n [3]: n [4]:	<pre>Import stat Importation des données countries_HDI = pd.read_csv('C:\\Users\\SCD UM\\Downloads\\stat\\countries.HDI.csv', encoding='latin1', header=0) effec1_quest_compil = pd.read_csv('C:\\Users\\SCD UM\\Downloads\\stat\\effec1.quest.compil.csv', encoding='latin1') effec2_quest_compil = pd.read_csv('C:\\Users\\SCD UM\\Downloads\\stat\\effec2.quest.compil.csv', encoding='latin1') effec3_quest_compil = pd.read_csv('C:\\Users\\SCD UM\\Downloads\\stat\\effec3.quest.compil.csv', encoding='latin1')</pre>
	<pre>effec2 = pd.merge(effec2_quest_compil, usages_effec2, on='Student_ID') effec3 = pd.merge(effec3_quest_compil, usages_effec3, on='Student_ID') # Ajoutons une nouvelle itération pour chaque DataFrame: effec1['itération'] = 1</pre>
	<pre>effec2['itération'] = 2 effec3['itération'] = 3 # Renommer les colonnes de countries_HDI countries_HDI = countries_HDI.rename(columns={"Norvège": "Country", "TH": "HDI", "1": "Identifier"}) Pour rassembler les données des différentes itérations ("row bind"), nous utiliserons la fonction concat() de pandas : final_df = pd.concat([effec1, effec2, effec3], ignore_index=True) final_df2 = pd.merge(final_df, countries_HDI, on="Country", how="left")</pre>
n [9]:	Enfin, pour conserver seulement certaines colonnes (par exemple 'HDI' et 'genre') #Extraction des colonnes sélectionnées dans un DataFrame final selected_columns_df = final_df2[['HDI', 'Gender', 'itération']] Pour quantifier le nombre de vidéos visionnées et de quiz réalisés par chaque apprenant, on peut créér deux nouvelles métriques quantitatives qui nous permettent de mesurer l'engagement des apprenants de manière plus précise. # Création des colonnes pour le nombre total de vidéos visionnées et de quiz réalisés video_columns = [col for col in final_df2.columns if 'S' in col] quiz_columns = [col for col in final_df2.columns if 'Quizz' in col and 'bin' in col]
	# Convertir les données non numériques en NaN et les remplacer par zéro final_df2[video_columns] = final_df2[video_columns].apply(pd.to_numeric, errors='coerce').fillna(0) final_df2[quiz_columns] = final_df2[quiz_columns].apply(pd.to_numeric, errors='coerce').fillna(0) Le nombre total de vidéos visionnées et de quiz réalisés par chaque étudiant final_df2['Total_Videos_visionnees'] = final_df2[video_columns].apply(sum, axis=1) final_df2['Total_Quizzes_realises'] = final_df2[quiz_columns].apply(sum, axis=1) final_df2['Total_Videos_visionnees'] = final_df2['Total_Videos_visionnees'].astype(int) # Compter le nombre d'étudiants dans chaque catégorie HDI HDI courts = final_df2['HDT'] value courts()
	HDI_counts = final_df2['HDI'].value_counts() print(HDI_counts) # Créer une nouvelle colonne HDI où 'M' et 'H' sont regroupées en 'I' final_df2['HDI_Grouped'] = final_df2['HDI'].replace({'M': 'I', 'H': 'I'}) # Verifier le resultat final_df2[['HDI', 'HDI_Grouped', 'Total_Videos_visionnees', 'Total_Quizzes_realises']].head() B 13292 TH 7260 M 354 H 313 Name: HDI, dtype: int64
t[11]: [12]:	HDI HDI_Grouped Total_Videos_visionnees Total_Quizzes_realises 0 B B 221 0 1 B B 221 0 2 TH TH 19179 0 3 TH TH 1116 4 4 TH TH 1949 0 # Compter le nombre d'étudiants dans chaque catégorie HDI_Grouped HDI_Grouped_counts = final_df2['HDI_Grouped'].value_counts() print(HDI_Grouped_counts)
