



Warsztaty modelowania

06 – optymalizacja i automatyzacja

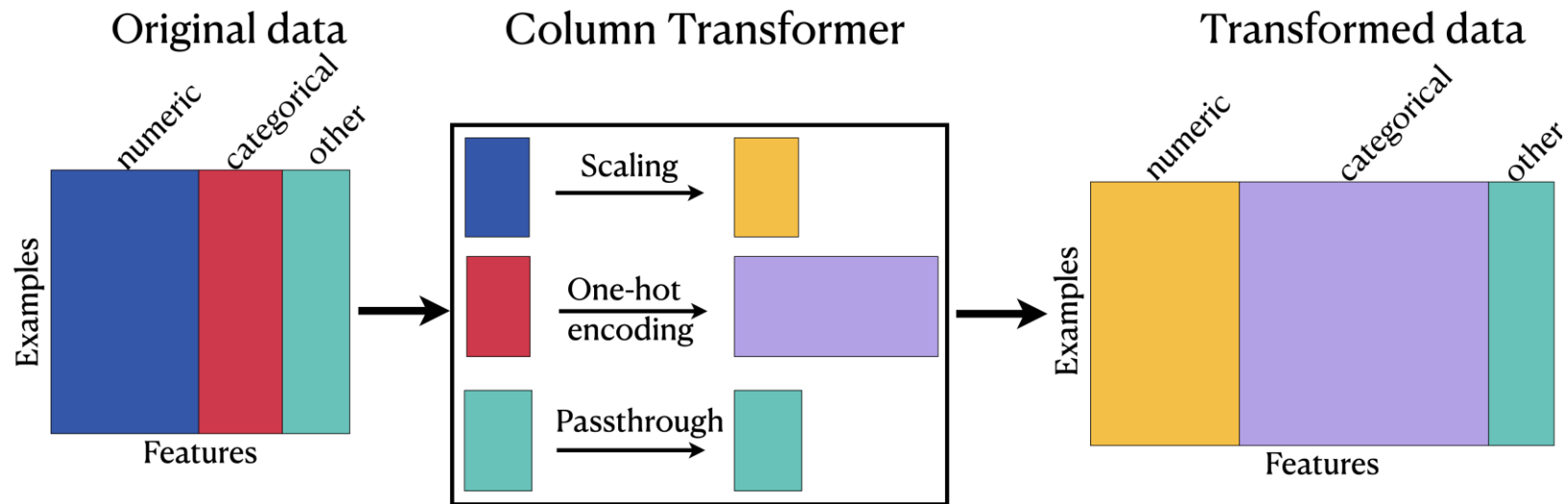
opracowała Patrycja Naumczyk

Optymalizacja i automatyzacja

1. Automatyzacja cz. I:
 - a) przekształcanie zmiennych
2. Optymalizacja:
 - a) walidacja krzyżowa
 - b) selekcja zmiennych wyjaśniających
3. Automatyzacja cz.II
 - a) pipeline'y

Automatyzacja cz. I – przekształcanie zmiennych

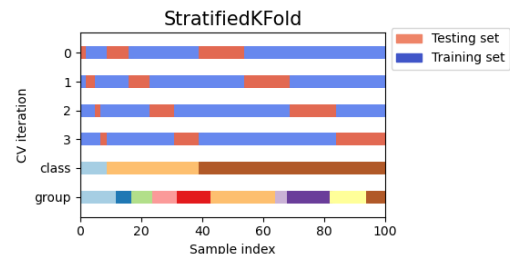
1. Zbiorowe kodowanie i skalowanie – [ColumnTransformer\(\)](#)
2. Ważne parametry:
 - a) transformers – lista tupli:
 - i. nazwa własna (str)
 - ii. transformer - estymator (ew. ,drop' lub ,passthrough')
 - iii. kolumny (np. lista)
 - b) remainder – estymator (ew. ,drop' lub ,passthrough')
 - c) sparse_threshold
 - d) verbose
3. Ważne metody:
 - a) .fit()
 - b) .transform()
 - c) .get_feature_names_out()
4. Ważne atrybuty:
 - a) .transformers_
 - b) .feature_names_in_



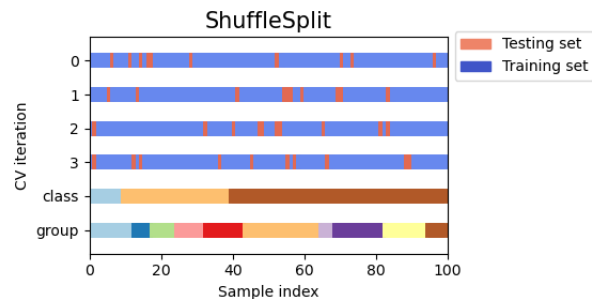
Optymalizacja – walidacja krzyżowa (cross-validation)

