



# Warsztaty modelowania

05 – uczenie i ewaluacja modelu

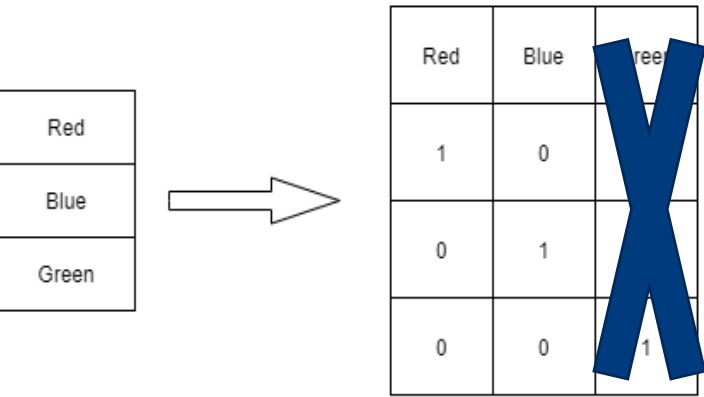
opracowała Patrycja Naumczyk

# Uczenie i ewaluacja – zasady ogólne

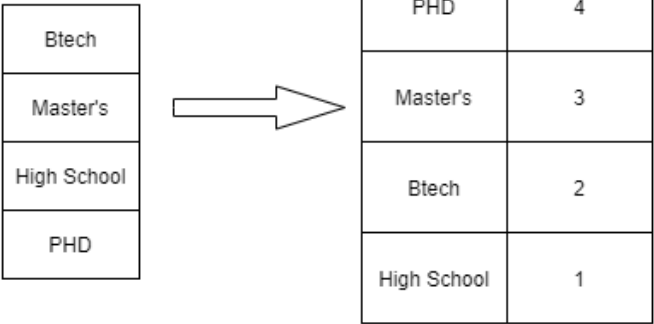
1. Uczenie
  - a) Przekształcenie zmiennych kategoryalnych
  - b) Regularyzować, czy nie regularyzować
2. Ewaluacja
  - a) Podział na set testowy i treningowy
  - b) Przeciekanie danych („data leakage”)
  - c) Właściwe metryki

# Uczenie

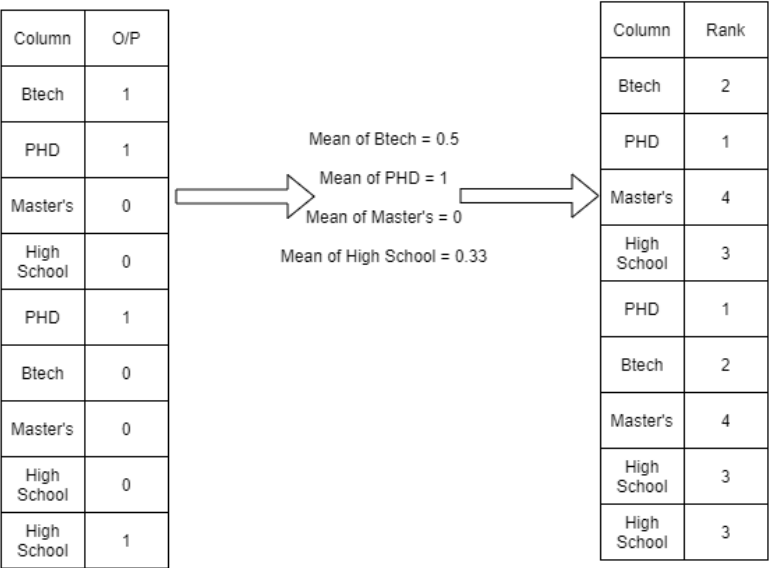
- 1. Przekształcenia zmiennych kategoryalnych
  - a) One-hot encoding
  - b) Ordinal encoding
  - c) Target encoding