	B 13292 TH 7260 I 667 Name: HDI_Grouped, dtype: int64 # Création des colonnes pour chaque type d'apprenant final_df2['Bystander'] = np.where((final_df2['Total_Videos_visionnees']==0) & (final_df2['Total_Quizzes_realises']==0), 1, 0) final_df2['Auditing'] = np.where((final_df2['Total_Videos_visionnees']>0) & (final_df2['Total_Quizzes_realises']==0), 1, 0) final_df2['Completer'] = np.where((final_df2['Total_Videos_visionnees']>0) & (final_df2['Total_Quizzes_realises']>0), 1, 0) final_df2['Disengaging'] = np.where((final_df2['Total_Videos_visionnees']==0) & (final_df2['Total_Quizzes_realises']==0), 1, 0) # Création de la table de synthèse
	<pre>summary_table = final_df2.groupby('itération')[['Bystander', 'Auditing', 'Completer', 'Disengaging']].sum() summary_table</pre>
[22]: [23]:	<pre>summary_table = summary_table.apply(lambda x: x / x.sum(), axis=1) print(summary_table)</pre>
	<pre># Set style of the plot sns.set_style("whitegrid") # Preparer les données df_plot = summary_table.reset_index().melt('itération') # Créer la bar plot plt.figure(figsize=(9, 7)) bar_plot = sns.barplot(x='itération', y='value', hue='variable', data=df_plot) # Placer la legend à l'extérieur plt.legend(bbox_to_anchor=(1.05, 1), loc=2, borderaxespad=0.) plt.title('Evolution des types d\'apprenants par itération') plt.xlabel('Itération')</pre>
	plt.ylabel('Proportion') # Sauvegarder la figure plt.savefig("Evolution des types d'apprenants par itération.png") plt.show() Evolution des types d'apprenants par itération Bystander Auditing Completer Disengaging
	0.4 5g 0.2 0.1
[25]:	1 2 3 4 Chi2 et mosaic plot:
	<pre># La Performance de test Chi-Square chi2, p, dof, expected = chi2_contingency(crosstab) print("Chi-Square value: ", chi2) print("p-value: ", p) Chi-Square value: 175.75674510524374 p-value: 7.273677317477502e-38 Chi-Square:</pre> Mosaic plot:
[26]:	from statsmodels.graphics.mosaicplot import mosaic import matplotlib.pyplot as plt plt.rcParams['font.size'] = 12 # change la taille de la police plt.figure(figsize=(12,8)) mosaic(final_df2, ['HDI', 'Gender'], title='Diagramme en mosaïque de IDH et Genre') plt.xlabel('IDH') # étiquette pour l'axe x plt.ylabel('Genre') # étiquette pour l'axe y plt.savefig("Diagramme en mosaïque de IDH et Genre.png") plt.show() <figure 0="" 864x576="" axes="" size="" with=""> Diagramme en mosaïque de IDH et Genre DH</figure>
	un homme - TH uncomme Bun homme un homme un homme un femme - TH une femme Bune femme une femme une femme une femme
[27]:	Mosaic Plot: import numpy as np # Calculate Cramer's V V = np.sqrt((chi2/(final_df2.shape[0]*(min(crosstab.shape)-1)))) print("Cramer's V: ", V) Cramer's V: 0.09081213724541108 5 Modèle linéaire, tests non paramétriques:
[28]:	<pre>**1. Test de Student : from scipy.stats import ttest_ind # Séparer les données en deux groupes group1 = final_df2[final_df2['Gender'] == 'une femme']['Total_Videos_visionnees'] group2 = final_df2[final_df2['Gender'] == 'un homme']['Total_Videos_visionnees'] # Effectuer le test t de Student t_stat, p_val = ttest_ind(group1, group2, equal_var=False) print(f"t-statistic: {t_stat}")</pre>
[29]:	<pre>print(f"p-value: {p_val}") t-statistic: 5.369331374528332 p-value: 8.209732654547305e-08 2. Test non paramétrique: from scipy.stats import mannwhitneyu # Effectuer le test U de Mann-Whitney u_stat, p_val = mannwhitneyu(group1, group2) print(f"U-statistic: {u_stat}") print(f"p-value: {p_val}") U-statistic: 10010895 0</pre>
[30]:	U-statistic: 10010895.0 p-value: 4.604958486880628e-11 3. Régression linéaire et tests de corrélation: from scipy.stats import pearsonr, spearmanr # Calculer la corrélation de Pearson pearson_corr, _ = pearsonr(final_df2['Total_Quizzes_realises'], final_df2['Total_Videos_visionnees']) print(f"Pearson correlation: {pearson_corr}") # Calculer la corrélation de Spearman spearman_corr, _ = spearmanr(final_df2['Total_Quizzes_realises'], final_df2['Total_Videos_visionnees']) print(f"La correlation de Spearman : {spearman_corr}")
