

# Opis Matematyczny Podziału Danych

Federated Learning z Stratyfikowaną 10-Fold Cross-Validation

## 1. Symbole matematyczne

| Symbol | Opis   |
|--------|--|
| D      | Całkowity dataset (zbiór wszystkich danych)        |
| D      | Liczba wierszy w datasecie D                       |
| D_tr   | Zbiór danych treningowych (80%)                    |
| D_val  | Zbiór danych walidacyjnych (10%)                   |
| D_ts   | Zbiór danych testowych (10%)                       |
| F_i    | i-ty fold treningowy ( $i = 1, 2, \dots, 10$ )     |
| C_k    | k-ty klient ( $k = 1, 2, \dots, K$ )               |
| K      | Liczba klientów (2, 3, 4 lub 5)                    |
| y      | Zmienna docelowa (klasy)                           |
| c      | Klasa ( $c$ należy do $\{c_1, c_2, \dots, c_m\}$ ) |
| p(c)   | Proporcja klasy $c$ w zbiorze                      |

## 2. Główny podział danych: Stratyfikowany 80/10/10

**Formuła podziału:**

$$D = D_{\text{tr}} \cup D_{\text{val}} \cup D_{\text{ts}}$$

gdzie:

- $|D_{\text{tr}}| = 0.80 \times |D|$  - dane treningowe
- $|D_{\text{val}}| = 0.10 \times |D|$  - dane walidacyjne
- $|D_{\text{ts}}| = 0.10 \times |D|$  - dane testowe

**Warunek stratyfikacji:**

Dla każdej klasy  $c$  proporcje są zachowane we wszystkich zbiorach:

$$p(c | D_{\text{tr}}) = p(c | D_{\text{val}}) = p(c | D_{\text{ts}}) = p(c | D)$$

**Właściwości:**

- $D_{\text{tr}}, D_{\text{val}}, D_{\text{ts}}$  są rozłączne (żadne dane się nie powtarzają)
- Stratyfikacja gwarantuje identyczne proporcje klas we wszystkich zbiorach

### 3. Stratyfikowana 10-Fold Cross-Validation

**Podział danych treningowych na 10 foldów:**

$$D_{\text{tr}} = F_1 \cup F_2 \cup F_3 \cup \dots \cup F_{10}$$

**Rozmiary foldów:**

$$|F_i| = |D_{\text{tr}}| / 10 \text{ dla } i = 1, 2, \dots, 10$$

**Warunek stratyfikacji foldów:**

Dla każdej klasy  $c$  i każdego foldu  $F_i$ :

$$p(c | F_i) = p(c | D_{\text{tr}}) = p(c | D)$$

*Każdy fold ma identyczny rozkład klas w przybliżeniu*

### 4. Mechanizm Cross-Validation

Dla każdej iteracji CV ( $i = 1, 2, \dots, 10$ ):

- **Dane treningowe w iteracji  $i$ :**  $D_{\text{train}}(i) = \text{wszystkie foldy OPRÓCZ } F_i$  (9 foldów)
- **Dane walidacyjne lokalne w iteracji  $i$ :**  $D_{\text{val\_local}}(i) = F_i$  (1 fold)

| Iteracja | Trening (9 foldów)                             | Walidacja lokalna |
|----------|--|-------------------|
| $i = 1$  | $F_2 \cup F_3 \cup \dots \cup F_{10}$          | $F_1$             |
| $i = 2$  | $F_1 \cup F_3 \cup \dots \cup F_{10}$          | $F_2$             |
| $i = 3$  | $F_1 \cup F_2 \cup F_4 \cup \dots \cup F_{10}$ | $F_3$             |
| ...      | ...  | ...               |
| $i = 10$ | $F_1 \cup F_2 \cup \dots \cup F_9$             | $F_{10}$          |

## 5. Podział na klientów (Federated Learning)

### Rotacja foldów:

Każdy klient otrzymuje te same 10 foldów, ale w różnej kolejności:

$$\text{start\_fold}(k) = ((k - 1) \times \text{floor}(10 / K)) \text{ MOD } 10$$

### Przykład dla K=2 klientów:

| Klient   | f1 | f2 | f3 | f4 | f5  | f6 | f7 | f8 | f9 | f10 |
|----------|----|----|----|----|-----|----|----|----|----|-----|
| Client 1 | F1 | F2 | F3 | F4 | F5  | F6 | F7 | F8 | F9 | F10 |
| Client 2 | F6 | F7 | F8 | F9 | F10 | F1 | F2 | F3 | F4 | F5  |

