
A case study for US Census Bureau:

Predict person's income is more or less
than 50000 dollar

Xi Chen

Overview

- Background
- Usage of model considerations
- Data overview
- Exploratory data analysis (EDA) & Insights
- Modeling & Performance
- Conclusion
- Next step
- Backup slides

Background

What is US Census Bureau

- A government agency responsible for collecting and analyzing demographic and economic data
- Conduct a nationwide census every 10 years to guide policy and funding decision

Why Census Data important

- Helps allocate billions of dollar for public services
- Provide valuable demographic insights for economic and social planning

Objective

- Analyze US Census data to identify factors influencing income levels
- Develop a data analysis and modeling pipeline to predict whether an individual earns more or less than \$50000 per year

Usage of model considerations

It is important to understand the usage of the model, considering it has a big impact on the analysis:

- **Impact on predicting target:**
 - There might be better target to predict other than 50k depends on the usage. For example, to predict income bucket, or directly predict whether to give out benefit(Social benefit) or not instead of purely based on income is larger than 50k or not.
- **Impact on the evaluation approach:**
 - We might use different way to evaluate the model depends on the objective. i.e, for tax fraud might focusing on the >50k.
 - To help to understand how good the model's performance need to reach.
- **Impact on the Data:**
 - There might be more features can be gathered for specific target.
 - Data regulation.
 - Model bias based on race/gender.
 - The potential data sampling approach, data drifting issue

Data overview

- Data size: Total data size is 299285 records. 199534 (66.66%) for training and 99762 (33.34 %) for testing.
- Data collected year : 94 and 95
- Feature size: In metadata file, there are 45 different features in total, but in dataset only contains 42 features.
- Main feature groups:
 - Demographics: Age, Gender, Race, Citizenship
 - Employment: Occupation, Industry, Work Class, Hours Worked
 - Financial: Capital Gains/Losses, Dividends, Tax Status
 - Education & Family: Education Level, Marital Status, Household Composition
 - Geography & Migration: State, Region, Migration Patterns

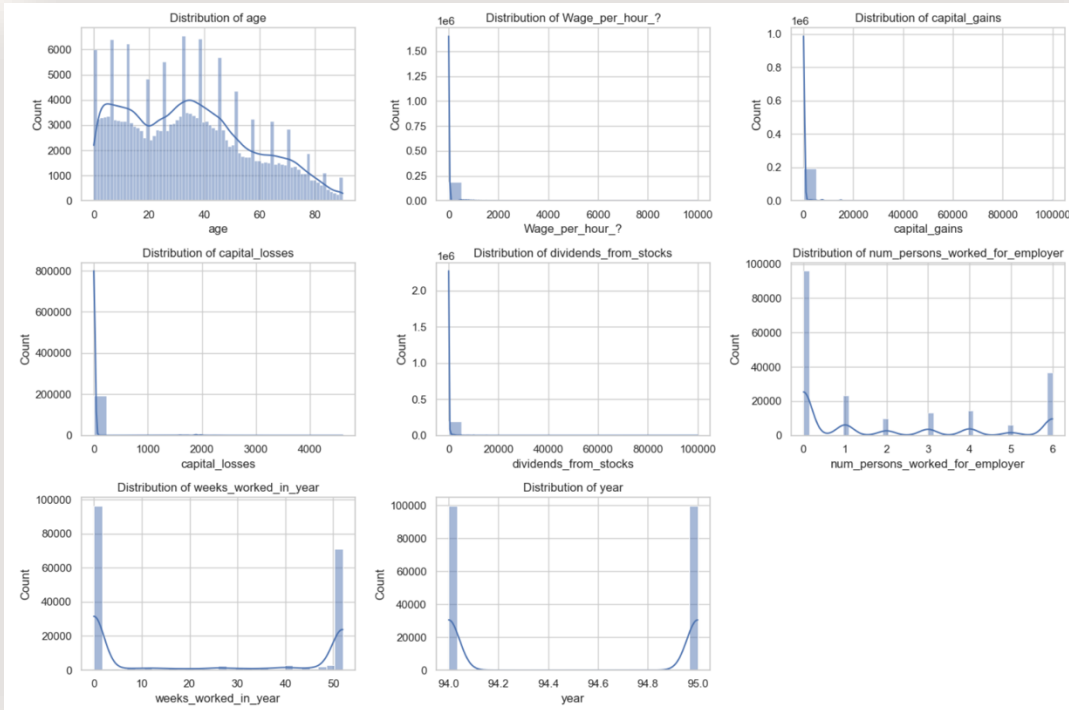
PS: There is no column name in the data, I matched it based on the feature order and the value report from metadata file.

