print("\nUnique values in Post_training_target:")
9 print(unique_post_training_targets)
```

Unique values in Pre\_training\_target:  
[0 1]

Unique values in Post\_training\_target:  
[1 0]

In [10]:

```
1 # Calculate the subscription renewal rate for each group
2 Pre_Target_achieved_rate = df1.groupby('Training_ID')['Pre_training_target'].sum() /
3 Post_Target_achieved_rate = df1.groupby('Training_ID')['Post_training_target'].sum()
4
5 # Print the subscription renewal rate for each group
6 print("Target Achievement Rate for pre Traning Program:")
7 print(Pre_Target_achieved_rate)
8
9 print("\nTarget Achievement Rate for post Traning Program:")
10 print(Post_Target_achieved_rate)
```

Target Achievement Rate for pre Traning Program:

Training\_ID

Training\_1 0.427273

Training\_2 0.487395

Training\_3 0.495868

Name: Pre\_training\_target, dtype: float64

Target Achievement Rate for post Traning Program:

Training\_ID

Training\_1 0.436364

Training\_2 0.596639

Training\_3 0.528926

Name: Post\_training\_target, dtype: float64

## Testing the hypothesis

The essence of this test is to check if the difference in the pre training and post training matrices are statistically significantly or due to random chance.

Alpha = 0.05

Null Hypothesis (H0): There is no significant difference between the target achievement rates before and after the training programs.

Alternative Hypothesis (H1): There is a significant difference between the target achievement rates before and after the training programs.

In [11]:

```
1 # Target Achievement Rate for pre Training Program
2 # Perform paired t-test
3 t_statistic, p_value = ttest_rel(Pre_Target_achieved_rate, Post_Target_achieved_rate)
4
5 # Print the results
6 print(f"Paired t-test results:")
7 print(f"t-statistic: {t_statistic}")
8 print(f"p-value: {p_value}")
```

Paired t-test results:

t-statistic: -1.6713772807696066

p-value: 0.23660723904680867

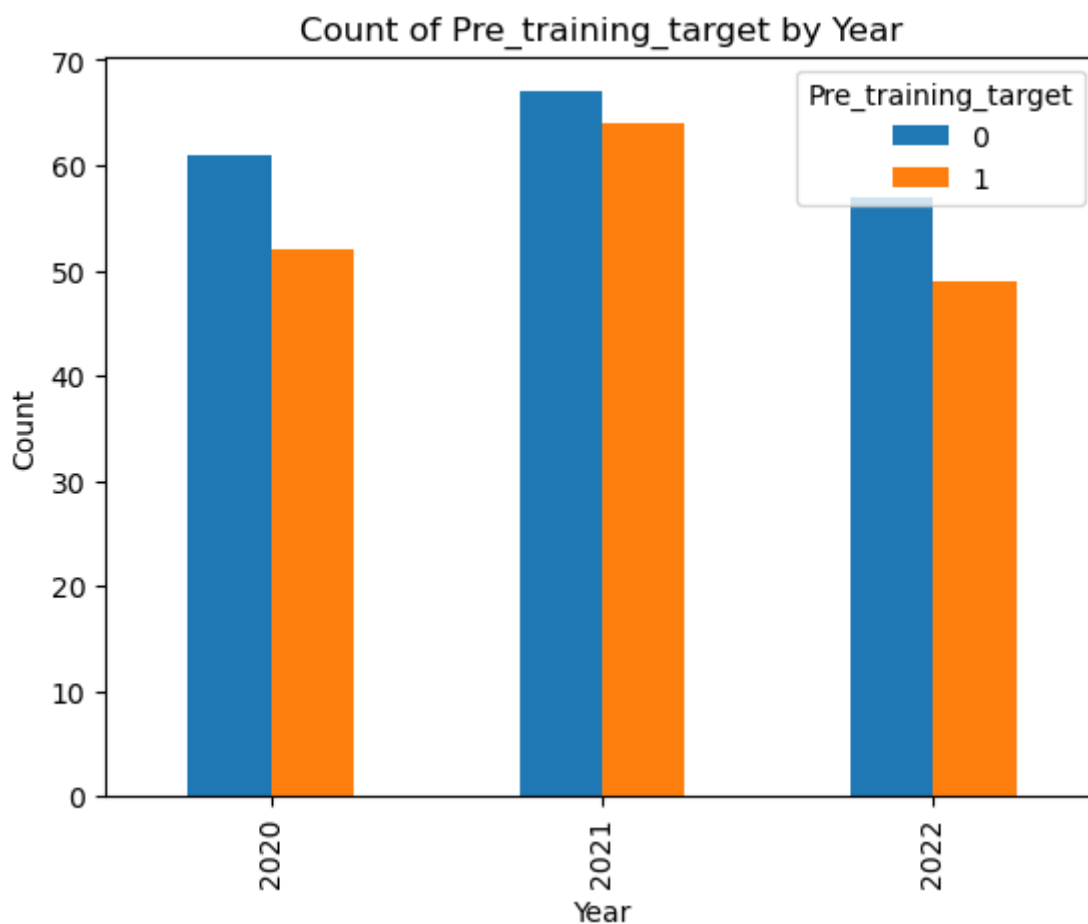
Interpretation:

With a p-value of 0.2366, which is greater than the typical significance level of 0.05, we do not have strong evidence to reject the null hypothesis. This means that we don't have sufficient evidence to conclude that There is a significant difference between the target achievement rates before and after the training programs.

## Understanding how the matrices changed year on year

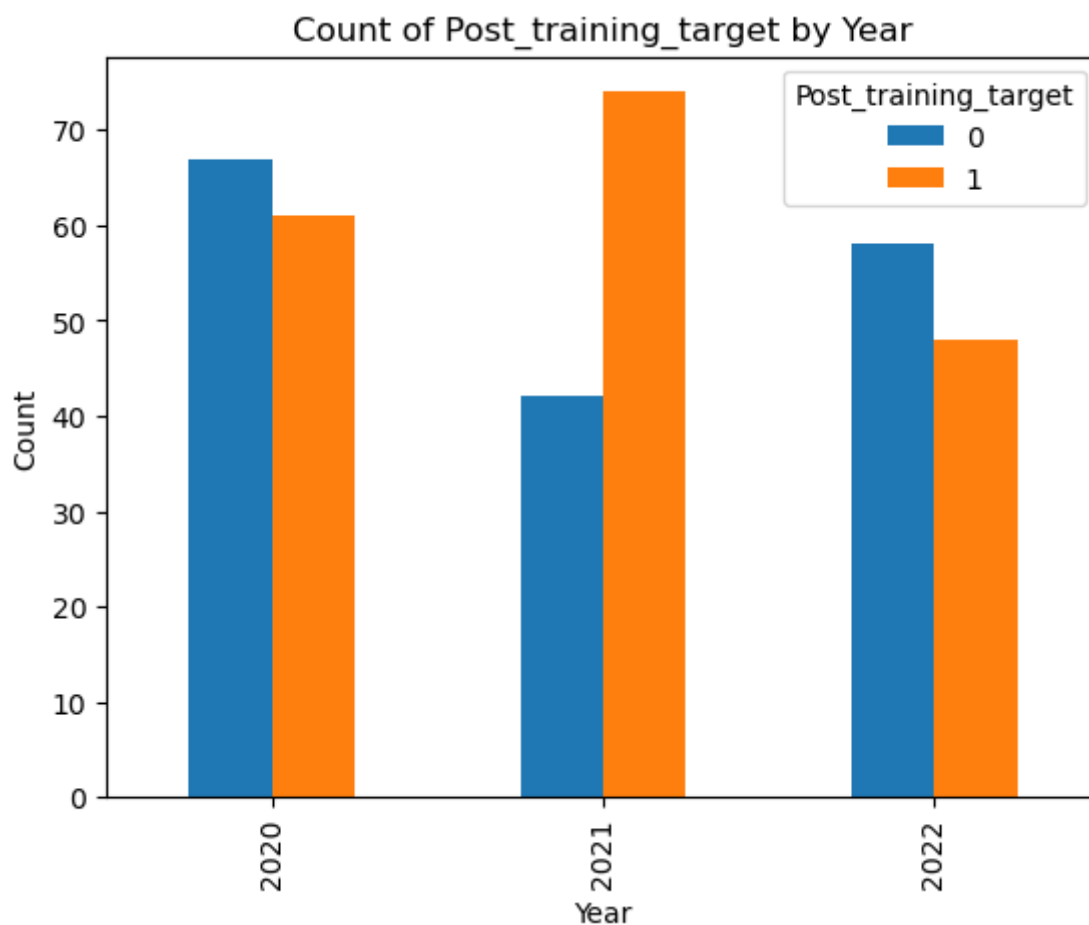
In [12]:

