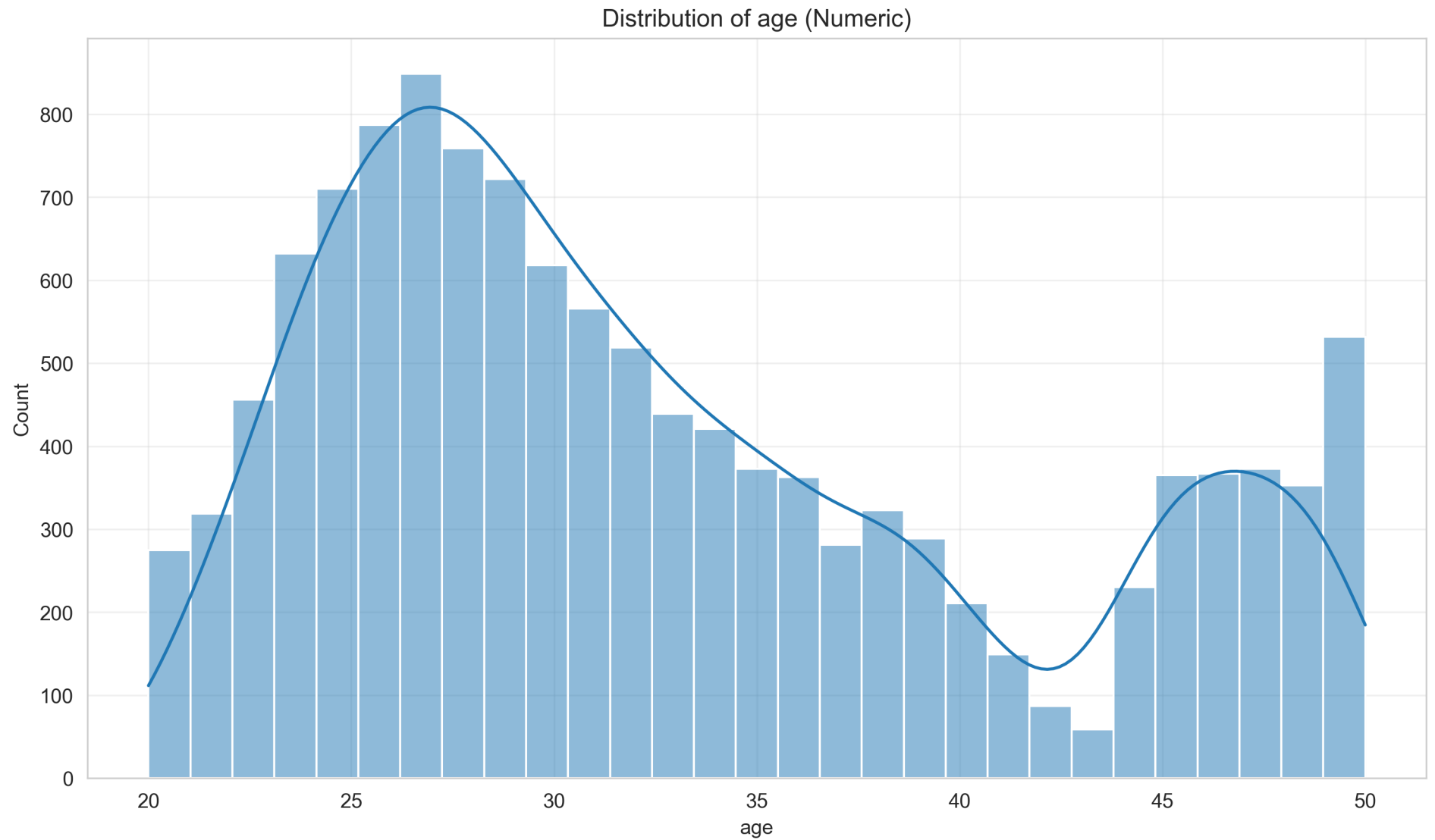


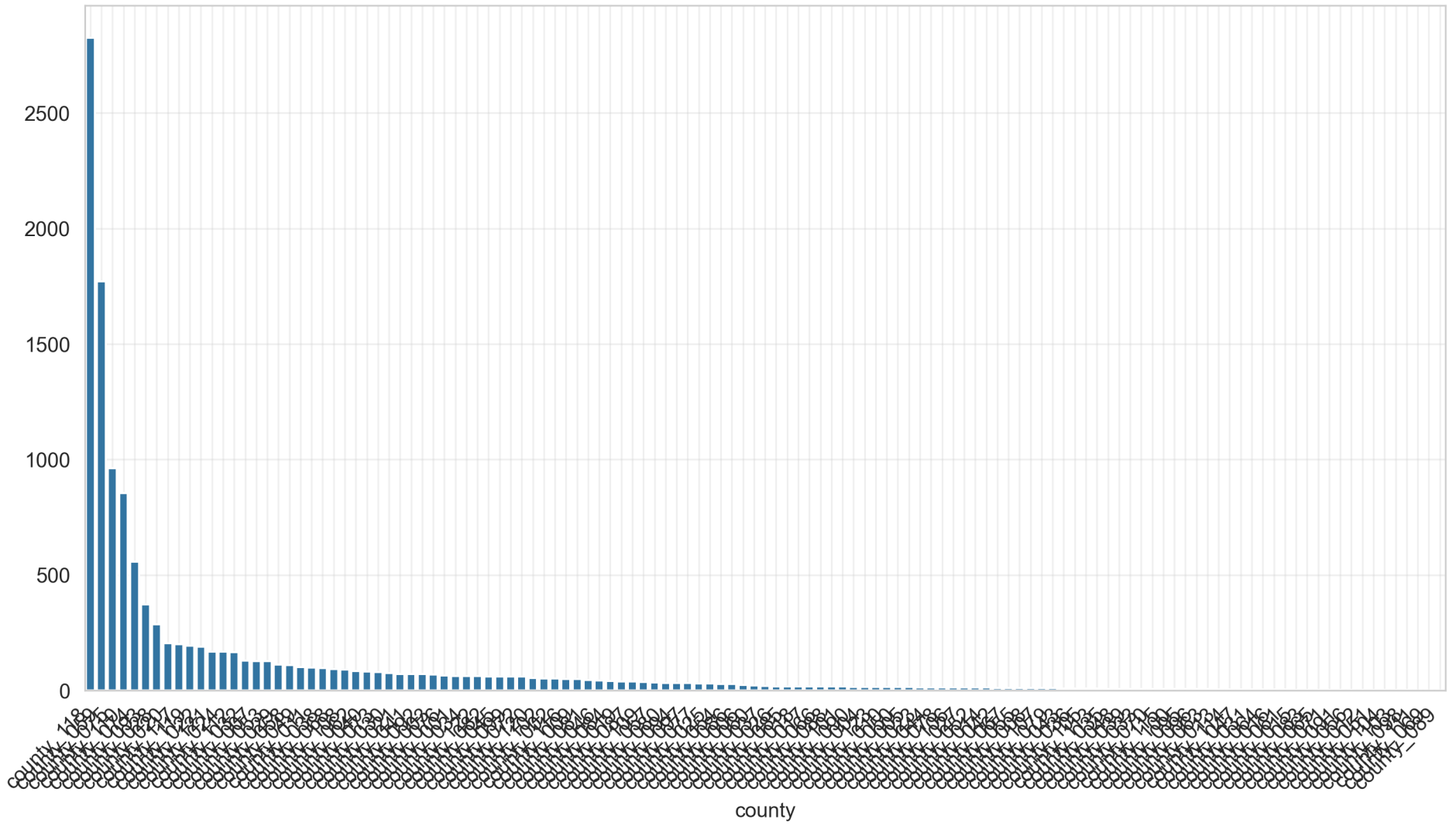
Job Change Prediction Model Analysis

EDA

