

In [ ]:

```
In [1]: import pandas as pd
import matplotlib.pyplot as plt
import numpy as np
import seaborn as sns
```

```
In [2]: df = pd.read_csv('C:\\\\Users\\\\lenovo\\\\Desktop\\\\courses\\\\Task3\\\\bank-additional.csv')
```

In [3]: df

Out[3]:

	age	job	marital	education	default	housing	loan	contact	month
0	30	blue-collar	married	basic.9y	no	yes	no	cellular	may
1	39	services	single	high.school	no	no	no	telephone	may
2	25	services	married	high.school	no	yes	no	telephone	jun
3	38	services	married	basic.9y	no	unknown	unknown	telephone	jun
4	47	admin.	married	university.degree	no	yes	no	cellular	nov
...	...	...	...	...	...	...	...	...	...
4114	30	admin.	married	basic.6y	no	yes	yes	cellular	jul
4115	39	admin.	married	high.school	no	yes	no	telephone	jul
4116	27	student	single	high.school	no	no	no	cellular	may
4117	58	admin.	married	high.school	no	no	no	cellular	aug
4118	34	management	single	high.school	no	yes	no	cellular	nov

4119 rows × 21 columns



In [4]: df.head()

Out[4]:

	age	job	marital	education	default	housing	loan	contact	month	day_of_
0	30	blue-collar	married	basic.9y	no	yes	no	cellular	may	
1	39	services	single	high.school	no	no	no	telephone	may	
2	25	services	married	high.school	no	yes	no	telephone	jun	
3	38	services	married	basic.9y	no	unknown	unknown	telephone	jun	
4	47	admin.	married	university.degree	no	yes	no	cellular	nov	

5 rows × 21 columns



In [5]: df.columns

```
Out[5]: Index(['age', 'job', 'marital', 'education', 'default', 'housing', 'loan',
   'contact', 'month', 'day_of_week', 'duration', 'campaign', 'pdays',
   'previous', 'poutcome', 'emp.var.rate', 'cons.price.idx',
   'cons.conf.idx', 'euribor3m', 'nr.employed', 'target'],
  dtype='object')
```

```
In [6]: df.dtypes
```

```
Out[6]: age           int64
job            object
marital        object
education      object
default         object
housing         object
loan            object
contact         object
month           object
day_of_week    object
duration        int64
campaign        int64
pdays           int64
previous        int64
poutcome        object
emp.var.rate    float64
cons.price.idx float64
cons.conf.idx  float64
euribor3m      float64
nr.employed    float64
target          object
dtype: object
```

```
In [7]: missing_values = df.isnull().sum()
print("Valeurs manquantes dans les colonnes :")
print(missing_values)
```

Valeurs manquantes dans les colonnes :

age	0
job	0
marital	0
education	0
default	0
housing	0
loan	0
contact	0
month	0
day_of_week	0
duration	0
campaign	0
pdays	0
previous	0
poutcome	0
emp.var.rate	0
cons.price.idx	0
cons.conf.idx	0
euribor3m	0
nr.employed	0
target	0

dtype: int64

```
In [8]: df.duplicated().sum()
```

```
Out[8]: 0
```

```
In [9]: df.isna().sum()
```

```
Out[9]: age          0  
        job          0  
       marital      0  
     education      0  
    default        0  
   housing         0  
    loan          0  
  contact         0  
month          0  
day_of_week     0  
duration        0  
campaign        0  
pdays           0  
previous        0  
poutcome        0  
emp.var.rate    0  
cons.price.idx  0  
cons.conf.idx   0  
euribor3m       0  
nr.employed     0  
target          0  
dtype: int64
```

```
In [10]: categorical_columns = df.select_dtypes(include=['object']).columns  
  
numeric_columns = df.select_dtypes(include=['int64', 'float64']).columns  
  
df_categorical = df[categorical_columns]  
df_numeric = df[numeric_columns]  
  
print("Dataframe des colonnes catégorielles :")  
print(df_categorical.head())  
  
print("\nDataframe des colonnes numériques :")  
print(df_numeric.head())
```

Dataframe des colonnes catégorielles :

	job	marital	education	default	housing	loan	\
0	blue-collar	married	basic.9y	no	yes	no	
1	services	single	high.school	no	no	no	
2	services	married	high.school	no	yes	no	
3	services	married	basic.9y	no	unknown	unknown	
4	admin.	married	university.degree	no	yes	no	

	contact	month	day_of_week	poutcome	target
0	cellular	may	fri	nonexistent	no
1	telephone	may	fri	nonexistent	no
2	telephone	jun	wed	nonexistent	no
3	telephone	jun	fri	nonexistent	no
4	cellular	nov	mon	nonexistent	no

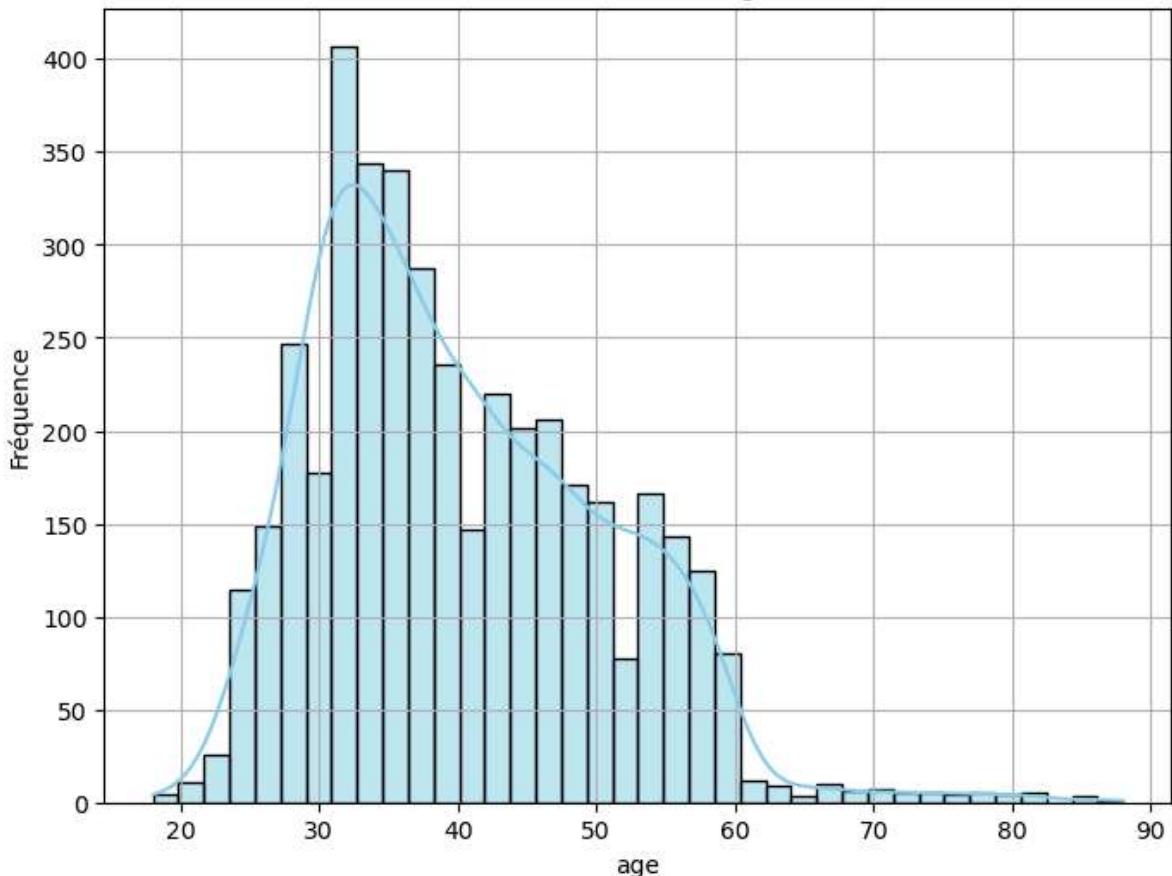
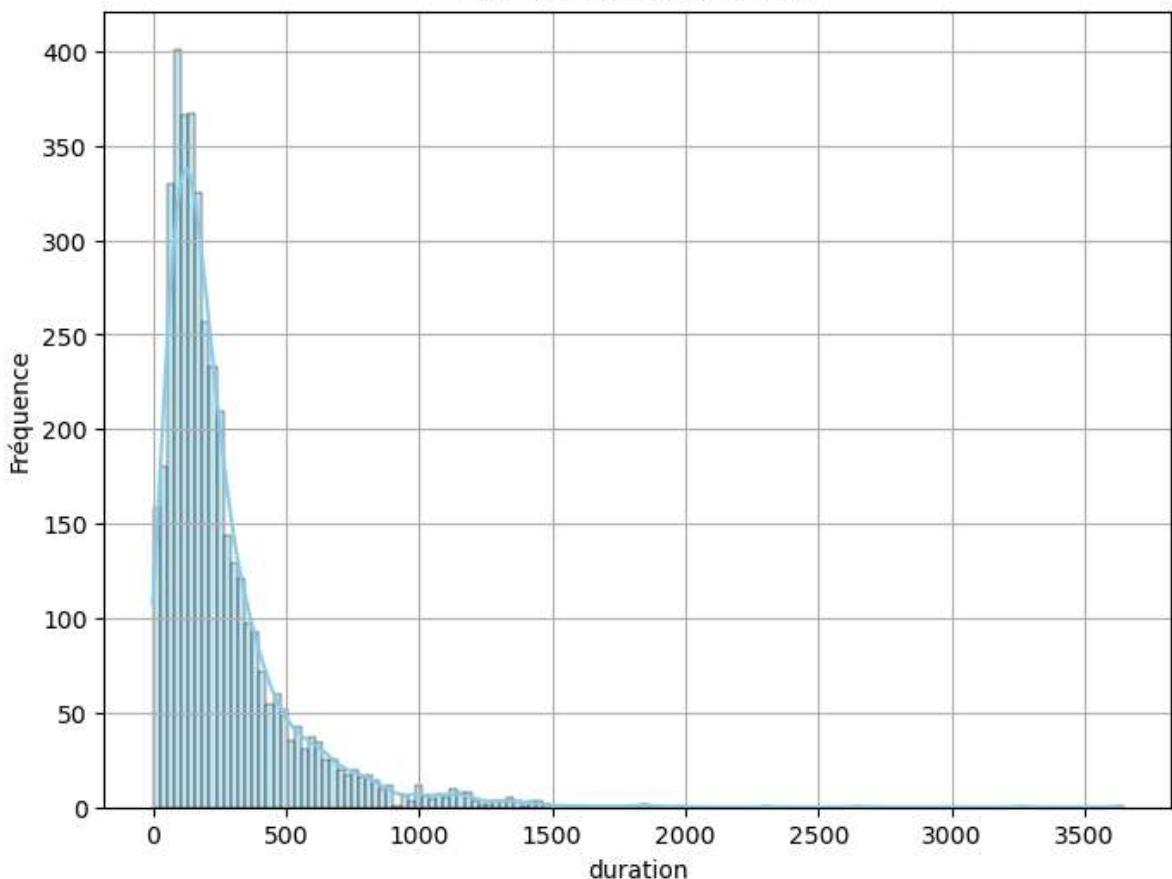
Dataframe des colonnes numériques :

