# Time-series analysis and prediction for the HICP dataset of Eurostat

*Capstone Project Defence*
*EPFL Extension School*

*Cornelia Blanke*
*24/01/2022*

# Overview of my Presentation

- Idea of the Project

  ➜ *Origin, Structure and Explanation of the HICP Data*

- Exploratory Data Analysis

  ➜ *Some selected Findings*

- Machine Learning Task 1

  ➜ *Analyze the Feature Importance for the Weights dataset*

- Machine Learning Task 2

  ➜ *Predict year 2019 out of previous years by Machine Learning*

# HICP = Harmonized Index of Consumer Prices

**Definition**

- The Harmonized Index of Consumer Prices (HICP) is an indicator that the member states of EU and EFTA calculate based on a harmonized method and that allows comparing inflation internationally.

- The European Statistical Office **Eurostat** publishes every month the results of all participating countries.
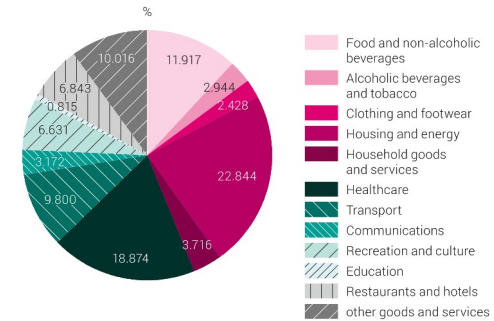
Source: FSO