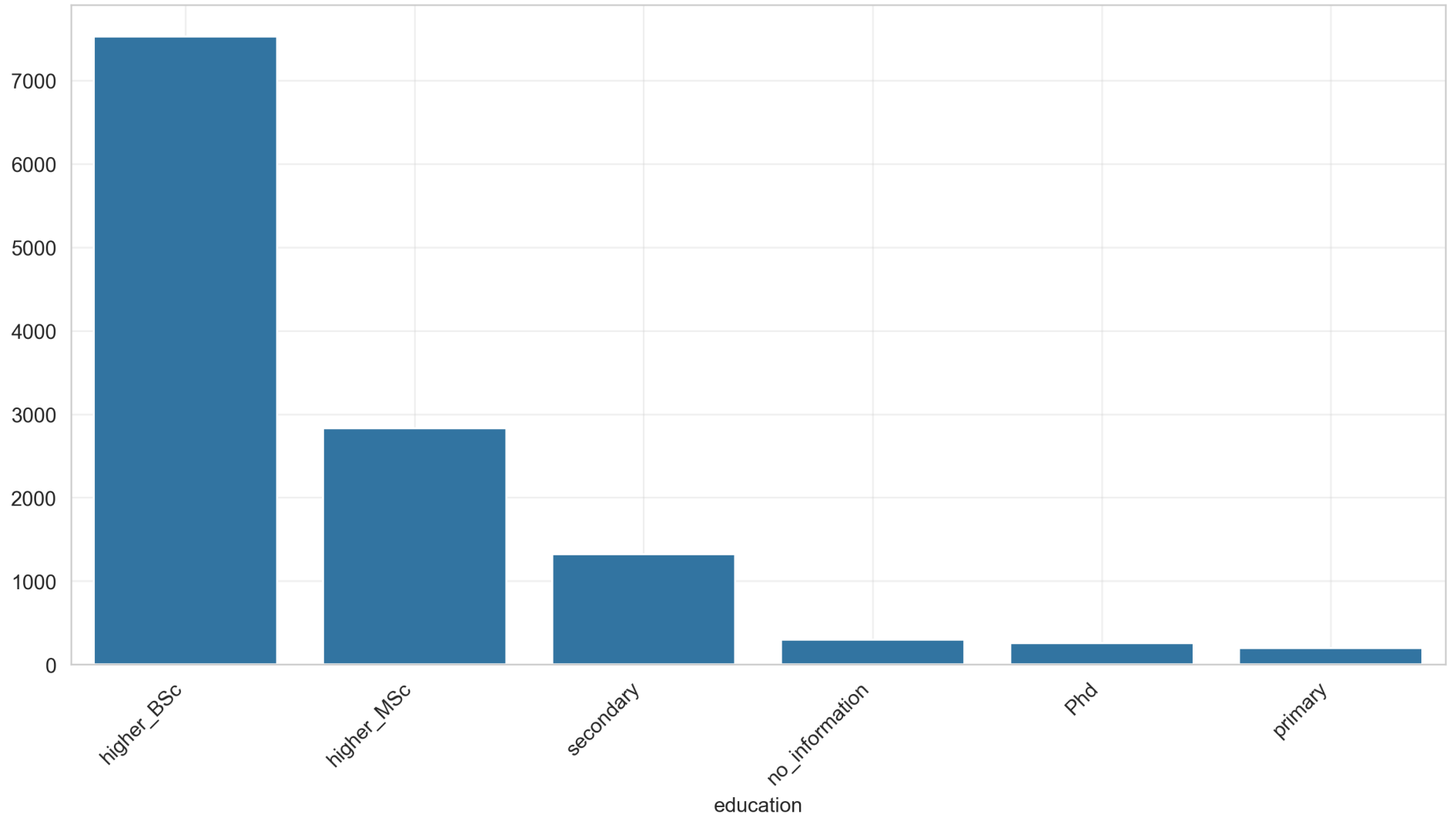
- No missing values
- Features histograms:



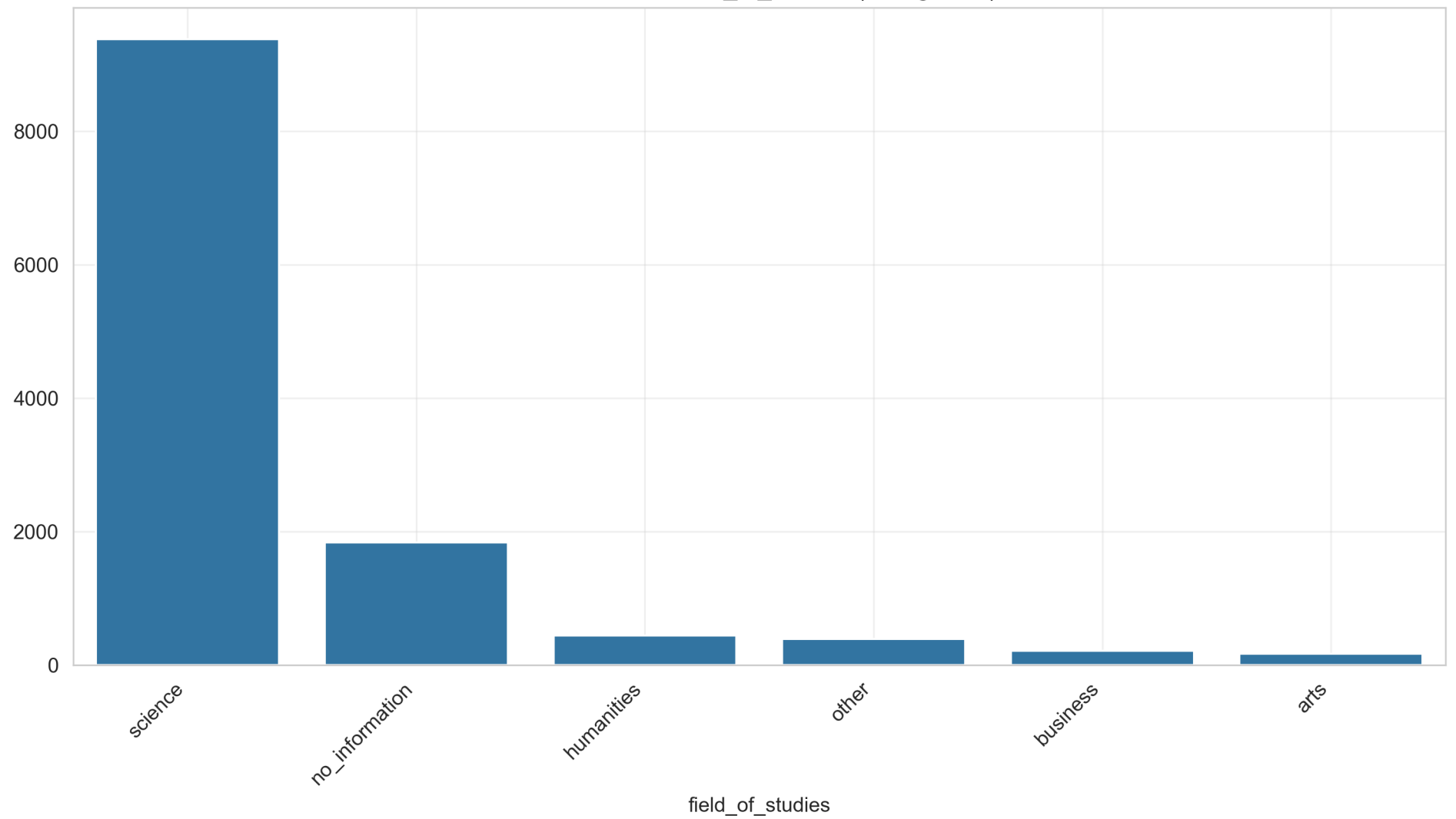
Distribution of county (Categorical)