Exploratory data analysis (EDA)

There are two groups of features: 8 Numeric features & 33 Categorical features.

We're mainly into the value distribution, frequency distribution to identify missing value, outliers, potential opportunity to extract more features for modeling purpose

EDA – Numeric feature



Feature

Findings

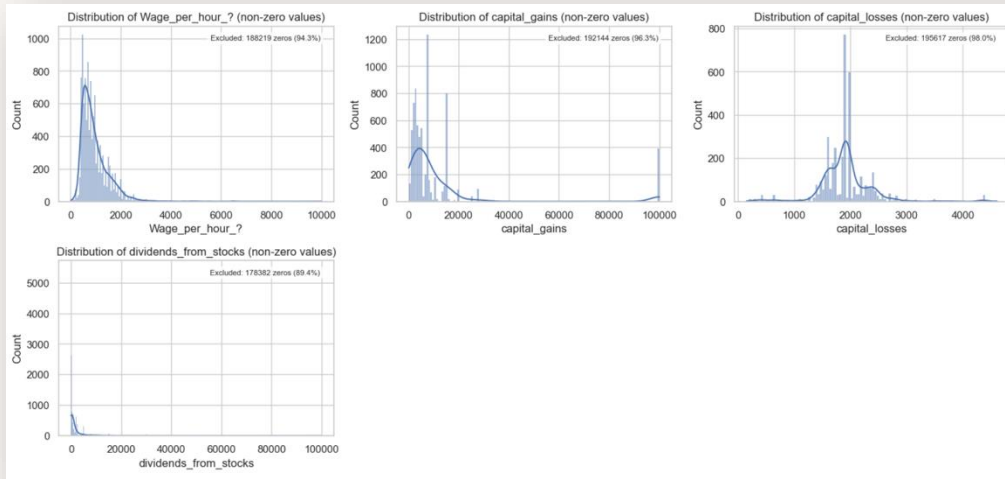
Age

There is a spike in each age group indicates there might have certain data sampling pattern.