# HICP = Harmonized Index of Consumer Prices

## The weights

- describe the composition of product groups in the consumer basket

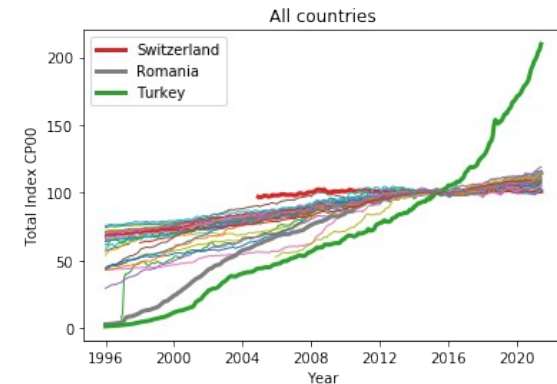- are different in each country

- are yearly updated

## The indices
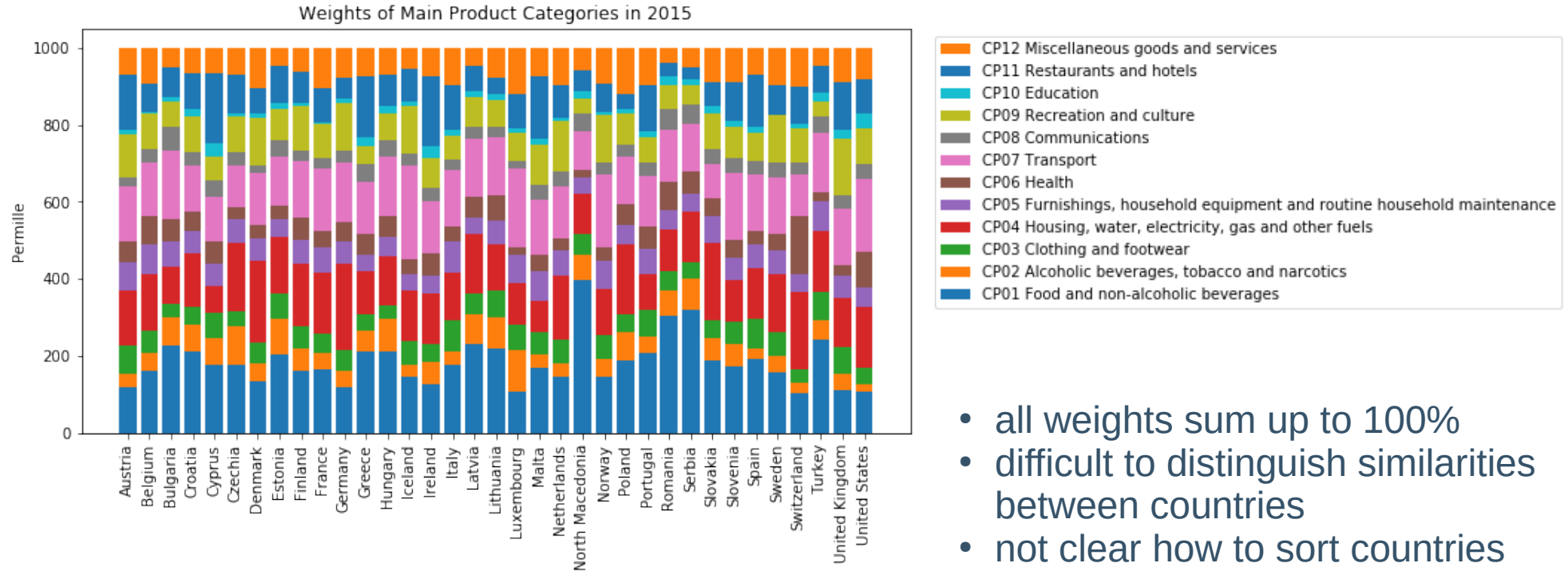
- describe the temporal development of the prices of product groups
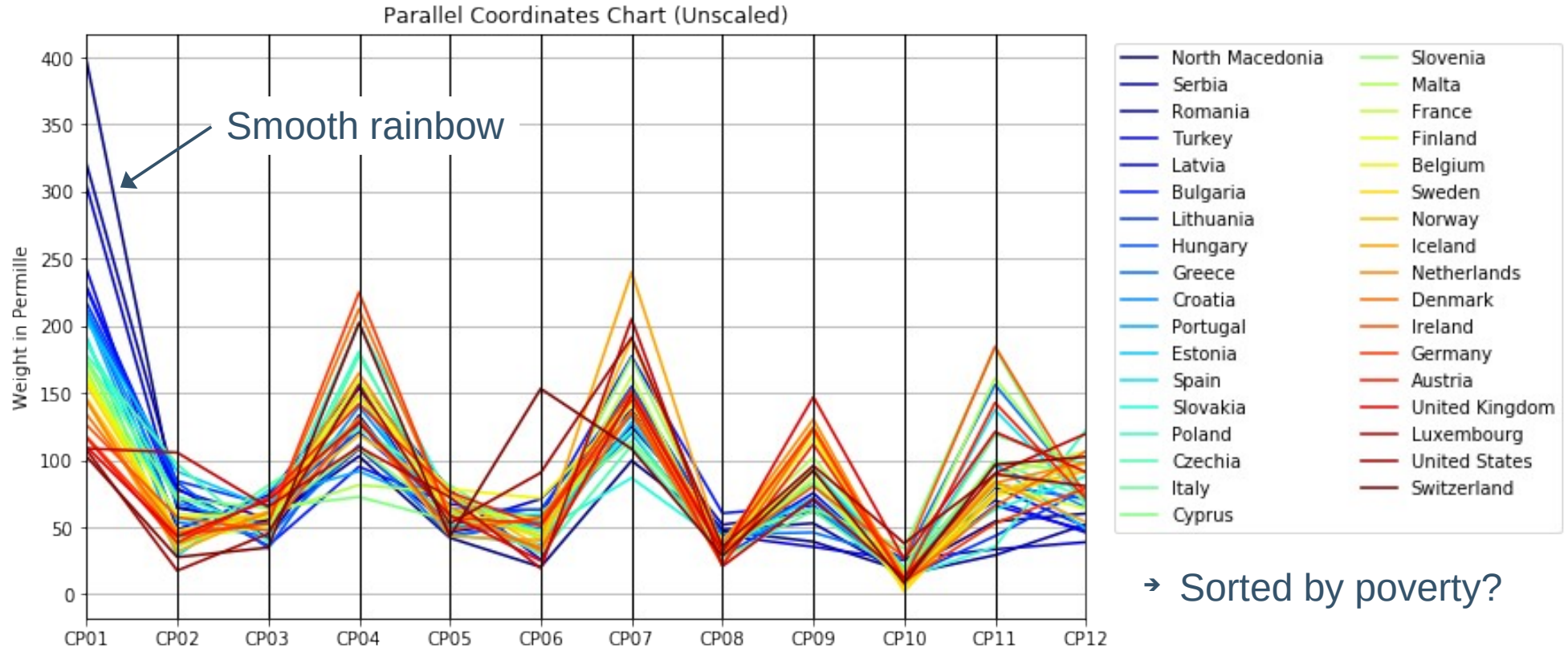
- are relative to 2015 = 100