### Współdzielone zbiory:

- D\_val - identyczny dla wszystkich klientów (SHARED)
- D\_ts - identyczny dla wszystkich klientów (SHARED)

## 6. Algorytm podziału danych ze stratyfikacją

### ALGORYTM: StratifiedFederatedDataSplit(D, K, y)

WEJŚCIE: D (dataset), K (liczba klientów), y (zmienna docelowa)

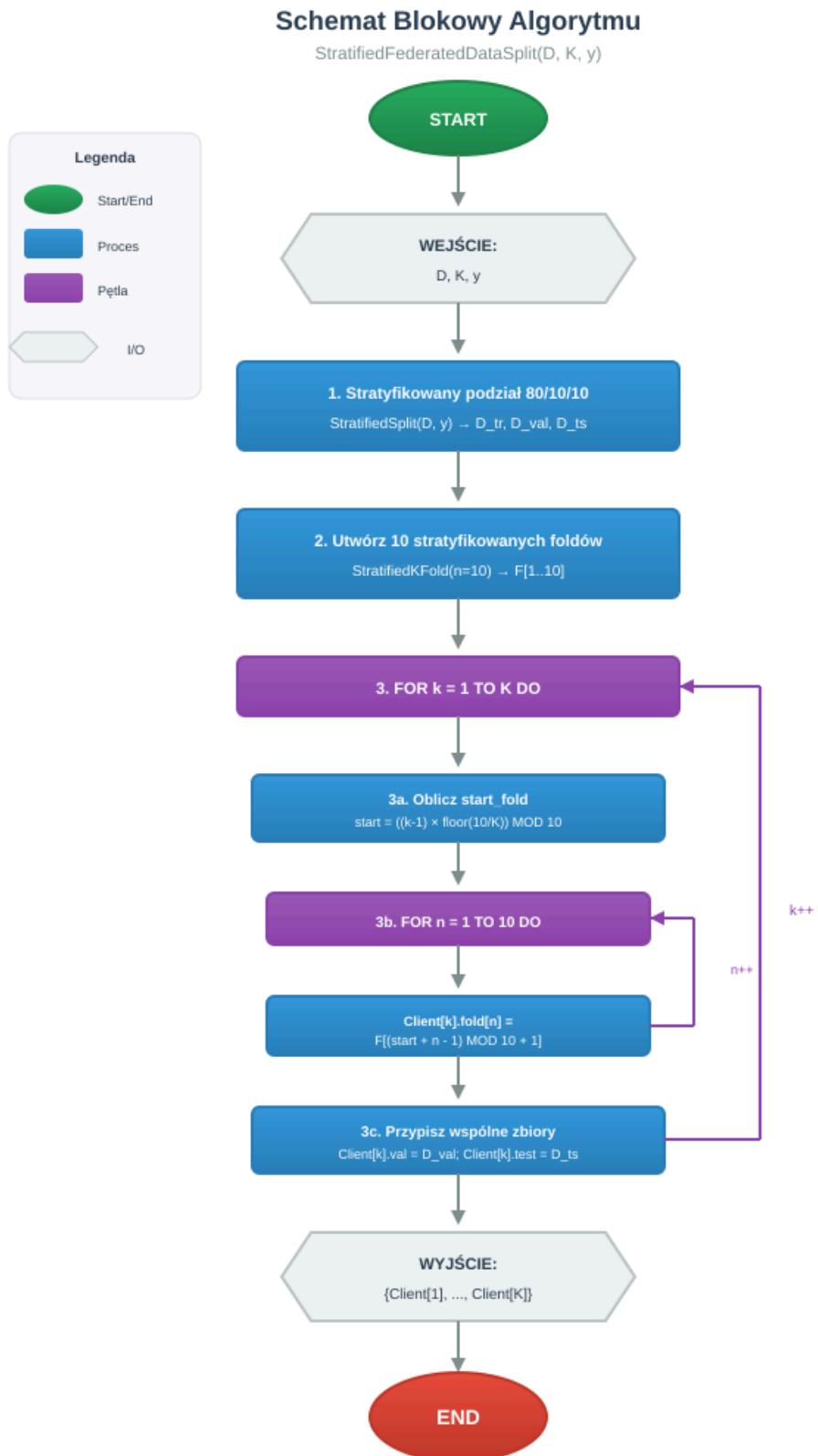
WYJŚCIE: K klientów z foldami, val, test