One Hot Encoding



Label Encoding



Target Encoding

# Uczenie

1. Regularyzacja:
  - a) L1 – Lasso
  - b) L2 – Ridge
2. Konieczne skalowanie zmiennych ciągłych (!!!)
  - a) Min-Max
  - b) Normalizacja (standardization / Z-score)

$$X_{\text{new}} = (X - X_{\text{min}}) / (X_{\text{max}} - X_{\text{min}})$$

$$X_{\text{new}} = (X - \text{mean}) / \text{Std}$$

$$\text{RSS} = \sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2.$$

$$\text{RSS} + \lambda \sum_{j=1}^p \beta_j^2$$

**RIDGE**

$$\text{RSS} + \lambda \sum_{j=1}^p |\beta_j|.$$

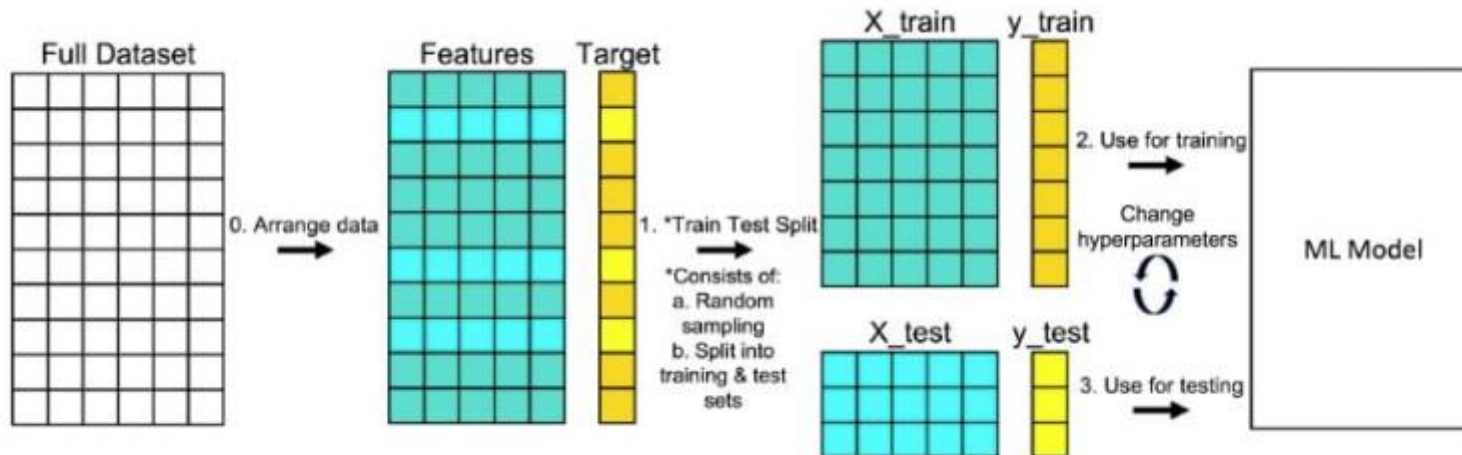
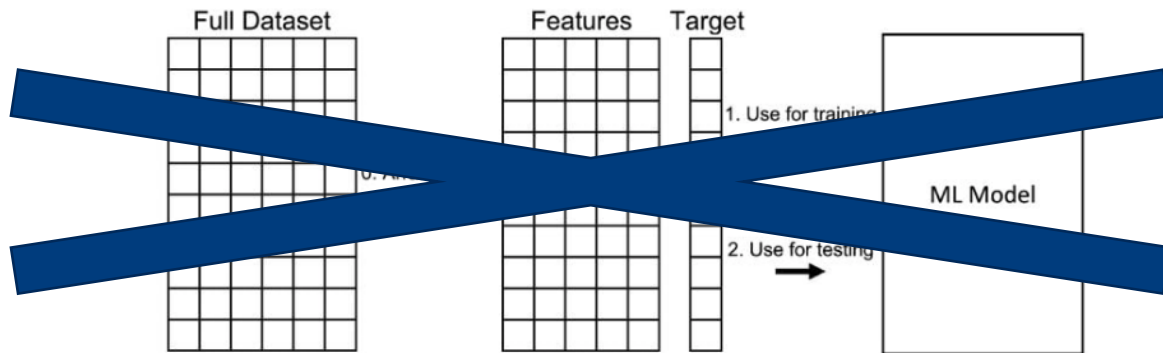
**LASSO**

$$\frac{\sum_{i=1}^n (y_i - x_i^J \hat{\beta})^2}{2n} + \lambda \left( \frac{1 - \alpha}{2} \sum_{j=1}^m \hat{\beta}_j^2 + \alpha \sum_{j=1}^m |\hat{\beta}_j| \right)$$