1. Optymalizacja uczenia modelu na secie TRENINGOWYM
2. Podstawowe rodzaje CV w scikit-learn:

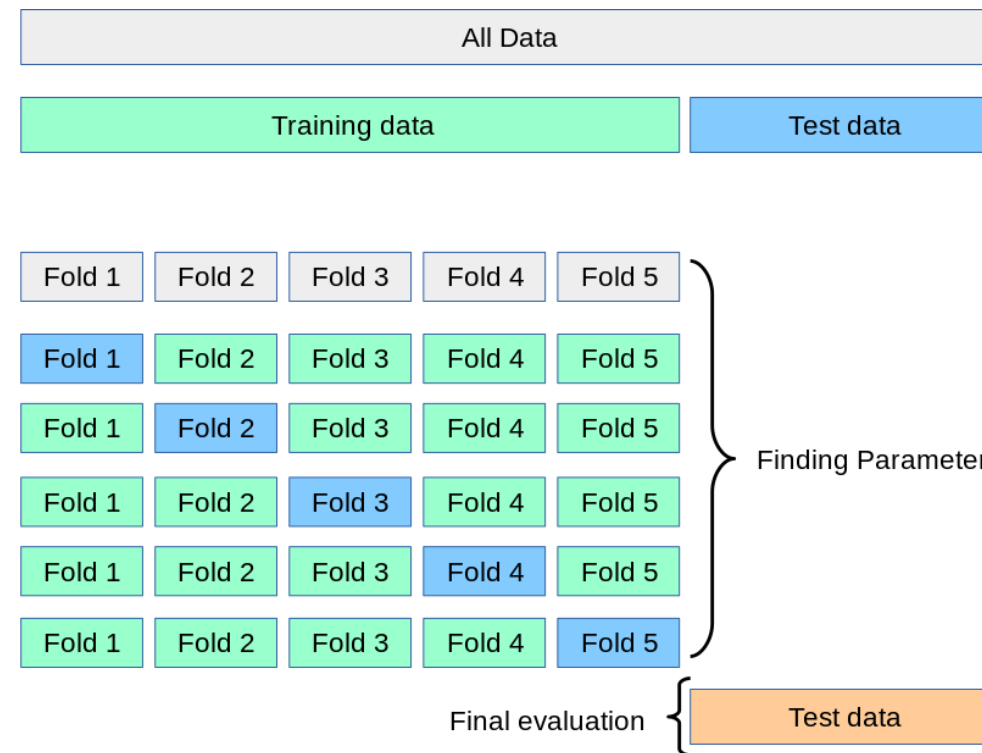
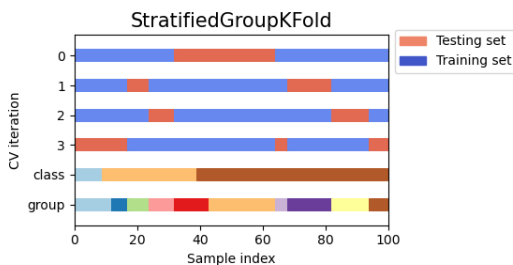
a) [StratifiedKFold\(\)](#)



b) [ShuffleSplit\(\)](#)



c) [StratifiedGroupKFold\(\)](#)



Optymalizacja – walidacja krzyżowa (cross-validation)