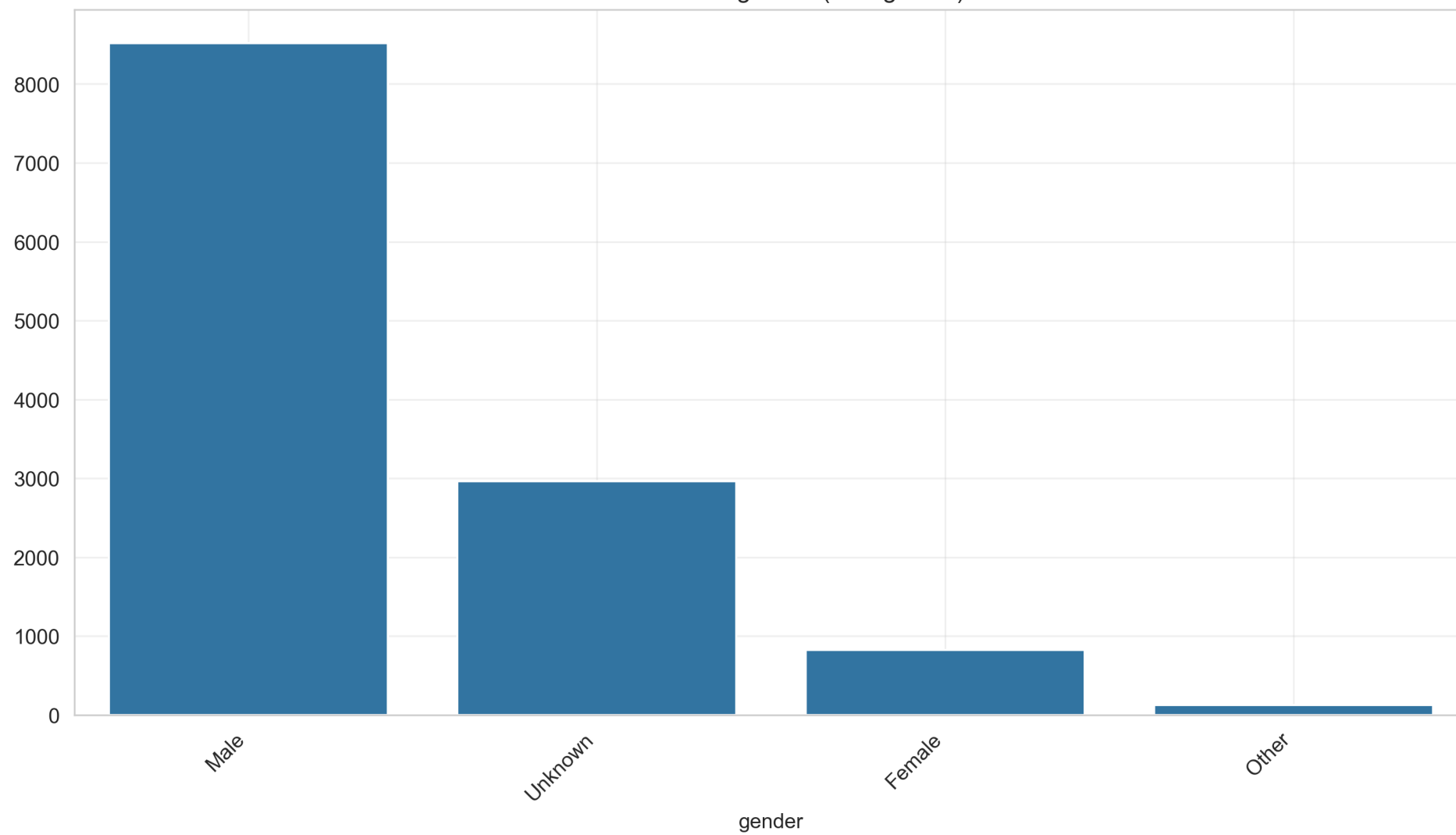
Distribution of education (Categorical)



