

*„Akademia Innowacyjnych Zastosowań Technologii Cyfrowych (AI Tech)”, projekt finansowany ze środków Programu Operacyjnego Polska Cyfrowa POPC.03.02.00-00-0001/20*

# Graniczna Analiza Danych I

## 1 Czym jest efektywność i jak ją ocenić?

**Przykład I:** Firma X produkuje telefony komórkowe w 5 fabrykach. W każdej z nich zatrudnia różną liczbę pracowników i produkuje różną liczbę telefonów dziennie (patrz tabela poniżej). Która fabryka pracuje efektywnie w porównaniu z pozostałymi?

Tabela 1: Tabela zawierająca dane o liczbie pracowników oraz liczbie produkowanych telefonów dla poszczególnych fabryk przykładowej firmy X

Fabryka	liczba pracowników	liczba produkowanych telefonów
A	50	250
B	80	760
C	15	63
D	95	650
E	70	550

**Odpowiedź:** Efektywność pracy fabryki można wyznaczyć jako średnią liczbę telefonów produkowanych przez jednego pracownika w danej firmie:

Tabela 2: Tabela pokazująca jak wyznaczyć efektywność poszczególnych fabryk – liczba produkowanych telefonów / liczba pracowników

Fabryka	produkcja/pracownicy
A	$250 / 50 = 5.00$
B	$760 / 80 = 9.50$
C	$63 / 15 = 4.20$
D	$650 / 95 = 6.84$
E	$550 / 70 = 7.86$

Zatem najlepsza wydajność obserwowana jest dla fabryki B (9.50), natomiast najgorsza dla C (4.20).

**Przykład II:** Ta sama firma X wprowadziła do swojej oferty drugi model telefonu (produkując model X1 oraz X2). Równocześnie oprócz liczby pracowników na efektywność pracy może wpływać dzienna liczba godzin pracy fabryki. Dane dotyczące poszczególnych fabryk w tym przykładzie przedstawione są w poniższej tabeli. Czy można teraz ocenić która z fabryk działa efektywnie? A może kilka z nich?

Tabela 3: Tabela zawierająca rozszerzone dane dotyczące fabryk firmy X – liczbę pracowników, czas pracy fabryki oraz dzienną produkcję telefonów modeli X1 oraz X2

Fabryka	l. pracowników	czas pracy (h)	produkcja X1 (szt/24h)	produkcja X2 (szt/24h)
A	50	12	250	140
B	80	8	416	108
C	15	4	29	19
D	95	6	342	151
E	70	18	550	230

**Odpowiedź:** Graniczna Analiza Danych

## 2 Graniczna Analiza Danych

**Graniczna Analiza Danych** (ang. Data Envelopment Analysis, DEA) jest bezparametryczną metodą oceny efektywności działania jednostek decyzyjnych (ang. decision making unit, DMU), które pobierają wiele nakładów (wejść) i produkują wiele efektów (wyjść). Metoda ta posiada zastosowanie w wielu dziedzinach, m.in.:

- ochrona zdrowia (szpitale, przychodnie, lekarze, ...),

- bankowość (całe banki, poszczególne oddziały, ...),
- szkolnictwo (szkoły, uczelnie, wydziały, systemy edukacyjne państw, ...),
- rolnictwo (gospodarstwa rolne, przedsiębiorstwa produkujące nasiona lub maszyny rolnicze, ...),
- ...

## 2.1 Efektywność

Model efektywności zaproponowany w pierwszym modelu Granicznej Analizy Danych (Charnes, Cooper i Rhodes, 1978) opiera się na intuicyjnej definicji efektywności, będącej ilorazem osiąganych efektów oraz pobieranych nakładów. Ze względu na to, że zarówno nakładów jak i efektów może być wiele zastosowano sumę ważoną do agregacji nakładów i efektów:

$$\text{efektywność} = \frac{\text{ważona suma efektów}}{\text{ważona suma nakładów}}$$

### Główne założenia:

- Efektywność może być mierzone jako stosunek wyjść i wejść danej jednostki,
- jednostki z najlepszym stosunkiem są uważane jako efektywne,
- granica efektywności – linia wyznaczona przez jednostki efektywne pokazująca najlepsze praktyki dla jednostek nieefektywnych,
- jednostki efektywne leżą na granicy efektywności (względna efektywność = 1), nieefektywne leżą poniżej granicy efektywności (względna efektywność mniejsza od 1),
- Hipotetyczna Jednostka Porównawcza (ang. Hypothetical Comparison Unit, HCU) dla jednostki nieefektywnej – sztuczna efektywna jednostka, która leży najbliżej analizowanej jednostki nieefektywnej, pokazuje jak powinna poprawić się analizowana jednostka, aby stać się efektywną.

