# **Country House Price Index (HPI)**
## BI and data-science solution for analyses and visualization

**LEO ANUCHIN, CKDU36**
anuchin.leonid@gmail.com
for Business Intelligence Laboratory
[Semester: Spring 2023/2024]

2024/05/12

**Link to GitHub:**
https://github.com/Leo4CDP/HPI_INDEX
**Link to video:**
https://drive.google.com/file/d/18V4iQqaYlJ5G4yNyZr8_TDl9waUI_4Xv/view?usp=drive_link
**Link to Google Colab:**
https://colab.research.google.com/drive/1FVyBKiqb4eGsSNfORimxIGstJJig7Ec7?usp=drive_link

# OUTLINE

**Aim of the work:** the work is to introduce BI and data science methods to work with house price index (HPI) data on the country level in chosen EU countries.

**Data sources:** API - Eurostat data (7 indicators).
**Main indicators to be displayed:** HPI growth, key macroeconomic indicators

**Report features:** during the project the database of indicators was gathered with use of API + ETL data transformations. The data was saved in database (SQL), and presented via Power BI.

**Notebook:** ML methods revealed the dependence of HPI from other macroeconomic indicators (p-value), and there was an attempt to forecasting its dynamics (ARIMA+regressors).

# USED APPLICATIONS

| TASK | TECHNOLOGY |
|------|------------|
| **ETL (Extract, transform, load)** | |
| 1. Data parsing via API | python script |
| 2. ETL engine | VS – SSIS + python scripts |
| 3. Database(s), special data storage solutions | MS SQL (Server name: MSI, database: CKDU36) |
| **Reporting engine** | |
| 4. Reporting engine | PowerBI Desktop |
| 5. Unique data visualization | tables, diagrams, geo-spatial data |
| **Data final preparation and analyses** | |
| 6. Technologies used for data analysis/data science functions | Google Colab – Jupiter Notebook python libraries: pandas, sklearn, mathplotlib, ipywidgets |

# WHAT IS HPI?

**HPI - House Price Index**. Indicator of house price trends. The HPI also functions as an analytical tool for estimating changes in the rates of mortgage defaults, prepayments, and housing affordability. [Investopedia, https://www.investopedia.com/terms/h/house-price-index-hpi.asp]

This **indicator shows the state of the residential real estate market at the country level** (however, in different parts of the country - cities and regions, the situation may vary).

The indicator is useful for **assess the economic situation**, the pace of its development, **monitoring the affordability and effectiveness of investments in residential real estate**.

# WHAT OTHER INDICATORS WE USE AND WHY?

**KEY MACROECONOMIC INDICATORS:**
- **I10GDP** – GDP growth data (index, quarterly)
- **I05HICP_T** – Inflation (customer price index) data (index, monthly – recalculated to quarterly)
- **I05HICP_R** – Housing rent prices data (index, monthly – recalculated to quarterly)
- **UNEMPLOYMENT** – Level of unemployment (15-72,%, quarterly)
- **I10LABOUR** – Spendings on labour (index, quarterly)

**KEY NATIONAL INDICATORS:**
- **POPULATION** – Total population (people, annual)

The indicators are partly based on report of BIS Determinants of house prices in central and eastern Europe (September 2007, ISSN 1682-7678), https://www.bis.org/publ/work236.pdf. The report considers this factors among connected to HPI.

# ETL PROCESS

**Link to GitHub:**
https://github.com/Leo4CDP/HPI_INDEX

# PREPARATION

1. Work folder is C:\SriptsHPI

2. When  first start the project, **firstly the SQL table - CKDU36.sql is deployed**.

3. Some steps (python scripts) **need connection to Python program**. No extra packages needed.

4. The first launch occurs using the **EurostatDownloads.csv** database. The file is used to monitor current updates and store API request codes.