BEGIN

1. // STRATYFIKOWANY PODZIAŁ 80/10/10  
[D\_train, D\_temp] = StratifiedSplit(D, y, test\_size=0.20)  
[D\_val, D\_test] = StratifiedSplit(D\_temp, y\_temp, test\_size=0.50)
2. // STRATYFIKOWANE FOLDY (StratifiedKFold)  
skf = StratifiedKFold(n\_splits=10, shuffle=True)  
FOR i = 1 TO 10 DO Fold[i] = D\_train[skf.indices[i]]
3. // PRZYPISANIE DO KLIENTÓW  
FOR k = 1 TO K DO  
    start\_fold = ((k-1) × floor(10/K)) MOD 10  
    Client[k].fold[n] = Fold[(start\_fold + n - 1) MOD 10 + 1]  
    Client[k].val = D\_val; Client[k].test = D\_test
4. RETURN {Client[1], ..., Client[K]}

END

## 7. Schemat blokowy algorytmu



## 8. Przykład: RT-IoT2022

**Dane:**  $|D| = 123,117$  wierszy,  $K = 2$  klientów

**Obliczenia:**

- $|D_{\text{tr}}| = 0.80 \times 123,117 = 98,493$  wierszy
- $|D_{\text{val}}| = 0.10 \times 123,117 = 12,312$  wierszy
- $|D_{\text{ts}}| = 123,117 - 98,493 - 12,312 = 12,312$  wierszy

**Weryfikacja:**  $98,493 + 12,312 + 12,312 = 123,117$

**Rozmiar folda:**  $|F_i| = 98,493 / 10 \approx 9,849$  wierszy

## 9. Podsumowanie wszystkich datasetów

| Dataset            | Total     | Train<br>(80%) | Val<br>(10%) | Test<br>(10%) | Folds | Stratyfikacja |
|--------------------|-----------|----------------|--------------|---------------|-------|---------------|
| RT-IoT2022         | 123,117   | 98,493         | 12,312       | 12,312        | 10    | TAK           |
| Letter Recognition | 20,000    | 16,000         | 2,000        | 2,000         | 10    | TAK           |
| Electric Power     | 2,049,279 | 1,639,423      | 204,928      | 204,928       | 10    | NIE*          |
| MIMIC-IV-ED        | 196       | 156            | 20           | 20            | 10    | TAK           |

\*Electric Power to dane czasowe (regresja) - stratyfikacja nie ma zastosowania.

## 10. Kluczowe koncepcje

| Koncepcja        | Opis   |
|------------------|--|
| Stratyfikacja    | Zachowanie identycznych proporcji klas w każdym zbiorze i foldzie  |
| Disjointness     | Zbiory $D_{\text{tr}}$ , $D_{\text{val}}$ , $D_{\text{ts}}$ są rozłączne - żaden wiersz się nie powtarza |
| Completeness     | $D_{\text{tr}} \cup D_{\text{val}} \cup D_{\text{ts}} = D$ - każdy wiersz użyty dokładnie raz            |
| Fold Rotation    | Klienci mają te same dane, foldy w różnej kolejności   |
| Shared Val/Test  | Wszystkie klienci używają identycznych $D_{\text{val}}$ i $D_{\text{ts}}$                                |
| Cross-Validation | 10 iteracji: fold $F_i$ jako walidacja lokalna, pozostałe 9 jako trening                                 |

## 11. Wzory matematyczne

**Całkowity rozmiar danych dla K klientów:**

$$\text{Total} = K \times |D_{\text{tr}}| + |D_{\text{val}}| + |D_{\text{ts}}| = (0.80K + 0.20) \times |D|$$

**Warunek stratyfikacji (dla każdej klasy c):**

$$\text{Dla każdego } c: |p(c|F_i) - p(c|D)| < \varepsilon, \text{ gdzie } \varepsilon \rightarrow 0$$

**Rozmiar danych w jednej iteracji CV:**

$$|D_{\text{train}(i)}| = |D_{\text{tr}}| - |F_i| = 0.9 \times |D_{\text{tr}}|$$

$$|D_{\text{val\_local}(i)}| = |F_i| = 0.1 \times |D_{\text{tr}}|$$