## 2.2 Notacja

Dane (wartości stałe):

- $K$  – liczba jednostek, ( $k = 1, 2, \dots, K$ ),

- $M$  – liczba nakładów (wejść) w analizowanym problemie ( $m = 1, 2, \dots, M$ ),
- $N$  – liczba efektów (wyjść) w analizowanym problemie ( $n = 1, 2, \dots, N$ ),
- $x_{mk}$  – wartość  $m$ -go wejścia  $k$ -ej jednostki,
- $y_{nk}$  – wartość  $n$ -go wyjścia  $k$ -ej jednostki.

Zmienne:

- $\nu_m$  – waga  $m$ -go nakładu,
- $\mu_n$  – waga  $n$ -go efektu.

Efektywność jednostki  $k$  wyznacza się więc wzorem:

$$E_k = \frac{\sum_{n=1}^N \mu_n \cdot y_{nk}}{\sum_{m=1}^M \nu_m \cdot x_{mk}}$$

**Kombinacja liniowa** – dla skończonej liczby wektorów  $z_1, z_2, \dots, z_h$ , kombinacja liniowa jest zdefiniowana jako:

$$z = \sum_{i=1}^h \lambda_i \cdot z_i = \lambda_1 \cdot z_1 + \lambda_2 \cdot z_2 + \dots + \lambda_h \cdot z_h$$

- kombinacja stożkowa – wszystkie wagi ( $\lambda_i$ ) muszą być nieujemne ( $\lambda_i \geq 0$ ),
- kombinacja wypukła – kombinacja stożkowa, w której dodatkowo suma wag ( $\lambda_i$ ) jest równa 1 ( $\lambda_i \geq 0, \sum_{i=1}^h \lambda_i = 1$ ).

## 2.3 Model CCR

**CCR** – Charnes, Cooper, Rhodes (1978). Miara efektywności bazuje na stałym efekcie skali (ang. constant return to scale, CRS).

Założenia:

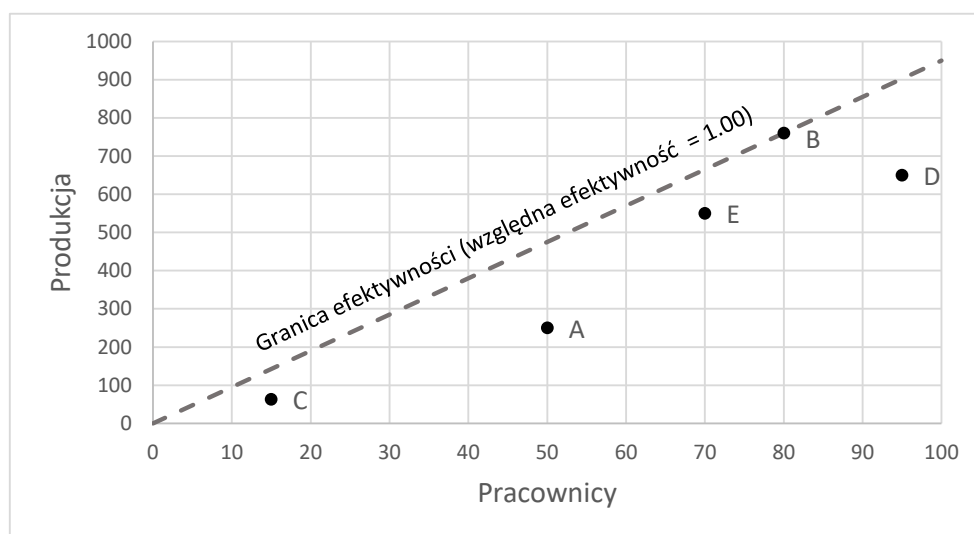
- plany produkcyjne mogą być dowolnie skalowane proporcjonalnie – realne jednostki mogą być utworzone poprzez pomnożenie nakładów i efektów przez tę samą liczbę (wzrost nakładów o  $n\%$  pociąga za sobą wzrost efektów również o  $n\%$ ),
- granica efektywności jest wyznaczana przez jednostkę o najlepszej efektywności (z najlepszym stosunkiem wyjść/wejść) $m$

- względna efektywność jest mierzona jako  $E_k/E_{max}$ , gdzie  $E_{max}$  jest jednostką o najwyższej efektywności.

