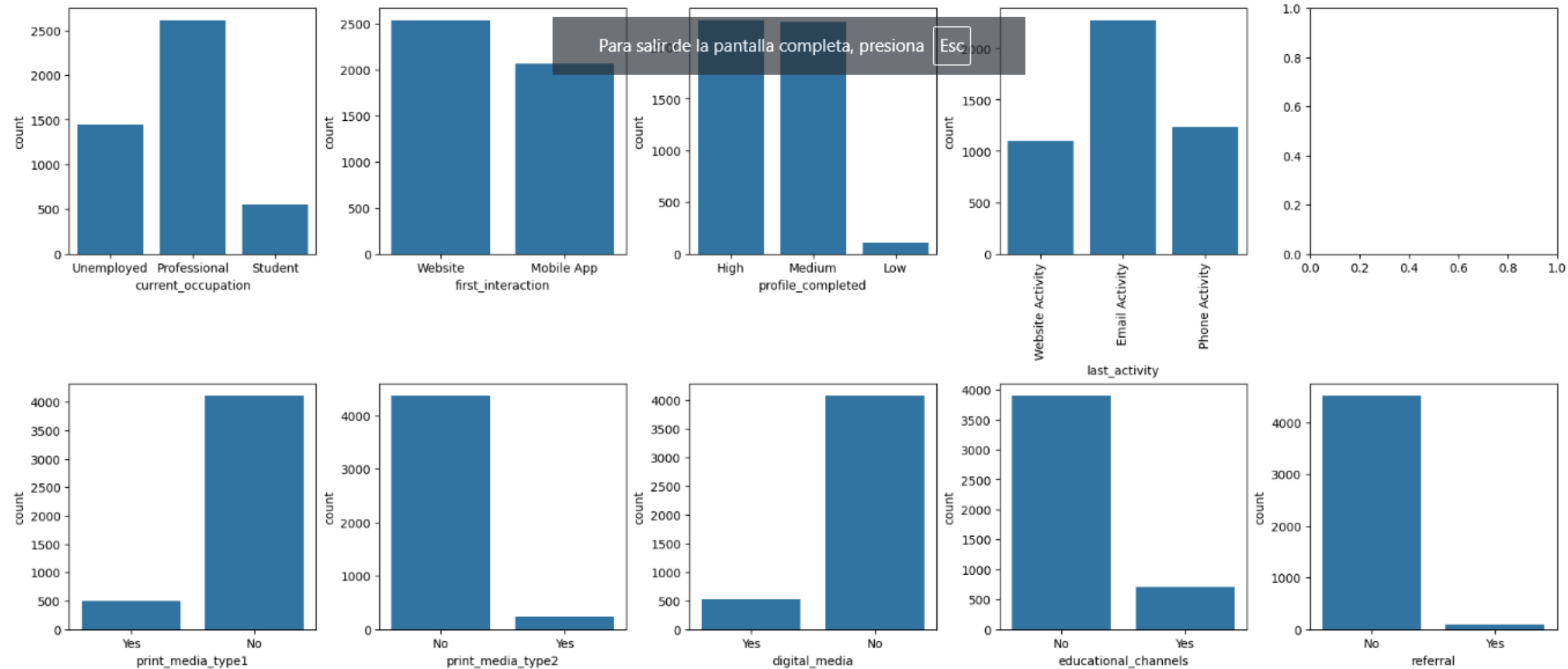
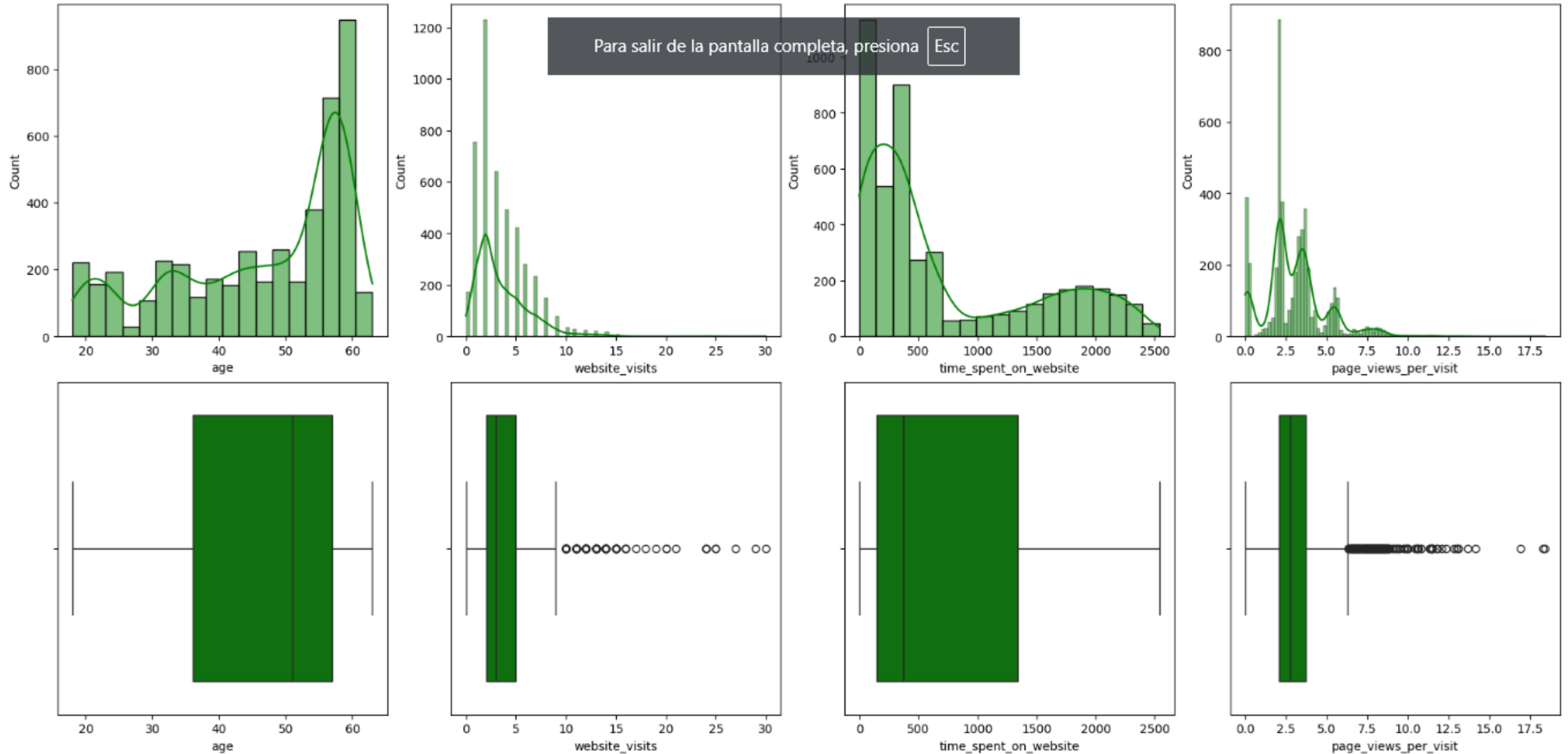