3. Optymalizacja parametrów – [GridSearchCV\(\)](#)

4. Ważne parametry:

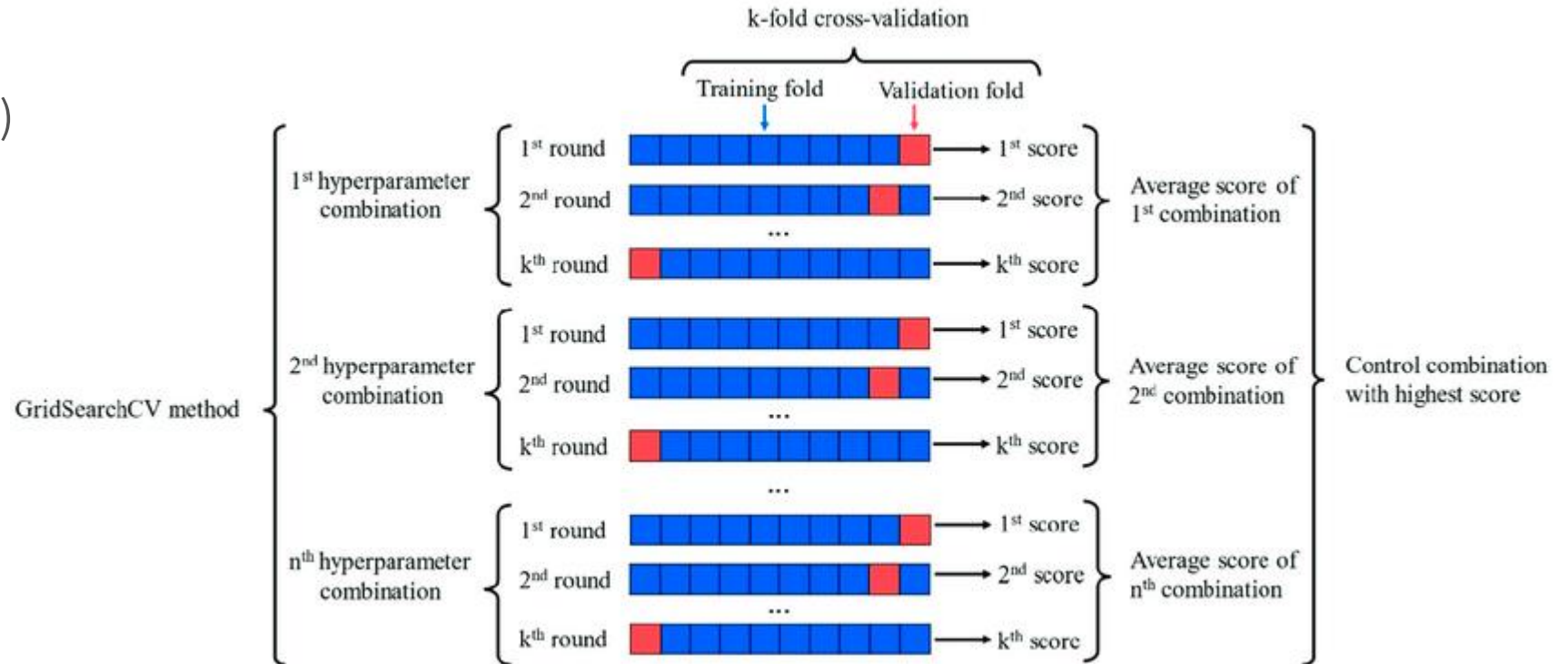
- a) estimator
- b) param_grid (słownik !)
- c) scoring
- d) cv

5. Ważne atrybuty:

- a) .cv_results_
- b) .best_estimator_
- c) .best_params_
- d) .feature_names_in_

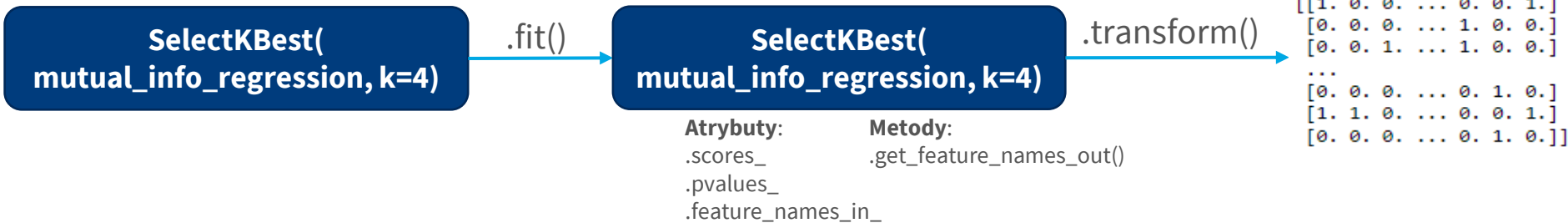
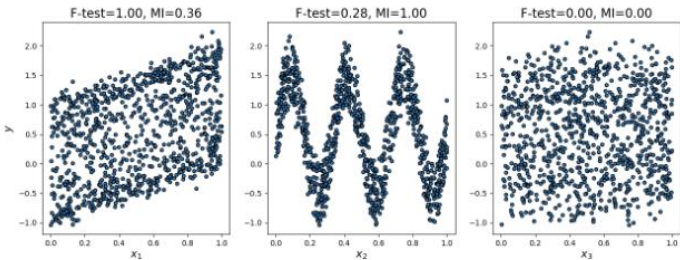
6. Ważne metody:

- a) .fit()
- b) .transform()
- c) .predict()