```
1 # Grouping by year and Pre_training_target and counting
2 grouped = df1.groupby(["Pre_training_data_year", "Pre_training_target"]).size().unstack()
3
4 # Plotting the bar chart
5 grouped.plot(kind="bar")
6
7 plt.title("Count of Pre_training_target by Year")
8 plt.xlabel("Year")
9 plt.ylabel("Count")
10 plt.legend(title="Pre_training_target", loc="upper right")
11
12 plt.show()
```



In [13]:

```
1 # Grouping by year and Post_training_target and counting
2 grouped = df1.groupby(["Post_training_data_year", "Post_training_target"]).size().un
3
4 # Plotting the bar chart
5 grouped.plot(kind="bar")
6
7 plt.title("Count of Post_training_target by Year")
8 plt.xlabel("Year")
9 plt.ylabel("Count")
10 plt.legend(title="Post_training_target", loc="upper right")
11
12 plt.show()
```



## The relationship between Training completion and target achievement rate

To start this, we will have to drop all training participants who did not actually complete the training.

In [14]:

```
1 # Drop rows where 'Training_completion_status' is 'No'
2 df2 = df1[df1['Training_completion_status'] != 'No']
3 df2.head()
```

Out[14]:

	Training_ID	Edu_background	Experience_level	Pre_training_Sales_revenue	Pre_training_d
0	Training_2	High school	senior	18046046	
1	Training_1	High school	senior	29646485	
2	Training_3	Masters	senior	30657715	
3	Training_3	College_degree	senior	5883971	
4	Training_2	College_degree	senior	15374421	

In [15]:

```
1 # using the unique command, Lets check if the changes where made.
2 unique_Training_completion_status = df2['Training_completion_status'].unique()
3 print(unique_Training_completion_status)
```

['Yes']

In [16]:

```
1 # Calculate the subscription renewal rate for each group
2 Pre_Target_achieved_rate1 = df2.groupby('Training_ID')['Pre_training_target'].sum()
3 Post_Target_achieved_rate1 = df2.groupby('Training_ID')['Post_training_target'].sum()
4
5 # Print the subscription renewal rate for each group
6 print("Target Achievement Rate for pre Traning Program:")
7 print(Pre_Target_achieved_rate1)
8
9 print("\nTarget Achievement Rate for post Traning Program:")
10 print(Post_Target_achieved_rate1)
```

Target Achievement Rate for pre Traning Program:

Training\_ID

Training\_1 0.430108

Training\_2 0.504950

Training\_3 0.495050

Name: Pre\_training\_target, dtype: float64

Target Achievement Rate for post Traning Program:

Training\_ID

Training\_1 0.451613

Training\_2 0.594059

Training\_3 0.544554

Name: Post\_training\_target, dtype: float64



In [17]:

```
1 # Target Achievement Rate for pre Training Program
2 # Perform paired t-test
3 t_statistic, p_value = ttest_rel(Pre_Target_achieved_rate1, Post_Target_achieved_rate1)
4
5 # Print the results
6 print(f"Paired t-test results:")
7 print(f"t-statistic: {t_statistic}")
8 print(f"p-value: {p_value}")
```

Paired t-test results:

t-statistic: -2.7215802288987496

p-value: 0.11264869218463308

Interpretation:

With a p-value of 0.1126, which is greater than the typical significance level of 0.05, we do not have strong evidence to reject the null hypothesis. This means that we don't have sufficient evidence to conclude that There is a significant difference between the target achievement rates before and after the training programs. However, when compared to the initial P-value of 0.23660 which included both those that completed the training as well as those who did not, we can see that those who completed the training has a lower p-value but thier outcome is due to random chance.

## Impact Analysis For Year 2021

In [18]:

```
1 # Filter the DataFrame based on the year
2 filtered_df = df2[df2['Pre_training_data_year'] == 2021]
3
4 # Calculate the subscription renewal rate for each group
5 Pre_Target_achieved_rate2 = filtered_df.groupby('Training_ID')['Pre_training_target']
6 Post_Target_achieved_rate2 = filtered_df.groupby('Training_ID')['Post_training_target']
7
8 # Print the subscription renewal rate for each group
9 print("Target Achievement Rate for pre Traning Program:")
10 print(Pre_Target_achieved_rate2)
11
12 print("\nTarget Achievement Rate for post Traning Program:")
13 print(Post_Target_achieved_rate2)
```

Target Achievement Rate for pre Traning Program:

Training\_ID

Training\_1      0.527778

Training\_2      0.515152

Training\_3      0.525000

Name: Pre\_training\_target, dtype: float64

Target Achievement Rate for post Traning Program:

Training\_ID

Training\_1      0.444444

Training\_2      0.666667

Training\_3      0.500000

Name: Post\_training\_target, dtype: float64

In [19]:

```
1 # Target Achievement Rate for pre Training Program
2 # Perform paired t-test
3 t_statistic, p_value = ttest_rel(Pre_Target_achieved_rate2, Post_Target_achieved_rate2)
4
5 # Print the results
6 print(f"Paired t-test results:")
7 print(f"t-statistic: {t_statistic}")
8 print(f"p-value: {p_value}")
```

Paired t-test results:

t-statistic: -0.20388503898940008

p-value: 0.8573067874378176

Interpretation:

With a p-value of 0.857, which is greater than the typical significance level of 0.05, we do not have strong evidence to reject the null hypothesis. This means that we don't have sufficient evidence to conclude that There is a significant difference between the target achievement rates before and after the training programs for the year 2021. Meaning that the Post training outcome of 2021 is due to random chance and not as a result of the trainings.

## Let's check the relationship between training evaluation score and Post training Sales revenue

In [20]:

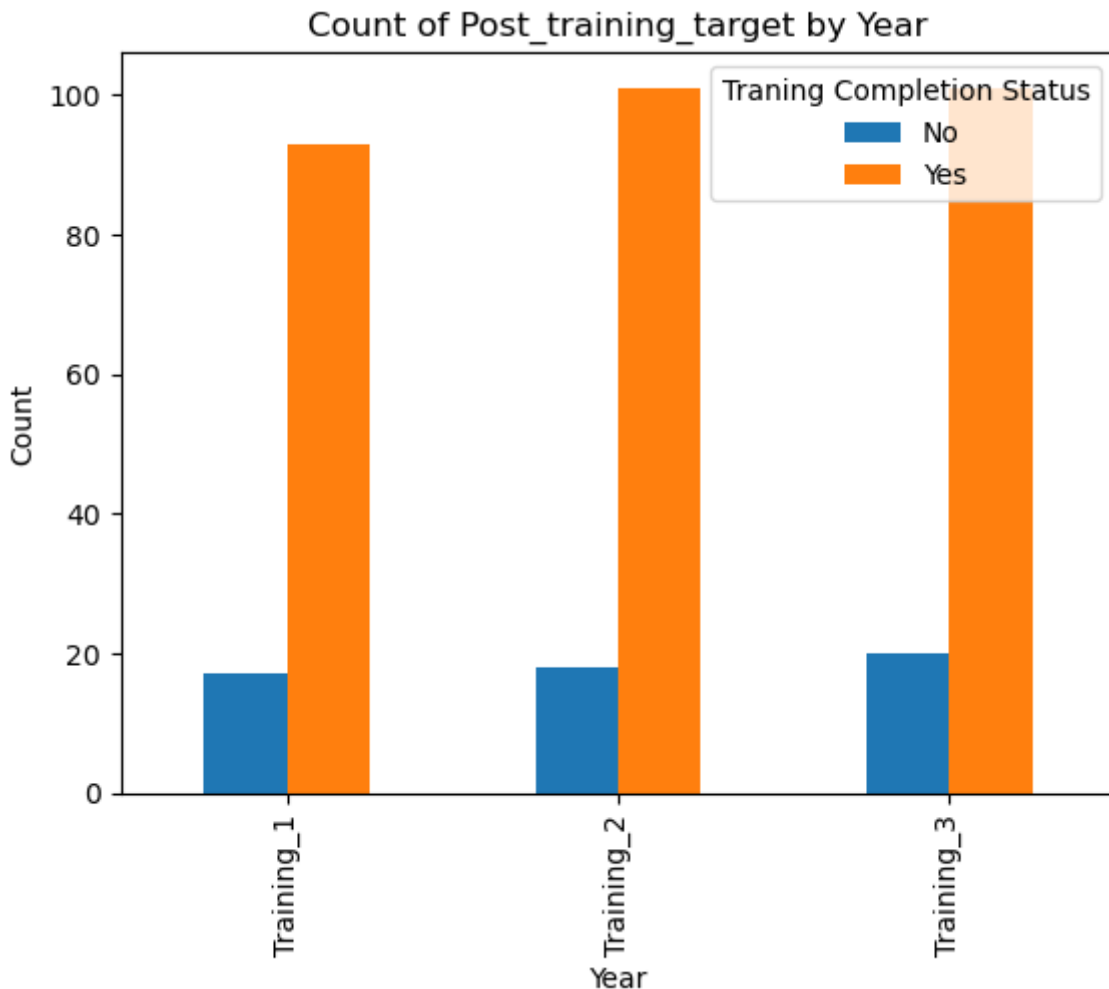
```
1 # Lets call in the needed variables
2 corr_data = df2[['Training_evaluation_score', 'Post_training_Sales_revenue']]
3
4 # Calculate correlation coefficient
5 correlation = corr_data['Training_evaluation_score'].corr(corr_data['Post_training_Sales_revenue'])
6
7 # Print the correlation coefficient
8 print("Correlation coefficient:", correlation)
```

Correlation coefficient: 0.01077660727173629

Since the correlation coefficient is close to zero(0.011), it implies that changes in the 'Training evaluation score' are not strongly associated with changes in the Post training Sales revenue.

In [21]:

```
1 # Grouping by year and Post_training_target and counting
2 grouped1 = df1.groupby(["Training_ID", "Training_completion_status"]).size().unstack
3
4 # Plotting the bar chart
5 grouped1.plot(kind="bar")
6
7 plt.title("Count of Post_training_target by Year")
8 plt.xlabel("Year")
9 plt.ylabel("Count")
10 plt.legend(title="Traning Completion Status", loc="upper right")
11
12 plt.show()
```



The diagram shows that a large population of their staffs actually completed the training program when compared to those who did not.

## Conclusion

After carefully examining the data in relation to the business objectives, the following insights was derived.

The overall best performing training program was Training\_2 after it recorded about a 10% increase in the post training target achievement rate. the second best was training\_3 after an increase of 3% and the worst performing was training\_1 with increase of just 0.9%. It is however important to know that the outcome of these training was noted not to be statistically significant but as a result of random chance. To carefully

understand this, employees who failed to complete the training was dropped. After carefully dropping employees that did not complete the program, training\_2 was still the best performing following also by training\_3 with a post training target achievement rate increase of 9% and 5% respectively. however, the result still failed to show that this result was statistically significant.

Year by Year examination of the data show that the trainees only achieved their post training targets in 2021 with traning\_2 still topping the chart after 70 trainees achieved their target while 40 failed. However, when the significance of this outcome was tested, it produced a p-value of over 80% strongly indicating that the outcome was due to random chance. This is true as 2021 is when Covid-19 restrictions and lockdown was eased leading to pending deals being finalized and closed which was then registered as part of their 2021 target.

## Recommendation

Based on the analysis conducted, the following recommendation were made:

- 1) The total efficiency of these training programs needs to be completely evaluated by completely taking a survey of the trainees who attended these programs not minding if they completed or not. This will help us understand if the trainees believe these programs had an effect on their efficiency or not. It's only when we have done this that we can fully know if any training program should be dropped or scraped all.
- 2) Conducting further analysis to explore potential factors that might contribute to the lack of statistical significance in the results. For example, we could explore demographic variables or other external factors that could impact training effectiveness.

In [ ]:

1	
---	--