- are monthly updated



HICP basket and weights, 2021

Food and non-alcoholic beverages
Alcoholic beverages and tobacco
Clothing and footwear
Housing and energy
Household goods and services
Healthcare
Transport
Communications
Recreation and culture
Education
Restaurants and hotels
other goods and services

Source: FSO – Harmonised Index of Consumer Prices (HICP)     © FSO 2021



All countries

Switzerland
Romania
Turkey

# Exploratory Data Analysis

*Some selected Findings*

# EDA – Analysis of the weights



Weights of Main Product Categories in 2015

- all weights sum up to 100%
- difficult to distinguish similarities between countries
- not clear how to sort countries

Parallel Coordinates Chart (Unscaled)

Smooth rainbow

➔ Sorted by poverty?

# EDA – Time Series Analysis

Indices of interest for ML task 2



2015 = 100%

Cornelia Blanke | Time-series analysis and prediction for the HICP dataset of Eurostat

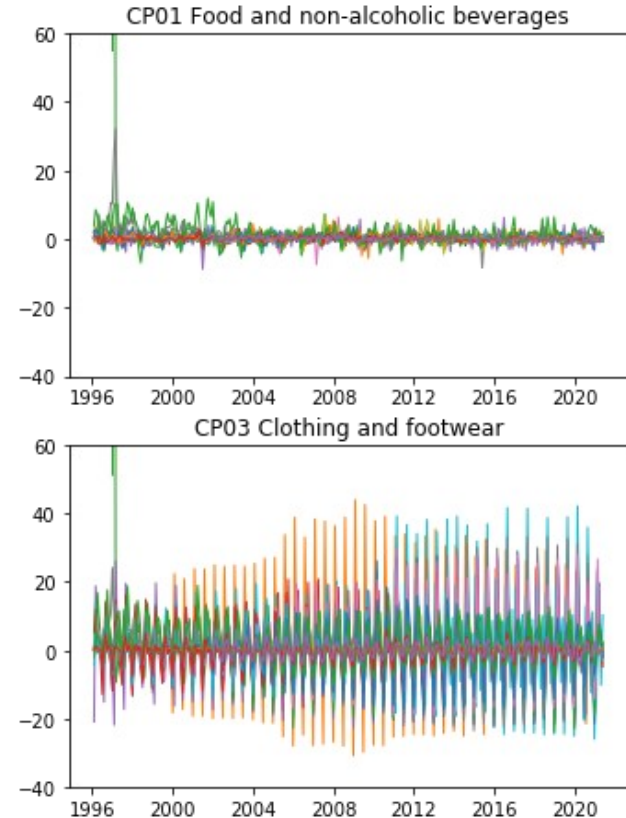# EDA – Compute the monthly rates of change

Often used in ML:

The (absolute) change

$$\Delta I(t_n) = I(t_n) - I(t_{n-1})$$

Here:

The (relative) rate of change

$$M(t_n) = \left( \frac{I(t_n) - I(t_{n-1})}{I(t_{n-1})} \right) \cdot 100\%$$



CP01 Food and non-alcoholic beverages



CP03 Clothing and footwear

# Machine Learning Task 1

*For each of the main product categories 1-12, which characteristics (= features) have the strongest impact on the weights dataset?*

*Use ML techniques to judge about **feature importance**.*

# ML1 – Data Preparation

|  | Feat. 1 | Feat. 2 | Feat. 3 | ... | Feat. 29 | Weight CP01 | Weight CP02 | ... | Weight CP12 |
|---|---|---|---|---|---|---|---|---|---|
| Country 1 |  |  |  |  |  |  |  |  |  |
| Country 2 |  |  |  |  |  |  |  |  |  |
| Country 3 |  |  |  |  |  |  |  |  |  |
| ⋮ |  |  |  |  |  |  |  |  |  |
| ⋮ |  |  |  |  |  |  |  |  |  |
| ⋮ |  |  |  |  |  |  |  |  |  |
| Country 31 |  |  |  |  |  |  |  |  |  |