**ELASTIC NET**

# Ewaluacja

1. Podział na sety treningowy i testowy
  - a) „Przeciekanie” danych
  - b) Podział proporcjonalny do zmiennej

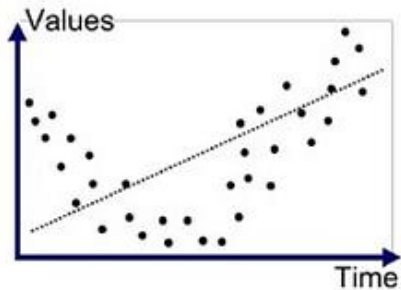


# Ewaluacja

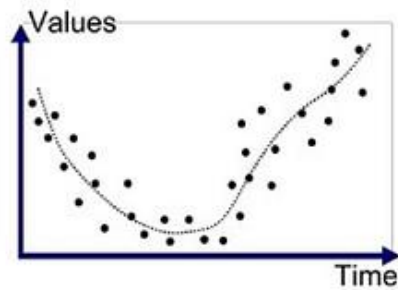
1. Podział na sety treningowy i testowy
  - a) „Przeciekanie” danych
  - b) Podział proporcjonalny do zmiennej
2. Metryki:
  - a)  $R^2$
  - b) Odchylenie

$$R^2 = 1 - \frac{SS_{\text{res}}}{SS_{\text{tot}}}$$

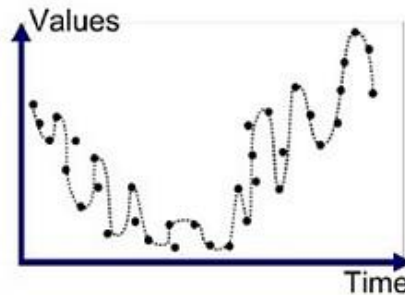
$$D(y, \hat{y}) = \frac{1}{n_{\text{samples}}} \sum_{i=0}^{n_{\text{samples}}-1} \begin{cases} (y_i - \hat{y}_i)^2, & \text{for } p = 0 \text{ (Normal)} \\ 2(y_i \log(y_i/\hat{y}_i) + \hat{y}_i - y_i), & \text{for } p = 1 \text{ (Poisson)} \\ 2(\log(\hat{y}_i/y_i) + y_i/\hat{y}_i - 1), & \text{for } p = 2 \text{ (Gamma)} \\ 2 \left( \frac{\max(y_i, 0)^{2-p}}{(1-p)(2-p)} - \frac{y_i \hat{y}_i^{1-p}}{1-p} + \frac{\hat{y}_i^{2-p}}{2-p} \right), & \text{otherwise} \end{cases}$$