**Przykład** Wracając do pierwszego przykładu firmy X względna efektywność dla każdej z fabryk przedstawiona jest poniżej:

Tabela 4: Tabela pokazująca jak przeliczyć efektywność fabryk na względną efektywność z zakresu 0-1 (względna efektywność = efektywność / maksymalna efektywność)

Fabryka	l. pracowników NAKŁAD	l. telefonów EFEKT	Efektywność	Względna efektywność
A	50	250	5.00	$5.00/9.50 = 0.53$
B	80	760	9.50	$9.50/9.50 = 1.00$
C	15	63	4.20	$4.20/9.50 = 0.44$
D	95	650	6.84	$6.84/9.50 = 0.72$
E	70	550	7.86	$7.86/9.50 = 0.83$

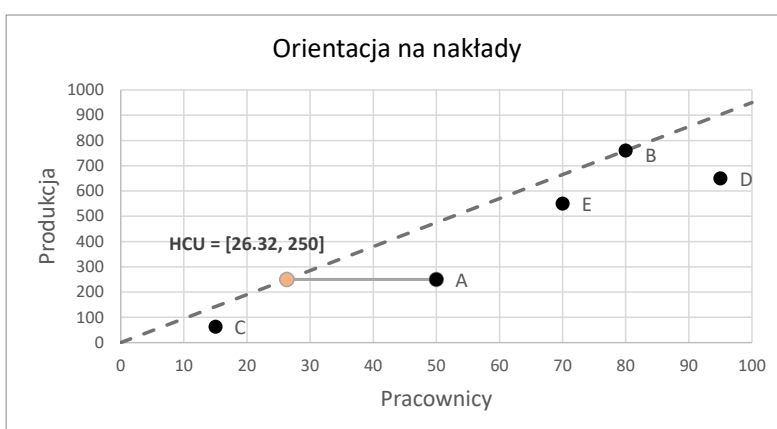


Rysunek 1: Granica efektywności dla przykładu fabryk firmy X

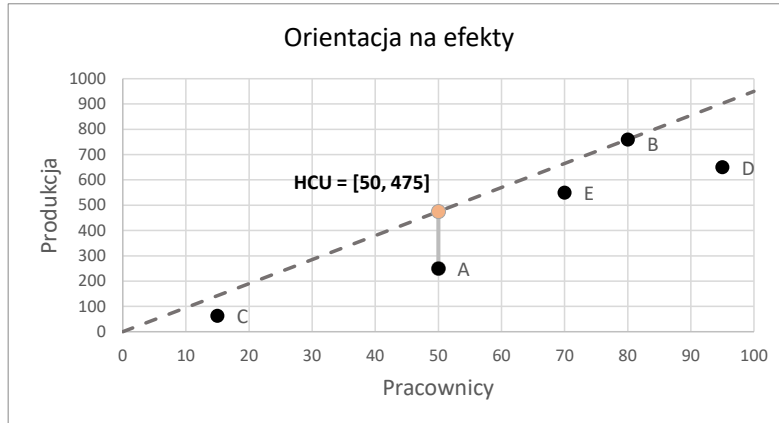
## 2.4 Orientacja na nakłady a orientacja na efekty

Poza wyznaczeniem wartości efektywności dla konkretnej jednostki, analityk może też mieć na celu ocenę jak należy poprawić działanie jednostek nieefektywnych, aby stały się efektywne. Można tu zastosować 2 polityki:

- orientacja na nakłady – o ile mniejsze nakłady powinna przetwarzać dana jednostka przy takich samych efektach, aby stać się efektywną,
- orientacja na efekty – o ile trzeba zwiększyć efekty jednostki nieefektywnej, przy takich niezmiennych nakładach, aby stać się efektywną



Rysunek 2: Hipotetyczna jednostka porównawcza dla modelu zorientowanego na nakłady



