

Coding occupations in the Population and Housing Census of Ecuador using machine learning techniques


National Institute
of Statistics and Census
of Ecuador

Méndez Diana, Espinoza Victor





Table of contents

- 01 ▶ **Introduction.**
 - 02 ▶ **Coding process.**
 - 03 ▶ **Methodology for machine learning projects.**
 - 04 ▶ **Developing machine learning models.**
 - 05 ▶ **Deployment in an integrated system.**
 - 06 ▶ **Key takeaways.**
- 

Introduction

International Standard Classification of
Occupations ISCO – 08
(4th level)



National Standard Classification of
Occupations
(6th level)

In your main job, what is the
occupation or tasks you
perform?



Coding occupations

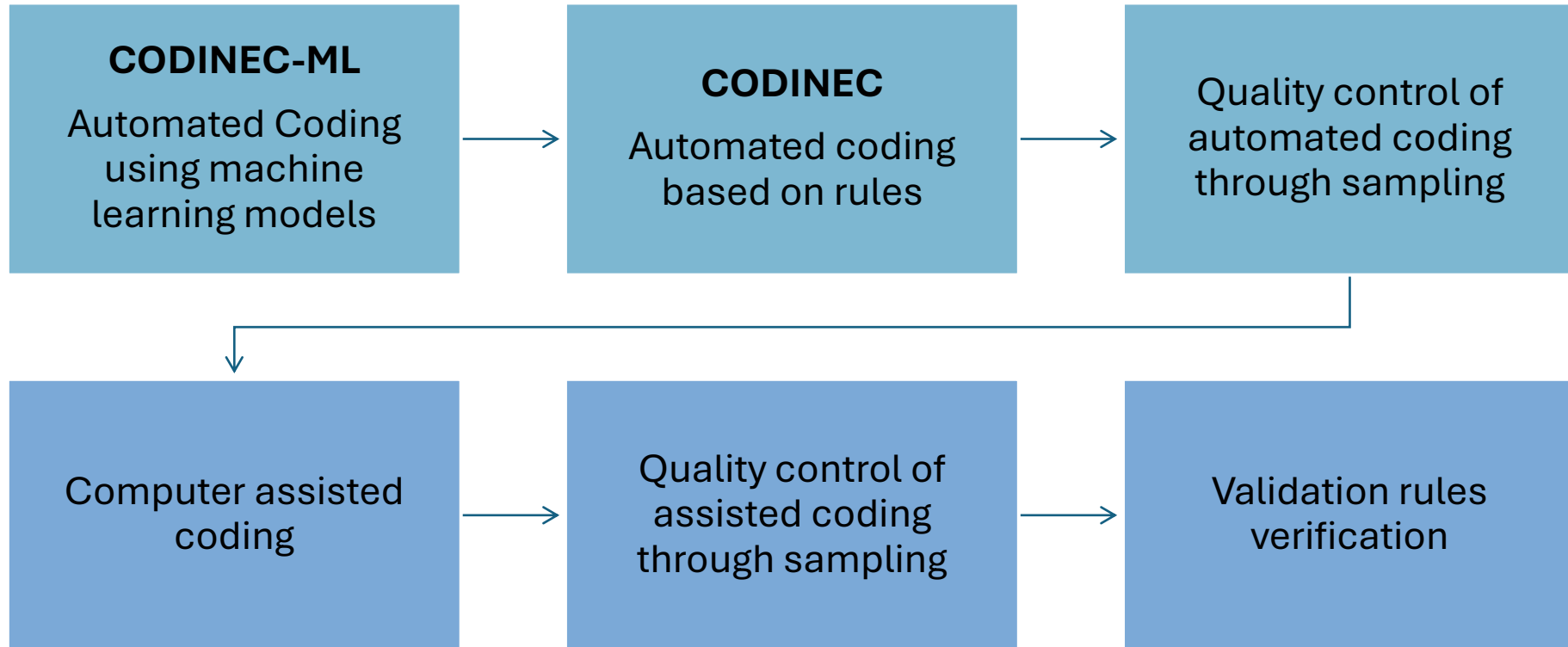


Occupation response

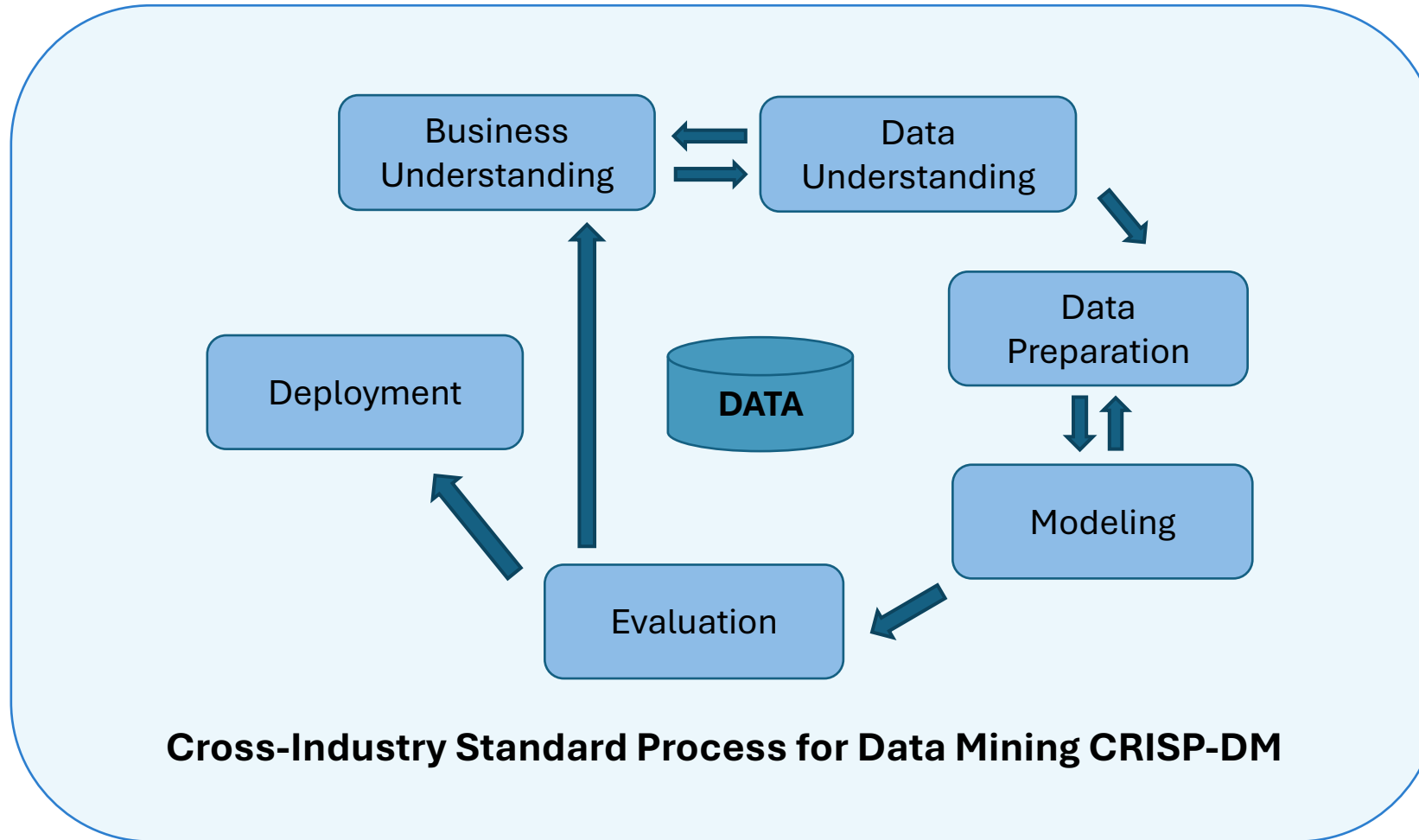
Ancillary information:

- Industry description
- Occupational category
- Age
- Educational level
- Social security

Coding process



Methodology for machine learning projects



Developing machine learning models

Business Understanding

Coding \approx 5.4 million of occupations

Main concerns:

Recruiting , training and managing hundreds of coders.

Ensuring sufficient financial resources for coding staff and operations.

Maintaining consistent quality control.

Solution:

Automating a part of the process using machine learning techniques.