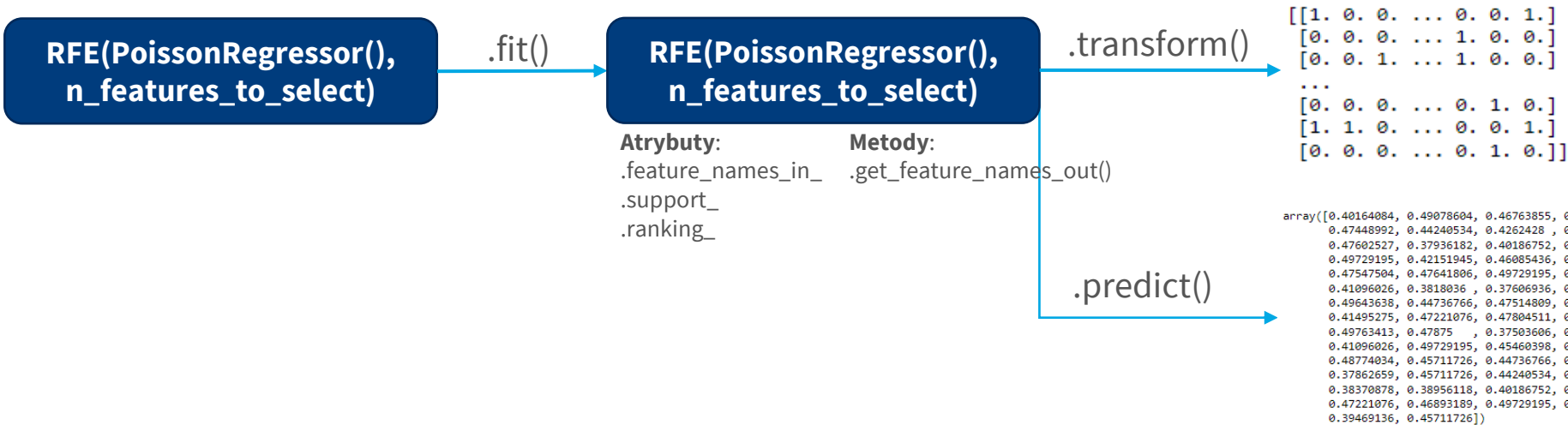


Optymalizacja – selekcja zmiennych (feature selection)

- 1. Moduł [feature_selection](#) w scikit-learn
- 2. Wybór jednoczynnikowy (związek ze zmienną wyjaśnianą)
 - a) [SelectKBest\(\)](#) i [SelectPercentile\(\)](#)
 - b) [r_regression\(\)](#), [f_regression\(\)](#), [mutual_info_regression\(\)](#)

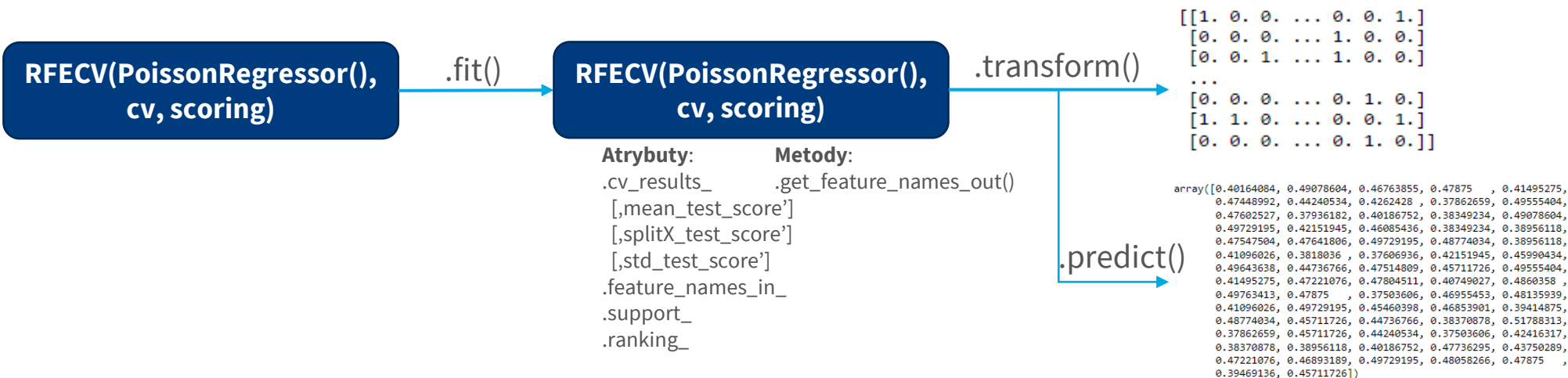
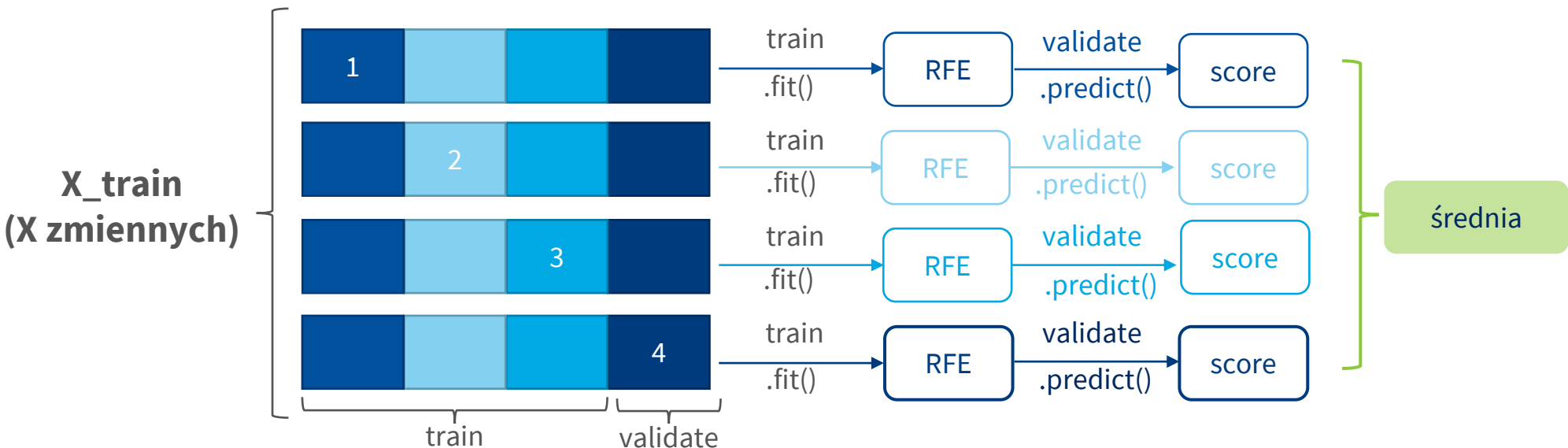