Rysunek 3: Hipotetyczna jednostka porównawcza dla modelu zorientowanego na efekty

#### 2.4.1 Model zorientowany na nakłady

**Przestrzeń efektywności** Aby zbadać efektywność jednostki  $DMU_o$  należy rozwiązać następujący problem programowania matematycznego.

$$\begin{aligned}
 \max \quad & E_o = \frac{\sum_{n=1}^N \mu_n \cdot y_{no}}{\sum_{m=1}^M \nu_m \cdot x_{mo}} \\
 \text{p.o.} \quad & \frac{\sum_{n=1}^N \mu_n \cdot y_{nk}}{\sum_{m=1}^M \nu_m \cdot x_{mk}} \leq 1, & k = 1, 2, \dots, K \\
 & \mu_n, \nu_n \geq 0 & m = 1, 2, \dots, M, n = 1, 2, \dots, N
 \end{aligned}$$

W tym problemie szukamy takiego zestawu wag nakładów i efektów ( $\mu_n$  oraz  $\nu_m$ ), dla którego względna efektywność badanej jednostki jest możliwie największa ograniczając przy tym efektywność wszystkich jednostek do wartości nie większych niż 1. W ten sposób znajdujemy najkorzystniejszy zestaw wag dla  $DMU_o$ . **Jeśli optymalna efektywność (funkcja celu) w tym modelu jest równa 1 to analizowana jednostka jest efektywna, w przeciwnym wypadku jest nieefektywna.**

Ze względu na to, że oryginalny model jest nieliniowy, nie może on być rozwiązany przy użyciu standardowych solverów problemów programowania liniowego. W prosty sposób można sprowadzić powyższy model do modelu liniowego (przy użyciu transformacji Charnesa-Coopera). W tym celu zakładamy, że mianownik

funkcji celu ma ustaloną wartość równą 1. Wtedy cały ułamek jest tym większy im większy jest licznik, więc wystarczy zmaksymalizować wartość licznika miary efektywności oraz dodać ograniczenie na stałą wartość mianownika. Ograniczenia na maksymalną efektywność wystarczy wymnożyć przez mianownik, aby otrzymać formę liniową (można to zrobić ze względu na nieujemność wszystkich elementów). Otrzymany model liniowy, przedstawiony poniżej można już rozwiązać przy pomocy dowolnego solwera liniowego.

$$\begin{aligned}
& \max \quad \sum_{n=1}^N \mu_n \cdot y_{no} \\
& \text{p.o.} \quad \sum_{m=1}^M \nu_m \cdot x_{mo} = 1 \\
& \quad \sum_{n=1}^N \mu_n \cdot y_{nk} \leq \sum_{m=1}^M \nu_m \cdot x_{mk}, \quad k = 1, 2, \dots, K \\
& \quad \mu_n, \nu_m \geq 0 \quad m = 1, 2, \dots, M, n = 1, 2, \dots, N
\end{aligned}$$

**Przestrzeń kombinacji jednostek** Alternatywnie do szukania maksymalnej efektywności, można zastosować model w przestrzeni kombinacji liniowych jednostek. Pod względem matematycznym jest to **model dualny** w stosunku do modelu w przestrzeni efektywności, więc optymalny wynik w obu przestrzeniach jest taki sam.

W tym przypadku szukamy takiej **kombinacji stożkowej** istniejących jednostek, która jest lepsza (lub przynajmniej niegorsza od badanej jednostki), tzn. ma co najmniej takie same nakłady i co najmniej takie same efekty. W przypadku orientacji na nakłady równocześnie model szuka najmniejszej wartości ( $\theta$ ), przez którą da się pomnożyć wszystkie wejścia badanej jednostki utrzymując nie gorsze wyjścia. Inaczej mówiąc szukana hipotetyczna jednostka musi mieć wejścia nie większe niż wejścia badanej jednostki, pomniejszone  $\theta$  razy oraz wyjścia co najmniej takie jak badana jednostka. Minimalna wartość zmiennej  $\theta$  jest równa efektywności badanej jednostki. W przypadku jednostek efektywnych szukaną kombinacją będzie badana jednostka ( $\lambda_o = 1$ , pozostałe  $\lambda_k = 0$ ).