Data Understanding

212.145
occupations
of Labour
Force Survey



114.672
occupations
of Census
Pilot Tests

326,817
occupations.
Ground Truth

Developing machine learning models

Data Preparation

Text preprocessing techniques

- Lowercasing
- Text cleansing: Removing punctuation marks, numbers, extra spaces and accents.
- Tokenization.
- Removing stop words
- Stemming.

Feature extraction techniques

- TF – IDF (Xgboost)
- Word Embeddings (Artificial Neural Networks)

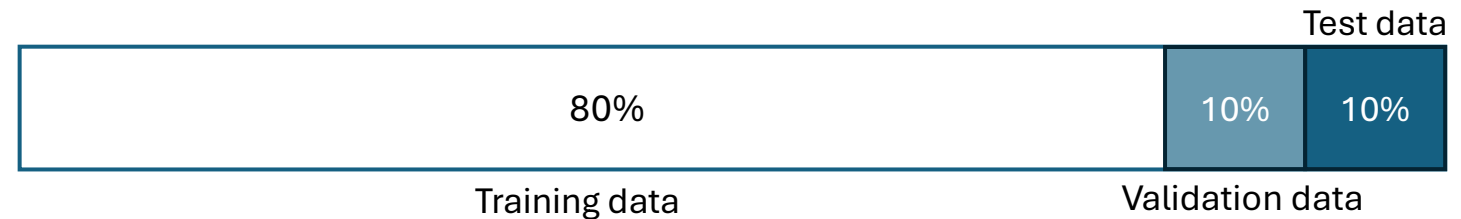
Ancilliary data preprocessing

- One Hot Encoding

Data balancing

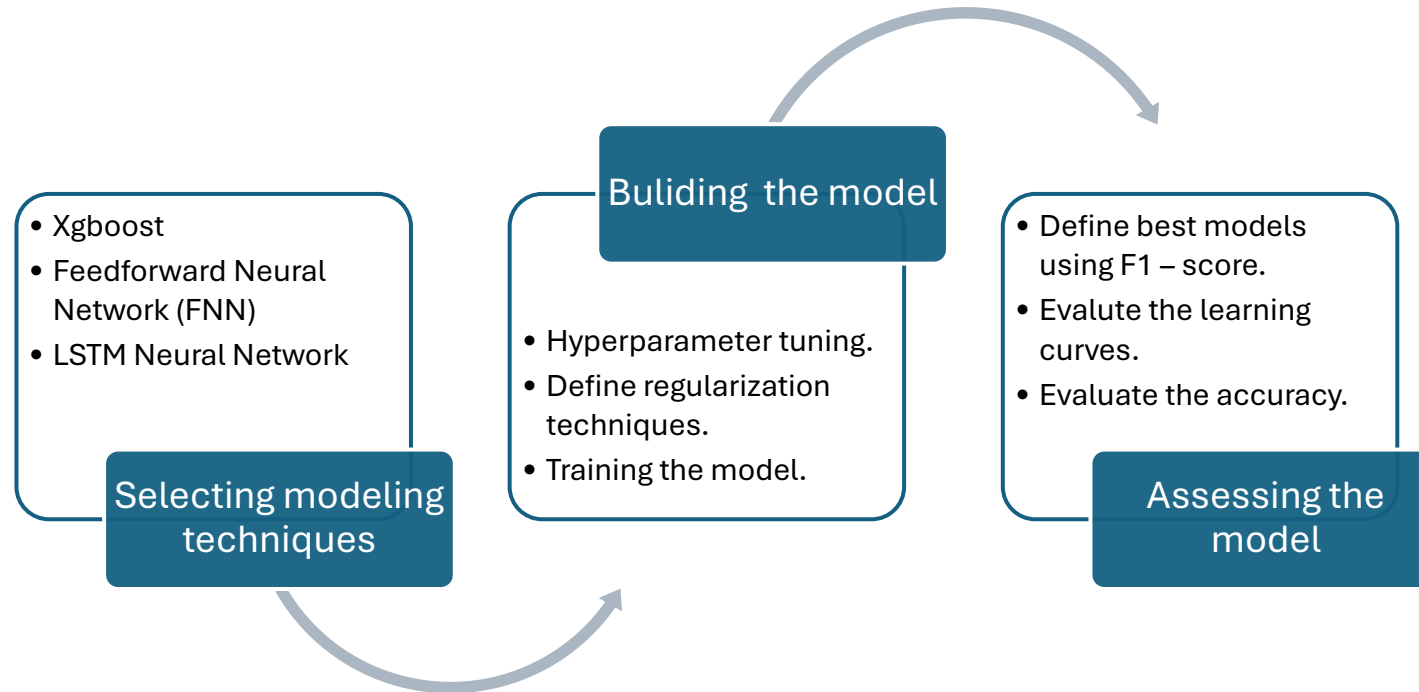
- Resample the training dataset (344 classes): under-sampling and over-sampling