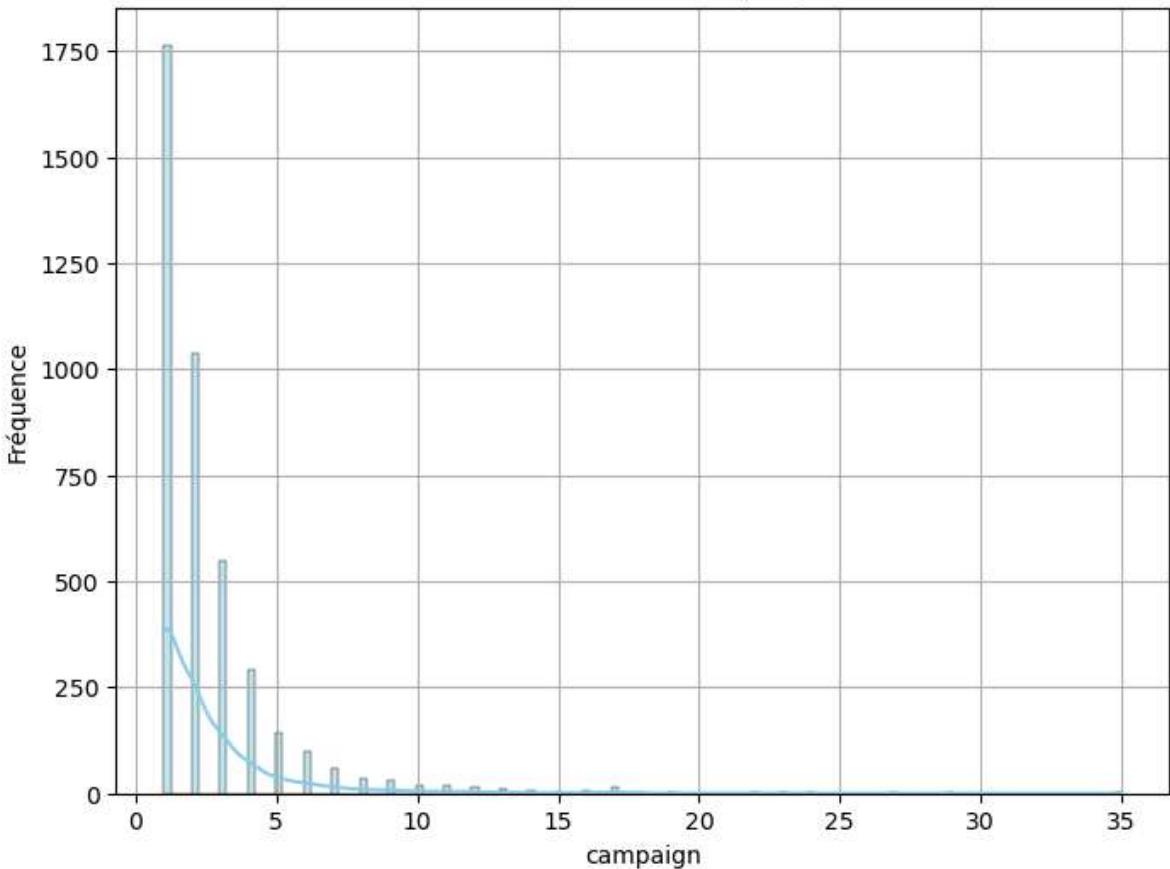
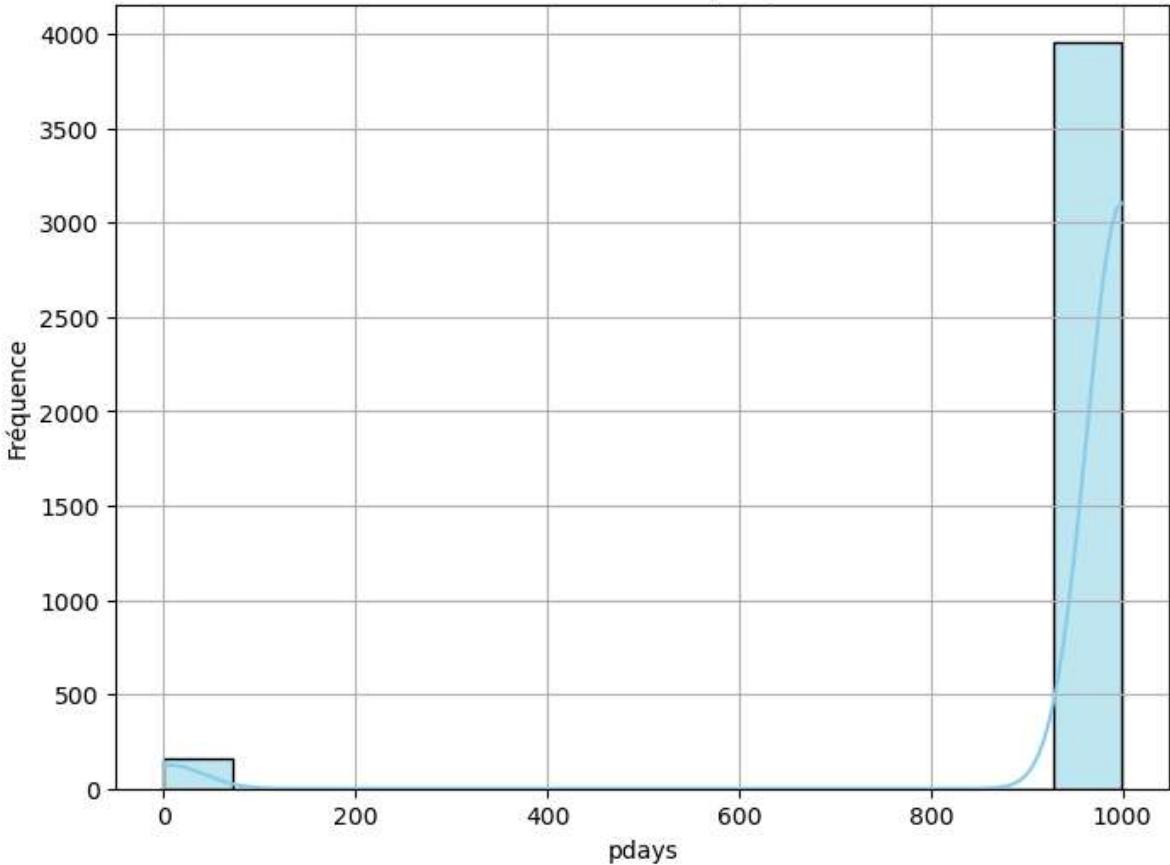
	age	duration	campaign	pdays	previous	emp.var.rate	cons.price.idx	\
0	30	487	2	999	0	-1.8	92.893	
1	39	346	4	999	0	1.1	93.994	
2	25	227	1	999	0	1.4	94.465	
3	38	17	3	999	0	1.4	94.465	
4	47	58	1	999	0	-0.1	93.200	