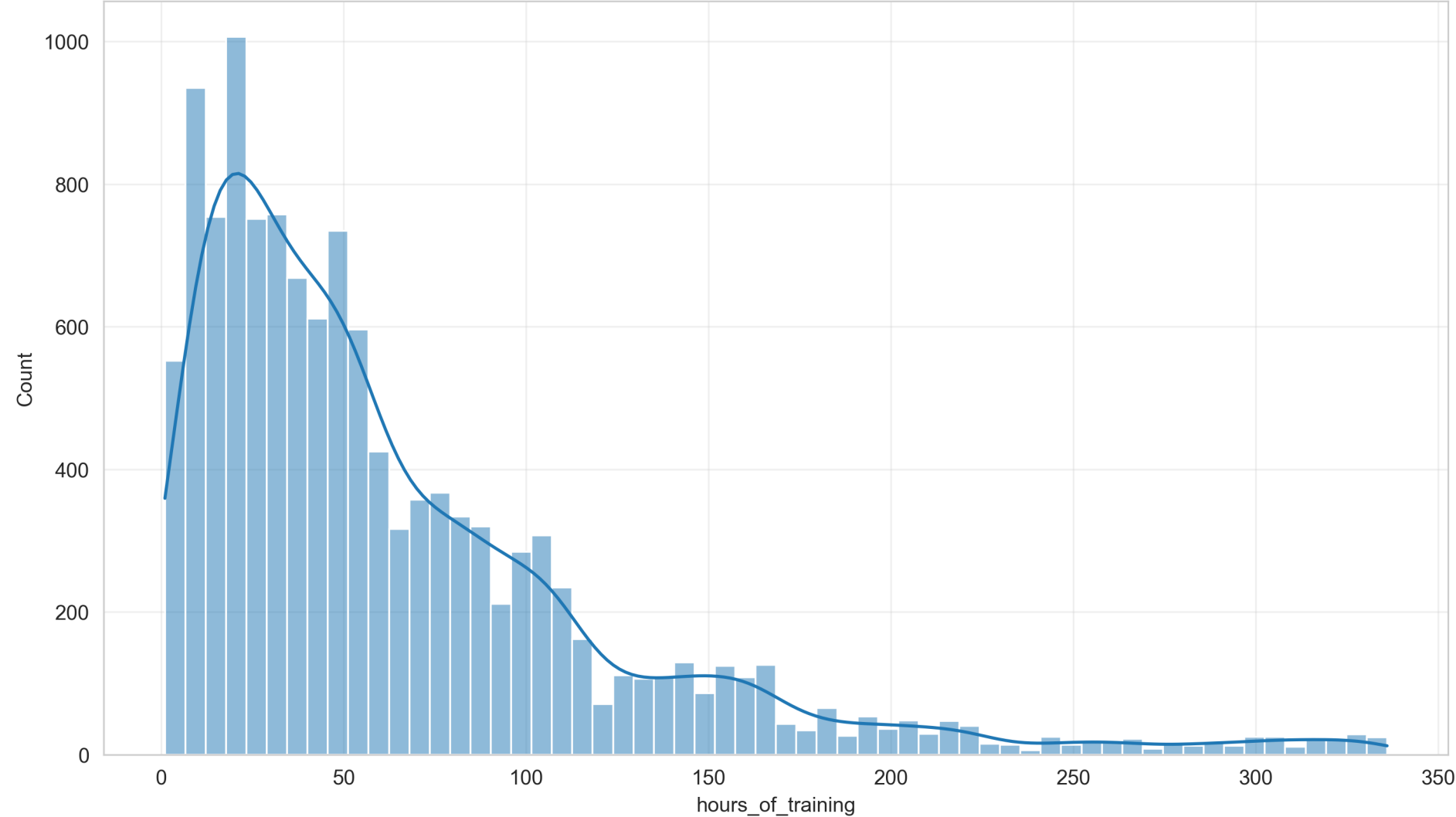
Distribution of field_of_studies (Categorical)



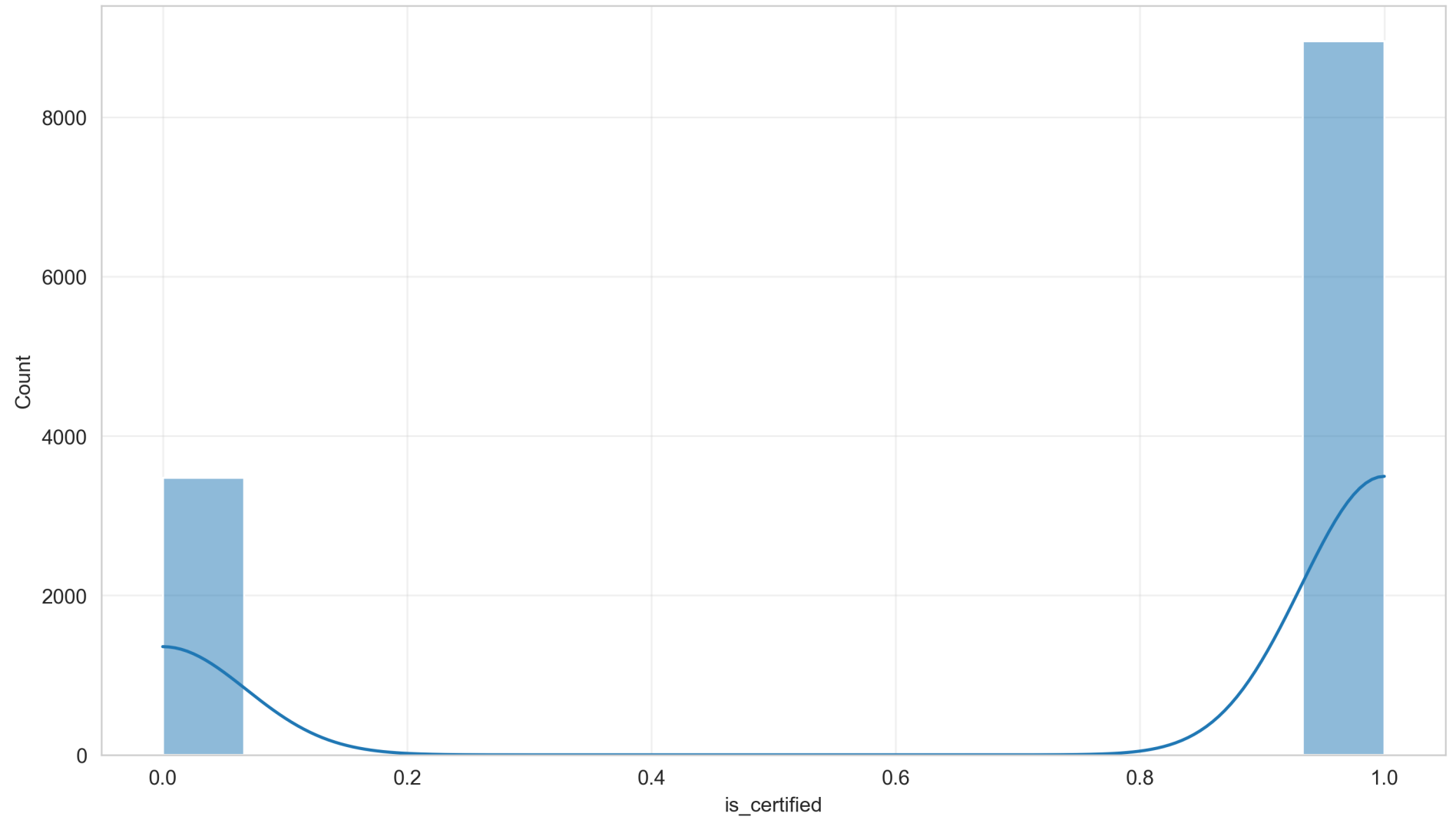
Distribution of gender (Categorical)