Features = Characteristics of Countries          Targets = Weights in 2015

# ML1 – Linear Regression

If all features are **scaled** to the same range, their coefficients in a **Linear Regression** model reflect their importance.
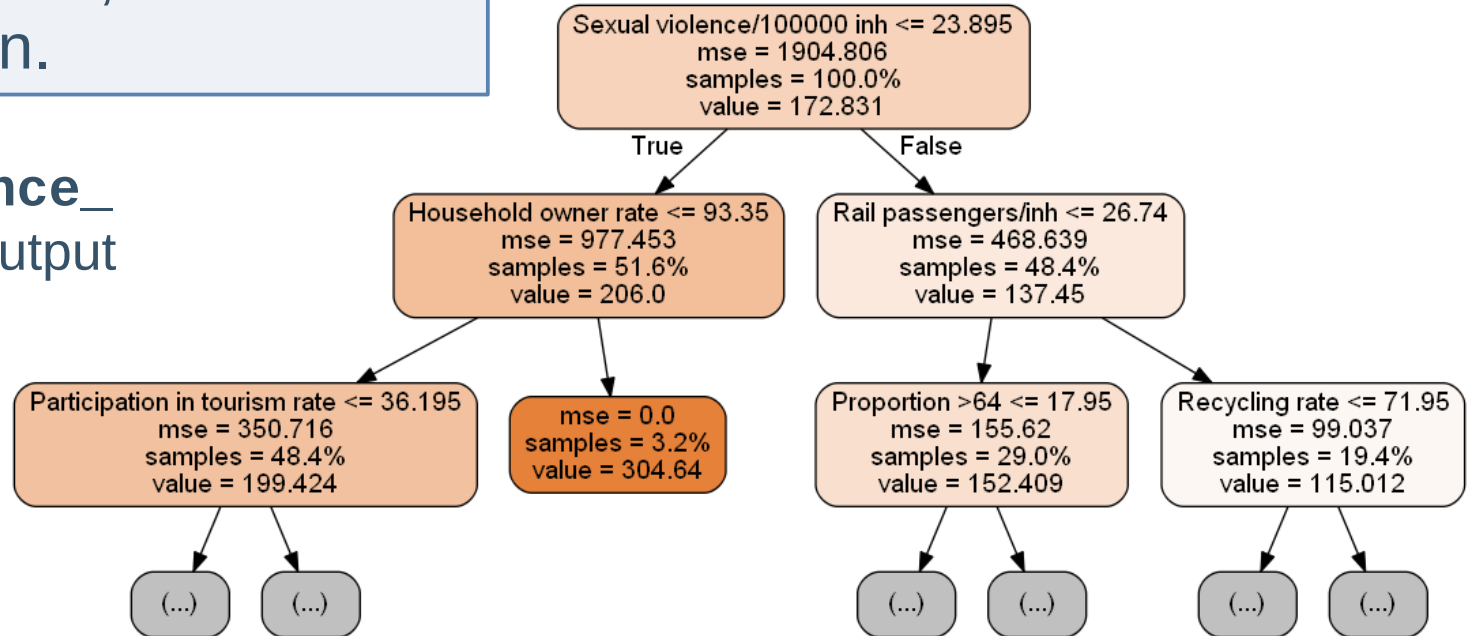
But:

- 29 features, 31 samples
- no validation or test set
- ➔ high risk of over-fitting
- ➔ only use 5 features



Features of 'df_task1_reduced'

# ML1 – Decision Tree Regressor

A **Decision Tree** provides an intuitive visualization about which feature matters most, which is the second, and so on.