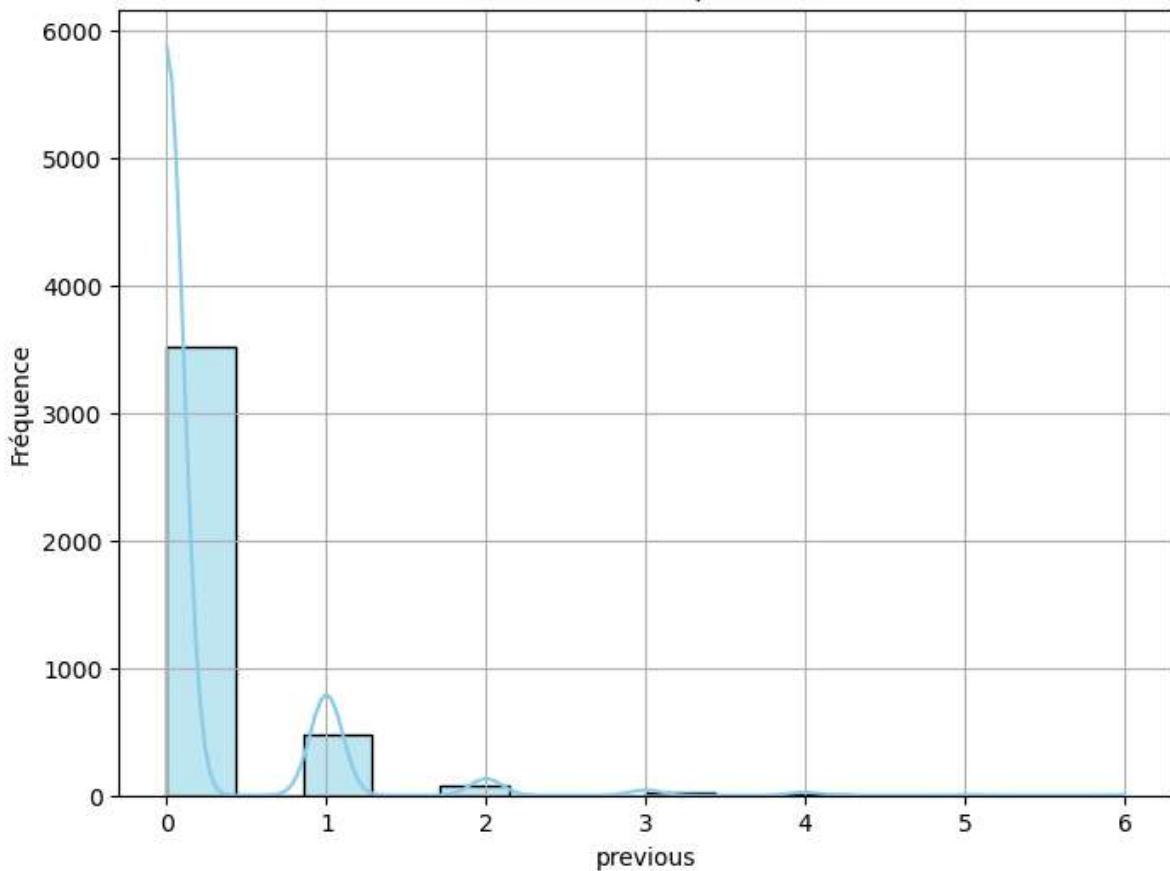
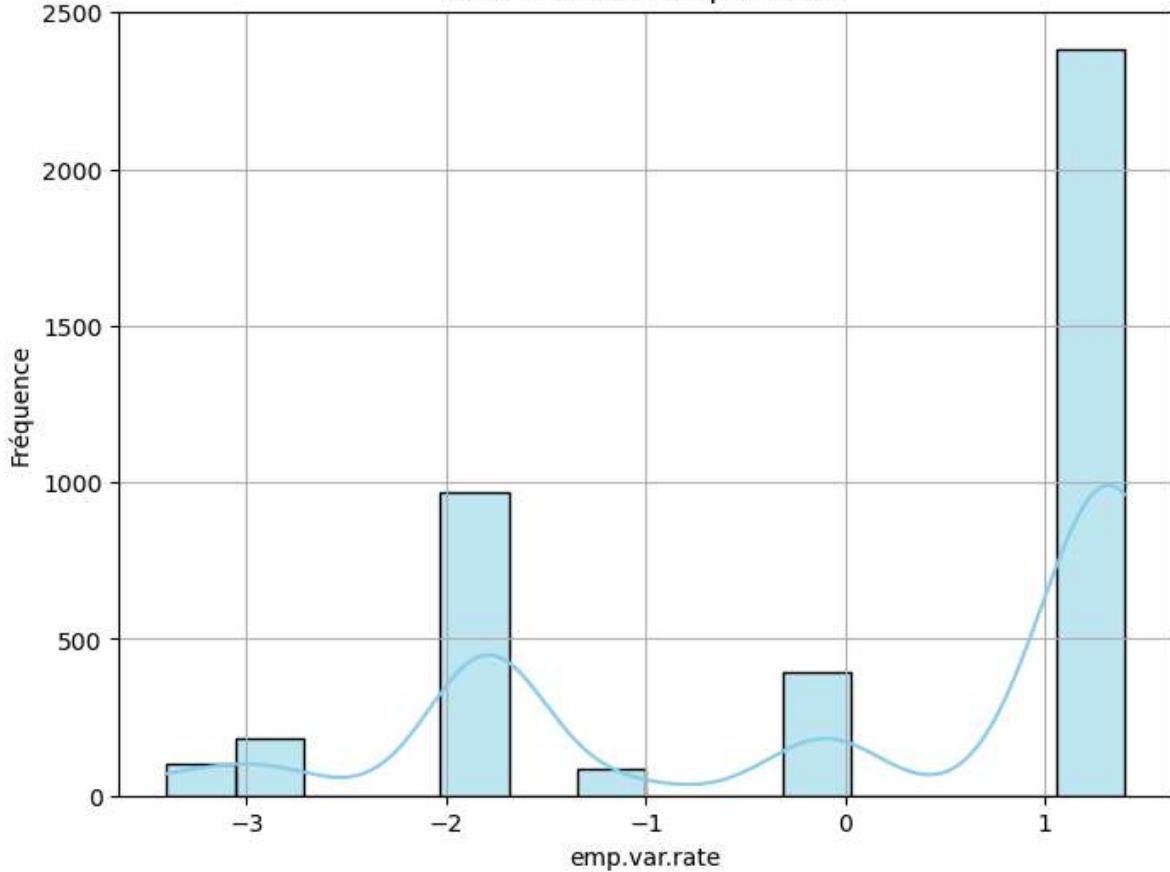
	cons.conf.idx	euribor3m	nr.employed
0	-46.2	1.313	5099.1
1	-36.4	4.855	5191.0
2	-41.8	4.962	5228.1
3	-41.8	4.959	5228.1
4	-42.0	4.191	5195.8