Underfitted



Good Fit/Robust



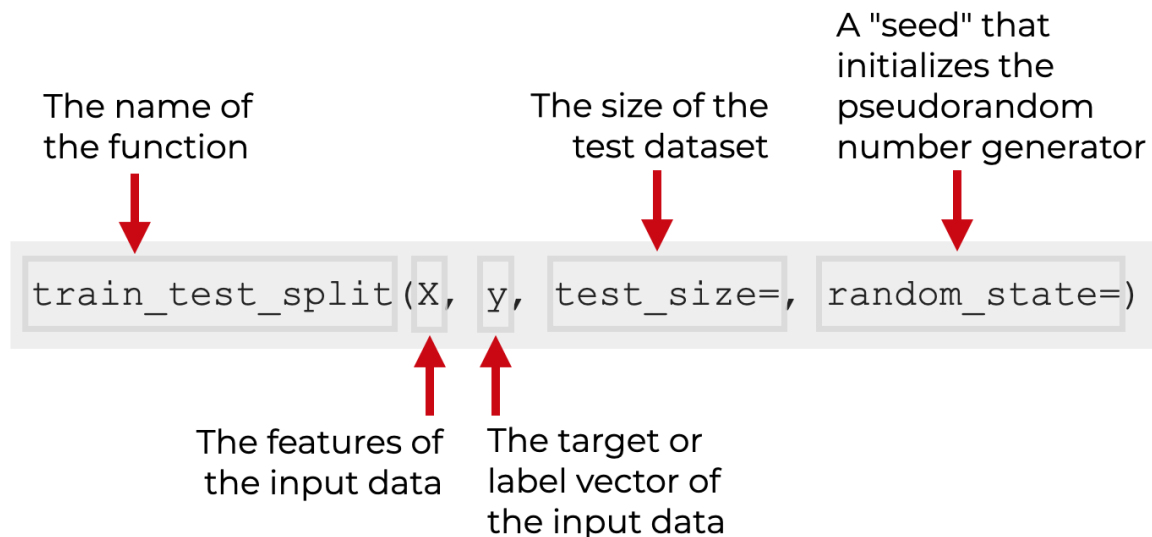
Overfitted

# Uczenie i ewaluacja – scikit-learn vs statsmodels

1. Podział na set treningowy i testowy – scikit-learn:
  - a) funkcja `train_test_split()`
2. Przekształcanie zmiennych – scikit-learn (metody `.fit()` `.transform()` `.fit_transform()`):
  - a) zmienne kategoryjne: kodowanie
    - i. obiekt `OneHotEncoder()`
    - ii. obiekt `OrdinalEncoder()`
  - b) zmienne ciągłe: skalowanie
    - i. obiekt `MinMaxScaler()`
    - ii. obiekt `StandardScaler()`
3. Uczenie modelu:
  - a) scikit-learn (metody `.fit()` `.predict()`):
    - i. obiekt `PoissonRegressor()`
    - ii. obiekt `GammaRegressor()`
  - b) statsmodels (metody `.fit()` `.fit_regularized()` `.predict()`):
    - i. obiekt `GLM( family=sm.families.Poisson() )`
    - ii. obiekt `GLM( family=sm.families.Gamma() )`
4. Ewaluacja modelu – scikit-learn:
  - a) funkcje poszczególnych metryk  
[https://scikit-learn.org/stable/modules/model\\_evaluation.html](https://scikit-learn.org/stable/modules/model_evaluation.html)

# Funkcja `train_test_split()`

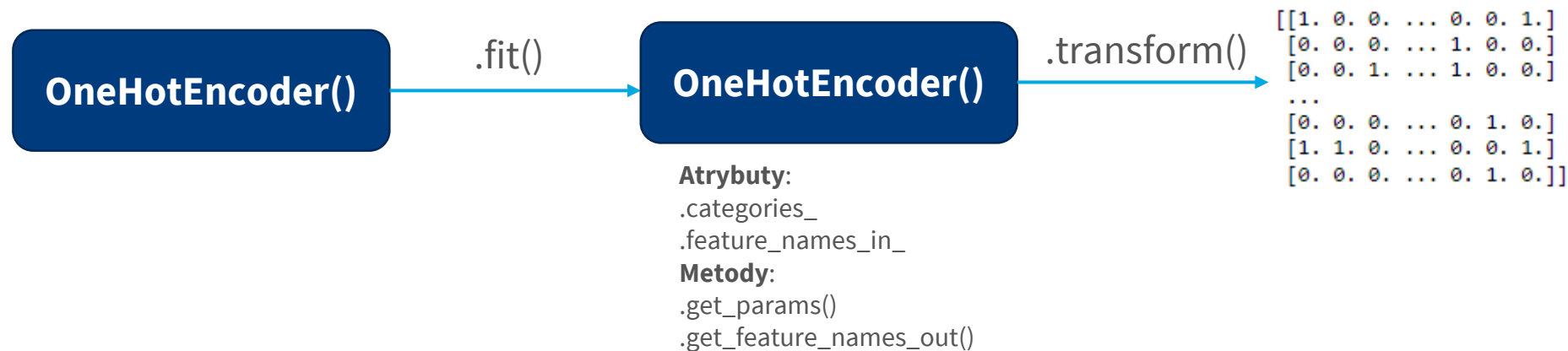
1. Podział na próbkę treningową i testową
2. Jeśli jest potrzeba wydzielenia osobno próbki „hold-out”, konieczne jest dwukrotne przeprowadzenie danych przez funkcję
3. Ważne parametry (patrz: [dokumentacja](#)):
  - a) **test\_size** – jaki odsetek zbioru ma być w próbce testowej (default: 0.25)
  - b) **random\_state** – zapamiętanie [stanu losowego](#), użytego przy podziale (powtarzalność podziałów!)
  - c) **stratify** – podział proporcjonalny do zmiennej
4. Zachowuje typ danych (dataframe -> dataframe, numpy array -> numpy array)