- 3. Recursive Feature Elimination ([RFE](#))
 - a) metoda eliminacji wstecznej na podstawie współczynników modelu (β)



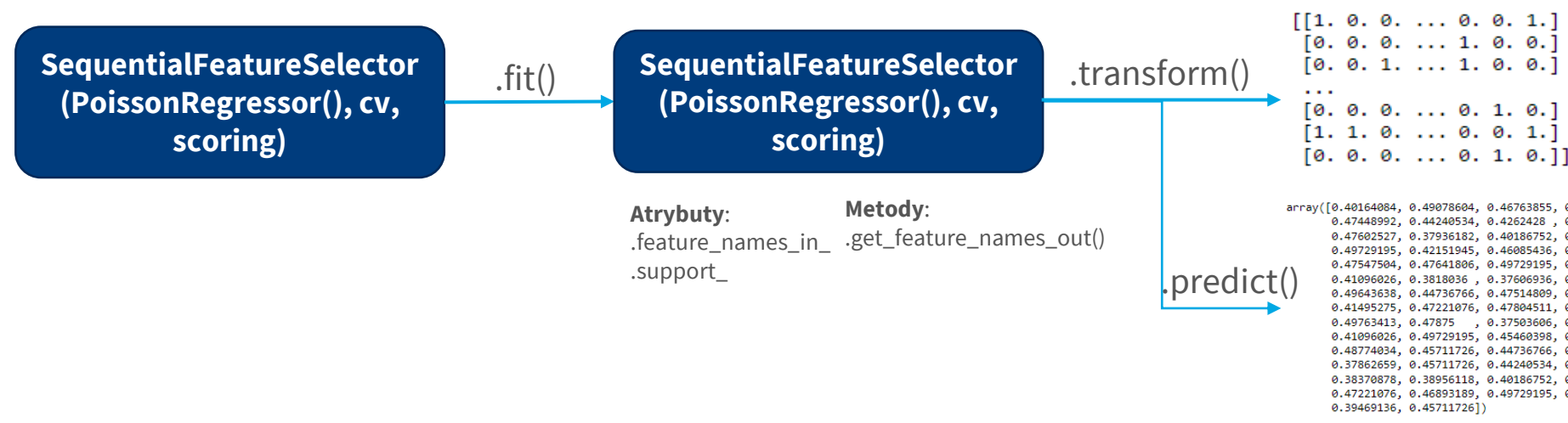
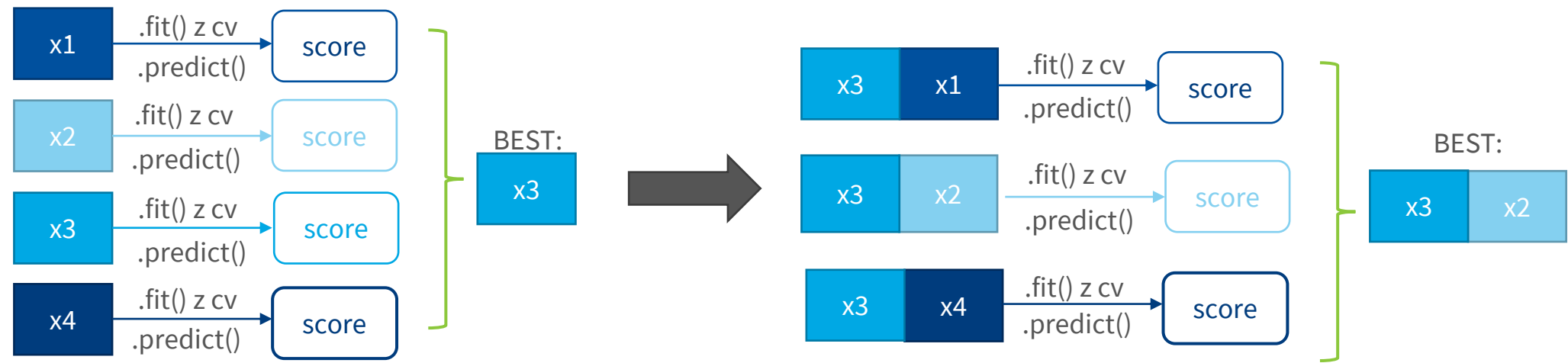
Optymalizacja – selekcja zmiennych (feature selection)