Wage per hour
Capital gains
Capital losses
Dividends

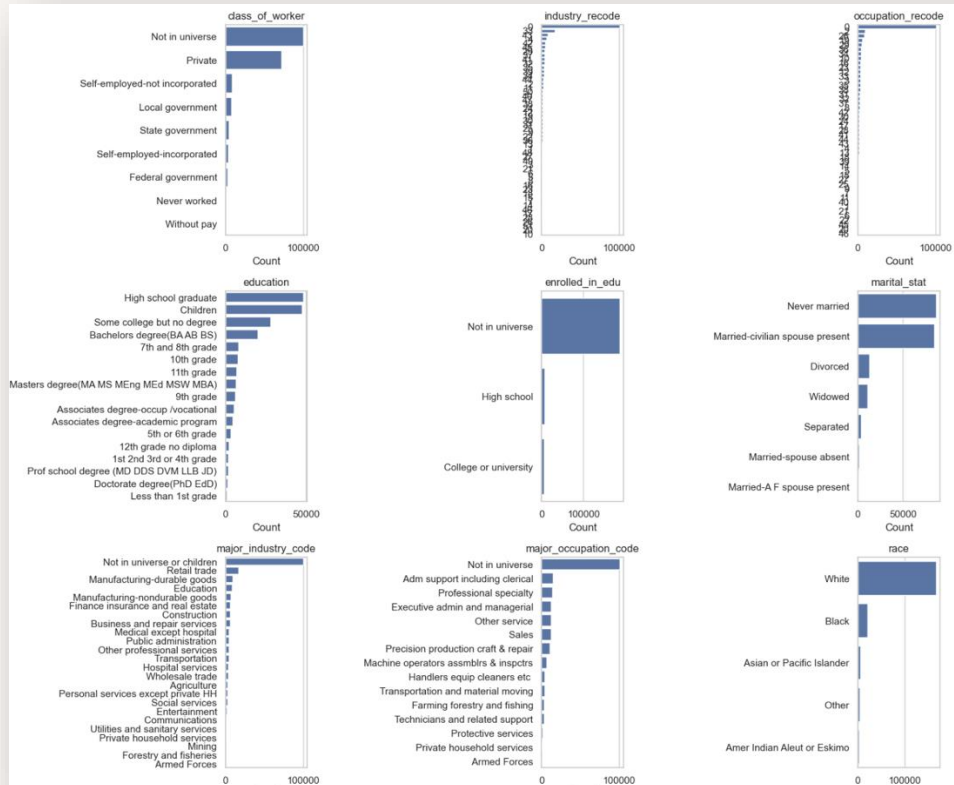
Majority are empty. Zoom in distribution(Remove 0) will show in next slide

EDA – Numeric feature



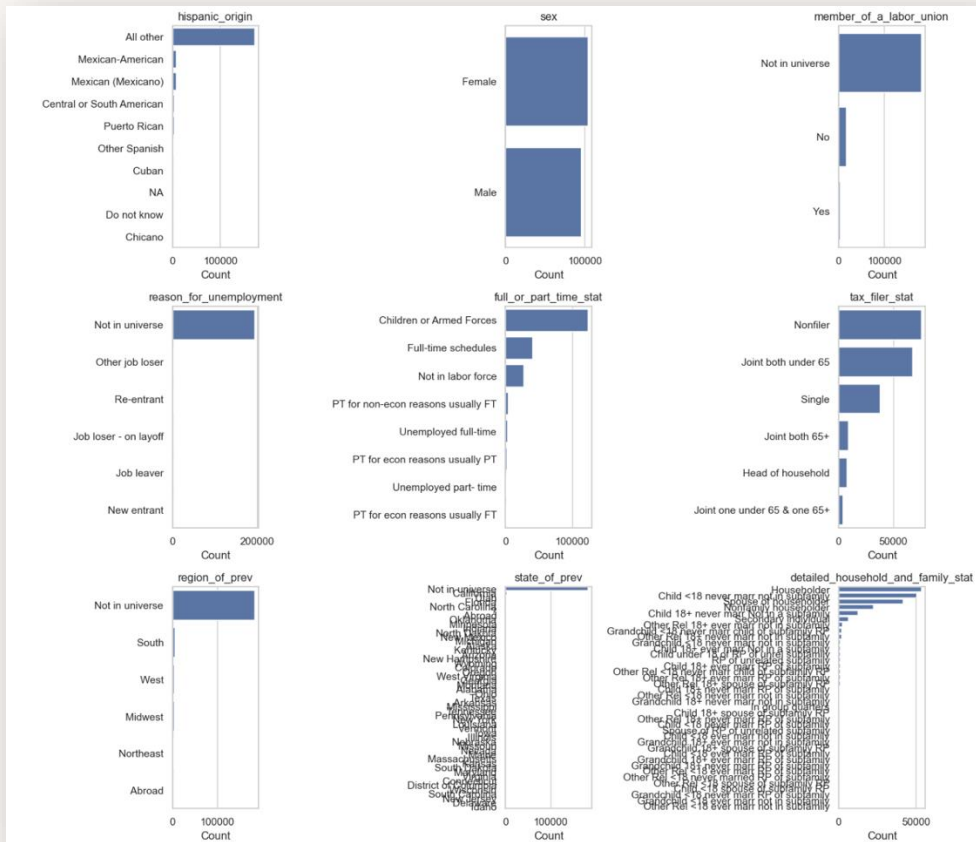
Feature	Findings
Wage per hour	This feature might be weekly income instead of hourly income. As the average is 800 dollar per hour, which is too high for hourly wage
Capital gains	There is a peak in 100000, there might be a cap on the declaration?
Capital losses	Compared with capital gain, the number of the capital losses is surprisingly small. I would expect the capital losses is also