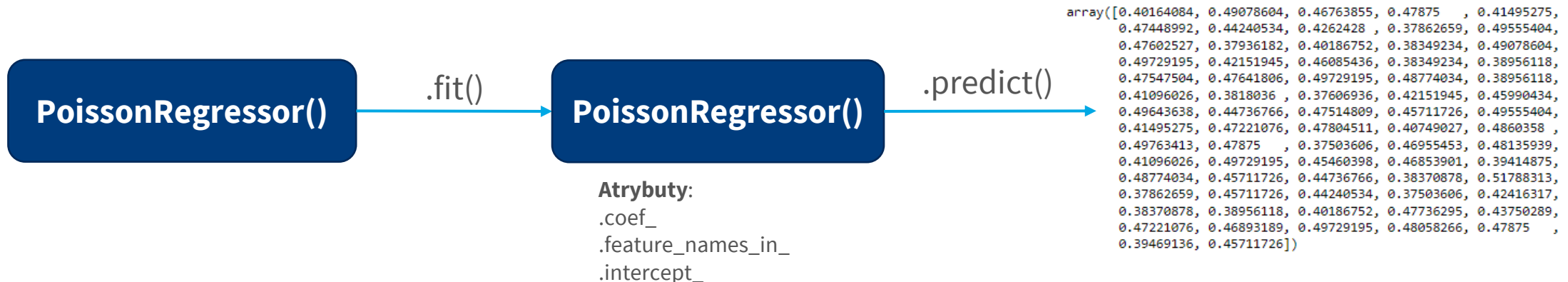
# Przekształcanie zmiennych

1. Zmienne kategoryjne -> kodowanie (np. [OneHotEncoder\(\)](#) [OrdinalEncoder\(\)](#))
2. Zmienne ciągłe -> skalowanie (np. [MinMaxScaler\(\)](#) [StandardScaler\(\)](#) )
3. Podstawowe metody na obiekcie:
  - a) **.fit()**
  - b) **.transform()**
  - c) **.fit\_transform()**
4. NIE zachowuje typu danych (!)
  - a) Wynikiem numpy array
  - b) Konieczność przekształcenia do dataframe