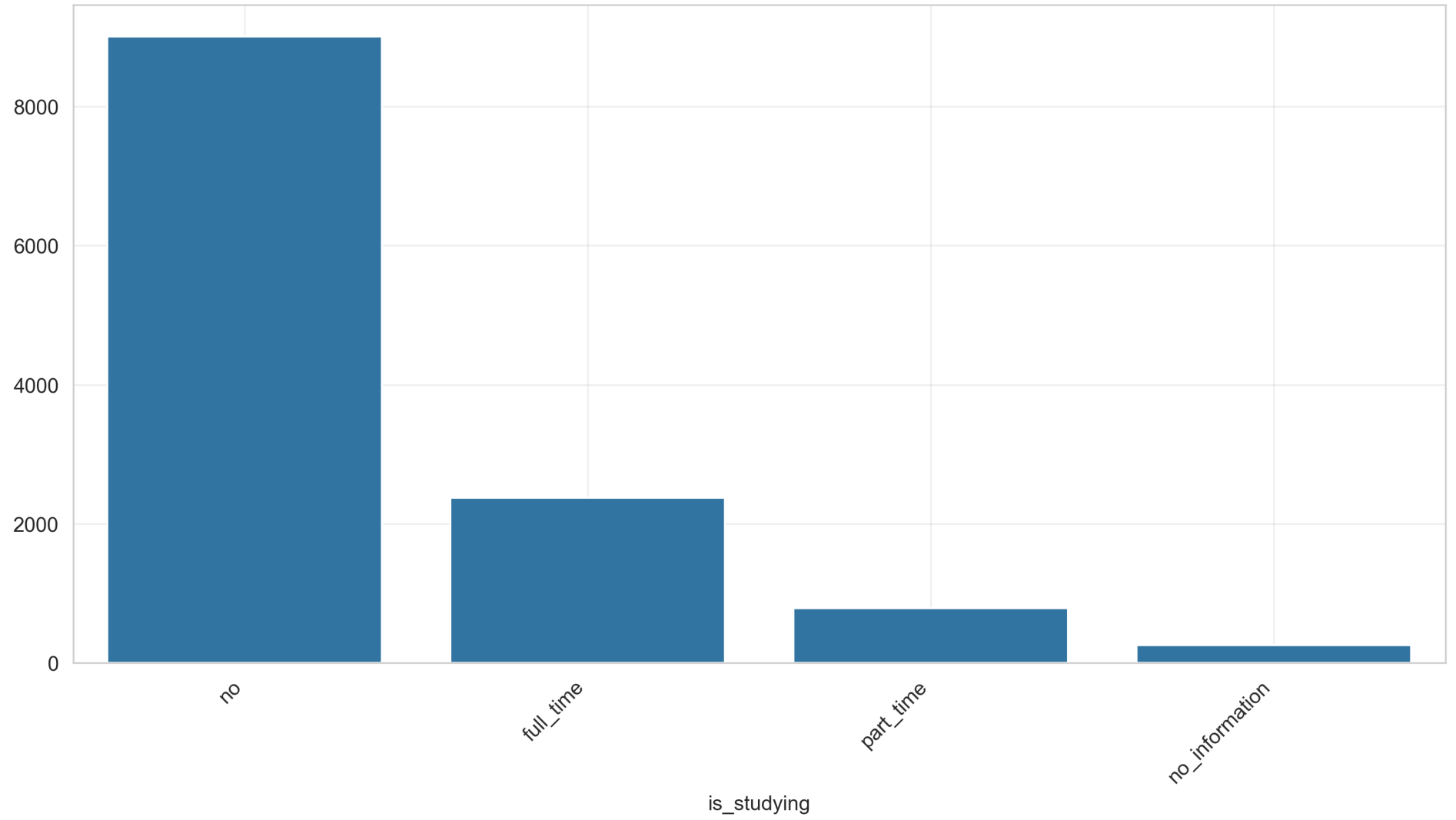
Distribution of hours_of_training (Numeric)



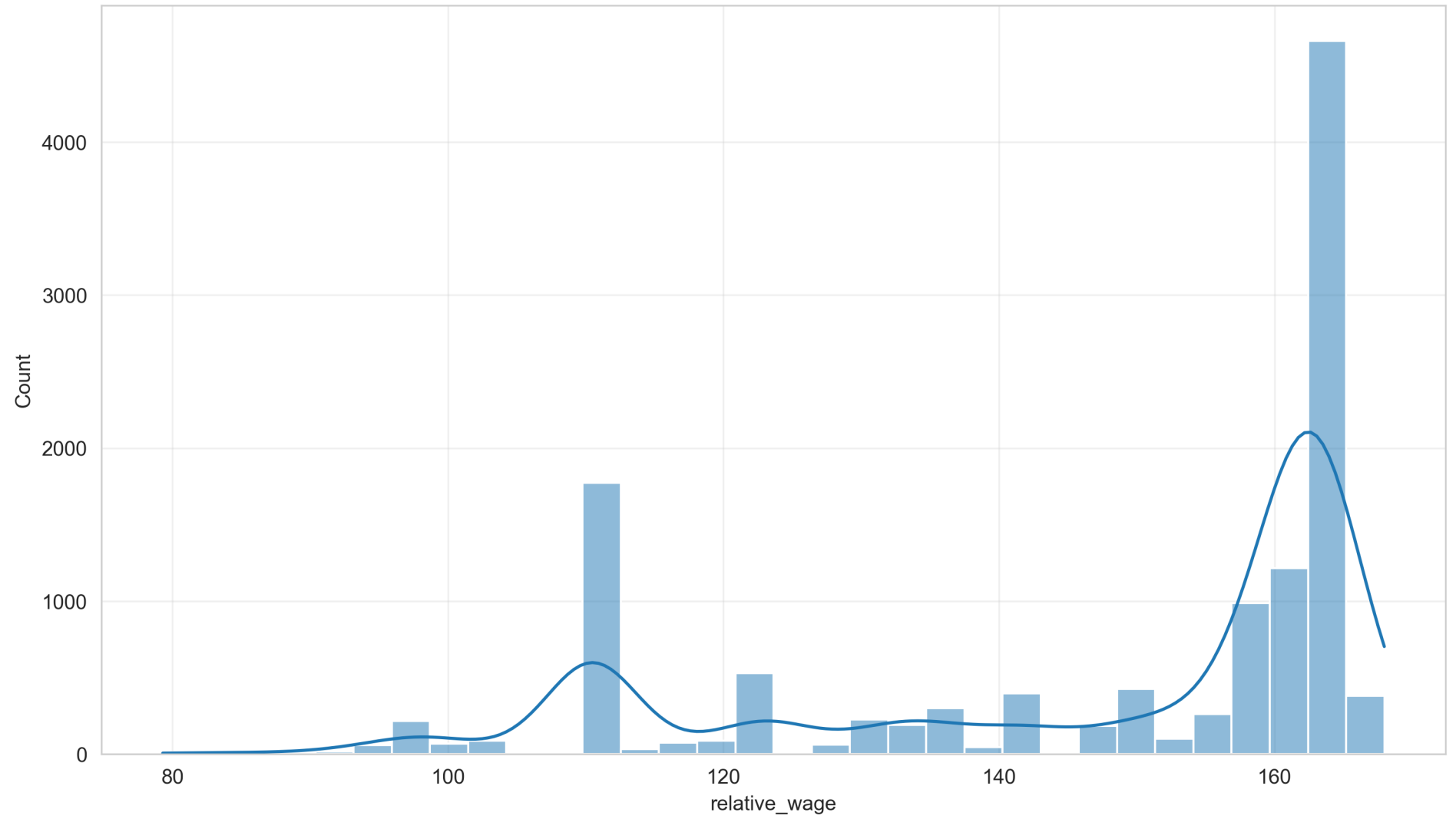
Distribution of is_certified (Numeric)