EDA – Numeric feature



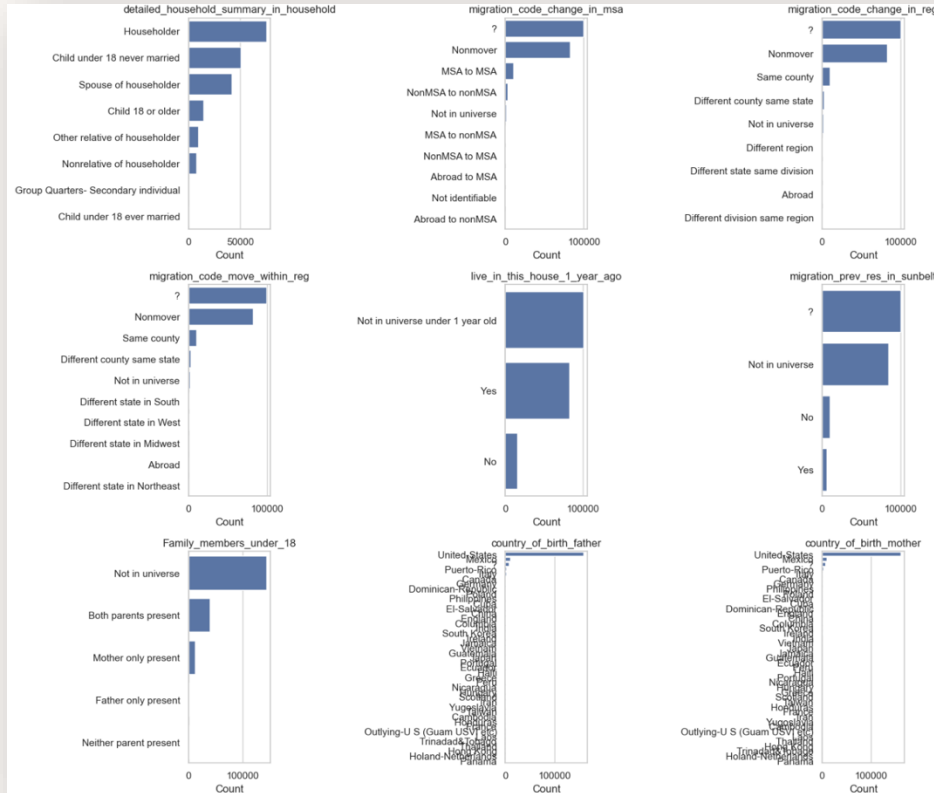
Feature	Findings
Class of worker	Class 0 might indicates students / children.
Industry recode	The class of worker have a lot of private might need pay a bit attention. Might be caused by the sampling approach
Occupation recode	
Major industry code	
Race	Majority are white, need to pay attention when it comes to race / gender bias if there is any ethical concerns might be raised by public