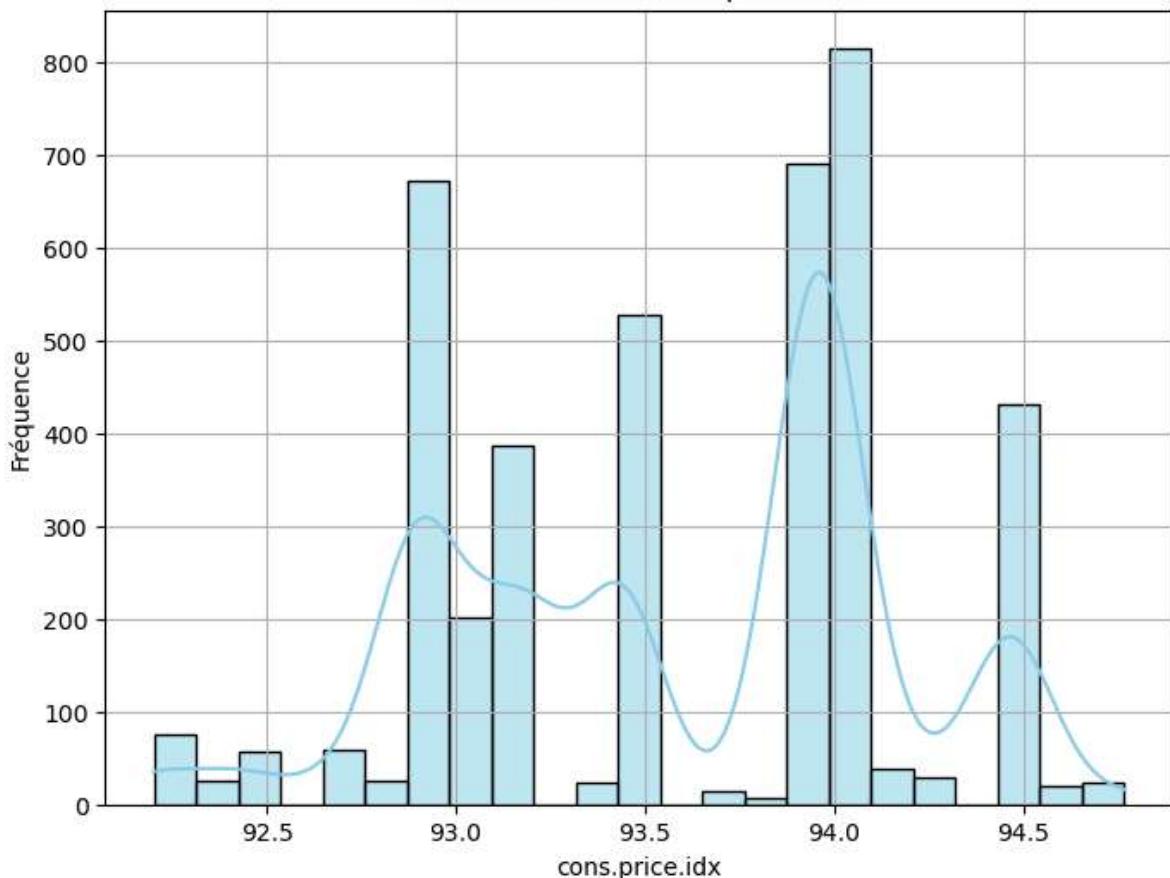
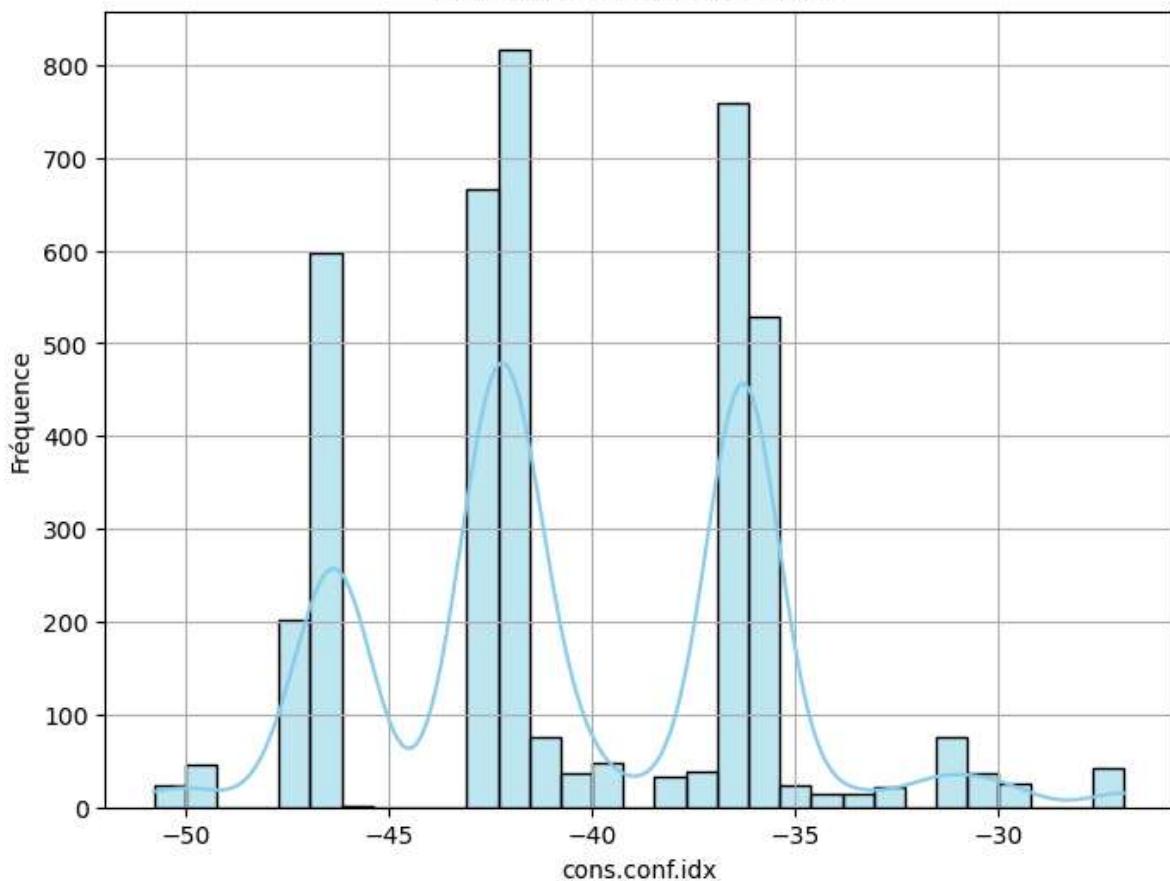
```
In [11]: significant_numeric_columns = ['age', 'duration', 'campaign', 'pdays', 'previous',
                                         'emp.var.rate', 'cons.price.idx', 'cons.conf.idx',
                                         'euribor3m', 'nr.employed']

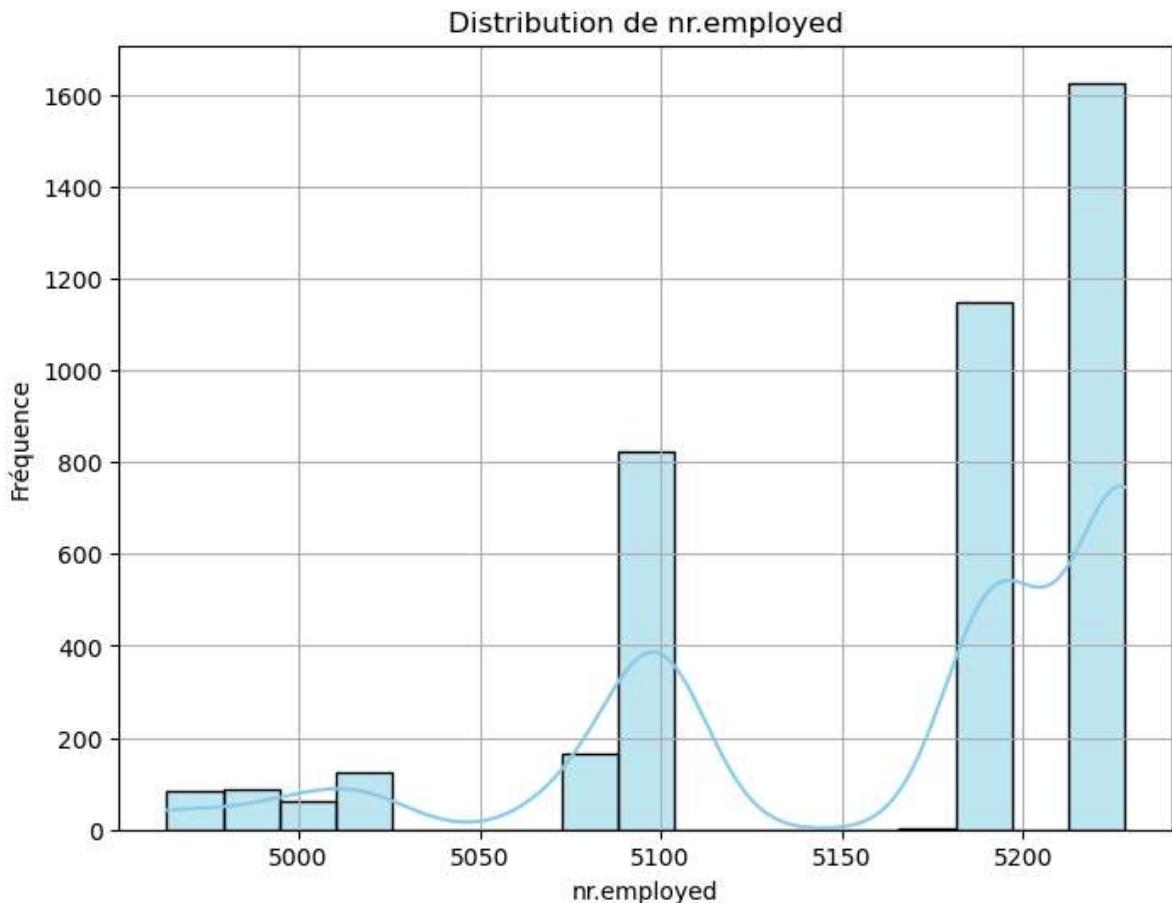
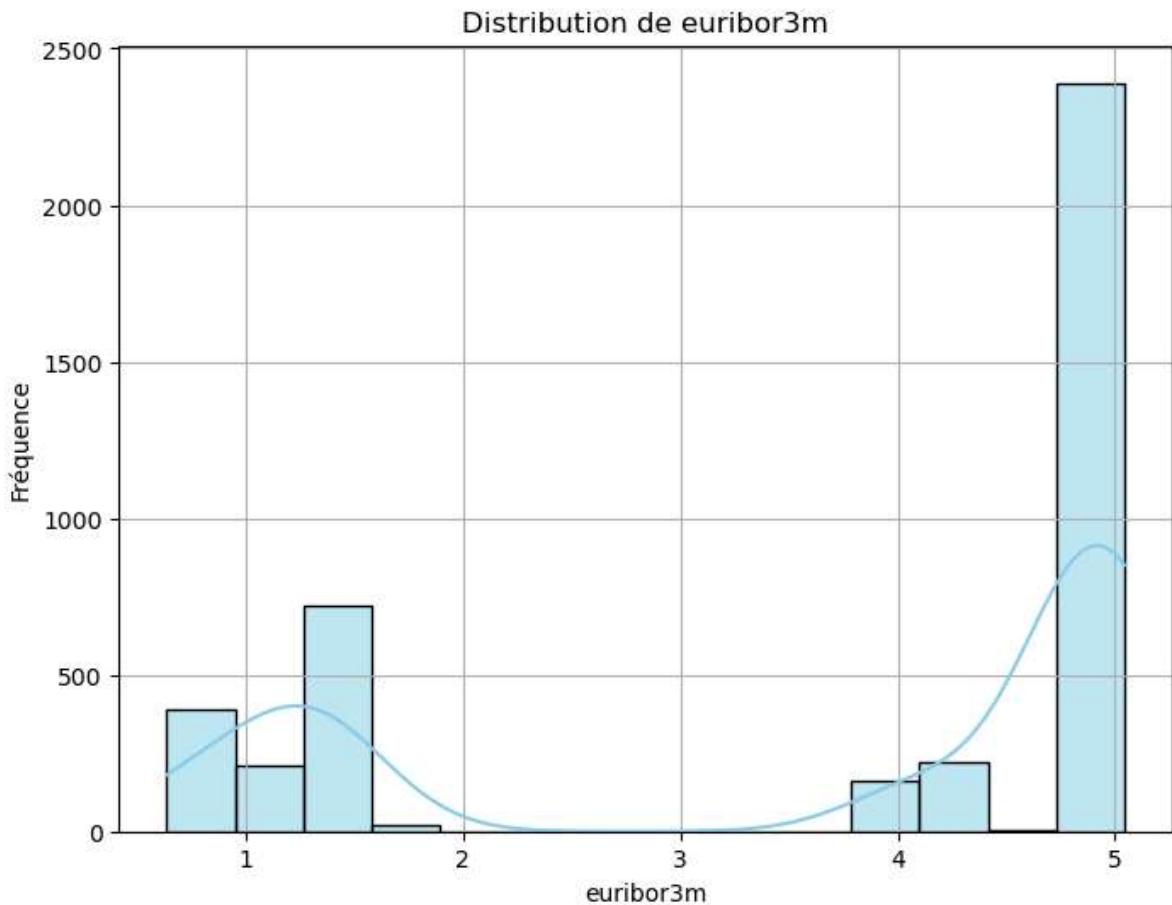
for column in significant_numeric_columns:
    plt.figure(figsize=(8, 6))
    sns.histplot(df_numeric[column], kde=True, color='skyblue', edgecolor='black')
    plt.title(f'Distribution de {column}')
    plt.xlabel(column)
    plt.ylabel('Fréquence')
    plt.grid(True)
    plt.show()
```