# STEPS DESCRIPTION

| # | Step | Description |
|---|------|-------------|
| 1 | NEW VERSION CHECK SCRIPT | SCRIPT1_Updates_checker.py is used.<br>It checks the last updates on Eurostat. If there are new updates it marks EurostatDownloads.csv rows with "1".<br>*If no updates are found if finishes SSIS (With Error). This may be handled by Send Email Task* |
| 2 | DOWNLOAD UPDATES | SCRIPT2_xmldownload.py is used.<br>The xml files are downloaded and stored to xml12 folder. This script used codes from EurostatDownloads.csv (flags + API keys) to facilitate download process (<1 MB instead of 50 MB) |
| 3 | CLEAN_DOWNLOAD_MARK | SCRIPT2_1_CLEAN6.py is used.<br>The script changes "1" in EurostatDownloads.csv to "0" to prevent unnecessary updates. |
| 4 | PARSE XML TO CSV | SCRIPT3_xml_parser.py is used.<br>It parses XML files to CSV on pre-elaborated chema. |
| 5 | SQL TR | truncate table STATDATA; truncate table I15HPI; truncate table UNEMPLOYMENT; truncate table I10LABOUR; truncate table POPULATION; truncate table I05HICP_T; truncate table I05HICP_R |
| 6 | IMPORT 5 CSV TO SQL SERVER | Import first 5 CSV to SQL (separate tables) They don't need further processing on this stage. |
| 7 | JOIN 3 TABLES ON BASE OF GDP SQL TABLE | We sort the rows (on KEY1) and using OLE BD JOIN the tables. |
| 8 | CLEAN SQL NA | In resulted STATDATA table we clean all rows with N/A (empty data) |

# STEPS DESCRIPTION

| # | Step | Description |
|---|------|-------------|
| 9 | CALCULATE NUMBER OF ROWS PER COUNTRY | We calculate number of rows per country. It is needed not to have countries that don't have enough information in Eurostat. For example the UK stopped providing data science 2020. Greece doesn't provide HPI. And etc. |
| 10 | KICK COUNTRIES WITH LOW ROWS | We kick countries with less then 56 rows (based on #9 results). |
| 11 | DICTIONAY TO TRANSFER MONTHLY CSV 2 QUATERLY CSV | Now, we can create dictionary of periods (YYYY-QX) and dictionary of countries. It may be used to parse xml of CPI and Rent prices data. These dictionaries help to avoid any mistakes on step 12 of process. |
| 12 | PARSE XML TO CSV 1 | SCRIPT4_M2Q.py is used. It parses CPI and Rent prices data from xml to csv format, based on Dictionaries from step 11. It also re-calculates monthly data to quarterly date. |
| 13 | IMPORT 2 CSV TO SQL SERVER 1 | Now we send csv received on step 12 to SQL tables. |
| 14 | JOIN 2 NEW TABLES TO GDP SQL TABLE | We sort the rows (on KEY1) and using OLE BD JOIN the tables. |
| 15 | CSV for Notebook | This step is export of STATDATA table to CSV4NOTEBOOK.csv. This file is needed to Google Colab (forecast) |

# RESULTS OF ETL PROCESS



**Both tables are accessible via Power BI.**

The data may be exported for different needs (for example in csv format to Jupiter Notebook)