EDA – Categorical feature



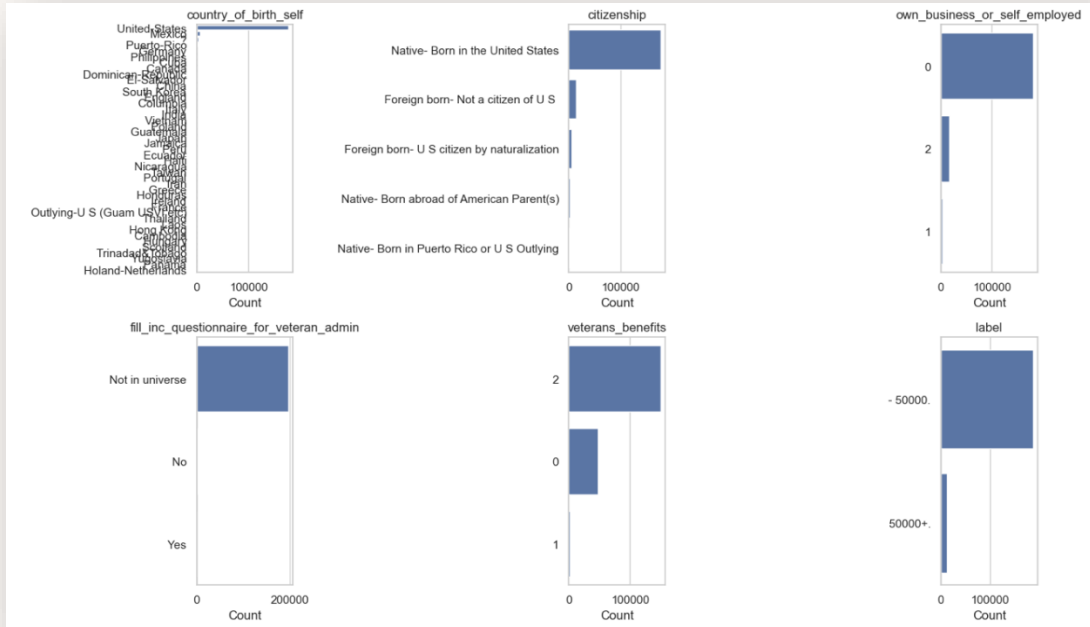
Feature	Findings
State_of_prev	Majority are the locals(who didn't move)
Household and family stat	This category might be able to break down to multiple categories

EDA – Categorical feature



Feature	Findings
Migration code change	What's the difference between "?" And Not in universe?
Country of birth father/mother/se If	Majority are native Americans

EDA – Categorical feature



Feature	Findings
Label	The distribution of income is highly imbalanced

EDA – Insights

- **Data correctness:**

1. No empty filed and obvious outlier
2. Potential conflicted records: (I,e age vs education: below 18 with PhD, age vs marital status: below 18 and “Widowed”, Occupation vs. Industry: Teacher in construction industry. To be consolidated with business.)

- **Data Sampling:**

1. There are spike in ages distribution. Need to find out why.
2. There are lots of people under 18. For predicting purpose, would be better to have more people above 18 if the predicting target is income > 50k or not.

- **Data preparation for modeling:**

1. Numeric feature: Some of the features distribution are skewed
2. Geo related feature: We can extract geo / embedding for data modeling
3. Feature engineering:
 - Expand the feature to multiple features such as household info
 - Statistics for the numeric value
4. Tackle imbalanced dataset

Modeling – results

Two different models (XGBoost) with accuracy of: 90.86% and 95.72%

		Predicted label		Predicted label	
		< 50k	> 50k	< 50k	> 50k
True label	< 50k	85686 (91.6%)	7890 (8.4%)	99.0 (99%)	907 (1.0%)
	> 50k	1226 (19.8%)	4960 (80.2%)	3366 (54.4%)	2860 (45.6%)
Model A: Accuray 90.86%			Model B: Accuray 95.72%		

Even the model B has overall higher accuracy. But depends on the objective Model A might be a better choice. i.e Identify someone is under-declare income for tax benefit.