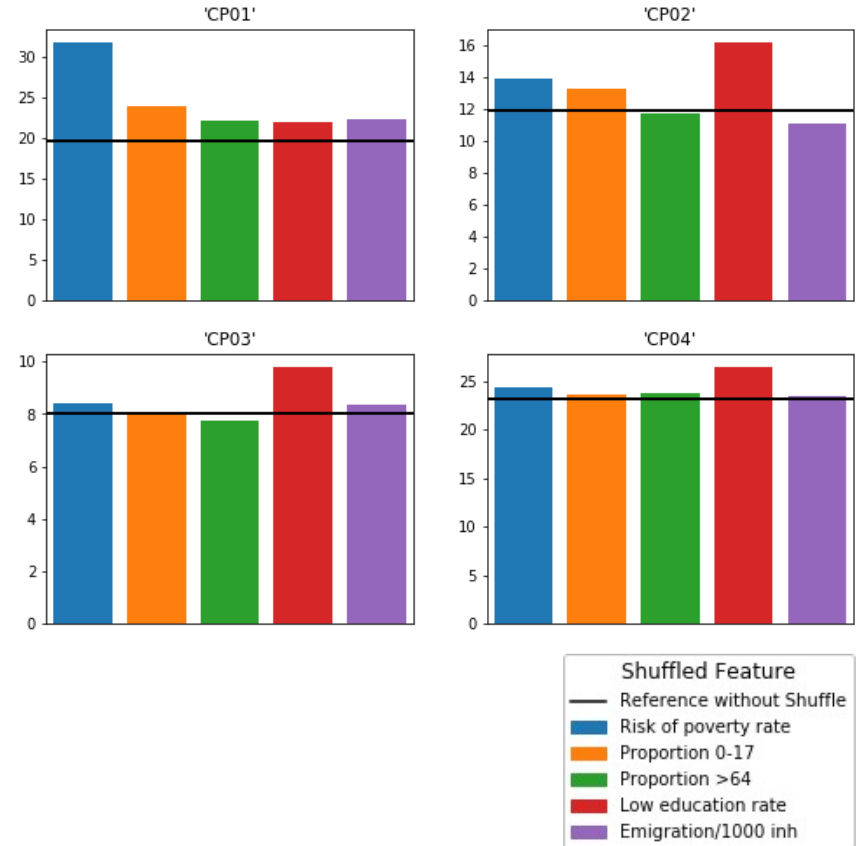`feature_importance_`
provides numerical output

Idea of **feature permutation**:
If it does not matter if a feature were disturbed, then that feature was not an important one.

Algorithm:

1. Shuffle each feature randomly

2. Repeat it five times

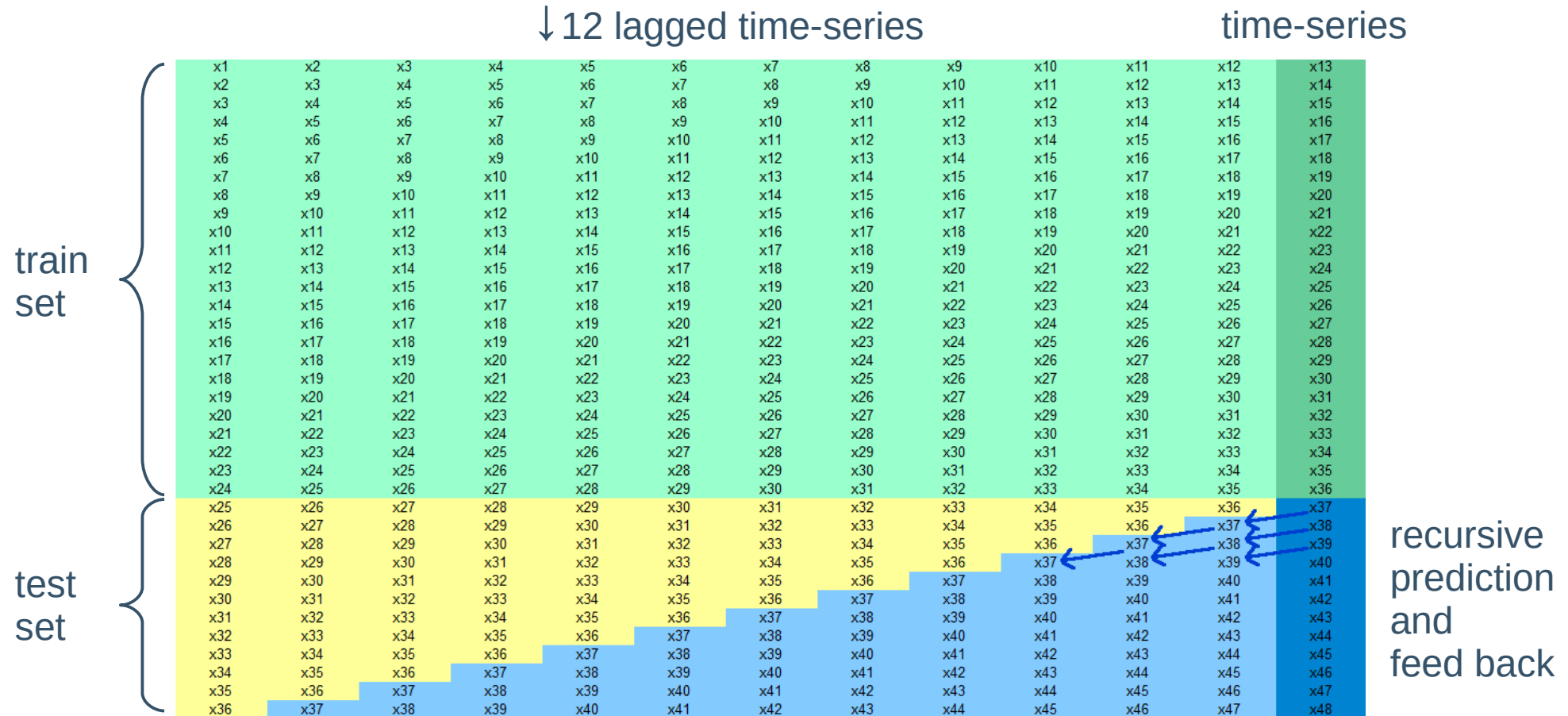3. Compare MAEs (e.g. in kNN)

