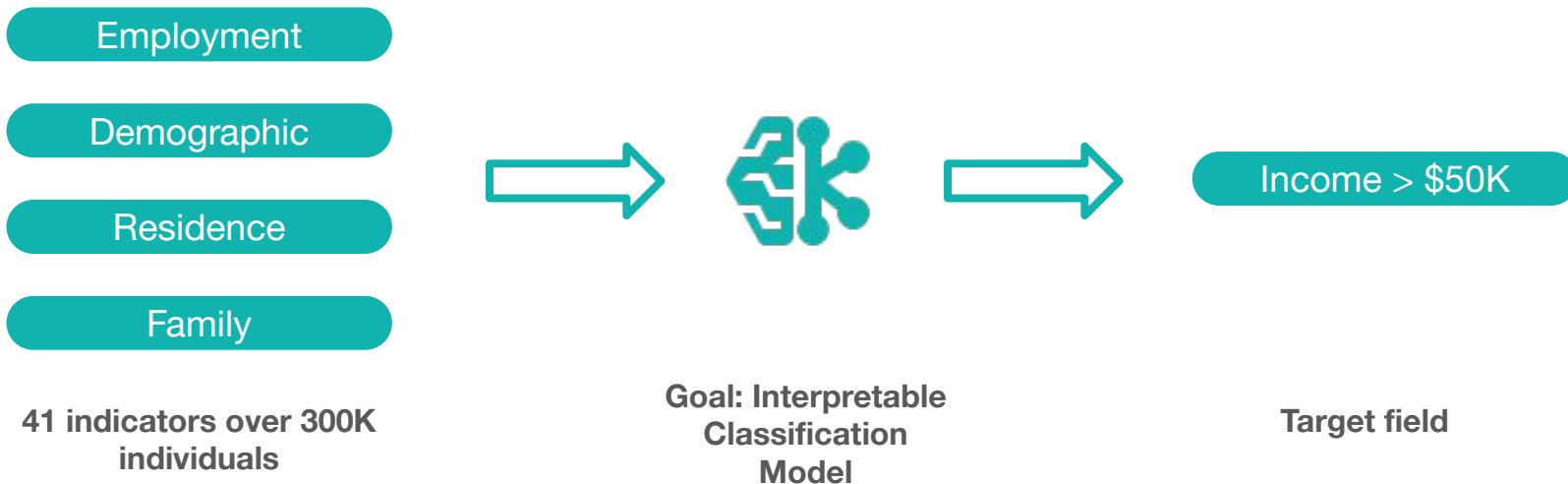# US Census
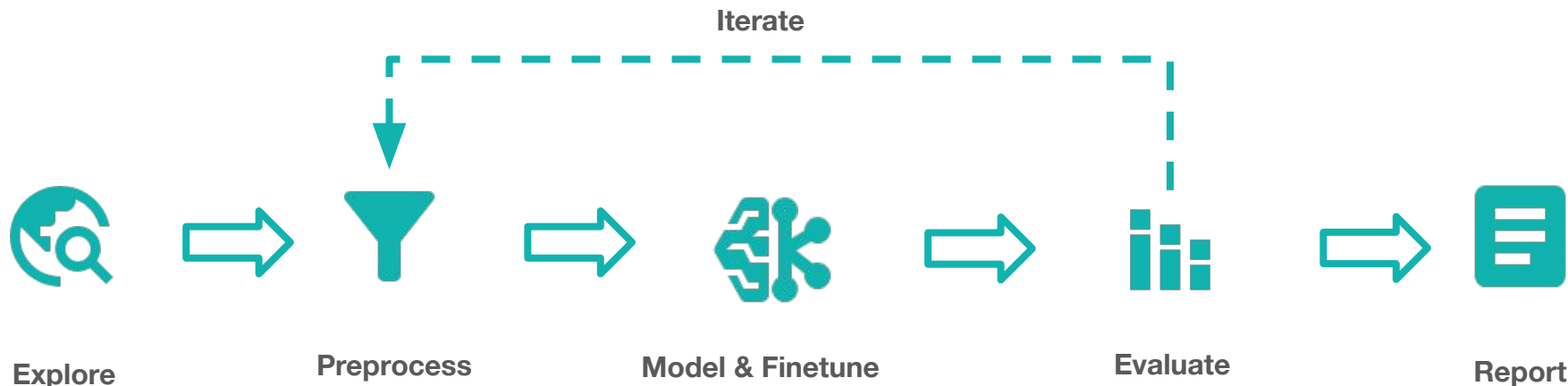## Income Prediction

Daniel Soukup

2025.11.03

# The Challenge

Our goal is finding identifying characteristics that distinguish between a person making more or less than $50,000 per year based on a sample dataset from the US Census archive.