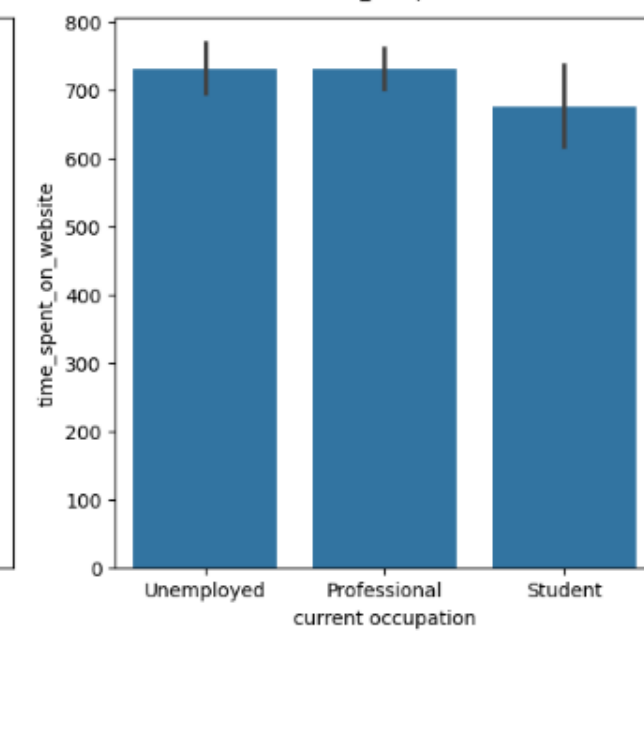
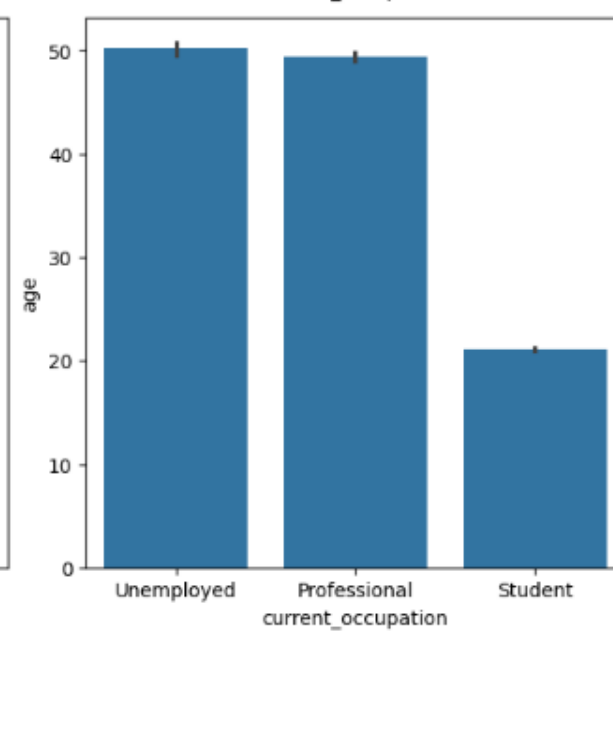
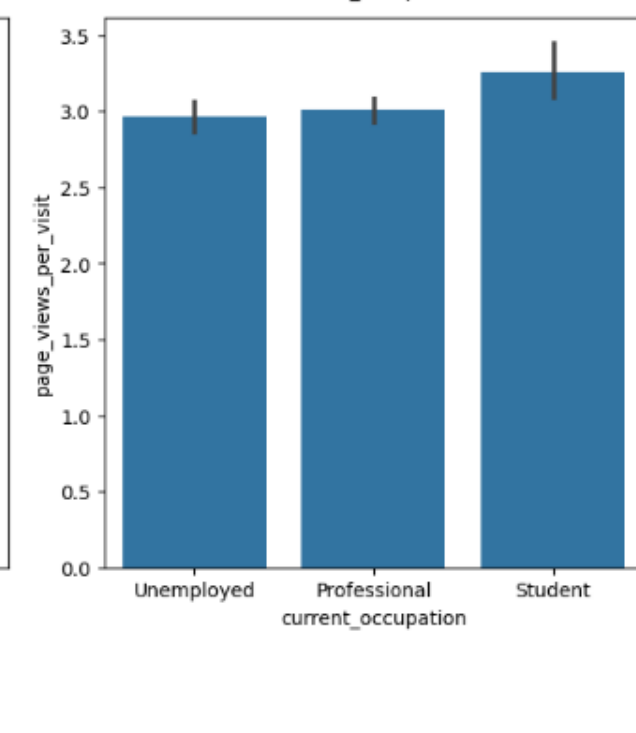
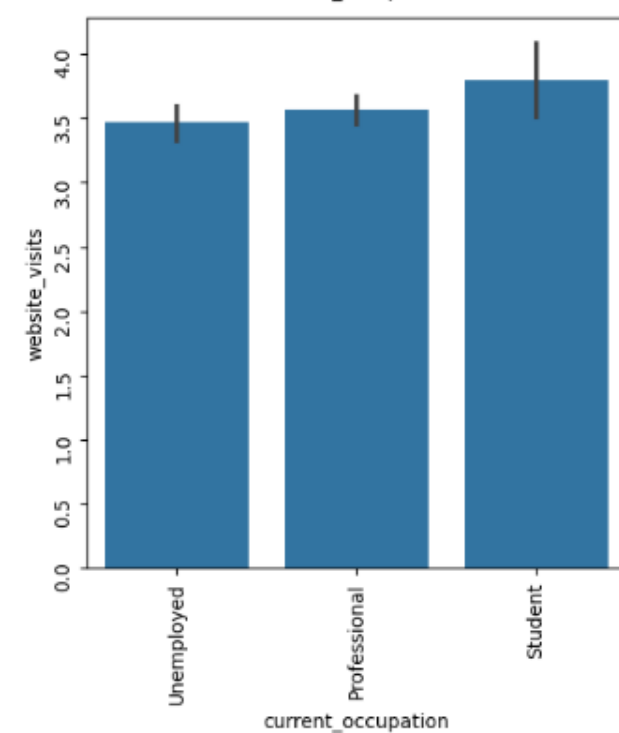
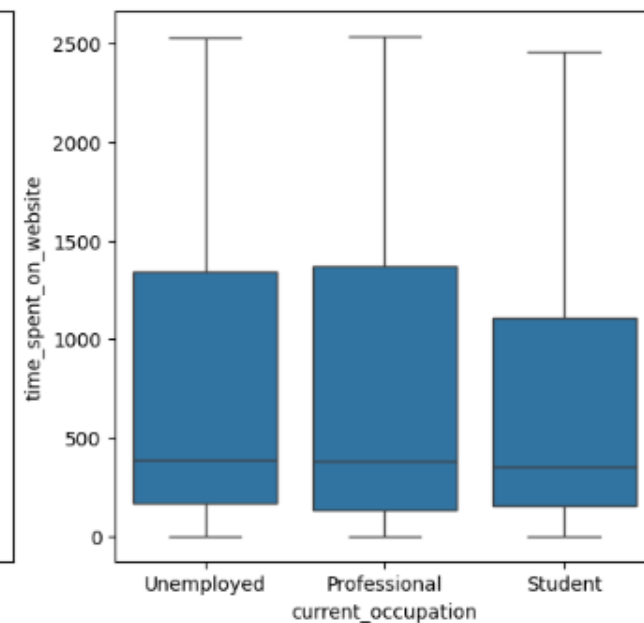
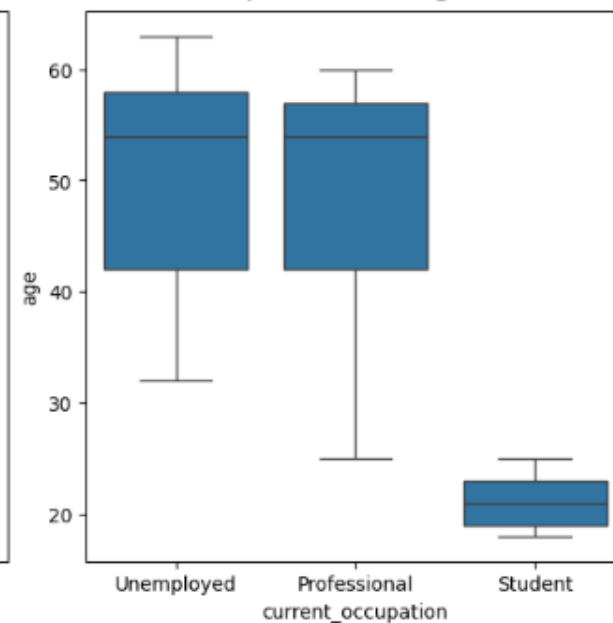
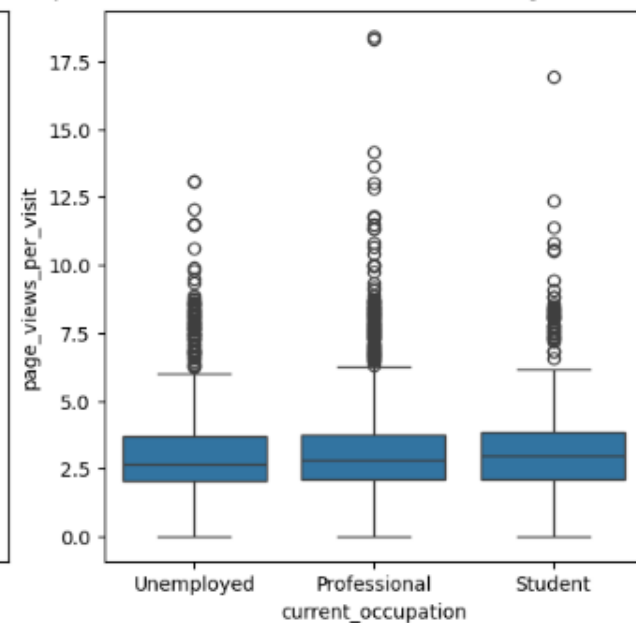
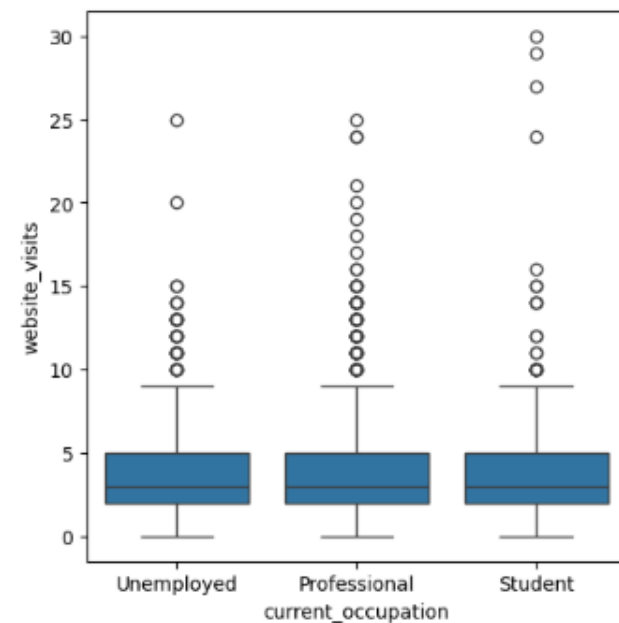
Countplots for categorical variables