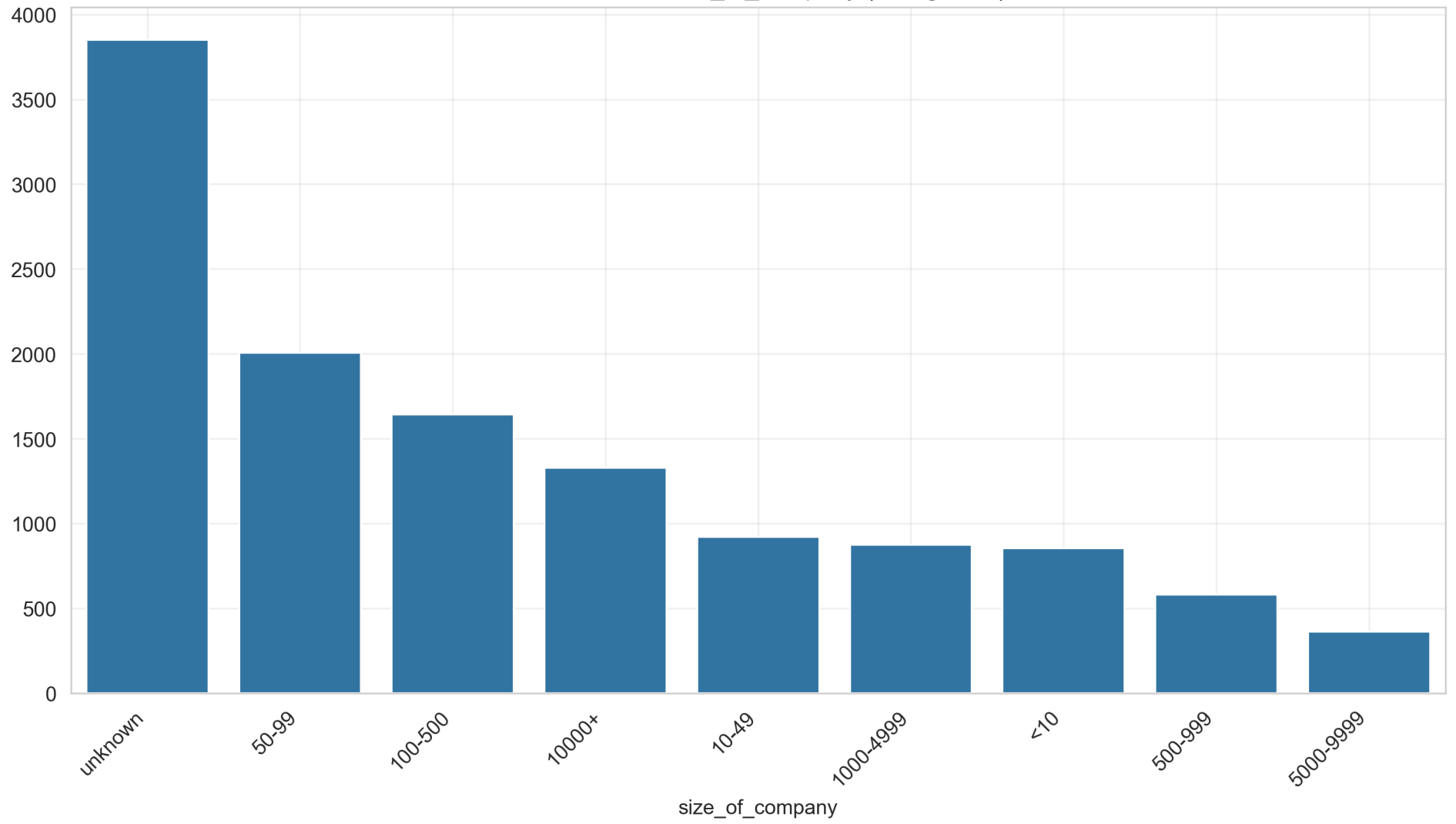
Distribution of is_studying (Categorical)



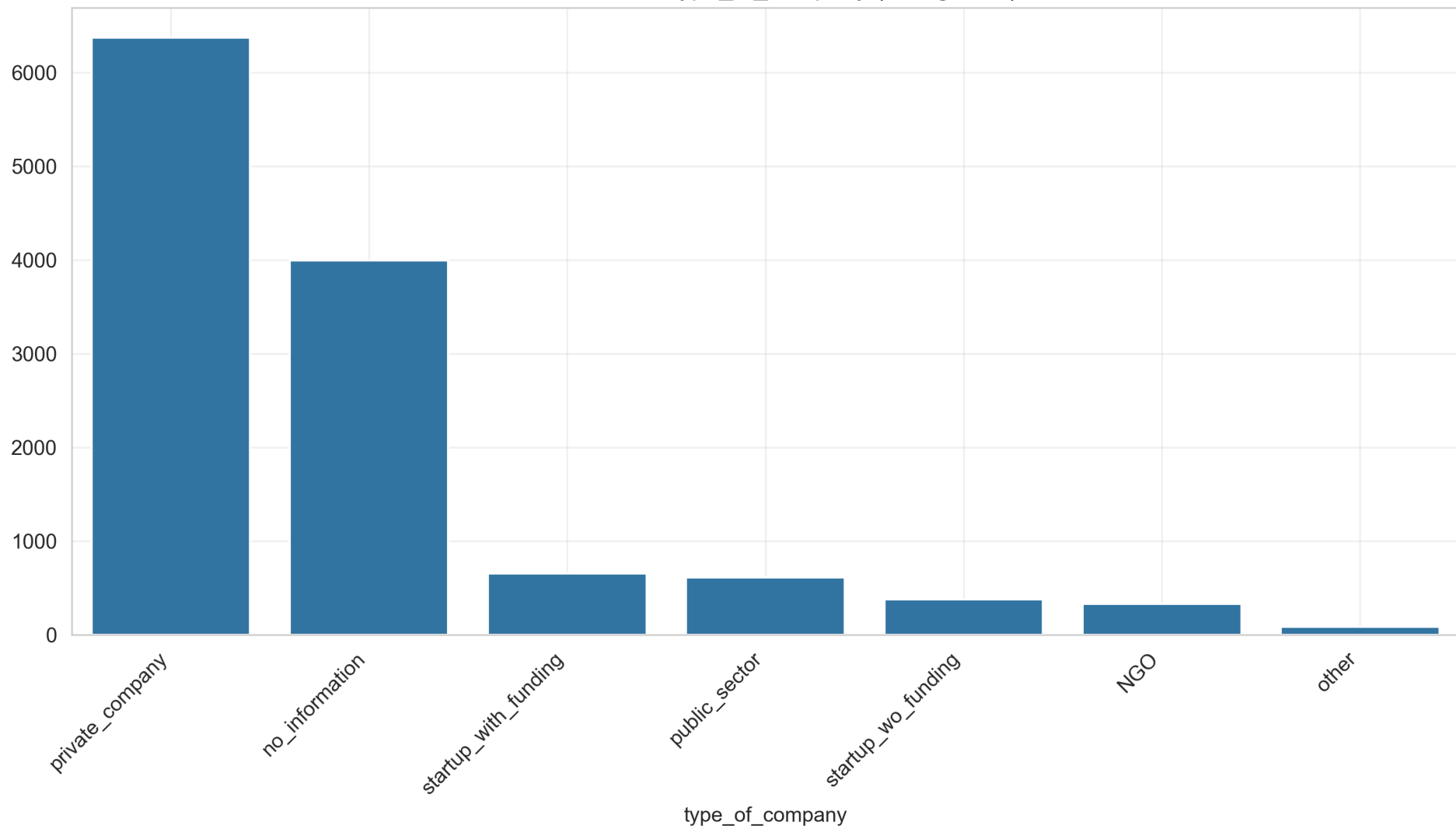
Distribution of relative_wage (Numeric)