Splitting data (Holdout Validation)



Developing machine learning models

Modeling



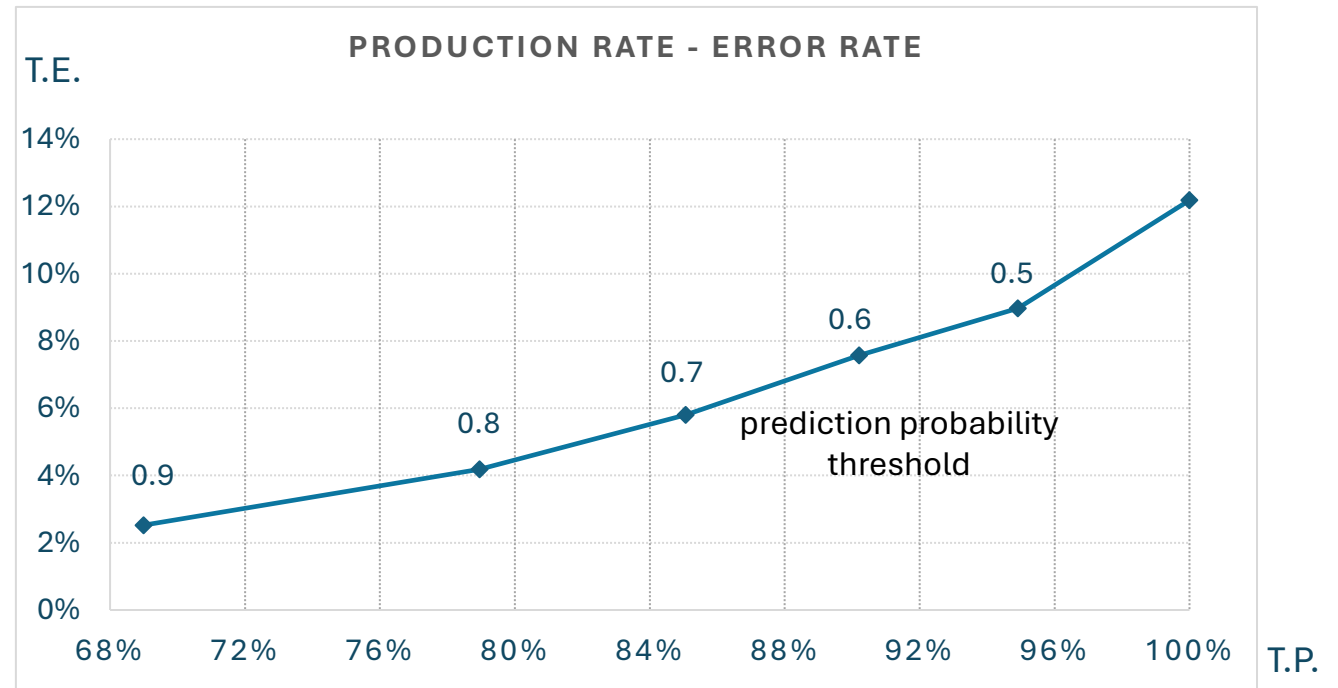
Developing machine learning models

Evaluation

Best model:
**Feedforward
Neural Network**

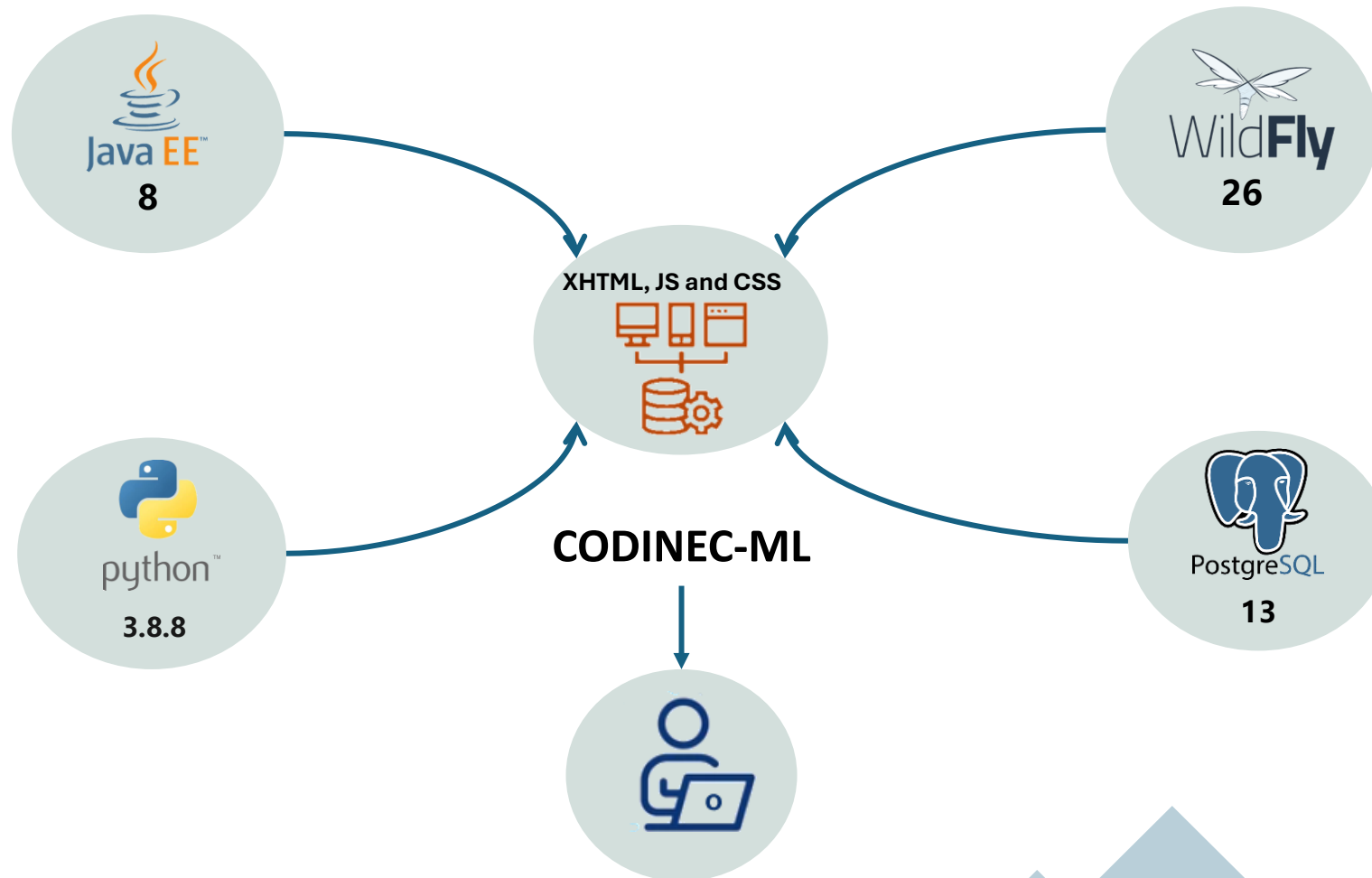
Accuracy: **87.82%**

Training time: **≈ 12
minutes**



Deployment in an integrated system

Web development tools and server



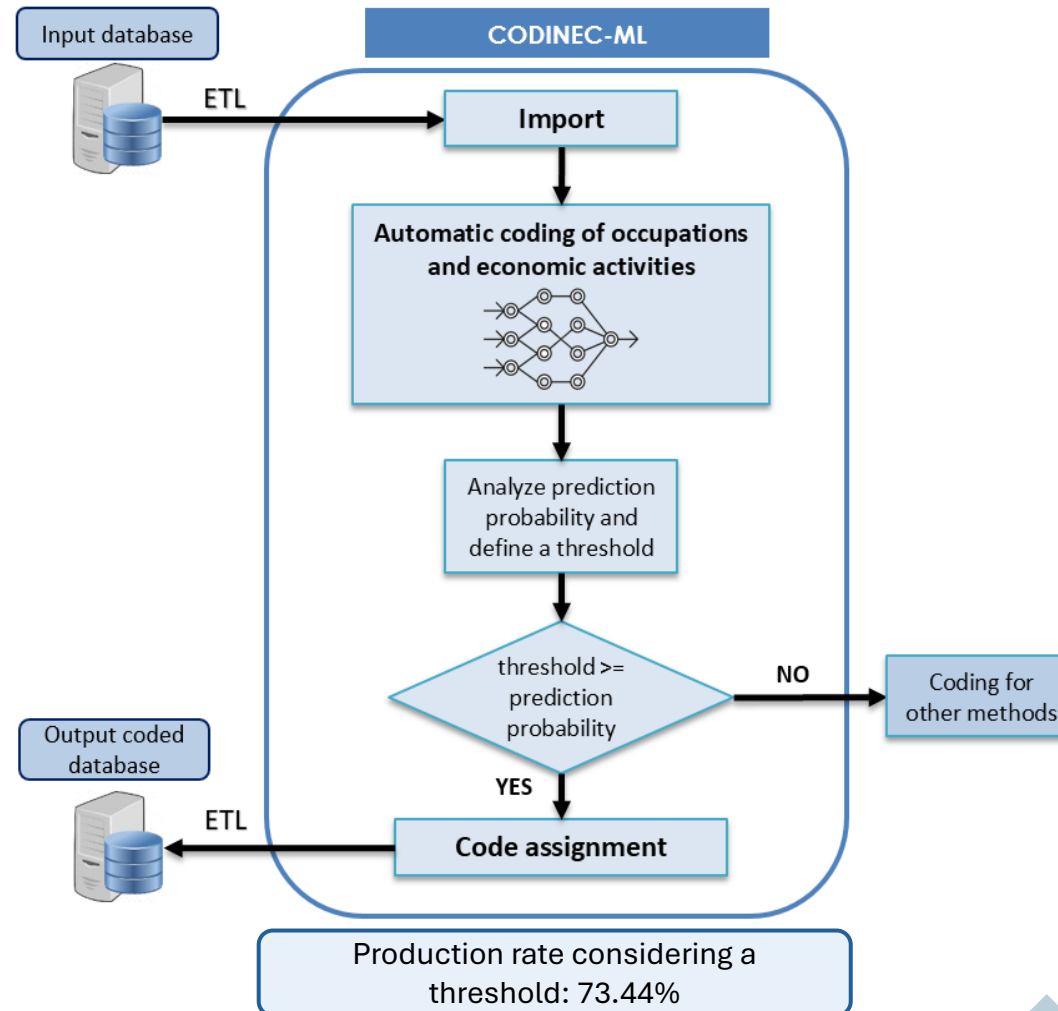
Deployment in an integrated system

Infrastructure Requirements

Feature	Database Server	Application Server	Web Server
Operating System	CentOS 7	CentOS 7	CentOS 7
CPU	24 cores	24 cores	24 cores
Memory	64 GB	32 GB	32 GB
Storage	1 TB	500 GB	500 GB

Deployment in an integrated system

Functionality



Key Takeaways

1

The base of any machine learning model lies in the training data. Ensuring the accuracy of labels is essential, as poor-quality ground truth can lead to unreliable predictions.

2

No model is perfect. However, before deploying a model in production, it must be considered that the error rate of the model is lower than of manual coding.

3

The production rate must be carefully calibrated, balancing available resources with the margin of expected error.

4

During the coding process, having a dedicated quality control team is non-negotiable. They play a critical role in monitoring the model's performance due to ambiguous and vague inputs that can affect the predictions.

The slide features a light blue background with decorative geometric shapes in the corners. The top-left corner has a dark blue square overlapping a medium blue square. The top-right corner has a medium blue square overlapping a light blue square. The bottom-right corner has a medium blue square overlapping a dark blue square. The text is centered in a bold, dark blue font.

**Thank you for your
attention!**