# POWER BI

**Link to GitHub:**
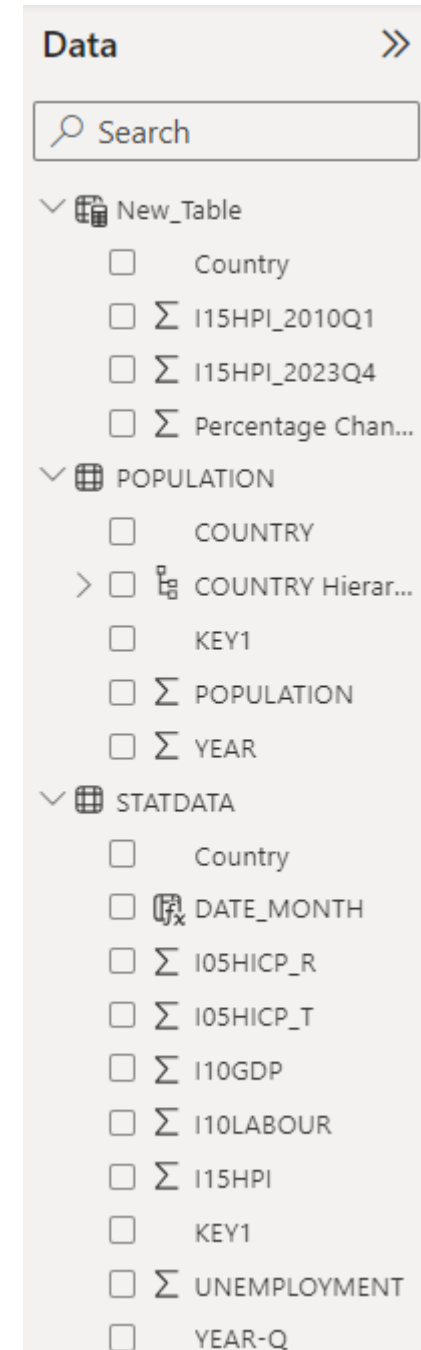https://github.com/Leo4CDP/HPI_INDEX

# ABOUT BI

**Power BI connect to MS SQL.**

1. It connects to STATDATA table (with all indicators in quarterly format.

As supplementary Power BI changes XXXX-QY format to date format to build diagrams (New column).

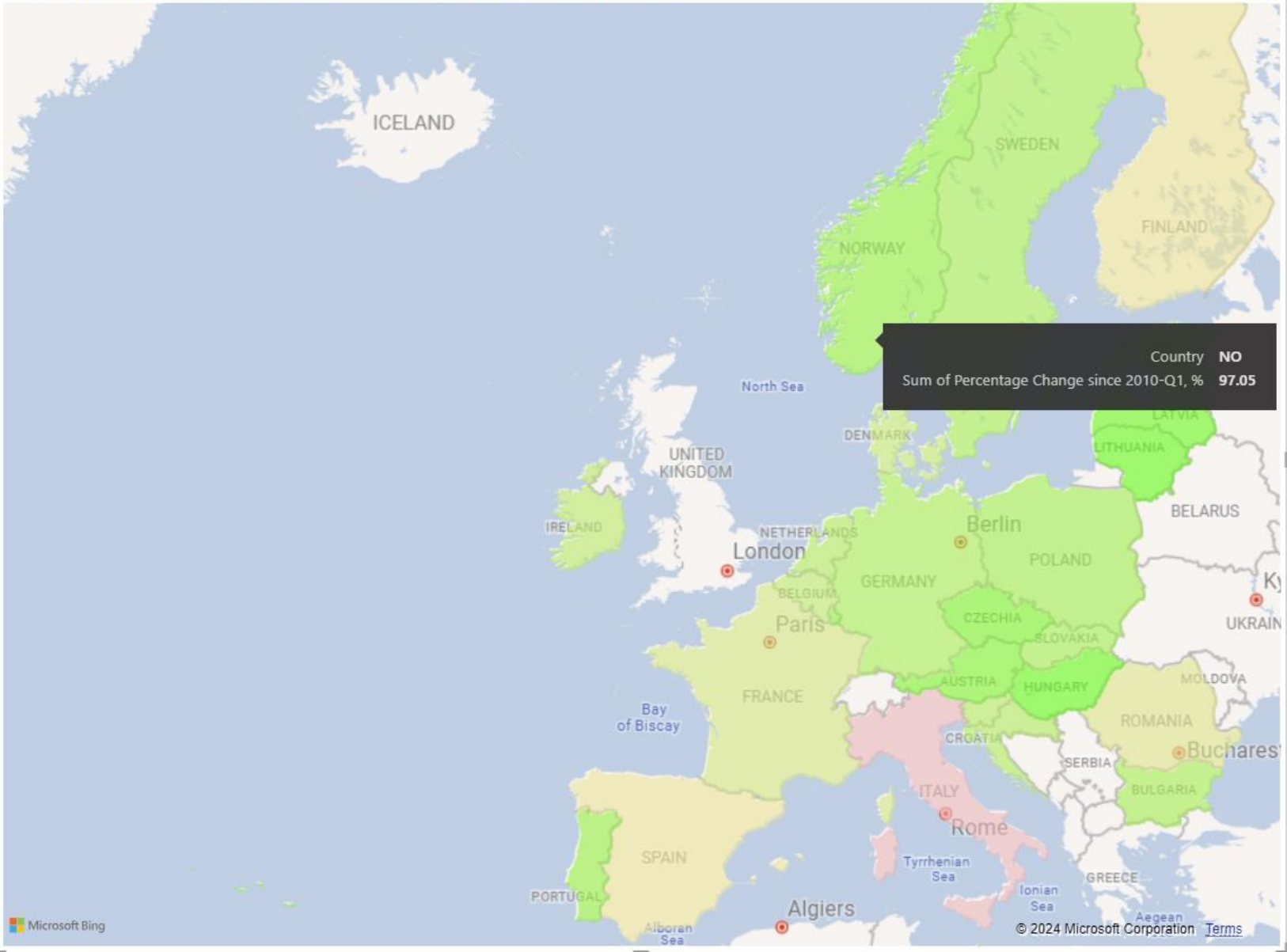2. It connects POPULATION database and is able to show population changes (annual basis).