Employment

Demographic

Residence

Family

**41 indicators over 300K individuals**

**Goal: Interpretable Classification Model**

Income > $50K

**Target field**

# Solution Framework

In addressing this problem, we followed an iterative process. Guided by our initial data exploration, the data was preprocessed and used to create competing classification models. Informed by subsequent evaluation, we refined the preprocessing and modeling steps to achieve higher performance and in-depth insights.

**Iterate**

**Explore**  **Preprocess**  **Model & Finetune**  **Evaluate**  **Report**

# Data Exploration

Our first step was to thoroughly review the data set, gaining statistical insights to the distribution of each field as well as uncovering any data quality concerns. This step informs the following preprocessing and feature engineering before modeling.

## Data Quality

Key concerns uncovered during EDA:

- **Duplicate** rows and **missing** values
- Highly **skewed distributions** with **extreme outliers**
- **High cardinality** categorical columns

## 8%
## high income

The **high class imbalance** of the target will affect model tuning and evaluation.

## Statistical Indicators

Confirming intuition, we found that **education & occupation** fields have significant relationship with having high income.

# Preprocessing

This step in our solution ensures that the raw data is best prepared for the modeling step. We addressed the data quality issues and encoded the categorical features adequately for our selected tree-ensemble model.

**Data Pipelines**

We created a **single, modular data processing pipeline** which captures multiple, custom data transformation steps. This approach ensures

- **consistency**: the same operation applied to train & test data sets,
- **no data leakage:** information from the test set is not used to create the pipeline**,** and
- **reduces tech debt** in the long run.

Data type mapping

Deduplication

Null Imputation

One-Hot Encoding

# Modeling & Fine Tuning

After the preprocessing of a mix of numeric and categorical features, we optimized an **XGBoost binary classification** model fine-tuning over a number of hyperparameters over 40 model variations. This model type, which successively combines simple decision trees to form an ensemble, is know to handle high dimensional, complex data sets well, while being **computationally efficient** and **robust to data skew** and outliers.

**Tuning Objectives**

- **Cross validation** to reduce variance of loss estimates.

- Increase **boosting rounds** to create complex models but limit tree depth and lower column/row sampling rates to reduce overfitting.

- Optimize area **under the precision-recall curve** to account for class imbalance.

| | |
|---|---|
| 4 | max depth |
| 30 | estimators |
| .3 | col sample |
| .8 | row sample |
| 50 | class weight |

# Evaluation & Interpretation

Analysing feature importance highlighted **sex**, **education** and **employment** indicators as most significant, the former reaffirming a well-know **bias** in this historical dataset. Our predictions correctly **identified 35% of high income** earners (recall) in the test set while overall the high income **predictions** being **77% of the time correct** (precision).

**.35** Recall

**.77** Precision



**Trade-offs**

Our model significantly outperformed naive baselines.

The **recall/precision trade-off** can be further tuned by adjusting the predicted probability cutoff threshold.

# Implementation on Dataiku DSS

The complete solution framework was implemented on Dataiku DSS Cloud using Python recipes ran on a custom project environment with 2-CPU cores and 16Gb memory.



The recipe notebooks are synced with Github.

Data stored on DSS managed AWS S3 storage.

# Next Steps

Our next steps will be ultimately informed by the key (business) drivers that initiated the project e.g. to understand the effect of individual predictors vs to optimize predictive power.

### Data Enrichment

Improve the input data quality & quantity:

- **Balance the classes** using up/down sampling or synthetic data generation.

- **Feature engineering &** alternative encodings.

- Adjust **outliers** and **normalize.**

### MLOps & Evaluation

More in-depth **experiment logging** and refined **analysis of the model errors** could reveal model limitations and inform future iterations.

**Feature importances and effects on individual samples** can be better understood by model agnostic methods such as SHAP or LIME.

### Alternative Models

A more thorough search through a **larger hyperparameter space** can further improve model performance.

**Alternative model types** can be easily tested, such as CatBoost which are built for high-dimensional categorical data.

# Q & A

Project: [Github](Github)

Contact: [LinkedIn](LinkedIn)

# Appendix: Bias Mitigation

Previously, I've implemented privacy and fairness risk mitigation with bias-corrected synthetic data on the Adult Census dataset: this allows for privacy-preserving data sharing and also aids the fair treatment of customers (data subjects) in downstream analysis and machine learning tasks. See more here and in our paper accepted at an ICLR 2021 workshop.

# Appendix: End-to-End Flow Build

# Appendix: Experiment Tracking

# Appendix: HP Tuning Analysis