**Distribution de age****Distribution de duration**

**Distribution de campaign****Distribution de pdays**

**Distribution de previous****Distribution de emp.var.rate**

**Distribution de cons.price.idx****Distribution de cons.conf.idx**



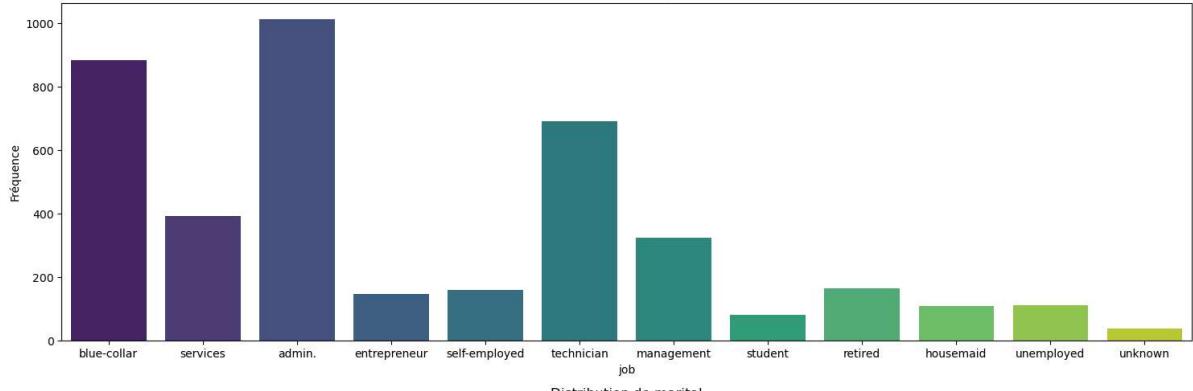
```
In [12]: num_plots = len(df_categorical.columns)
plt.figure(figsize=(15, 5*num_plots))

for i, column in enumerate(df_categorical.columns, 1):
    plt.subplot(num_plots, 1, i)
    sns.countplot(data=df_categorical, x=column, palette='viridis')
```

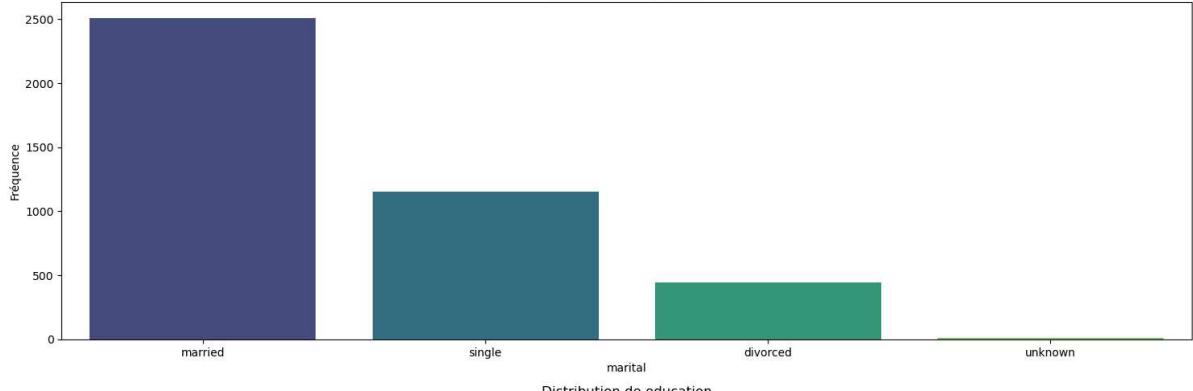
```
plt.title(f'Distribution de {column}')
plt.xlabel(column)
plt.ylabel('Fréquence')

plt.tight_layout()
plt.show()
```