	70000 60000 10000 10000 0 1 2 3 4 5 Total_Quizzes_realises
[31]:	# Fit the model model = ols('Total_Videos_visionnees ~ C(HDI) + C(Gender)', data=final_df2).fit() # Perform ANOVA and print the table anova_table = sm.stats.anova_lm(model, typ=2) print(anova_table)
[32]:	sum_sq df F PR(>F) C(HDI) 3.455544e+10 3.0 39.106792 4.221582e-25 C(Gender) 4.616417e+09 1.0 15.673357 7.585347e-05 Residual 2.664990e+12 9048.0 NaN NaN import statsmodels.api as sm from statsmodels.formula.api import ols # Fit the model model = ols('Total_Videos_visionnees ~ C(HDI) + C(Gender)', data=final_df2).fit() # Perform ANOVA and print the table anova_table = sm.stats.anova_lm(model, typ=2)
	<pre># Add significance level stars def add_stars(p): if p < 0.001: return '***' elif p < 0.01: return '**' elif p < 0.05; return '*' else: return '.' anova_table['stars'] = [add_stars(p) for p in anova_table['PR(>F)']] print(anova_table)</pre>
n []: [34]:	sum_sq df F PR(>F) stars C(HDI) 3.455544e+10 3.0 39.106792 4.221582e-25 *** C(Gender) 4.616417e+09 1.0 15.673357 7.585347e-05 *** Residual 2.664990e+12 9048.0 NaN NaN . Pour obtenir les statistiques inférentielles, on peut utiliser la fonction summaryde la bibliothèque statsmodels. Cette fonction affiche les coefficients de la régression, l'erreur standard, les valeurs de tet de p, l'intervalle de confiance, et d'autres statistiques applicables. # Affiche le tableau de résumé print(model.summary())
	OLS Regression Results Dep. Variable: Total_Videos_visionnees Model: OLS Adj. R-squared: 0.015 Model: Least Squares F-statistic: 36.55 Date: Tue, 04 Jul 2023 Prob (F-statistic): 2.34e-30 Time: 23:10:46 Log-Likelihood: -1.0111e+05 No. Observations: 9053 AIC: 2.022e+05 Df Residuals: 9048 BIC: 2.023e+05 Df Model: 4 Covariance Type: nonrobust ====================================
	C(HDI)[T.H] 3121.5829 1101.389 2.834 0.005 962.612 5280.554 C(HDI)[T.M] 2616.9782 1047.104 2.499 0.012 564.417 4669.540 C(HDI)[T.TH] 5737.8907 552.637 10.383 0.000 4654.597 6821.184 C(Gender)[T.une femme] 1537.7852 388.432 3.959 0.000 776.371 2299.199 ================================
[35]:	<pre># Obtenir le sommaire sous forme de dataframe summary_df = pd.read_html(model.summary().tables[1].as_html(), header=0, index_col=0)[0] # Définir une fonction pour mapper les p-valeurs à des étoiles def significance(p): if p < 0.001: return '***' elif p < 0.01: return '**' elif p < 0.05: return '*' else: return ''' </pre>
	# Appliquer la fonction à la colonne de p-valeur summary_df['P> t '] = summary_df['P> t '].apply(significance) # Afficher le tableau print(summary_df)
[38]:	Intercept 19200.000 C(HDI)[T.H] 5280.554 C(HDI)[T.M] 4669.540 C(HDI)[T.TH] 6821.184 C(Gender)[T.une femme] 2299.199 import statsmodels.api as sm import statsmodels.formula.api as smf # Créez le modèle avec le terme d'interaction model_interaction = smf.ols('Total_Videos_visionnees ~ C(HDI) * C(Gender)', data=final_df2).fit() # Affichez le résumé du modèle print(model_interaction.summary())
	Dep. Variable: Total_Videos_visionnees R-squared: 0.016
	Intercept 1.824e+04 555.378 32.845 0.000 1.72e+04 1.93e+04