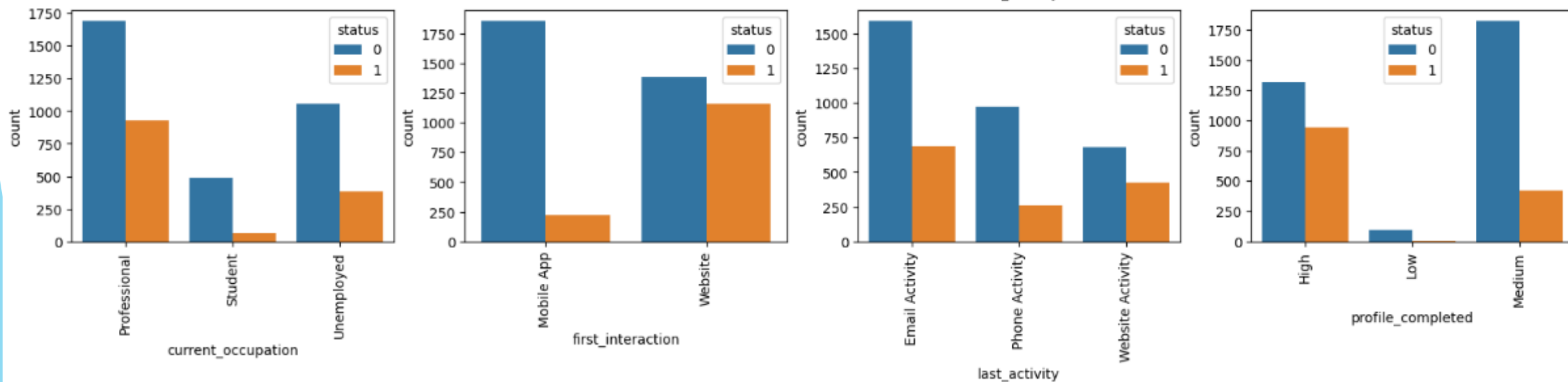
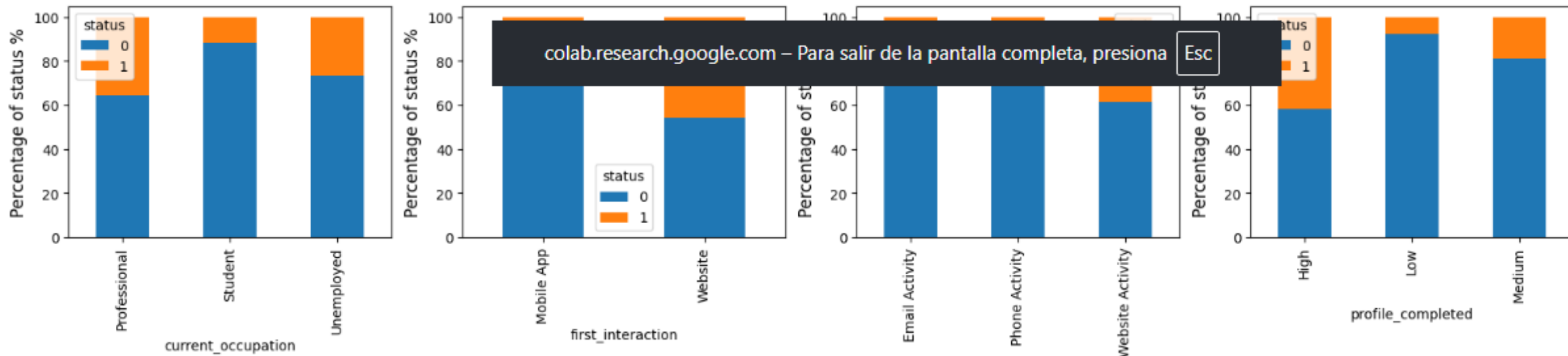
Histograms and boxplots for the different numerical variables