# Machine Learning Task 2

*For 20 pre-selected indices,*
*predict all months (i.e. 12) of year 2019 out of years 2005-2018.*

*Use different approaches and ML models and compare.*

# ML2 – Recursive prediction of one time-series

1. Data preparation: create 12 lagged features

2. Split into train and test set

3. Learn a ML model that is able to predict the next month out of 12 previous months

4. Use ML model to predict the next month

5. Feed this result to the known values and predict therefrom the next month

Recursion

# ML2 – Recursive prediction of one time-series

# ML2 – Recursive prediction of one time-series

```python
from sklearn.linear_model import LinearRegression
lr1 = LinearRegression()

# initialize y_pred
y_lr1_pred = np.empty((20,12))

# outer loop over all indices
for i in np.arange(df_task2.shape[1]):

    # fit LR model for index i
    lr1.fit(X_train[i], y_train[i])

    # create a local copy
    X_pred = X_test[i].copy()

    # inner loop over all months
    for month in np.arange(12):

        # predict month
        y_lr1_pred[i,month] = lr1.predict(X_pred)

        # shift data by 1 to the left, append y_pred, reshape to (1,12)
        X_pred = np.append(X_pred[:,1:], y_lr1_pred[i,month]).reshape(1,12)
```
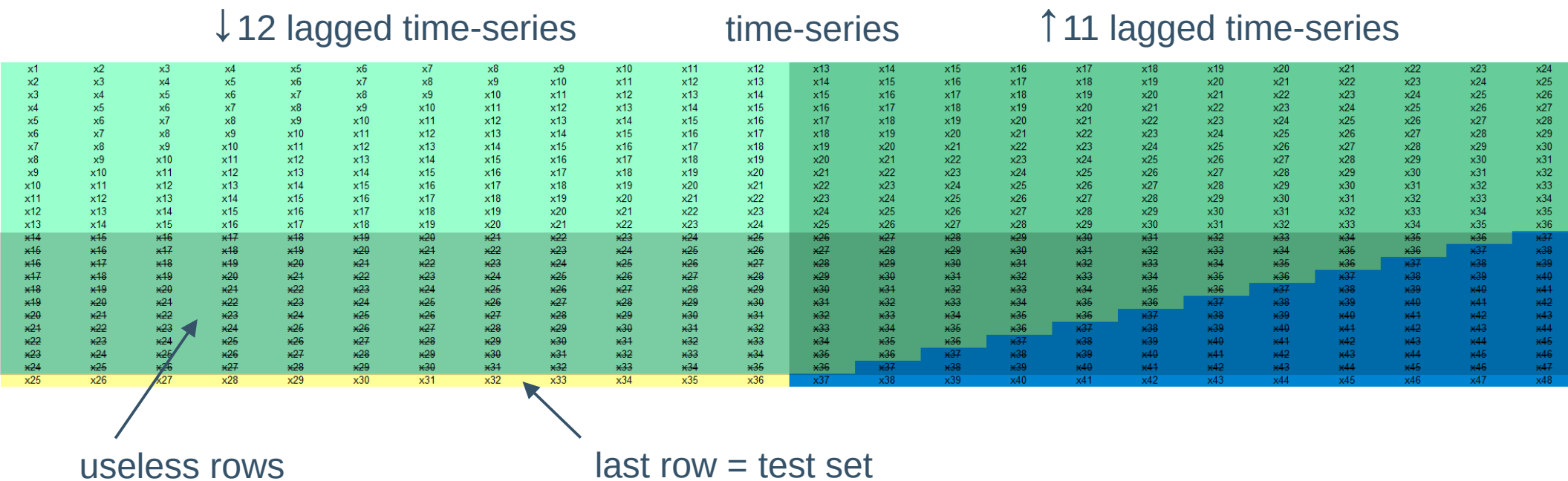
# ML2 – One-step multi-output prediction of one time-series