- 4. Recursive Feature Elimination with Cross-Validation ([RFECV](#))
 - a) selekcja optymalnej liczby zmiennych na bazie RFE



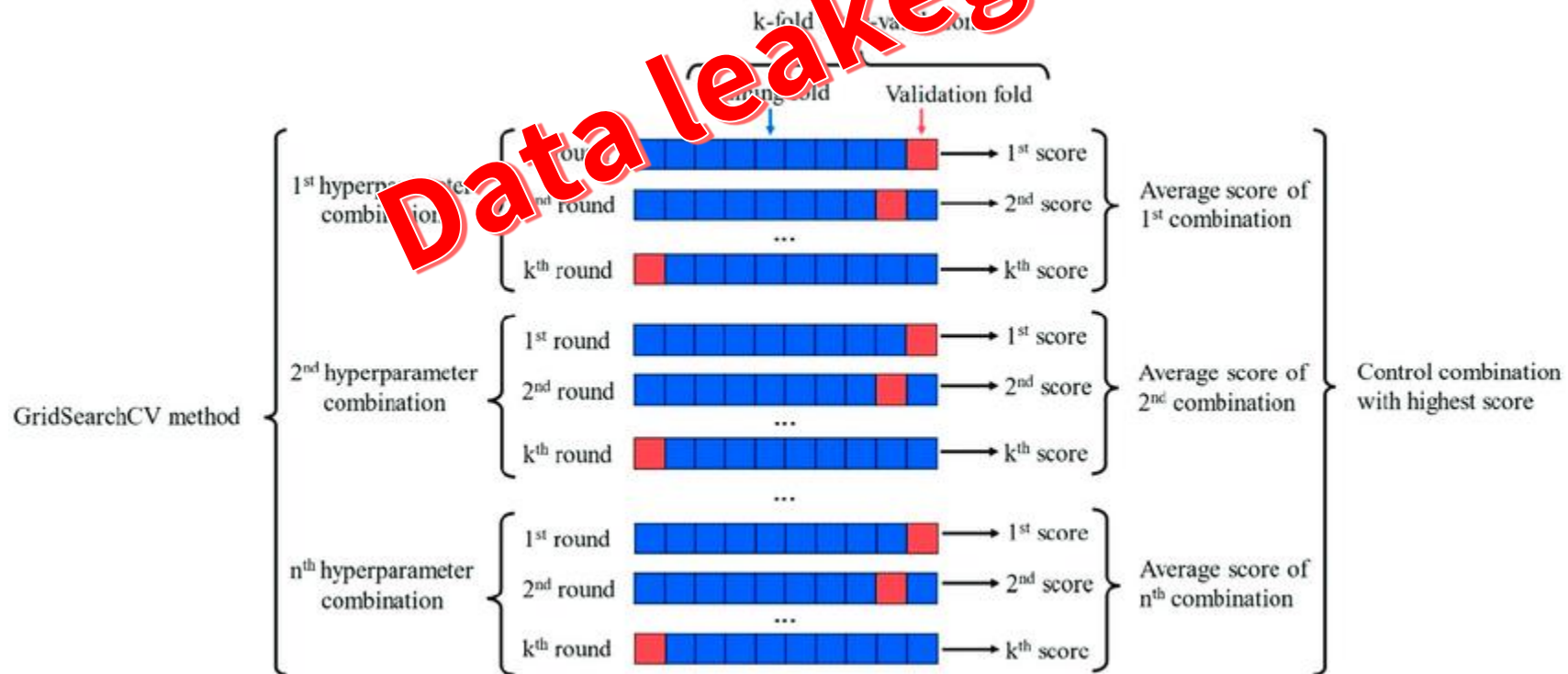
Optymalizacja – selekcja zmiennych (feature selection)

- 5. Sequential Feature Selector ([SFS](#))
 - a) regresja z eliminacją wsteczną lub **wprzód**