$$\begin{aligned}
& \min \quad \theta \\
& \text{p.o.} \quad \sum_{k=1}^K \lambda_k \cdot x_{nk} \leq \theta \cdot x_{no}, & n = 1, 2, \dots, N \\
& \quad \sum_{k=1}^K \lambda_k \cdot y_{nk} \geq y_{mo}, & m = 1, 2, \dots, M \\
& \quad \theta \geq 0, \\
& \quad \lambda_k \geq 0, & k = 1, 2, \dots, K
\end{aligned}$$

Zmienne  $\lambda_k$  optymalnego rozwiązania pozwalają na znalezienie jednostki będącej hipotetyczną jednostką porównawczą (HCU) badanej. Na podstawie HCU możemy określić potrzebną poprawę badanej jednostki aby uzyskać efektywność:

$$\begin{aligned}
x_{m,HCU} &= \sum_{k=1}^K \lambda_k \cdot x_{mk} \\
y_{n,HCU} &= \sum_{k=1}^K \lambda_k \cdot y_{nk}
\end{aligned}$$

Poprawki potrzebne do osiągnięcia efektywności to różnica pomiędzy wartościami wejść i wyjść jednostki HCU oraz badanej:

$$\begin{aligned}
\Delta x_m &= x_{mo} - x_{m,HCU} \\
\Delta y_y &= y_{n,HCU} - y_{no}
\end{aligned}$$

#### 2.4.2 Model zorientowany na efekty

Aby wyznaczyć efektywność jednostki  $DMU_o$  w przestrzeni efektywności przy orientacji na efekty należy rozwiązać model przedstawiony poniżej. W tym przypadku **minimalizowana** jest ważona suma wejść przy ustalonej wartości ważonej sumy wyjść (równiej 1). Pozostałe ograniczenia (względna efektywność nie większa niż 1 i nieujemne wagi) są takie same jak w modelu zorientowanym na nakłady.

$$\begin{aligned}
\min \quad & \sum_{m=1}^M \nu_n \cdot x_{mo} \\
\text{p.o.} \quad & \sum_{n=1}^N \mu_n \cdot y_{no} = 1 \\
& \sum_{n=1}^N \mu_n \cdot y_{nk} \leq \sum_{m=1}^M \nu_n \cdot x_{mk}, & k = 1, 2, \dots, K \\
& \mu_n, \nu_n \geq 0 & m = 1, 2, \dots, M, n = 1, 2, \dots, N
\end{aligned}$$

Ostateczna efektywność badanej jednostki jest **odwrotnością** wartości funkcji celu uzyskanej w tym modelu.

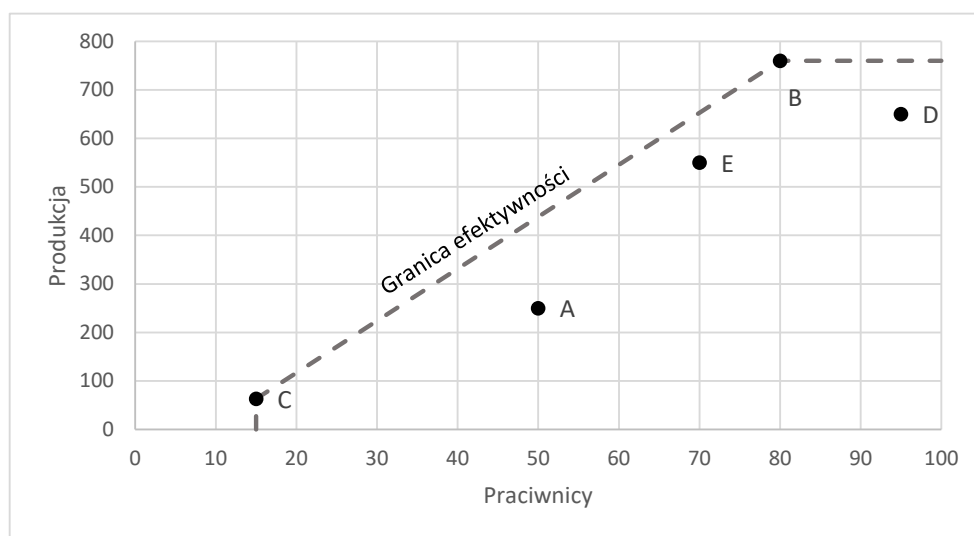
W przestrzeni kombinacji jednostek model zorientowany na efekty jest podobny do zorientowanego na nakłady, jednak tym razem wartość zmiennej  $\theta$  jest maksymalizowana oraz przez tę zmienną mnożone są **wyjścia** badanej jednostki (po prawej stronie ograniczeń). Wejścia pozostają wartościami wejściowymi. Podobnie jak w przestrzeni efektywności wartość efektywności badanej jednostki jest odwrotnością optymalnej wartości  $\theta$ .