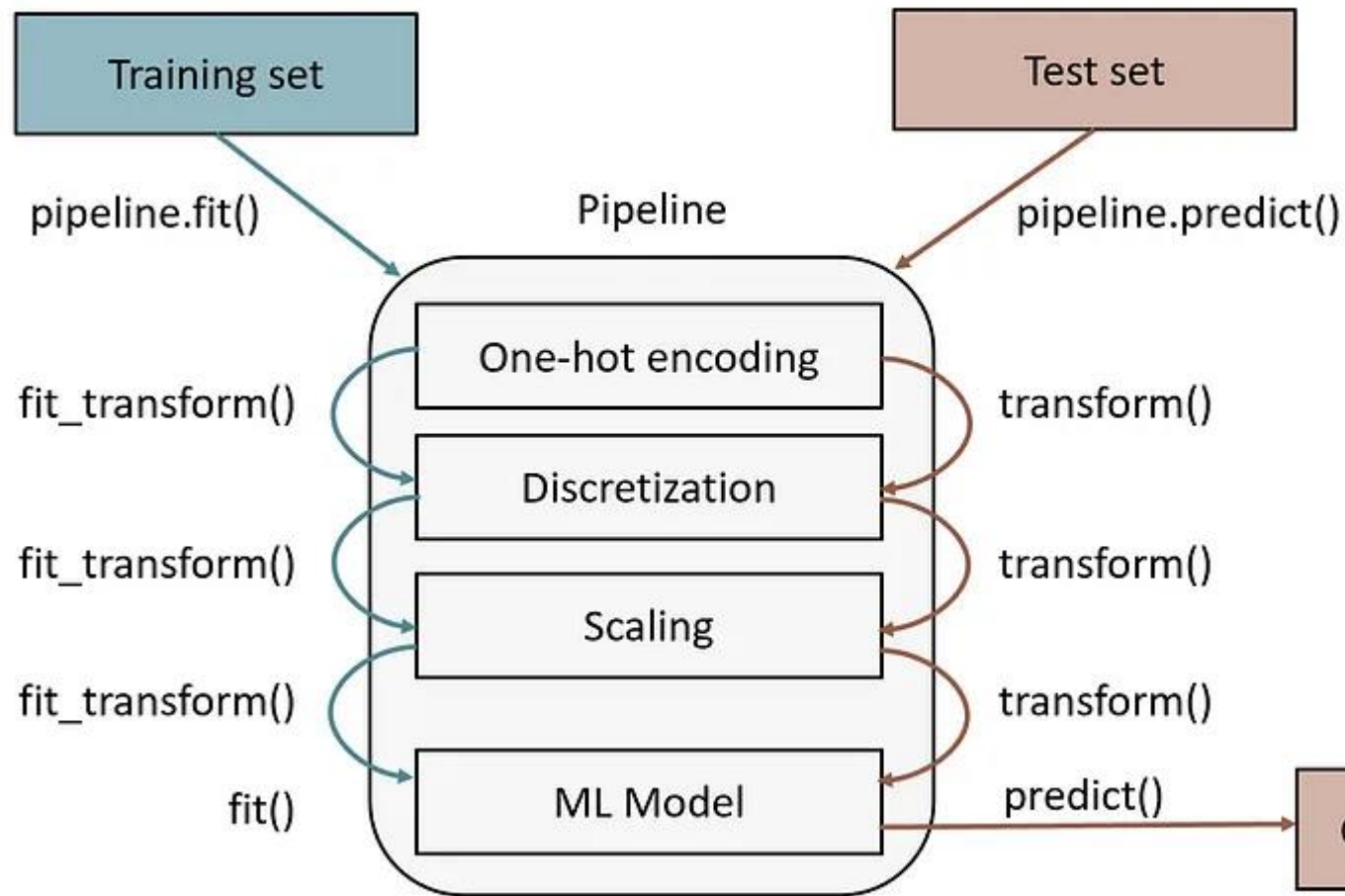


# Uczenie modelu – scikit-learn

1. Model częstości -> [PoissonRegressor\(\)](#)
2. Model średniej szkody -> [GammaRegressor\(\)](#)
3. Podstawowe metody na obiekcie:
  - a) **.fit()**
  - b) **.predict()**
4. NIE zachowuje typu danych (!)
  - a) Wynikiem numpy array
  - b) Konieczność przekształcenia do dataframe
5. NIE wylicza prawdopodobieństwa poszczególnych zmiennych wyjaśniających, ani F całego modelu
6. Ocena modelu tylko indywidualna przez użytkownika (domyślnie D<sup>2</sup>)
7. Regularyzacja tylko L2 (parametr  $\alpha \in [0.0, \infty)$  domyślnie wartość 1.0)



## Co na którym zbiorze?

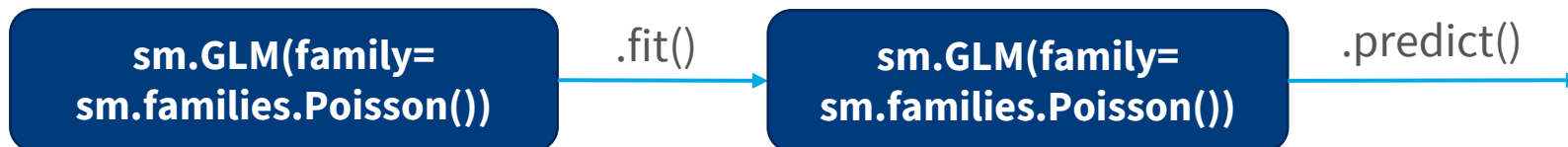


# Uczenie modelu – statsmodels

1. Równanie modelu zapisane tekstowo (notacja R, [biblioteka patsy](#)):  
zmienna\_wyjaśniana ~ zmienna\_wyjaśniająca1 + zmienna\_wyjaśniająca2 + ...