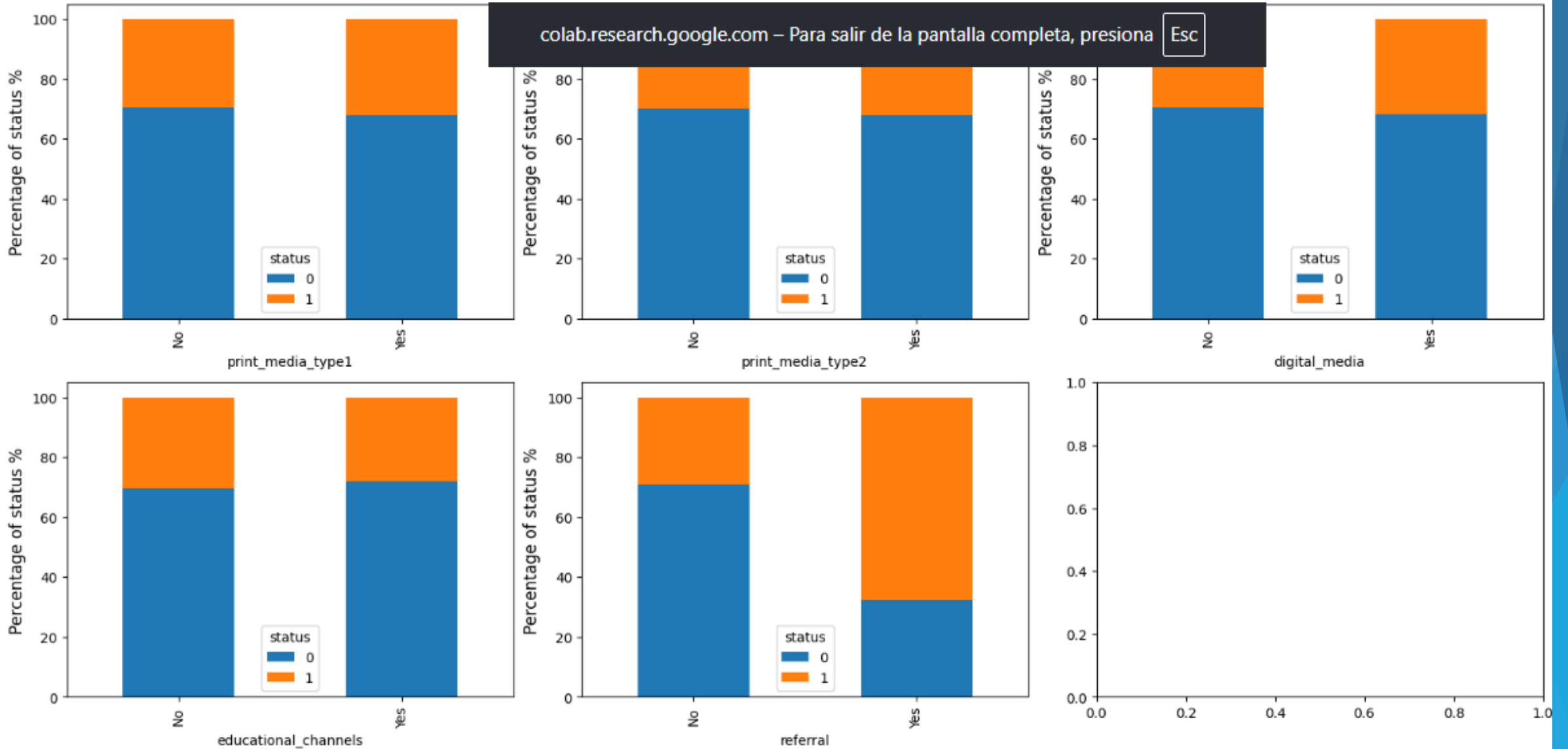
Bivariate boxplots, for numerical variables (y axis) vs current occupation (categorical variable in x axis)