3. New_Table was created to check the HPI change since 2010. This table is used to build geospatial map.

**Data** »

Search

- ∨ New_Table
  - ☐ Country
  - ☐ Σ I15HPI_2010Q1
  - ☐ Σ I15HPI_2023Q4
  - ☐ Σ Percentage Chan...
- ∨ POPULATION
  - ☐ COUNTRY
  - > ☐ COUNTRY Hierar...
  - ☐ KEY1
  - ☐ Σ POPULATION
  - ☐ Σ YEAR
- ∨ STATDATA
  - ☐ Country
  - ☐ DATE_MONTH
  - ☐ Σ I05HICP_R
  - ☐ Σ I05HICP_T
  - ☐ Σ I10GDP
  - ☐ Σ I10LABOUR
  - ☐ Σ I15HPI
  - ☐ KEY1
  - ☐ Σ UNEMPLOYMENT
  - ☐ YEAR-Q

REPORT ON HPI AND CONNECTED MACROECONOMIC INDICATORS (EUROSTAT DATA, 2010-2023)

# COUNTRY REPORT (1)



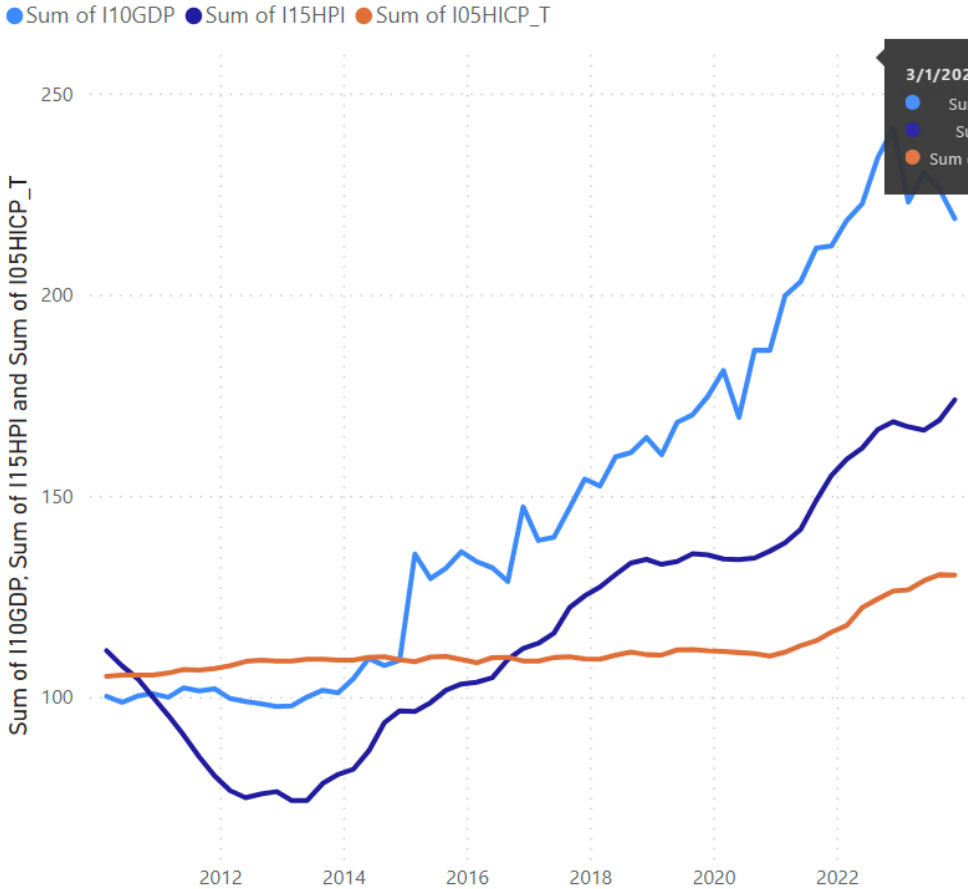Change of House prices in different countries since 2010-Q1