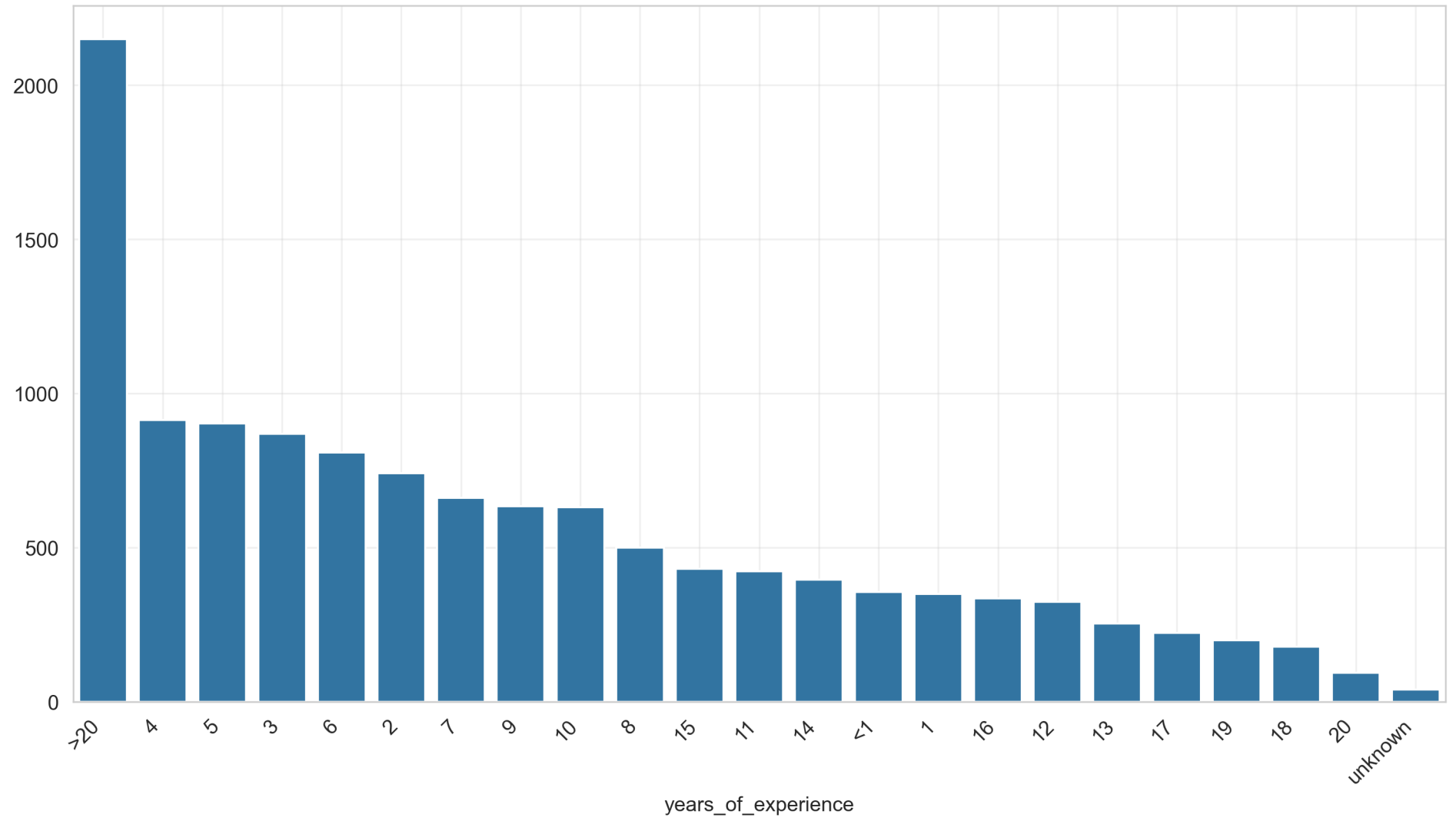
Distribution of size_of_company (Categorical)



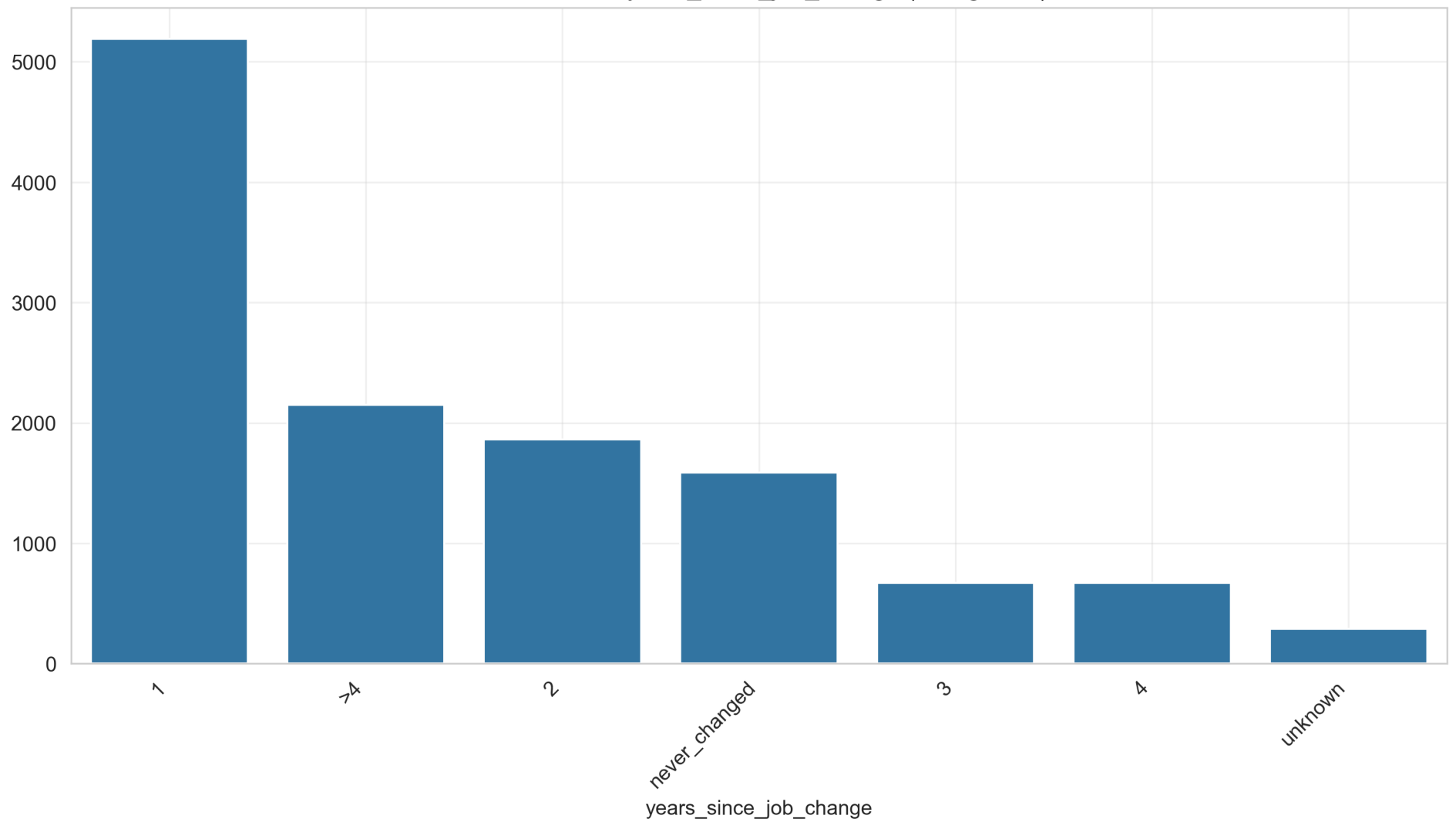
Distribution of type_of_company (Categorical)