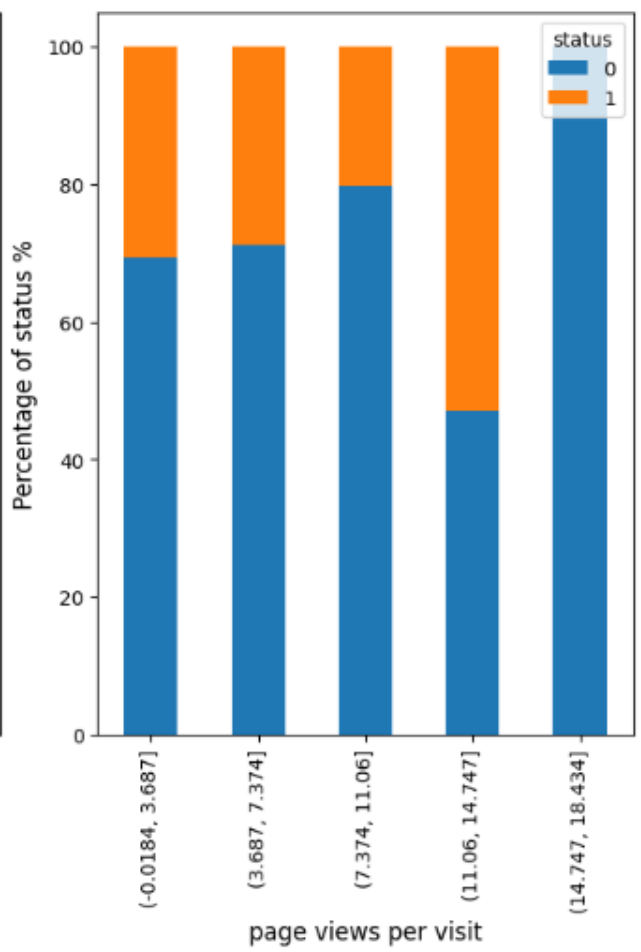
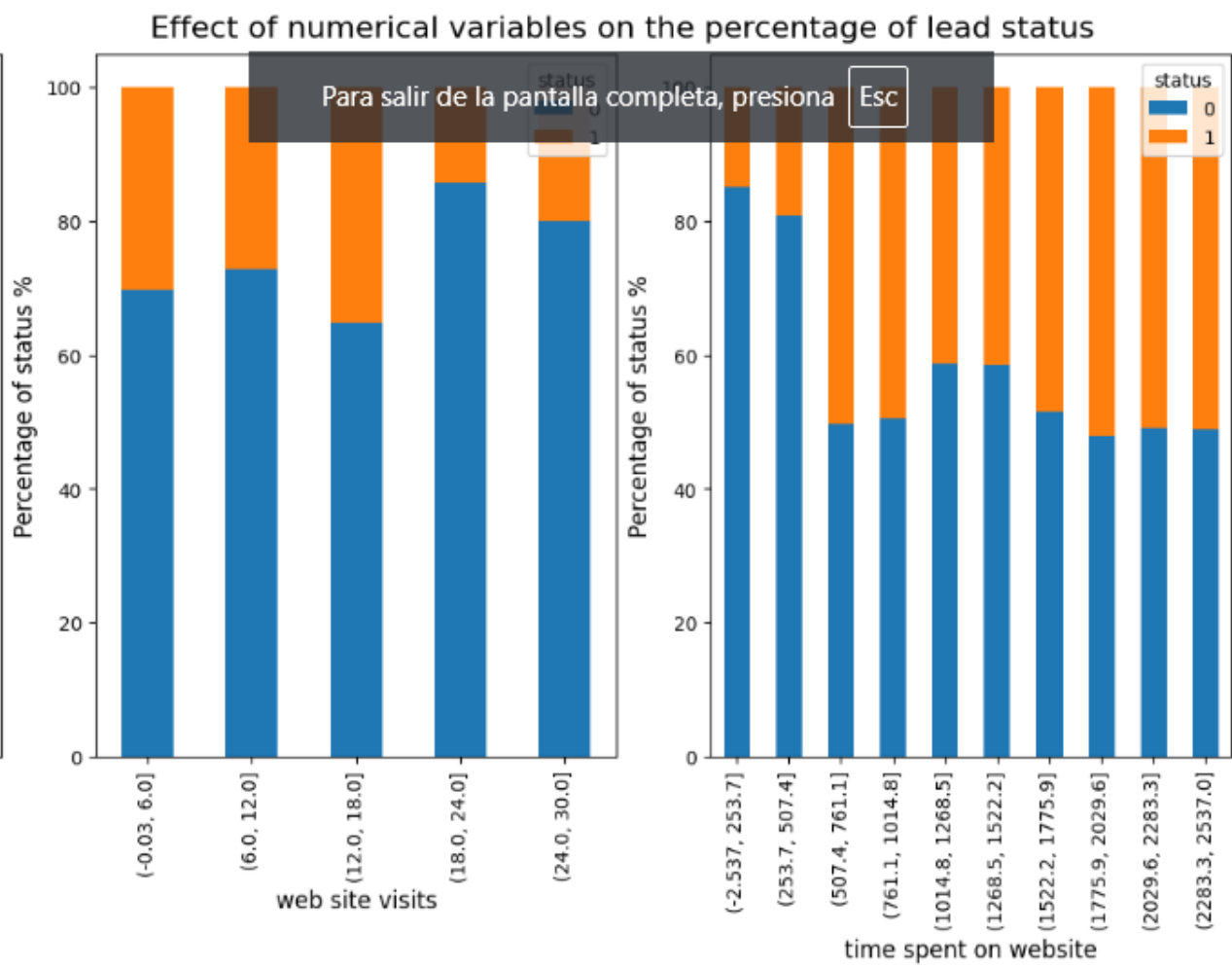
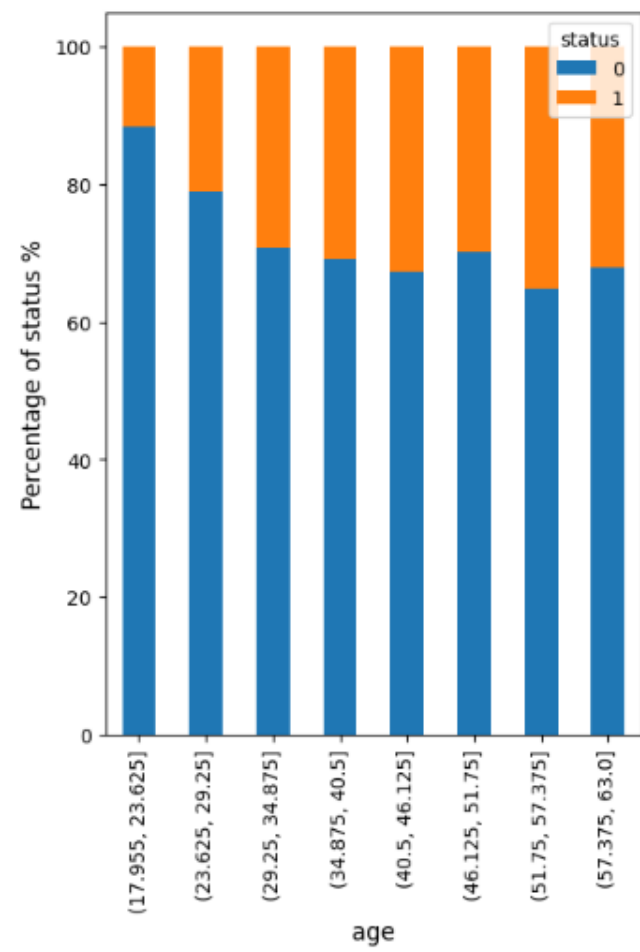
Effect of current\_occupation, first\_interaction, profile\_completed and last\_activity, on lead status



Effect of channels (print media, digital media, educational channels, referrals) on percentage of lead status