| Country | Percentage Change since 2010-Q1, % |
|---|---|
| EE | 222.09 |
| HU | 181.84 |
| LT | 166.34 |
| LV | 140.12 |
| CZ | 123.00 |
| AT | 117.67 |
| LU | 102.36 |
| NO | 97.05 |
| PT | 96.09 |
| PL | 83.29 |
| SK | 81.72 |
| BG | 79.87 |
| SE | 78.89 |
| DE | 75.30 |
| MT | 69.67 |
| NL | 66.70 |
| HR | 66.18 |
| SI | 62.62 |
| IE | 55.90 |
| BE | 54.34 |
| DK | 48.08 |
| FR | 33.19 |
| RO | 22.88 |
| FI | 11.85 |
| ES | 11.22 |
| CY | -3.86 |
| IT | -6.93 |

# COUNTRY REPORT (2)

**MAIN MACROECONOMIC INDICATORS VS HPI**



GDP, CPI (Inflation) and HPI

●Sum of I10GDP ●Sum of I15HPI ●Sum of I05HICP_T

LABOUR, RENTS vs HPI (income and prices)

●Sum of I15HPI ●Sum of I10LABOUR ●Sum of I05HICP_R

3/1/2022 12:00:00 AM
Sum of I10GDP  218.44
Sum of I15HPI  159.07
Sum of I05HICP_T  117.77

Search

Filters on this page

Country
is IE

Search

☑ IE    56
☐ IT    56
☐ LT    56
☐ LU    56
☐ LV    56
☐ MT    56
☐ NL    56

☑ Require single selection

Add data fields here

Filters on all pages

Add data fields here

# COUNTRY REPORT (3)

# COUNTRY REPORT (4)

**SCATTER CHART: RELATIVITY OF SOME INDICATORS CHANGE**

# **FORECASTING**
# **(FIRST ATTEMPT)**

(Jupiter Notebook python)

**Link to Google Colab:**
https://colab.research.google.com/drive/1FVyBKiq
b4eGsSNfORimxIGstJJig7Ec7?usp=drive_link

**Link to CSV:**
https://drive.google.com/file/d/1qd3XR6Gz8ACCRp
jY0QBW7n1WpTIrZv3h/view?usp=drive_link

# ABOUT HPI FORECASTING

**WHY IS IT IMPORTANT TO FORECAST HPI:**
- The HPI measures the cost and affordability of housing for residents. Its prediction allows making decisions that affect the availability of demand (mortgages, preferential programs) and supply (construction volumes) in the medium term.
- the indicator lags behind other macroeconomic indicators - it is prepared six months later than data on inflation and economic growth.

**INPUTS**

To solve the problem of HPI forecasting at national level, we will use a previously generated data array (based on Eurostat data). **We use CSV4NOTEBOOK.csv – one of products of ETL process. It is to be put to Content folder of Notebook.**

A preliminary analysis of the data showed that the following **restrictions have a significant impact on the results of predictions**:

> 1) the **economic crisis of 2008**, its effects are observed until the end of 2013-2015. Intervals before 2010 were completely excluded from analyses. (~Dummy-variable);
> 2) **Covid effect**. In almost all countries, the macroeconomic indicators studied (for example, GDP, inflation) responded to the pandemic faster than real estate prices, which also somewhat reduced the quality of predictions.

# MODEL DESCRIPTION

Preliminary tests showed that **ARIMA (1,1,1) + regressors** and **Linear regression (=ARIMA (0,0,0))** show the most alike results to the expected. The sample for each country **(with ipywidgets)** is divided on 2 parts (parameters may be adjusted according to macroeconomic situation in national economic):

- test data (after 2014-01-01)
- validation data (after 2021-01-01)

**Used forecasting models\*:**

- ARIMA (1,1,1) with regressors
- ARIMA (0,0,0) with regressors (=linear regression)

**Regressors:**

- **I10GDP** – GDP growth data (index)
- **I05HICP_T** – Inflation (customer price index) data (index)
- **I05HICP_R** – Housing rent prices data (index)
- **UNEMPLOYMENT** – Level of unemployment (15-72,%)
- **I10LABOUR** – Spendings on labour (index)