1. Data preparation:
   create 12 lagged features and 1 un-lagged and 11 lagged targets

2. Split into train and test set

3. Learn a multi-target ML model that is able to predict the next 12 months out of 12 previous months

4. Use ML model to predict the next 12 months

# ML2 – One-step multi-output prediction of one time-series



↓12 lagged time-series        time-series        ↑11 lagged time-series

useless rows

last row = test set

# ML2 – One-step multi-output prediction of one time-series

```python
lr2 = LinearRegression()

# initialize y_pred
Y_lr2_pred = np.empty((20,1,12))

# outer loop over all indices
for i in np.arange(df_task2.shape[1]):

    # fit LR model for index i
    lr2.fit(X_train[i], Y_train[i])

    # predict next year
    Y_lr2_pred[i] = lr2.predict(X_test[i])
```
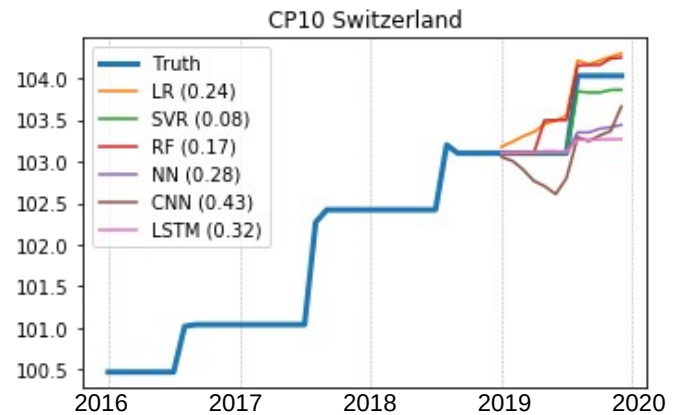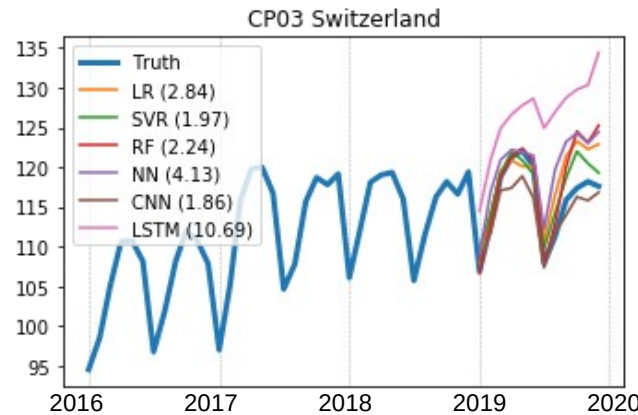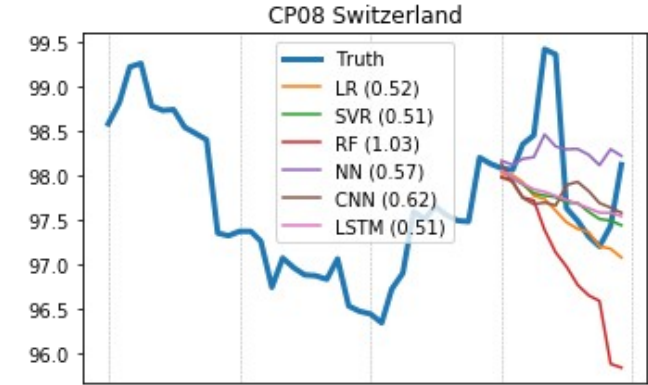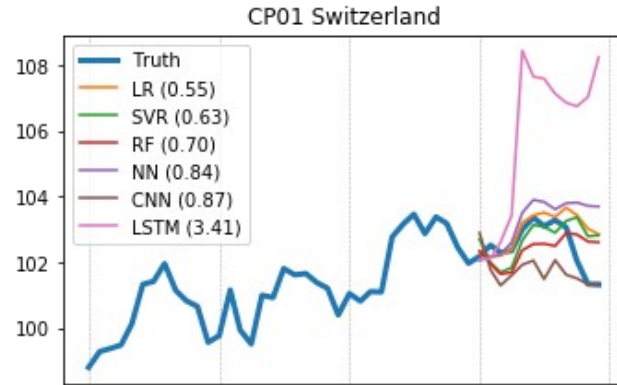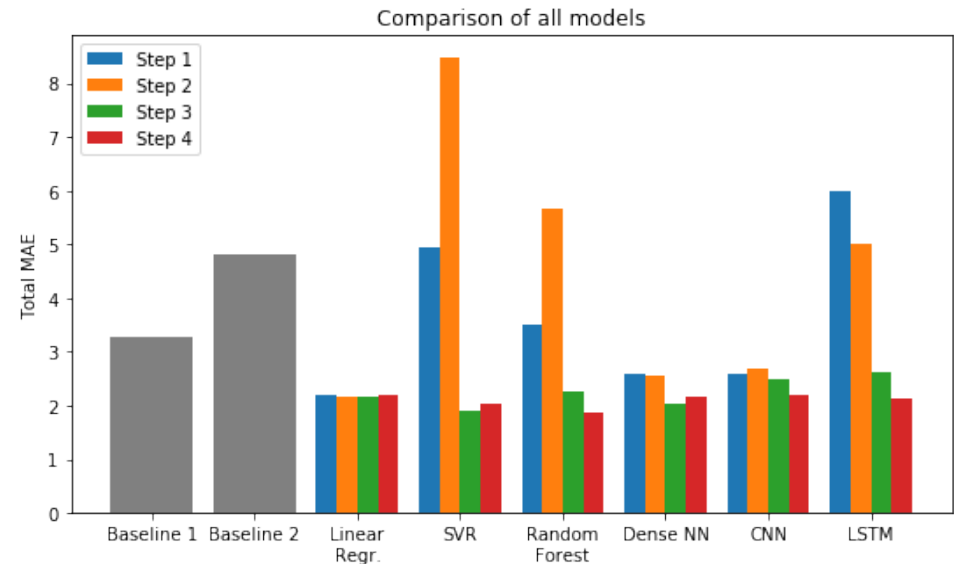
# ML2 – Comparison and Results

Step 3:
Recursive approach based on
rates of change

# ML2 – Comparison and Results

- Rather good performance with Linear Regression

- Several models perform better with stationary rates of change

- Deep learning is not superior here

- Approach number 4 runs faster

- Problem: dataset is too small and contains too much randomness

Keep it simple!



Comparison of all models

# Summary

**Exploratory Data Analysis**

- Data preparation straightforward

- Insights into main product groups and sub-categories

  - for the weights

  - for the indices

- Find similarities between countries

- Computation of rates of change

- Discover seasonality

**Machine Learning Task 1**

- Use ML to judge about feature importance

- Features are strongly correlated

**Machine Learning Task 2**

- Time-series prediction using different approaches and ML models (scikit-learn, tensorflow)

- Good results with basic models

- No enhancement with more effort