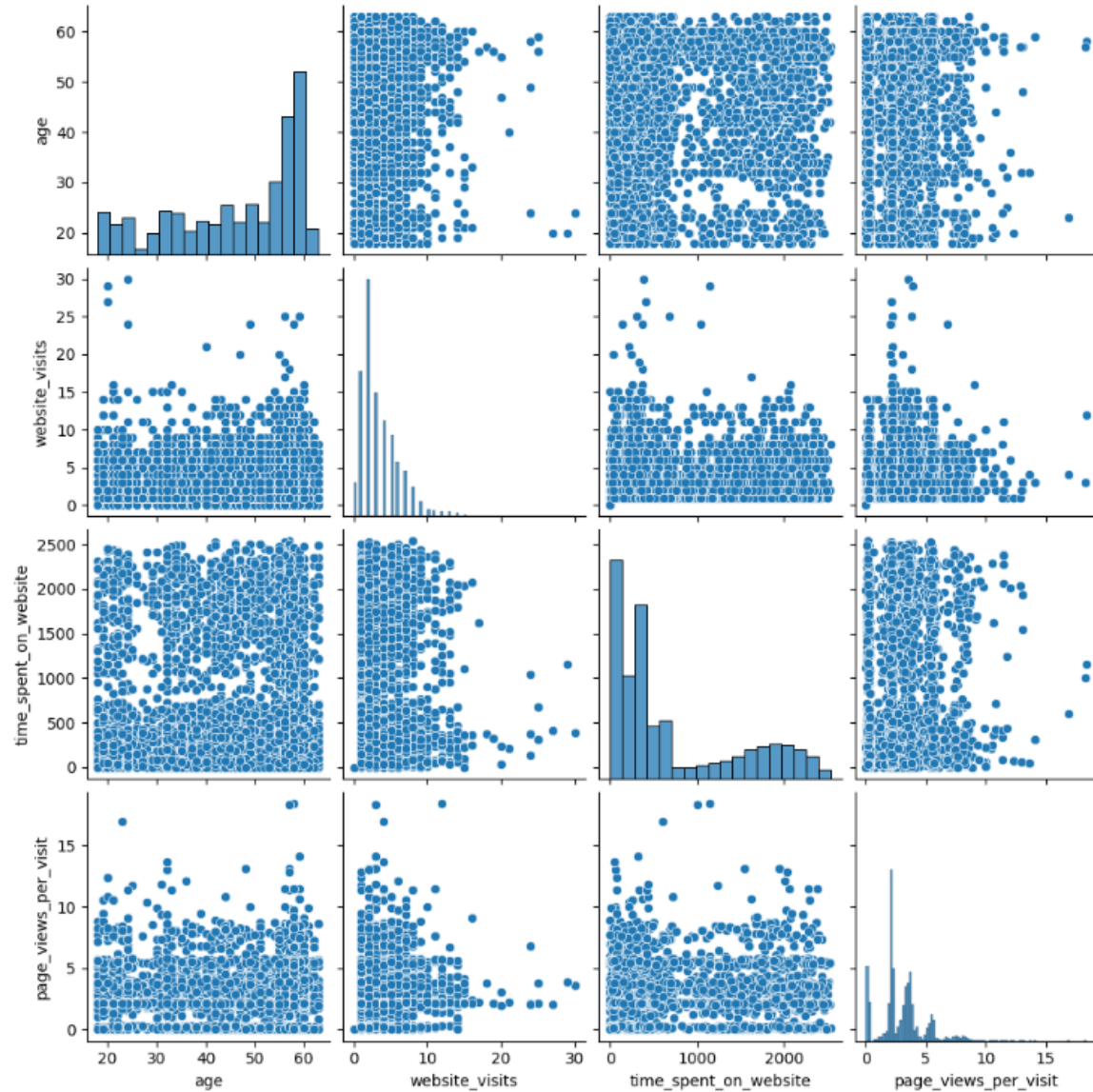


# Effect of numerical variables on percentage of lead status



## Pairplot for numerical variables

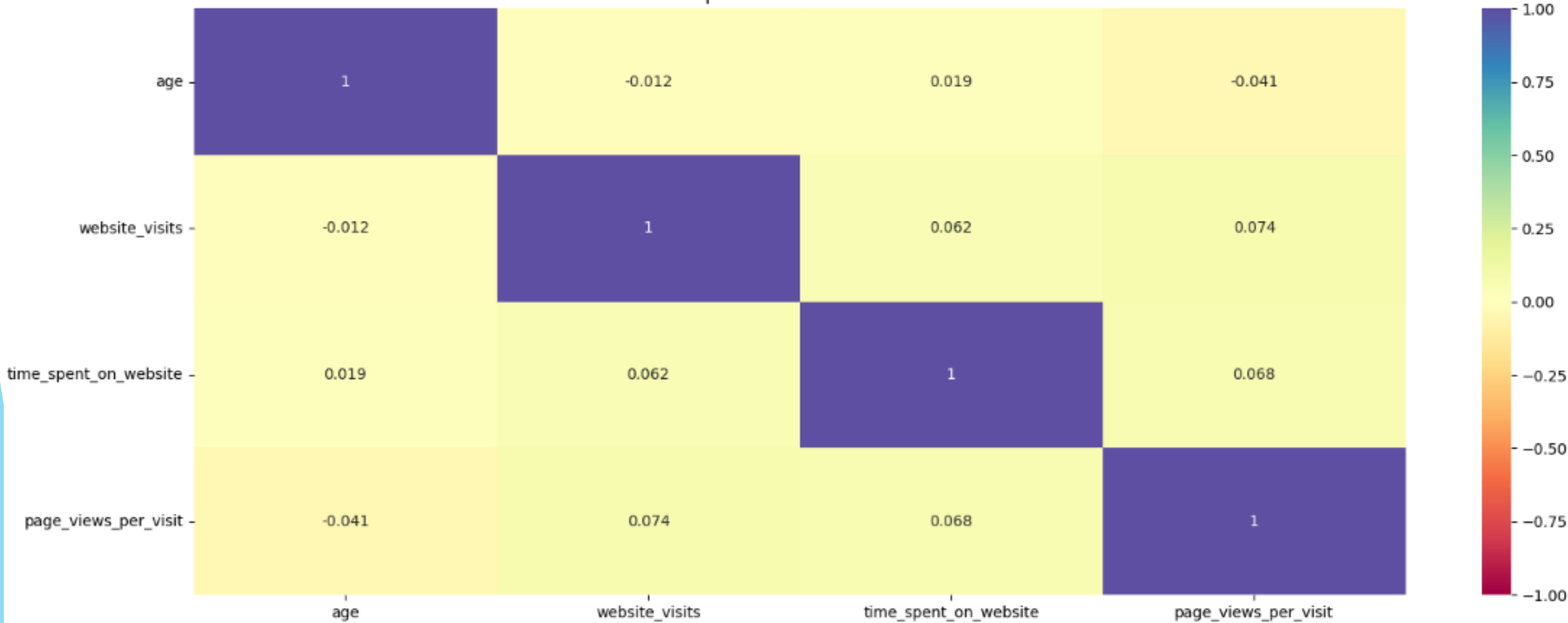
Pairplot for the numerical variables



Pairplot for the numerical variables, segmented by current occupation

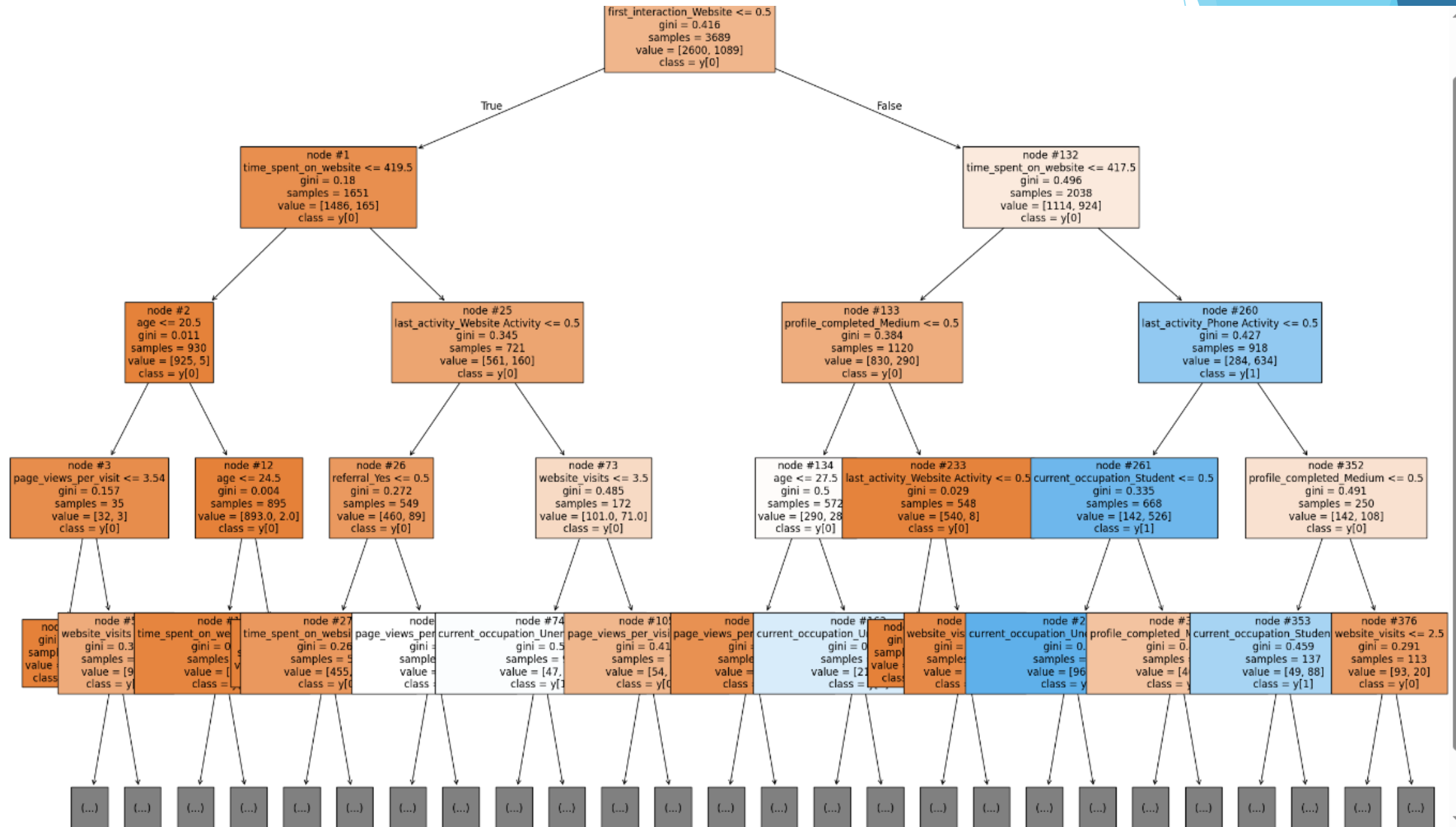


Heatmap for the numerical variables





# Decision tree



# Tuning of random forest

## Tuning of the random forest:

## Metrics of tuned random forest on training data:

## Metrics of tuned random forest on test data:

```
✓ 0s
# Tuned Random Forest Regressor

Note: we optimize the parameters:\
• max_depth = max number of levels in each decision tree
• max_features = max number of features considered for splitting a node
• n_estimators = number of trees in the forest

Running the code below takes some time.

# Model Performance on the train data:
y_pred_train_rf_tuned = rf_tuned_regressor.predict(x_train)
metrics_score(y_train, y_pred_train_rf_tuned)
```

```
rf_tuned = RandomForestClassifier(random_state = 1)

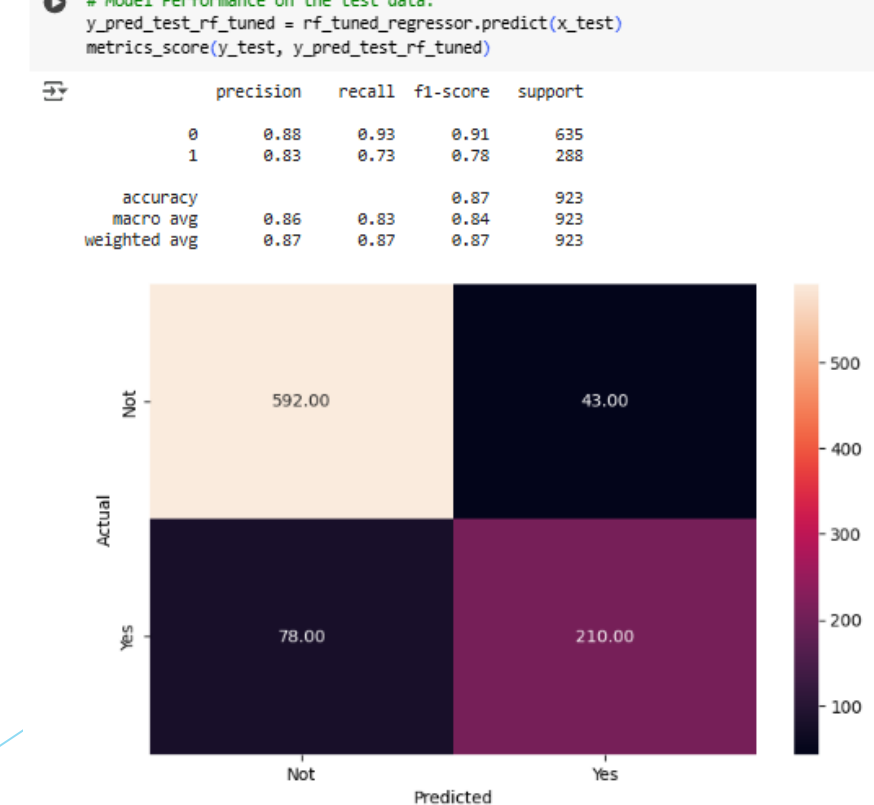
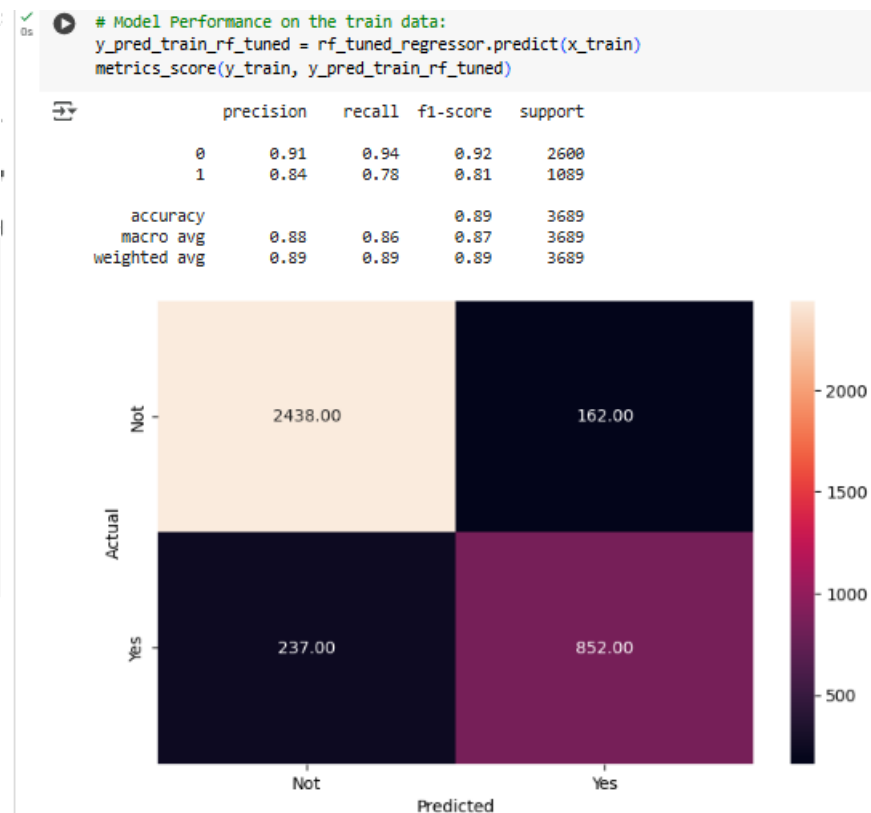
# Grid of parameters to choose from
rf_parameters = {"n_estimators": [100, 110, 120],
                 "max_depth": [5, 7, 15, 20, 30, 50],
                 "max_features": [0.8, 1, 2]}

# Run the grid search
rf_grid_obj = GridSearchCV(rf_tuned, rf_parameters, scoring = 'neg_mean_squared_error', cv = 5)
rf_grid_obj = rf_grid_obj.fit(x_train, y_train)