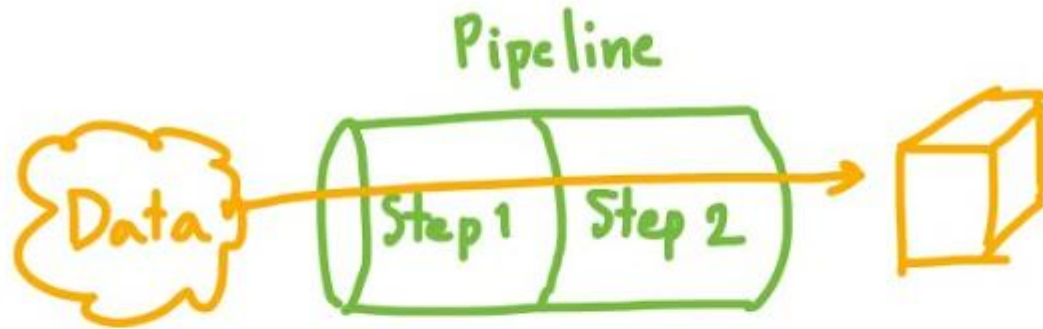


Automatyzacja – ścieżki przekształceń (pipeline)

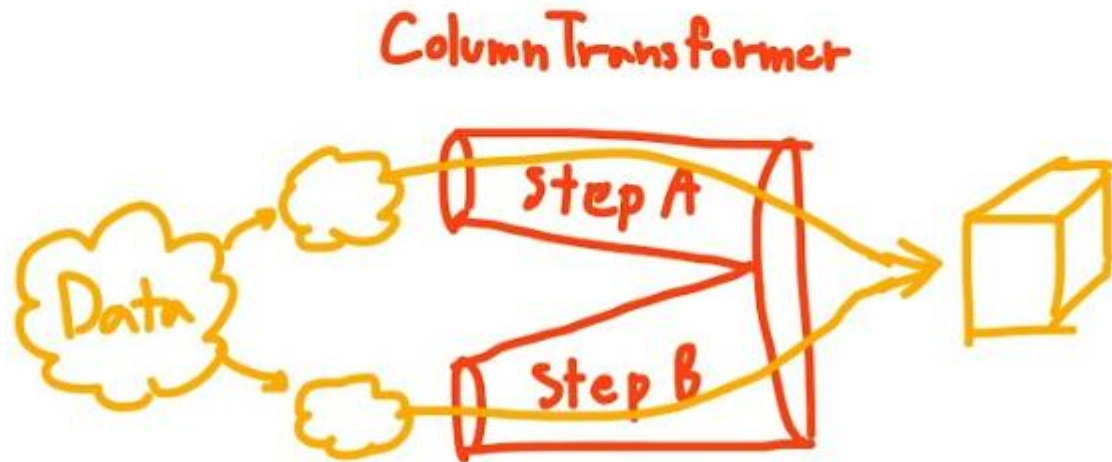
```
In [12]: 1 poisson_gscv = GridSearchCV(  
2         PoissonRegressor(),  
3         param_grid = params,  
4         scoring='neg_mean_poisson_deviance',  
5         cv=10  
6     ).fit(df_train_trans, df_train.czy_zkodowany)  
7     poisson_gscv
```



Automatyzacja – ścieżki przekształceń (pipeline)



```
Pipeline(steps=[  
    (nazwa_str, transformer),  
    (nazwa_str, transformer),  
    ...  
])
```