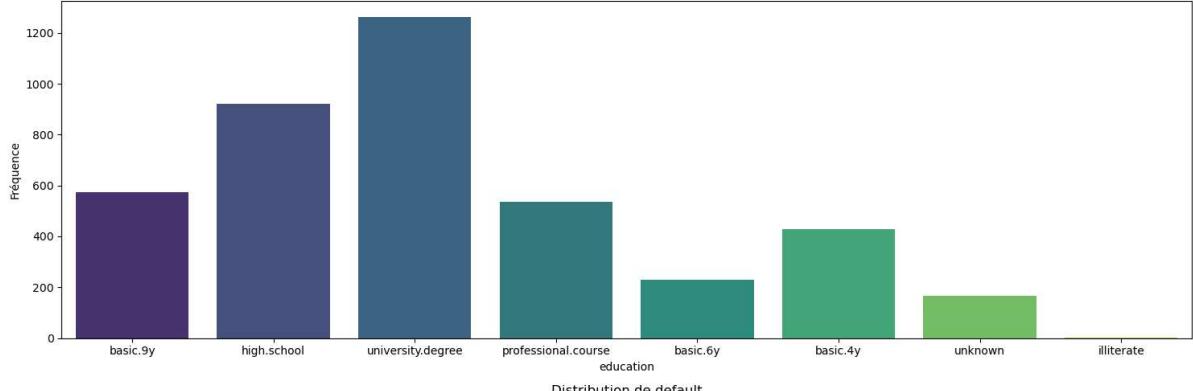
Distribution de job



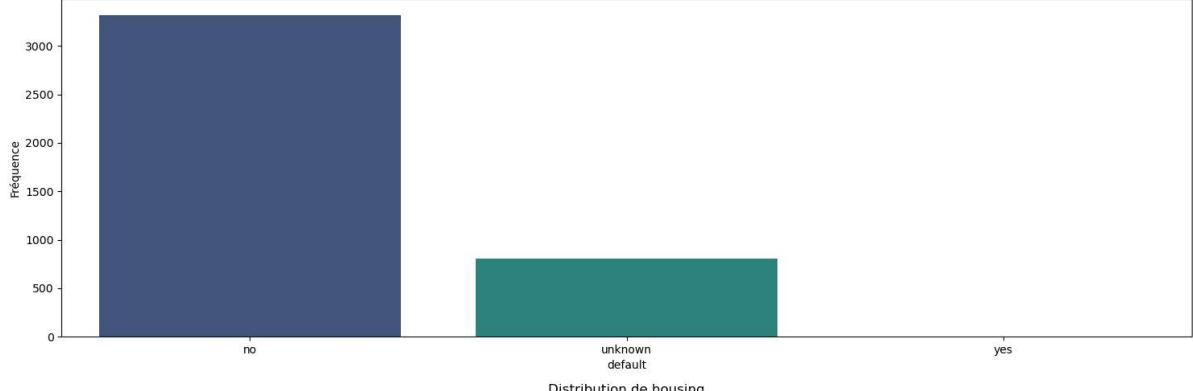
Distribution de marital



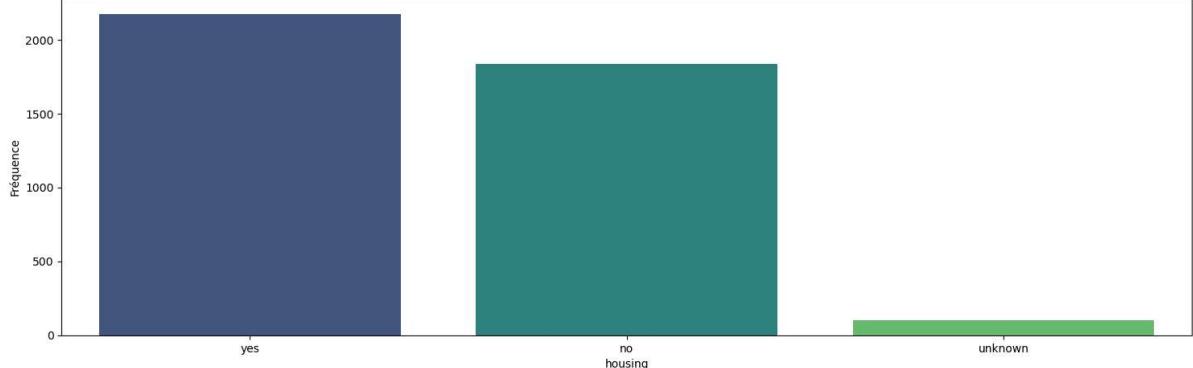
Distribution de education



Distribution de default

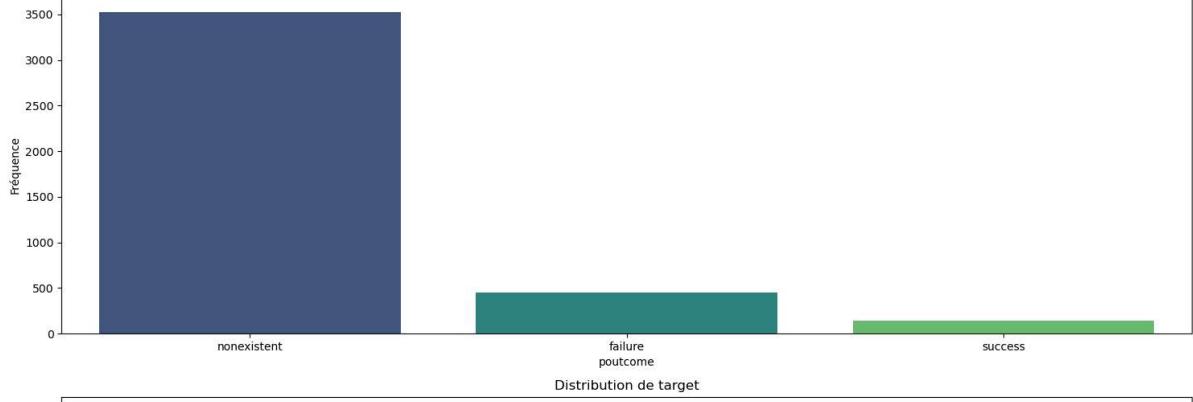
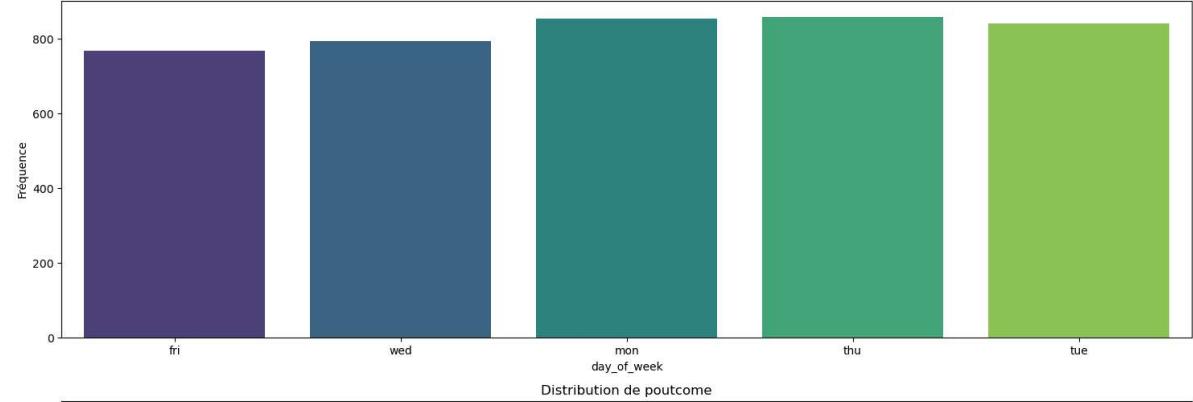
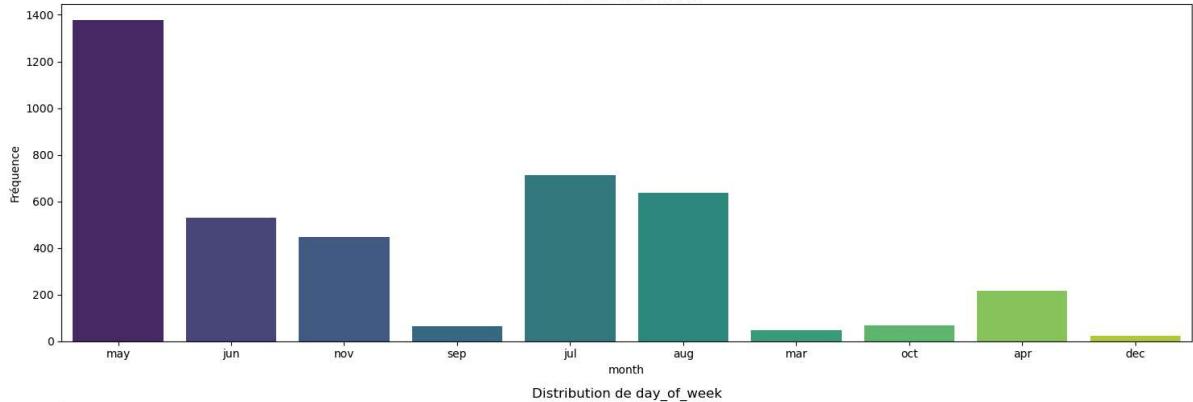
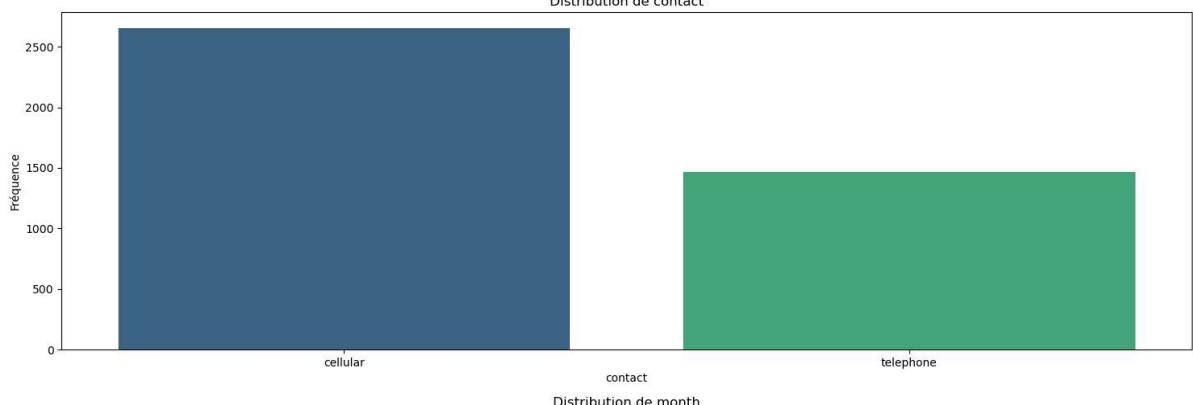
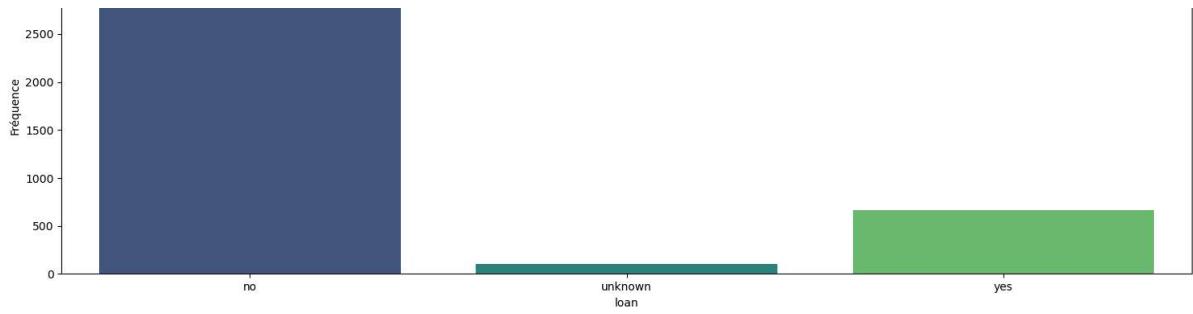