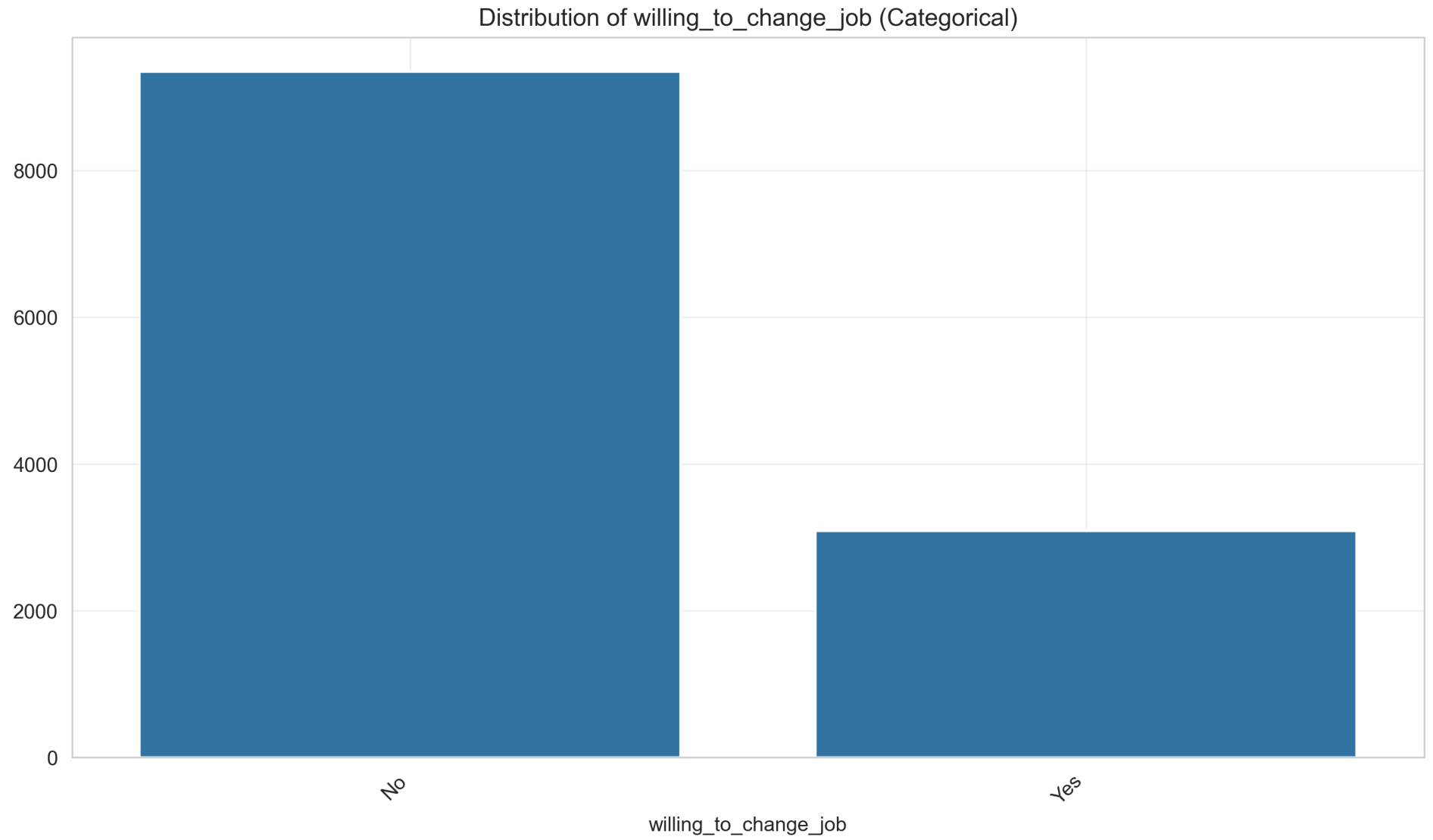
Distribution of years_of_experience (Categorical)



Distribution of years_since_job_change (Categorical)



- Target histogram:



Data Preprocessing Pipeline

1. **Categorical Variables** (One-hot encoding)

- gender
- education
- field_of_studies
- is_studying
- county
- years_since_job_change
- years_of_experience
- size_of_company
- type_of_company

2. **Numerical Variables** (Standardization)

- age
- relative_wage
- hours_of_training

3. **Special Handling**

- County grouping: Top 10 most frequent counties kept, rest grouped as "Other"
- Experience binning: Converted to categorical ranges
 - 0-3 years
 - 3-7 years
 - 7-15 years
 - 15+ years
- Special values handling:
 - '>20' years converted to 21
 - '<1' year converted to 0

Model Selection

After testing multiple algorithms including Random Forest, XGBoost, Gradient Boosting, and KNN, **Logistic Regression** was chosen.

Model	CV Score	CV Std	Test Score	Training Time (s)
XGBoost	0.7942	0.0095	0.6757	0.798
Logistic Regression	0.7602	0.0103	0.7554	0.197
Gradient Boosting	0.7021	0.0023	0.6908	7.985
Random Forest	0.6748	0.0115	0.6711	8.580
KNN	0.6617	0.0100	0.6639	0.593

Top 20 Feature Importances in Logistic Regression Model