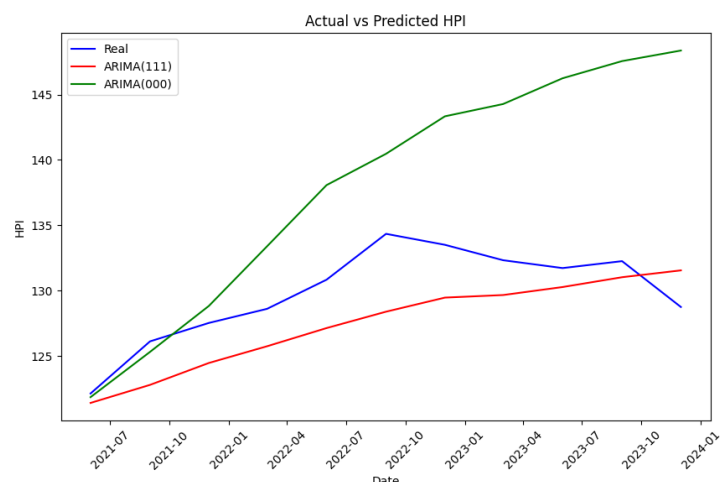**Jupiter Notebook calculates MSE for each model, and shows  p-value of every regressor.**

\* Test on neuro networks models (Keras) showed worser results, probably because of limited sample rows.
\*\* Population indicator was excluded from analyses on this stage, as for some countries (it shows jumps and falls due to statistical methodology and updates on census results)

# Examples of predictions: COVID-effect

Economic and historic factors influence forecast results. For example, **if we start forecast since 2020-01-01 Hungary macro indicators (GDP, CPI, unemployment and etc.) suffered less than in France, and for Hungary model shows better quality.** But if we take **later date (2021-06-01)**, we see, that **model based on macroeconomic regressors expected better HPI growth, than it occurred in fact.**
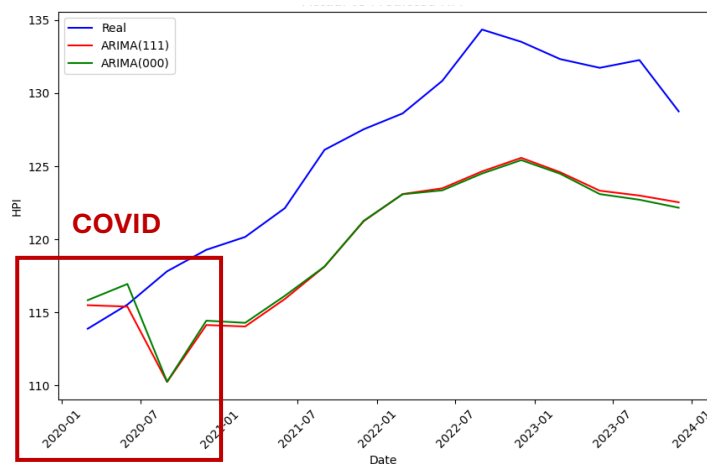
## FRANCE