[39]:	Notes: [1] Standard Errors assume that the covariance matrix of the errors is correctly specified. import statsmodels.api as sm import statsmodels.formula.api as smf # Créez le modèle avec le terme d'interaction model_interaction = smf.ols('Total_Videos_visionnees ~ C(HDI) * C(Gender)', data=final_df2).fit() # Obtenez le résumé du modèle sous forme de DataFrame summary_df = pd.read_html(model_interaction.summary().tables[1].as_html(), header=0, index_col=0)[0] # Ajoutez une colonne pour les étoiles d'importance summary_df['Significance'] = '' summary_df['Significance'] = '' summary_df.loc[summary_df['P> t '] <= 0.001, 'Significance'] = '***'
	<pre>summary_df.loc[(summary_df['P> t '] > 0.001) & (summary_df['P> t '] <= 0.01), 'Significance'] = '**' summary_df.loc[(summary_df['P> t '] > 0.001) & (summary_df['P> t '] <= 0.05), 'Significance'] = '*' summary_df.loc[(summary_df['P> t '] > 0.05) & (summary_df['P> t '] <= 0.1), 'Significance'] = '.' summary_df.loc[(summary_df['P> t '] > 0.1, 'Significance'] = '.' # Imprimez le DataFrame print(summary_df)</pre>
	C(HDI)[T.H]:C(Gender)[T.une femme]
[43]:	Pour la régression logistique, notre variable dépendante doit être transformée en une variable dichotomique (0 ou 1). Supposons que nous utilisions la variable "succès", qui indique si un utilisater a réussi à obtenir un certificat (1 pour le succès, 0 pour l'échec). final_df2['Success'] = final_df2['Completer'] model = smf.glm(formula="Success ~ C(Gender) + C(HDI)",
	Dep. Variable: Success No. Observations: 9053 Model: GLM Df Residuals: 9048 Model Family: Binomial Df Model: 4 Link Function: Logit Scale: 1.0000 Method: IRLS Log-Likelihood: -4524.8 Date: Tue, 04 Jul 2023 Deviance: 9049.7 Time: 23:28:13 Pearson chi2: 9.05e+03 No. Iterations: 4 Pseudo R-squ. (CS): 0.002517 Covariance Type: nonrobust
[44]:	
	Intercept 1.2019 0.071 17.001 0.000 1.063 1.340 C(Gender)[T.une femme] -0.1039 0.056 -1.847 0.065 -0.214 0.006 C(HDI)[T.H] 0.1194 0.155 0.772 0.440 -0.184 0.423 C(HDI)[T.M] -0.1281 0.140 -0.913 0.361 -0.403 0.147 C(HDI)[T.TH] 0.2753 0.077 3.577 0.000 0.124 0.426 ## Calculate Odds Ratios OR = np.exp(model.params) ## Calculate lower and upper confidence intervals CI_lower = np.exp(model.conf_int()[0]) CI_upper = np.exp(model.conf_int()[1])
	Intercept 1.2019 0.071 17.001 0.000 1.063 1.340 C(Gender)[T.une femme] -0.1039 0.056 -1.847 0.065 -0.214 0.006 C(HDI)[T.H] 0.1194 0.155 0.772 0.440 -0.184 0.423 C(HDI)[T.M] -0.1281 0.140 -0.913 0.361 -0.403 0.147 C(HDI)[T.TH] 0.2753 0.077 3.577 0.000 0.124 0.426 ===================================
[45]:	Intercept 1.2819 0.071 17.081 0.080 1.083 1.348
[45]:	Intercept 1.2019 0.071 17.001 0.000 1.003 1.340 C(Gender)[T.une femme] -0.1839 0.056 -1.847 0.065 -0.214 0.006 C(HD1)[T.H] 0.1194 0.155 0.772 0.440 -0.194 0.423 C(HD1)[T.H] 0.1201 0.104 -0.193 0.301 -0.403 0.147 C(HD1)[T.H] 0.2753 0.077 3.77 0.000 0.124 Calculate Older Ratios OR = np.exp(model.params) # Calculate Older Ratios OR = np.exp(model.conf.int()[0]) CLupper = np.exp(model.conf.int()[1]) # Create a DataFrame for the Odds Ratios and Confidence Intervals results.df = pd.DataFrame("Odds Ratio": OR,
[45]:	Tree
[45]:	Intercept (1.2019 0.071 17.000 1.00000 1.00000 1.00000 1.0000 1.0000 1.0000 1.0000 1.0000 1.00000 1.0000 1.0000 1.0000 1.
	Tricking Tricking 0.100
[46]:	Commerce
[47]:	Content Cont
[47]:	1
[47]:	Property
[47]:	### Committee 1985
[47]:	Security for the large of the control of the contro
[47]:	Company Comp
[47]: [51]:	Martin M
[47]: [51]:	March Marc
[47]: [52]:	March Marc
[46]: [51]: [52]:	March Marc
[47]: [51]: [52]:	March Marc
[46]: [51]: [52]:	
[47]: [51]:	
[47]: [51]: [52]:	The content of the
[46]: [51]: [52]:	March Marc