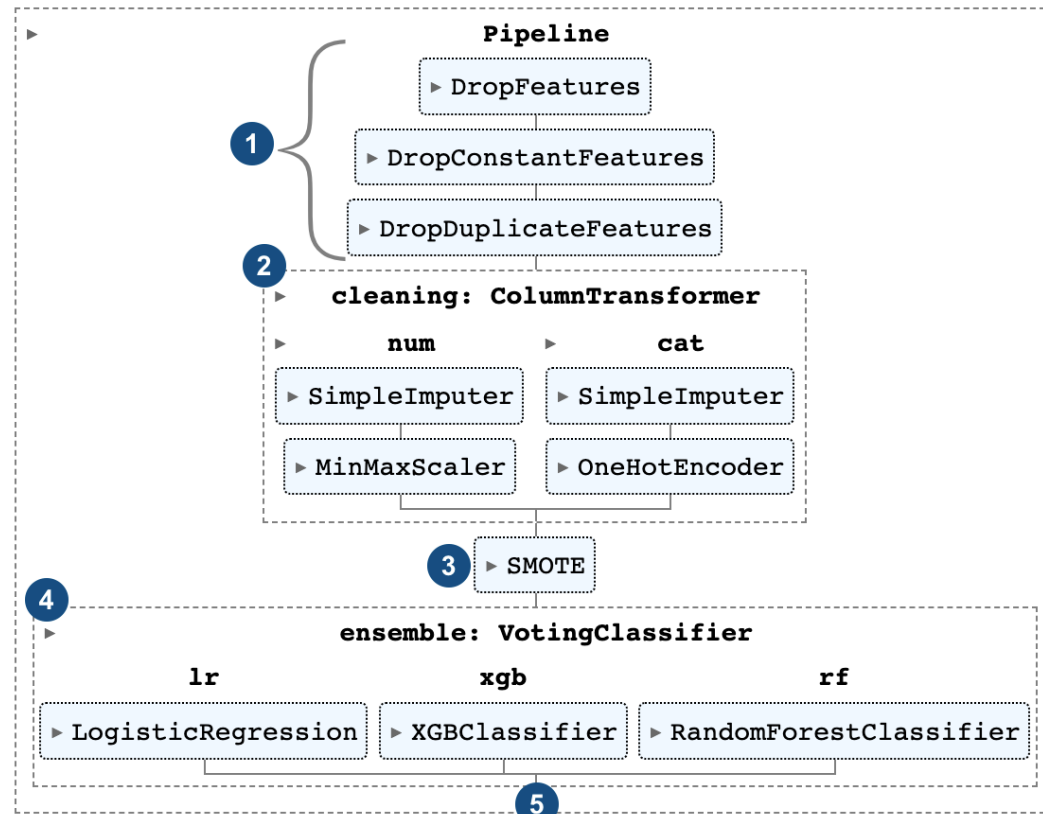
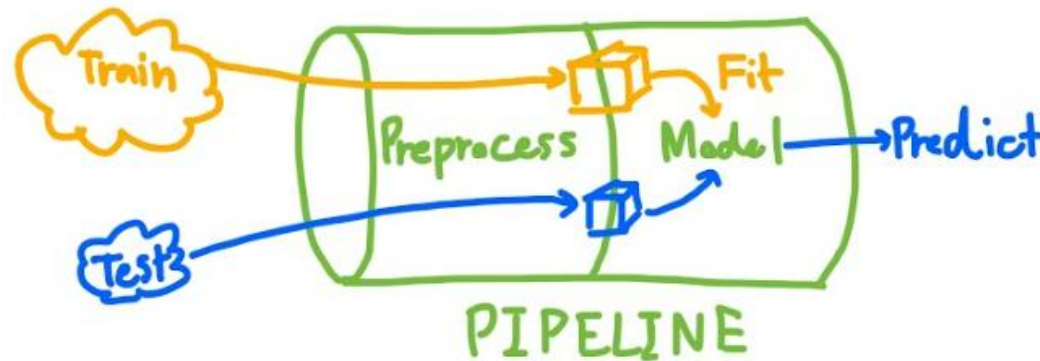


```
ColumnTransformer(transformers=[  
    (nazwa_str, transformer, kolumny),  
    (nazwa_str, transformer, kolumny),  
    ...  
])
```

Automatyzacja – ścieżki przekształceń (pipeline)

1. Ścieżki „krok po kroku” – [Pipeline\(\)](#)
2. Ważne parametry:
 - a) steps – lista tupli:
 - i. nazwa własna (str)
 - ii. transformer - estymator
 - b) verbose
3. Ważne metody:
 - a) .fit()
 - b) .predict()
 - c) .transform()
 - d) .get_params()
 - e) .get_feature_names_out()
4. Ważne atrybuty:
 - a) .feature_names_in_
5. „ładne” wyświetlanie:

```
from sklearn import set_config
set_config(display='diagram')
```



Automatyzacja – ścieżki przekształceń (pipeline)

1. [GridSearchCV\(\) z pipeline](#):

a) Optymalizacja parametrów bez data leakage

```
11 params = {
12     "alpha": list(np.logspace(-20, 1, num=100)),
13     # "solver" : ['lbfgs', 'newton-cholesky']
14 }
15
16 poisson_gscv_data_leakage = GridSearchCV(
17     PoissonRegressor(),
18     param_grid = params,
19     scoring=scores,
20     cv=10,
21     refit=False
22 ).fit(df_train_trans, df_train.czy_szkoda)
```



```
33 params = {
34     "model__alpha": list(np.logspace(-20, 1, num=100)),
35     # "solver" : ['lbfgs', 'newton-cholesky']
36 }
37
38 poisson_gscv_pipe = GridSearchCV(
39     pipeline,
40     param_grid = params,
41     scoring=scores,
42     cv=10,
43     refit=False
44 ).fit(df_train, df_train.czy_szkoda)
```

b) Weryfikacja optymalnej metody przetwarzania

```
grid_step_params = [{'col_trans__num_pipeline__minmax_scale': ['passthrough']},
                    {'col_trans__num_pipeline__std_scale': ['passthrough']}]
```

c) Oba ww. na raz:

```
merge_dict = {**dict_1, **dict_2}
```

```
43 step_params = [
44     {**params, **{"col_trans__minmax_numeric": ['passthrough']}},
45     {**params, **{"col_trans__standard_numeric": ['passthrough']}},
46 ]
```