**FRANCE**
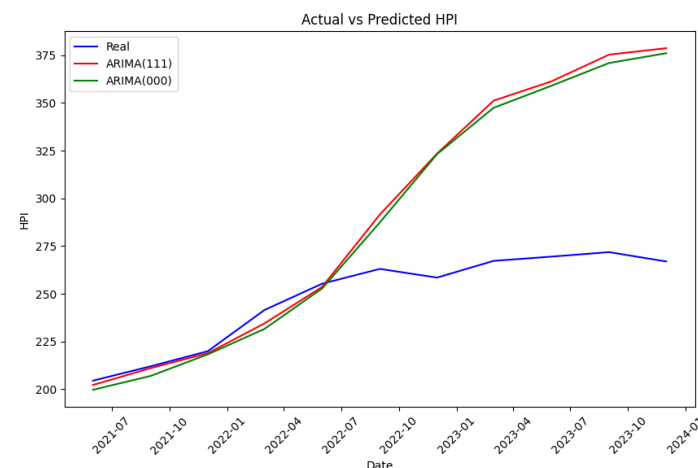train_start = "2013-01-01"
valid_start = "**2021-06-01**"

ARIMA 111 - MSE: 107.7782
**ARIMA 000 - MSE: 10.3012**



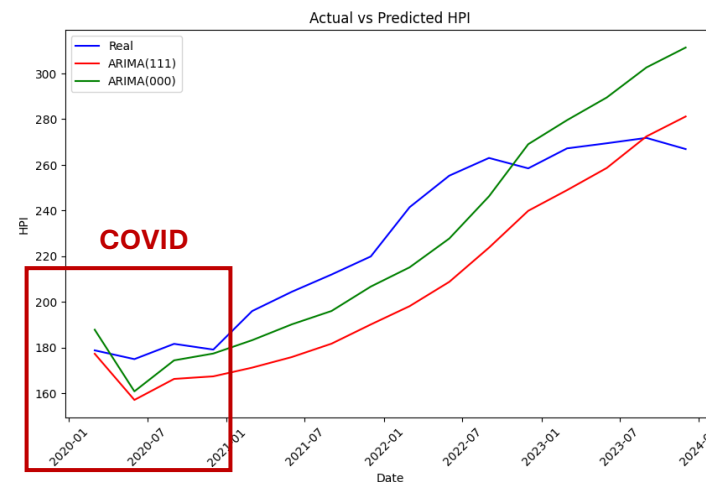train_start = "2013-01-01"
valid_start = "2020-01-01"

ARIMA 111 - MSE: 48.7467
ARIMA 000 - MSE: 47.6718

## HUNGARY



train_start = "2013-01-01"
valid_start = " 2021-06-01"

ARIMA 111 - MSE: 3740.1201
ARIMA 000 - MSE: 3981.8389



train_start = "2013-01-01"
valid_start = "**2020-01-01**"

**ARIMA 111 - MSE: 403.6217**
ARIMA 000 - MSE: 656.9420

# KEY PRELIMINARY FINDINGS

- The **model allows you to predict the level of HPI depending on other macroeconomic indicators at the national level**, even in conditions of significant turbulence in the macroeconomic situation currently observed. If we use available macroeconomic forecasts (like Central bank forecasts and etc, we may predict house price growth at national level).

- **Model quality results are different for each country**. The most strongly influenced indicators are those associated with a long test base - since countries unevenly emerged from the 2008 crisis (for example, for Romania the model gives the most accurate results if the test sample is taken from 2017; and for Hungary from 2015).

- The **model can be improved by adding additional data sources** (construction volumes, bank lending data) and clarifying HPI data based on real sector data. Also, it is advisable to analyze major economies to find conditions (test base, model parameters) under which the model shows greater accuracy.

# FUTHER STEPS

# FUTHER STEPS

- Improving the forecasting model (selection of parameters for large national economies, expanding the list of parameters by including information on new construction)
- Testing ETL solution scripts over longer periods
- Scientific publication based on p-value and forecasting results

# THANK YOU!

# ANNEX: P-VALUE

ARIMA(0,0,0)
print(model_fit_0.summary())

```
                               SARIMAX Results
==============================================================================
Dep. Variable:                        y   No. Observations:                   21
Model:                            ARIMA   Log Likelihood                 -32.307
Date:                  Sun, 12 May 2024   AIC                             78.614
Time:                          17:08:25   BIC                             85.926
Sample:                               0   HQIC                            80.201
                                   - 21
Covariance Type:                    opg
==============================================================================
                 coef    std err          z      P>|z|      [0.025      0.975]
------------------------------------------------------------------------------
const       -215.2462     28.173     -7.640      0.000    -270.465    -160.028
x1             1.4689      0.701      2.096      0.036       0.095       2.842
x2            -0.0677      0.347     -0.195      0.845      -0.747       0.612
x3             0.9971      0.144      6.901      0.000       0.714       1.280
x4            -1.6161      1.745     -0.926      0.354      -5.036       1.804
x5             0.2206      1.072      0.206      0.837      -1.881       2.322
sigma2         1.2698      0.761      1.668      0.095      -0.222       2.762
===================================================================================
Ljung-Box (L1) (Q):                   0.00   Jarque-Bera (JB):                1.53
Prob(Q):                              0.97   Prob(JB):                        0.46
Heteroskedasticity (H):               0.70   Skew:                            0.59
Prob(H) (two-sided):                  0.64   Kurtosis:                        2.39
===================================================================================
```

**Regressors:**
x1 - I10GDP
x2 - I05HICP_T)
x3 - I05HICP_R
x4 - UNEMPLOYMENT
x5 - I10LABOUR