$$\begin{aligned}
\max \quad & \theta \\
\text{p.o.} \quad & \sum_{k=1}^K \lambda_k \cdot x_{nk} \leq x_{no}, & n = 1, 2, \dots, N \\
& \sum_{k=1}^K \lambda_k \cdot y_{nk} \geq \theta \cdot y_{mo}, & m = 1, 2, \dots, M \\
& \theta \geq 0, \\
& \lambda_k \geq 0, & k = 1, 2, \dots, K
\end{aligned}$$

## 2.5 Model BCC

**Model BCC** (Banker, Charnes, Cooper, 1984) jest modelem granicznej analizy danych, w który, zakłada się **zmienny efekt skali** (ang. variable return to scale, VRS). Oznacza to, że brak założenia, że proporcjonalny wzrost nakładów przekłada się na proporcjonalny wzrost efektów. W tym przypadku zakłada się, że możliwe są tylko plany produkcyjne będące **kombinacją wypukłą** istniejących planów produkcyjnych.

Dla przykładu firmy X (z jednym nakładem i jednym efektem) granica efektywności dla modelu BCC przedstawiony jest na rysunku poniżej.



Rysunek 4: Granica efektywności dla przykładu fabryk firmy X w modelu BCC

W przestrzeni kombinacji liniowych jedyna różnica w modelach matematycznych w stosunku do modelu CCR polega na dodaniu dodatkowego ograniczenia gwarantującego sumę wag ( $\lambda_k$ ) równą 1.

### Orientacja na nakłady:

$$\begin{aligned}
& \min \quad \theta \\
& \text{p.o.} \quad \sum_{k=1}^K \lambda_k \cdot x_{nk} \leq \theta \cdot x_{no}, & n = 1, 2, \dots, N \\
& \quad \sum_{k=1}^K \lambda_k \cdot y_{nk} \geq y_{mo}, & m = 1, 2, \dots, M \\
& \quad \sum_{k=1}^K \lambda_k = 1 \\
& \quad \theta \geq 0, \\
& \quad \lambda_k \geq 0, & k = 1, 2, \dots, K
\end{aligned}$$

### Orientacja na efekty:

$$\begin{aligned}
& \max \quad \theta \\
& \text{p.o.} \quad \sum_{k=1}^K \lambda_k \cdot x_{nk} \leq x_{no}, & n = 1, 2, \dots, N \\
& \quad \sum_{k=1}^K \lambda_k \cdot y_{nk} \geq \theta \cdot y_{mo}, & m = 1, 2, \dots, M \\
& \quad \sum_{k=1}^K \lambda_k = 1 \\
& \quad \theta \geq 0, \\
& \quad \lambda_k \geq 0, & k = 1, 2, \dots, K
\end{aligned}$$

W przypadku przestrzeni miary efektywności model musi uwzględniać dodatkową zmienną ( $\mu_0$  lub  $\nu_0$ ) reprezentującą punkt, w którym granica efektywności przecina oś pionową. Zmienna ta jest uwzględniona zarówno w funkcji celu jak i w ograniczeniach. Przy modelu zorientowanym na nakłady dodatkowa zmienna występuje przy części dotyczącej wyjść, natomiast w modelu zorientowanym na efekty w części dotyczącej wejść. **Dodatkowa zmienna ( $\mu_0$  lub  $\nu_0$ ) nie ma ograniczenia na znak – może przyjmować również wartości ujemne.**