Conclusion

- The model can predict the income threshold of 50k with 95.72% accuracy
- Impact & benefit (Depends on the use case)
- Limitaions: The data sample is bit too old, and highly imbalanced data is challenging for modeling
- Domain knowledge: Some of the features, value and objective will need business to clarify

Next step

Business:

- Objective: Explore the other possible objectives for modeling target
- Data sampling: Evaluate if the sampling approach is representative or not
- Attributes: Identify if there is more attributes can be gathered
- Conflicting record's value: Verify the rules and consolidate it for data cleansing

Technical:

- Feature engineering: Review it together with business
- Model enhancement: Tryout deep learning models (i.e deep learning model).
- Explore ensemble methods

Backup slides

Modeling – setup (Backup)

Data preparation:

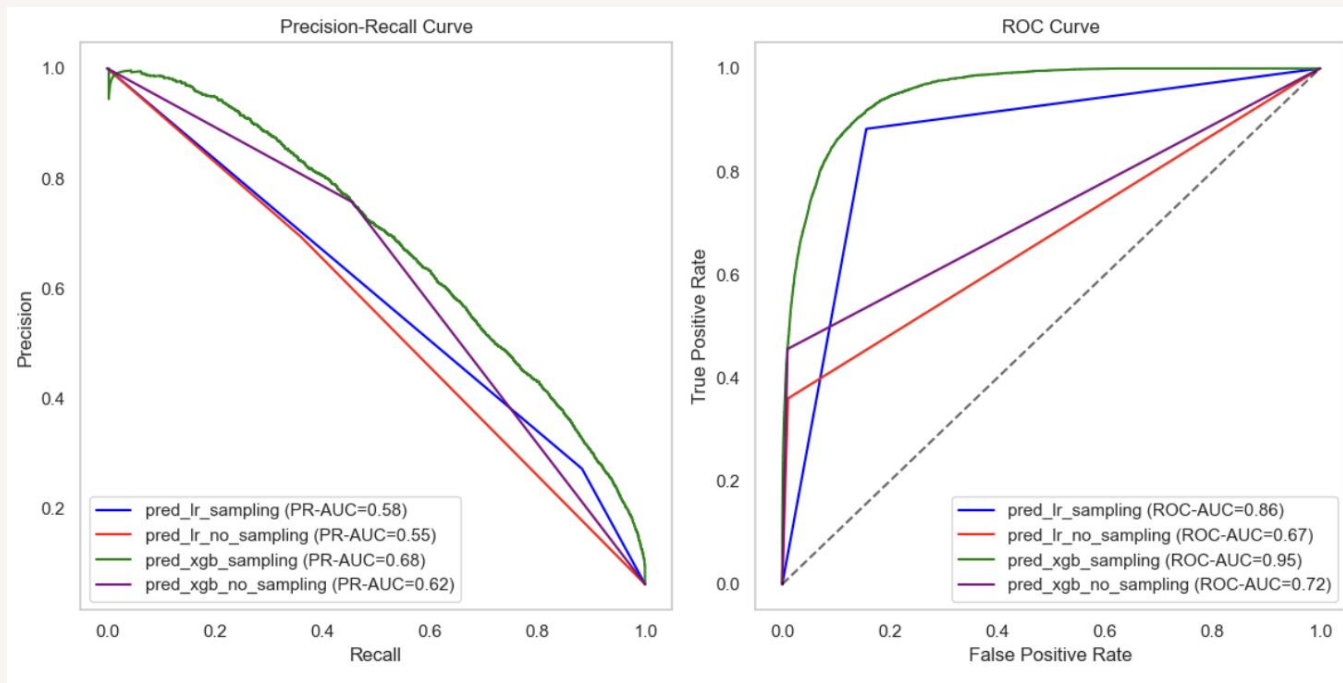
- Encoding
- Standardization
- Feature engineering

Hyperparameter tuning: Cross-validation

Model:

- LR
- XGBOOST
- Deep learning

Backup



Backup