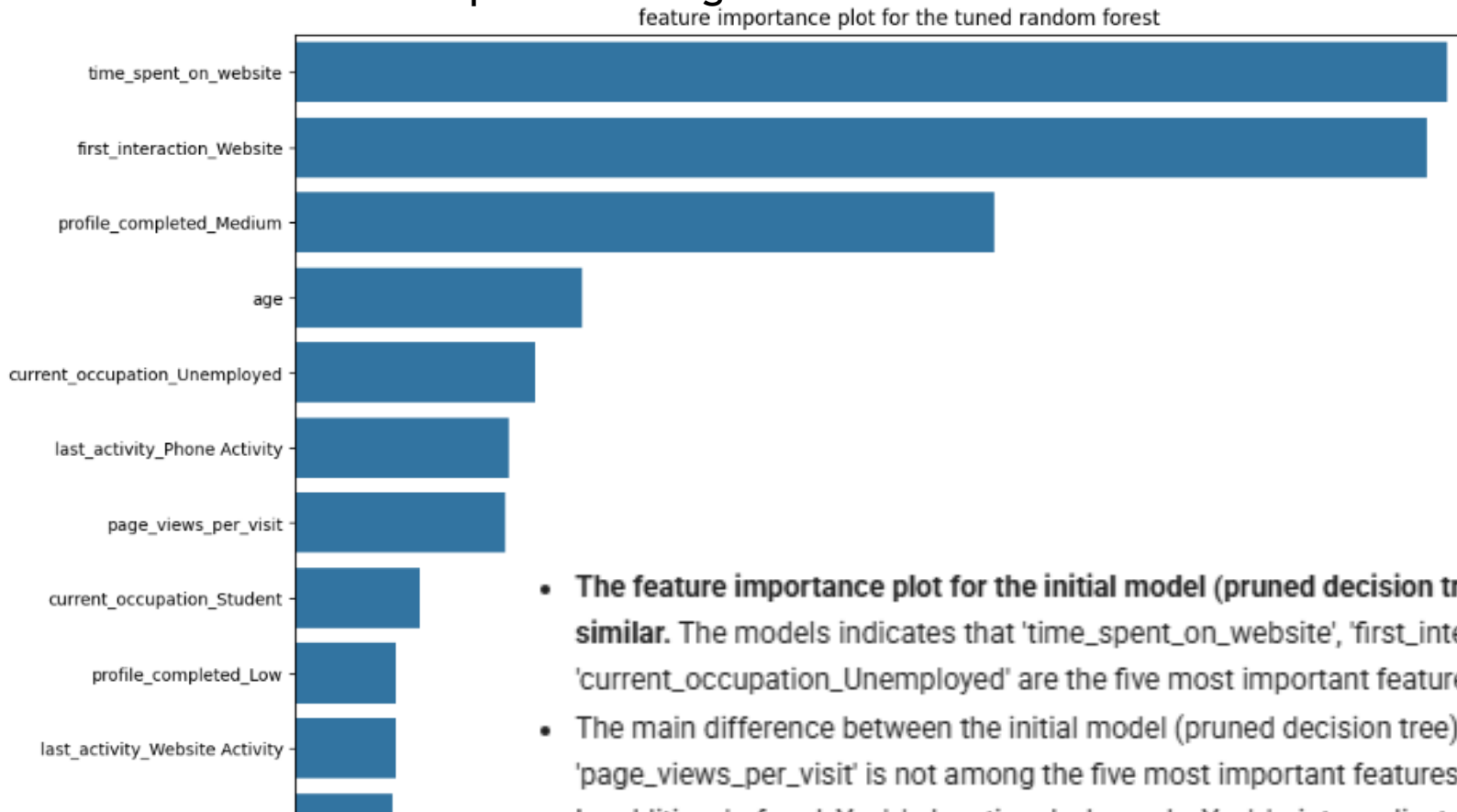
# Set the rf_tuned_regressor to the best combination of parameters
rf_tuned_regressor = rf_grid_obj.best_estimator_

rf_tuned_regressor.fit(x_train, y_train)
```

```
RandomForestClassifier
RandomForestClassifier(max_depth=7, max_features=0.8, random_state=1)
```



## ‘Feature importance’ figure for the tuned Random Forest model



- The feature importance plot for the initial model (pruned decision tree) and that for tuned model (tuned random forest) are quite **similar**. The models indicates that 'time\_spent\_on\_website', 'first\_interaction\_Website', 'profile\_completed\_Medium', 'age', 'current\_occupation\_Unemployed' are the five most important features.
- The main difference between the initial model (pruned decision tree) and tuned model (tuned random forest), is that 'page\_views\_per\_visit' is not among the five most important features in the tuned random forest.
- In addition, 'referral\_Yes', 'educational\_channels\_Yes', 'print\_media\_type1\_Yes', 'print\_media\_type2\_Yes' have a low importance.
- We can say that leads with higher likelihood to become paid customers are characterized by: spend more time on website, their first interaction is through website, have a High level of profile completion; they are older, and their current occupation is professional.
- The channels (print\_media\_type1, print\_media\_type2, digital\_media, educational\_channels, referral) have a low importance in the achievement of paid customers. Also, referral is the channel with higher importance.

## Actionable Insights and Recommendations

### Conclusions

- We have utilized the decision tree, pruned decision tree, random forest and tuned random forest.
- The tuning of random forest resulted in avoidance of overfitting.
- The obtained values of precision and recall compared between train and test are very similar, with a precision of 0.83-0.84 and a recall of 0.73-0.78.
- We have identified the key factors involved in the conversion of leads to paid customers, by means of the 'feature importance' diagram.
- It is possible to improve the tuning by including other parameters in the optimization, and by modifying the parameter values used as possible values in the optimization.

### Recommendations:

- The variables 'time\_spent\_on\_website' and 'first\_interaction\_Website' are the most important ones for identifying which leads are more likely to convert to paid customers.
- Then, to improve the conversion of unpaid to paid customers, the ExtraaLearn company should focus on: i) increasing the advertising that motivates people to look at the website; ii) reviewing the website (the information that is given to the reader) and improving it if possible.
- The variable 'current\_occupation\_Unemployed' has a significantly higher importance than 'current\_occupation\_Student'. Therefore, the ExtraaLearn company should focus its advertising and communication campaigns on unemployed people and Professional people rather than students.