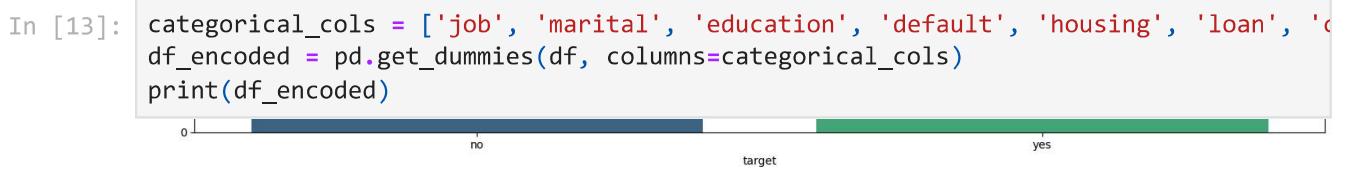
Distribution de housing



Distribution de loan







	age	duration	campaign	pdays	previous	emp.var.rate	cons.price.idx	\
0	30	487	2	999	0	-1.8	92.893	
1	39	346	4	999	0	1.1	93.994	
2	25	227	1	999	0	1.4	94.465	
3	38	17	3	999	0	1.4	94.465	
4	47	58	1	999	0	-0.1	93.200	
...	...	...	...	...	...	...	...	
4114	30	53	1	999	0	1.4	93.918	
4115	39	219	1	999	0	1.4	93.918	
4116	27	64	2	999	1	-1.8	92.893	
4117	58	528	1	999	0	1.4	93.444	
4118	34	175	1	999	0	-0.1	93.200	
	cons.conf.idx	euribor3m	nr.employed	...	month_oct	month_sep	\	
0	-46.2	1.313	5099.1	...	False	False		
1	-36.4	4.855	5191.0	...	False	False		
2	-41.8	4.962	5228.1	...	False	False		
3	-41.8	4.959	5228.1	...	False	False		
4	-42.0	4.191	5195.8	...	False	False		
...	...	...	...	...	...	...		
4114	-42.7	4.958	5228.1	...	False	False		
4115	-42.7	4.959	5228.1	...	False	False		
4116	-46.2	1.354	5099.1	...	False	False		
4117	-36.1	4.966	5228.1	...	False	False		
4118	-42.0	4.120	5195.8	...	False	False		
	day_of_week_fri	day_of_week_mon	day_of_week_thu	day_of_week_tue	\			
0	True	False	False	False				
1	True	False	False	False				
2	False	False	False	False				
3	True	False	False	False				
4	False	True	False	False				
...	...	...	...	...				
4114	False	False	True	False				
4115	True	False	False	False				
4116	False	True	False	False				
4117	True	False	False	False				
4118	False	False	False	False				
	day_of_week_wed	poutcome_failure	poutcome_nonexistent	\				
0	False	False	True					
1	False	False	True					
2	True	False	True					
3	False	False	True					
4	False	False	True					
...	...	...	...					
4114	False	False	True					
4115	False	False	True					
4116	False	True	False					
4117	False	False	True					
4118	True	False	True					
	poutcome_success							
0	False							
1	False							
2	False							
3	False							
4	False							
...	...							
4114	False							
4115	False							
4116	False							
4117	False							
4118	False							

[4119 rows x 64 columns]

```
In [14]: df_encoded['target'].value_counts()
```

```
Out[14]: target
no      3668
yes     451
Name: count, dtype: int64
```

```
In [27]: X = df_encoded.drop('target', axis=1)
Y = df_encoded['target']
print(x.shape)
print(y.shape)
print(type(x))
print(type(y))
```

```
(4119, 63)
(4119,)
<class 'pandas.core.frame.DataFrame'>
<class 'pandas.core.series.Series'>
```

```
In [28]: from sklearn.model_selection import train_test_split
```

```
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.25, random_st
print("Taille de l'ensemble d'entraînement:", len(X_train))
print("Taille de l'ensemble de test:", len(X_test))
```

```
Taille de l'ensemble d'entraînement: 3089
```

```
Taille de l'ensemble de test: 1030
```

```
In [29]: from sklearn.model_selection import train_test_split
from sklearn.tree import DecisionTreeClassifier
from sklearn.metrics import accuracy_score
```

```
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.25, random_st
clf = DecisionTreeClassifier()
clf.fit(X_train, y_train)
y_pred = clf.predict(X_test)
```

```
In [30]: accuracy = accuracy_score(y_test, y_pred)
print("Accuracy:", accuracy)
```

```
Accuracy: 0.8883495145631068
```

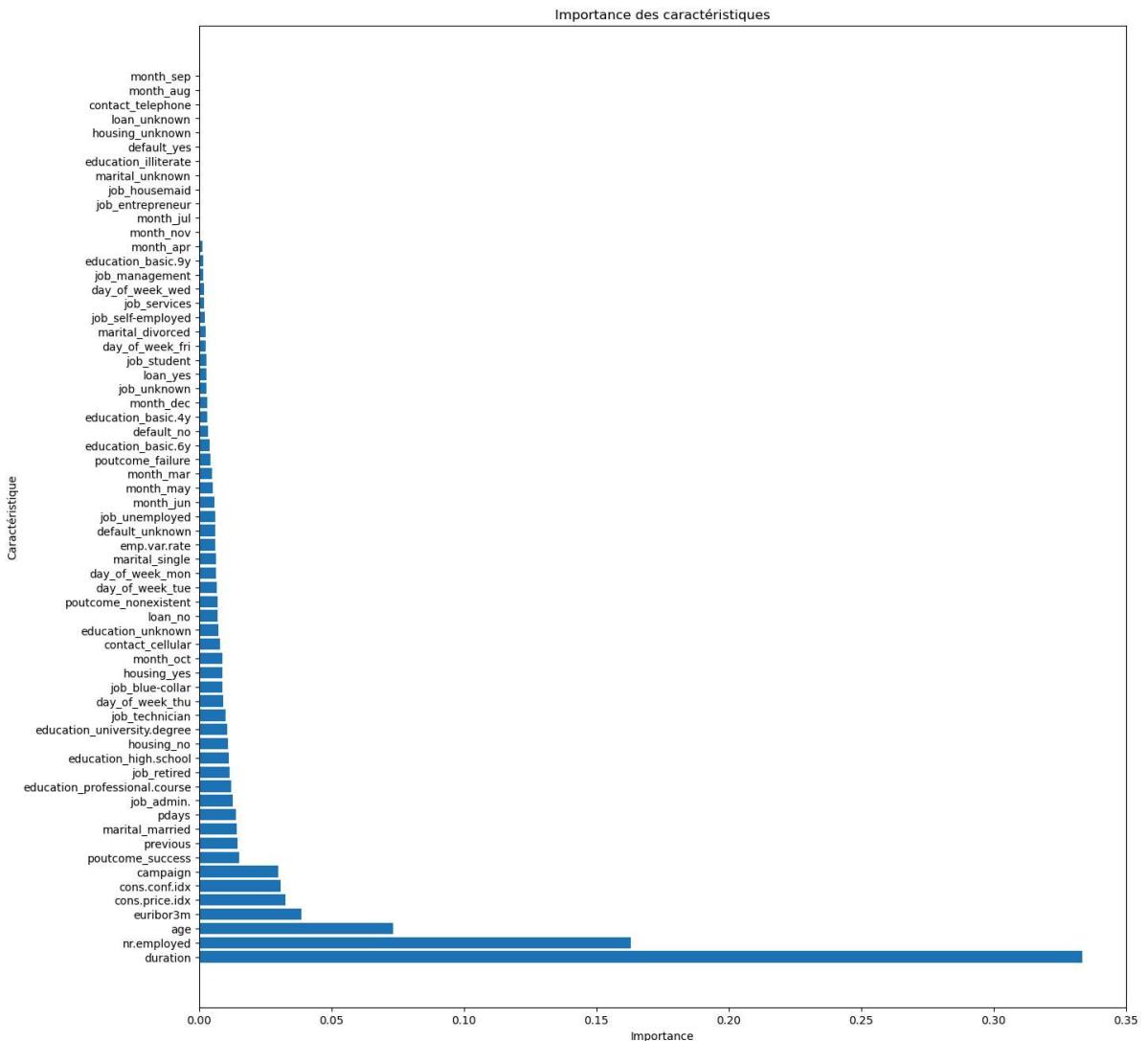
```
In [32]: feature_importances_ = clf.feature_importances_
feature_importance_dict = dict(zip(X.columns, feature_importances_))
sorted_features = sorted(feature_importance_dict.items(), key=lambda x: x[1], reverse=True)

print("Caractéristiques importantes :")
for feature, importance in sorted_features:
    print(f"{feature}: {importance}")
```

Caractéristiques importantes :

duration: 0.33350863560236477  
nr.employed: 0.16315978322217245  
age: 0.07332839905625593  
euribor3m: 0.03867597955363051  
cons.price.idx: 0.03249293067723582  
cons.conf.idx: 0.030890977085346515  
campaign: 0.02980647689928591  
poutcome\_success: 0.01502593277824243  
previous: 0.014644049611674412  
marital\_married: 0.014174976872401178  
pdays: 0.01396524260608575  
job\_admin.: 0.012838091360922113  
education\_professional.course: 0.012102253882318531  
job\_retired: 0.011570746019616454  
education\_high.school: 0.011194535683212329  
housing\_no: 0.010960033664227714  
education\_university.degree: 0.010798341307117376  
job\_technician: 0.010188910767095456  
day\_of\_week\_thu: 0.009198599682476832  
job\_blue-collar: 0.008975934309572824  
housing\_yes: 0.008829274178329329  
month\_oct: 0.008809901145028035  
contact\_cellular: 0.007909417006674451  
education\_unknown: 0.007403281325364322  
loan\_no: 0.007040363819696394  
poutcome\_nonexistent: 0.006928306994886804  
day\_of\_week\_tue: 0.006850931341941681  
day\_of\_week\_mon: 0.006320439581576696  
marital\_single: 0.006302019356930936  
emp.var.rate: 0.006246493067232722  
default\_unknown: 0.0060646361393186295  
job\_unemployed: 0.0060425999344159865  
month\_jun: 0.0057390721732683316  
month\_may: 0.005337059871276105  
month\_mar: 0.004995978862559062  
poutcome\_failure: 0.0043840134656910855  
education\_basic.6y: 0.0039273453963482655  
default\_no: 0.0032880100992683146  
education\_basic.4y: 0.0031514241335199957  
month\_dec: 0.003082509468064045  
job\_unknown: 0.002882965376894681  
loan\_yes: 0.002846761990708497  
job\_student: 0.0028182943708014133  
day\_of\_week\_fri: 0.0026304080794146514  
marital\_divorced: 0.0023693013950609916  
job\_self-employed: 0.0021048246468800937  
job\_services: 0.0018761113539083486  
day\_of\_week\_wed: 0.0018527856626935355  
job\_management: 0.0014796045446707406  
education\_basic.9y: 0.0014796045446707406  
month\_apr: 0.001196786608466346  
month\_nov: 0.0002463705046370558  
month\_jul: 6.227291854675453e-05  
job\_entrepreneur: 0.0  
job\_housemaid: 0.0  
marital\_unknown: 0.0  
education\_illiterate: 0.0  
default\_yes: 0.0  
housing\_unknown: 0.0  
loan\_unknown: 0.0  
contact\_telephone: 0.0  
month\_aug: 0.0  
month\_sep: 0.0

```
In [33]: plt.figure(figsize=(15,16))
plt.barh(range(len(sorted_features)), [val[1] for val in sorted_features], align='center')
plt.yticks(range(len(sorted_features)), [val[0] for val in sorted_features])
plt.xlabel('Importance')
plt.ylabel('Caractéristique')
plt.title('Importance des caractéristiques')
plt.show()
```



```
In [35]: from sklearn.model_selection import GridSearchCV

param_grid = {
    'max_depth': [None, 10, 20, 30],
    'min_samples_split': [2, 5, 10],
    'min_samples_leaf': [1, 2, 4]
}

clf = DecisionTreeClassifier()

grid_search = GridSearchCV(estimator=clf, param_grid=param_grid, cv=5)

grid_search.fit(X_train, y_train)

print("Meilleurs hyperparamètres : ", grid_search.best_params_)

best_clf = grid_search.best_estimator_

y_pred = best_clf.predict(X_test)
accuracy = accuracy_score(y_test, y_pred)
print("Accuracy du meilleur modèle:", accuracy)
```