Orientacja na nakłady:

$$\begin{aligned}
 \max \quad & \sum_{n=1}^N \mu_n \cdot y_{no} + \mu_0 \\
 \text{p.o.} \quad & \sum_{m=1}^M \nu_m \cdot x_{mo} = 1 \\
 & \sum_{n=1}^N \mu_n \cdot y_{nk} + \mu_0 \leq \sum_{m=1}^M \nu_m \cdot x_{mk}, & k = 1, 2, \dots, K \\
 & \mu_n, \nu_m \geq 0 & m = 1, 2, \dots, M, n = 1, 2, \dots, N \\
 & \mu_0 \text{ free}
 \end{aligned}$$

Orientacja na efekty:

$$\begin{aligned}
 \min \quad & \sum_{m=1}^M \nu_m \cdot x_{mo} + \nu_0 \\
 \text{p.o.} \quad & \sum_{n=1}^N \mu_n \cdot y_{no} = 1 \\
 & \sum_{n=1}^N \mu_n \cdot y_{nk} \leq \sum_{m=1}^M \nu_m \cdot x_{mk} + \nu_0, & k = 1, 2, \dots, K \\
 & \mu_n, \nu_m \geq 0 & m = 1, 2, \dots, M, n = 1, 2, \dots, N \\
 & \nu_0 \text{ free}
 \end{aligned}$$

### 3 Zadanie domowe - część I

Dany jest zbiór danych zawierający informacje o lotniskach w 11 polskich miastach:

- Warszawa (WAW),
- Kraków (KRK),
- Wrocław (WRO),
- Poznań (POZ),
- Łódź (LCJ),

- Gdańsk (GDN),
- Szczecin (SZZ),
- Bydgoszcz (BZG),
- Rzeszów (RZE),
- Zielona Góra (IEG)

ocenionych przy pomocy 4 wejść i 2 wyjść.

- wejścia:
  - i1: roczna przepustowość terminalu zdefiniowana jako przepływ pasażerów, który port lotniczy może obsłużyć bez poważnych niedogodności (w milionach pasażerów rocznie);
  - i2: maksymalna przepustowość zdefiniowana jako średnia liczba operacji (przylotów i/lub odlotów), które można wykonać na pasach startowych portu lotniczego (w liczbie operacji na godzinę);
  - i3: przepustowość płyty postojowej lotniska definiowana jako średnia liczba samolotów, które może obsłużyć lotnisko (w liczbie samolotów na godzinę);
  - i4: obszarciążenia lotniska zdefiniowany jako liczba mieszkańców mieszkających w promieniu 100 km od lotniska (w mln mieszkańców);
- wyjścia:
  - o1: ruch pasażerski mierzony całkowitą liczbą pasażerów obsługanych przez port (w mln pasażerów rocznie);
  - o2: liczba operacji statku powietrznego (jeden całkowity ruch to lądowanie lub start statku powietrznego) (w tysiącach ruchów rocznie).

Oceny poszczególnych lotnisk dostępne są odpowiednio w plikach inputs.csv oraz outptus.csv.

Dane zostały zaczerpnięte z: Kadziński M., Labijak A., Napieraj M., "Integrated Framework for Robustness Analysis Using Data Envelopment Model with Application to Efficiency Analysis of Polish Airports", Omega, Vol. 67, 1-18, 2017.

Przy użyciu biblioteki PULP napisz skrypt w pythonie, który dla każdej jednostki wyznaczy jej miarę efektywności. Dodatkowo dla każdej z jednostek nieefektywnych znajdź hipotetyczną jednostkę porównawczą oraz poprawki potrzebne do

osiągnięcia efektywności. Użyj modelu CCR granicznej analizy danych **zorientowanego na nakłady**.

Skrypt powinien wczytywać dane z plików CSV tak, aby po zmianie danego pliku na inny (o takiej samej strukturze) obliczał on efektywność jednostek (oraz inne wartości podane w treści zadania) dla nowych danych.

Na kolejnych zajęciach konieczne będzie uzupełnienie skryptu o dodatkowe funkcjonalności. Oprócz samego kodu źródłowego należy przygotować krótkie sprawozdanie (w formie dokumentu PDF) zawierające wyniki obu części zadania (z dzisiejszych oraz kolejnych zajęć).



Rzeczpospolita  
Polska

Unia Europejska  
Europejski Fundusz  
Rozwoju Regionalnego



*„Akademia Innowacyjnych Zastosowań Technologii Cyfrowych (AI Tech)”, projekt finansowany ze środków Programu Operacyjnego Polska Cyfrowa POPC.03.02.00-00-0001/20*