$$y \sim x + C(x1) + x:x1 + x2*x3$$

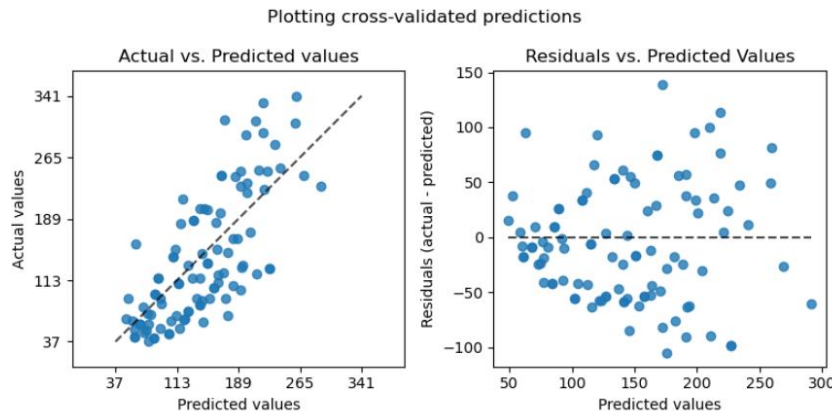
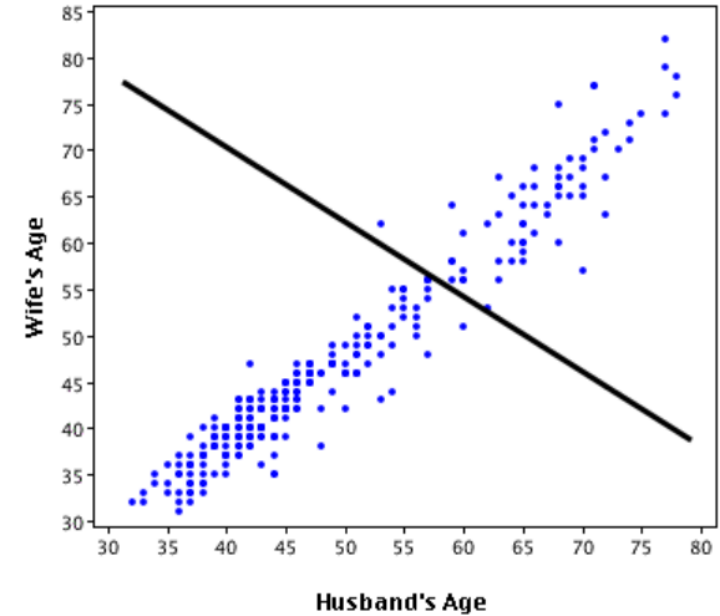
2. Formatowanie danych wejściowych (funkcja `dmatrices()`)
3. Jeden obiekt dla GLMs:
  - a) [statsmodels.api.GLM\(\)](#), wskazanie rodziny poprzez parametr *family*
  - b) dla rodziny Poisson i Gamma konieczne wskazanie zmiennej wagowej (exporyz dla częstości, liczba szkód dla średniej szkody)
4. Podstawowe metody na obiekcie:
  - a) **.fit()** *ALE! .fit() zwraca nowy obiekt(!)*
  - b) **.predict()**
  - c) regularyzacja poprzez metodę **fit\_regularized()**
    - i. parametr  $\alpha \in [0.0, \infty)$ , domyślnie 0
    - ii. parametr  $L1\_w \in [0.0, 1.0]$ ,  $L1 = L1\_w$ ,  $L2 = 1 - L1\_w$



Metody:  
.summary()

# Funkcje ewaluacyjne... i jedna metoda... i kilka wykresów

1. [Funkcje](#) ewaluacyjne scikit-learn
  - a) [r2\\_score\(\)](#)
  - b) [mean\\_poisson\\_deviance\(\)](#)
  - c) [mean\\_gamma\\_deviance\(\)](#)
2. Metoda `.score()`
  - a) Wykonywana na obiekcie z dopasowanym modelem
  - b) W zależności od modelu, różna metryka -> sprawdź dokumentację!
  - c) Dla regresji Poissona i Gamma:  $D^2$
3. [Wykresy](#) ewaluacyjne
  - a) Funkcjonalność dostępna od scikit-learn v.1.2
  - b) Zestawienia „actual-predicted” oraz „residuals-predicted”



$$D^2(y, \hat{y}) = 1 - \frac{\text{dev}(y, \hat{y})}{\text{dev}(y, y_{\text{null}})}$$