Meilleurs hyperparamètres : {'max\_depth': 10, 'min\_samples\_leaf': 1, 'min\_samples\_split': 5}  
 Accuracy du meilleur modèle: 0.9029126213592233

In [39]:

```
from sklearn.tree import DecisionTreeClassifier
dt = DecisionTreeClassifier(criterion='gini', max_depth=10, min_samples_split=5)
dt.fit(X_train, y_train)
```

Out[39]:

```
DecisionTreeClassifier(max_depth=10, min_samples_split=5)
```

In [45]:

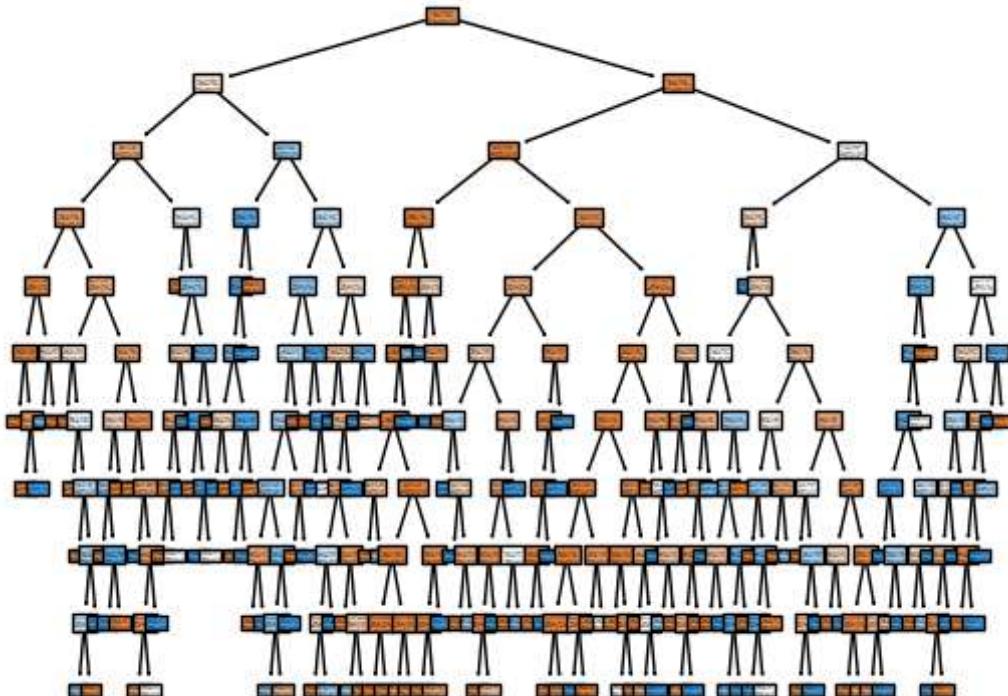
```
y_pred = dt.predict(X_test)

accuracy = accuracy_score(y_test, y_pred)
print("Accuracy du modèle sur les données de test:", accuracy)
```

Accuracy du modèle sur les données de test: 0.9

In [46]:

```
plot_tree(dt, class_names=cn, filled=True)
plt.show()
```



In [50]:

```
dt1 = DecisionTreeClassifier(criterion='entropy', max_depth=4, min_samples_split=15)
dt1.fit(X_train, y_train)
```

Out[50]:

```
DecisionTreeClassifier(criterion='entropy', max_depth=4, min_samples_split=15)
```

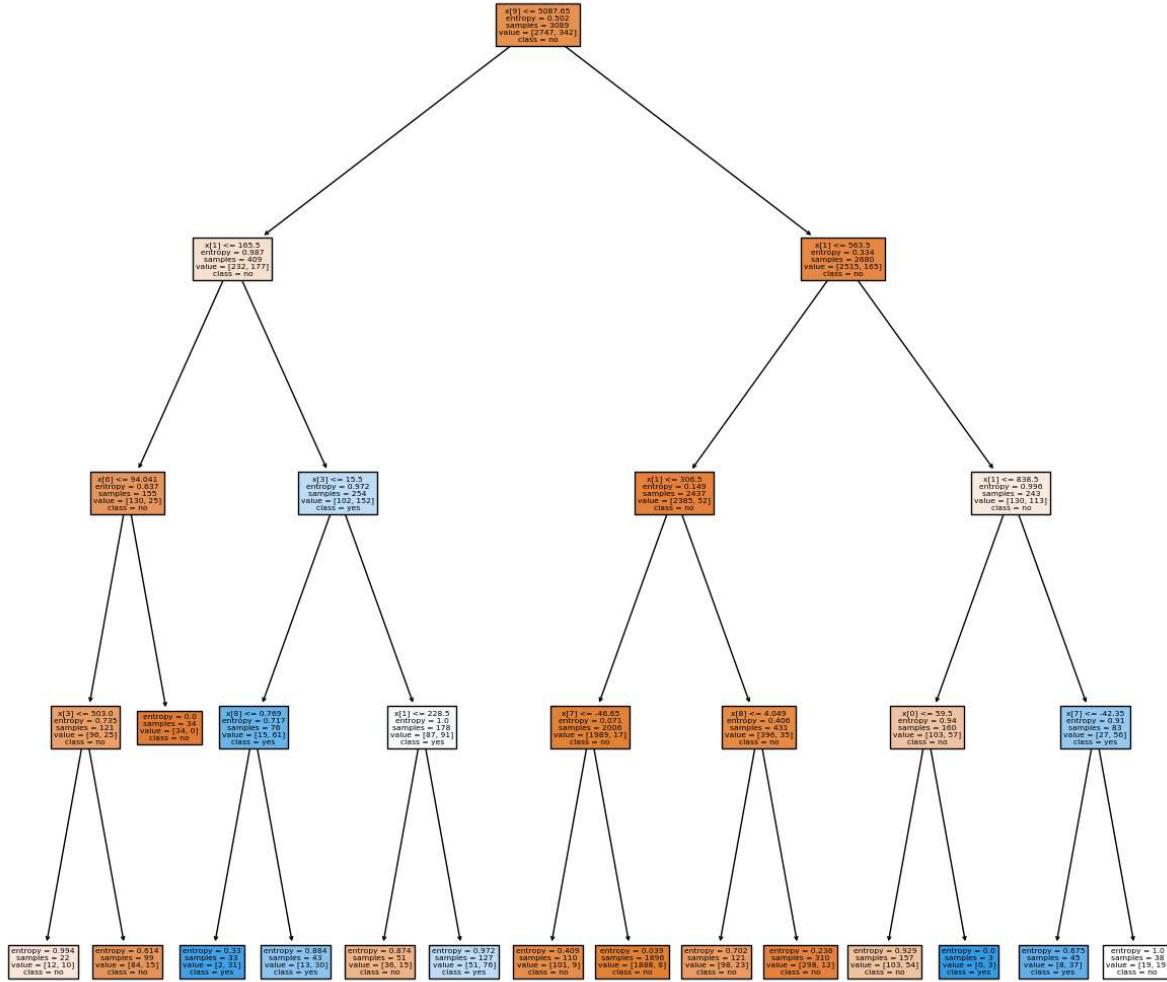
In [52]:

```
from sklearn.metrics import accuracy_score
y_pred = dt.predict(X_test)

accuracy = accuracy_score(y_test, y_pred)
print("Accuracy du modèle sur les données de test:", accuracy)
```

Accuracy du modèle sur les données de test: 0.9

```
In [54]: plt.figure(figsize=(15,15))
plot_tree(dt1,class_names=cn,filled=True)
plt.show()
```



In [ ]: