

Integracja usługi centralnej z rejestrami państwowymi (reuse)

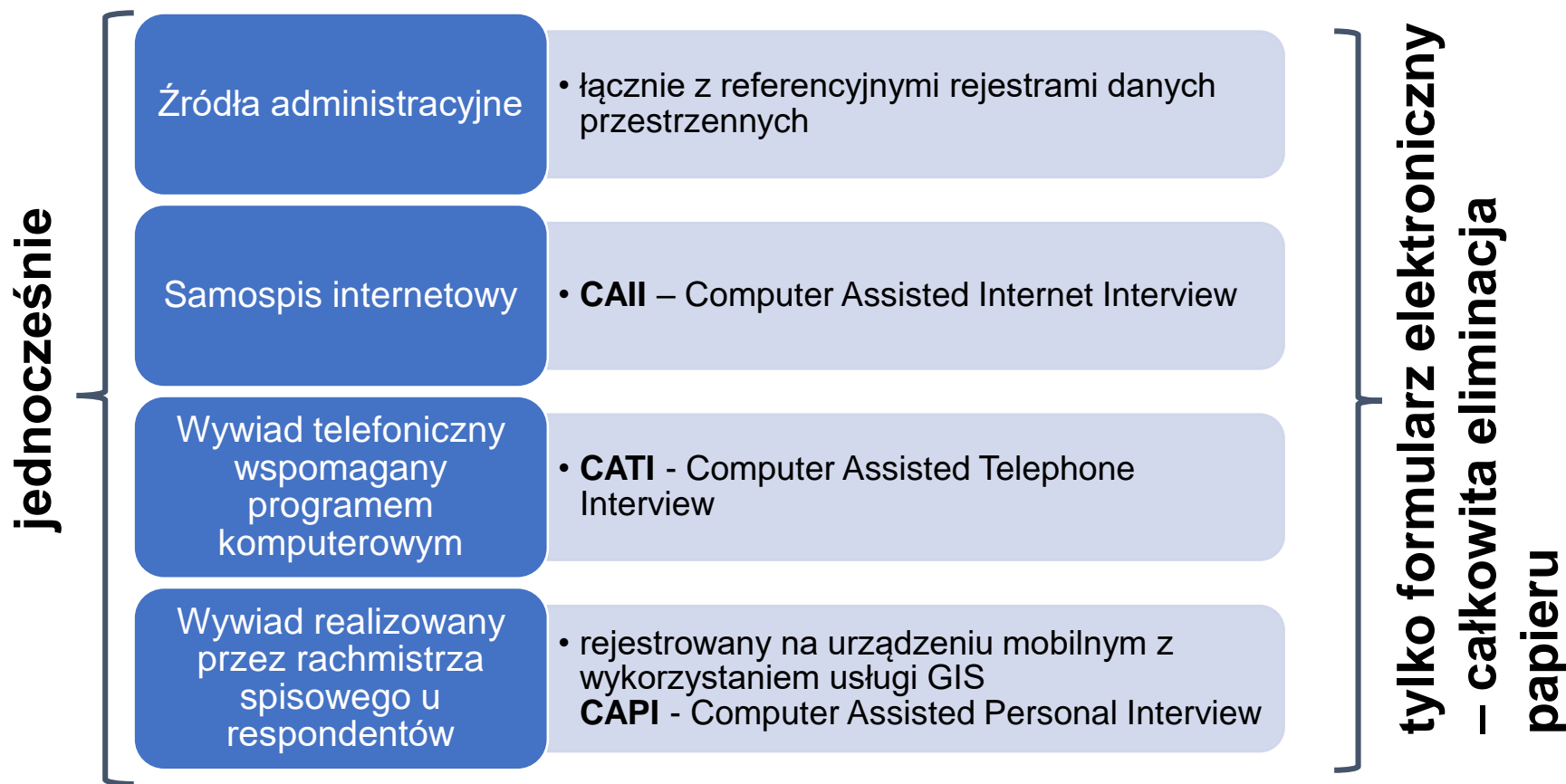
Wykorzystanie źródeł informacji
w badaniach statystycznych

dr inż. Janusz Dygaszewicz

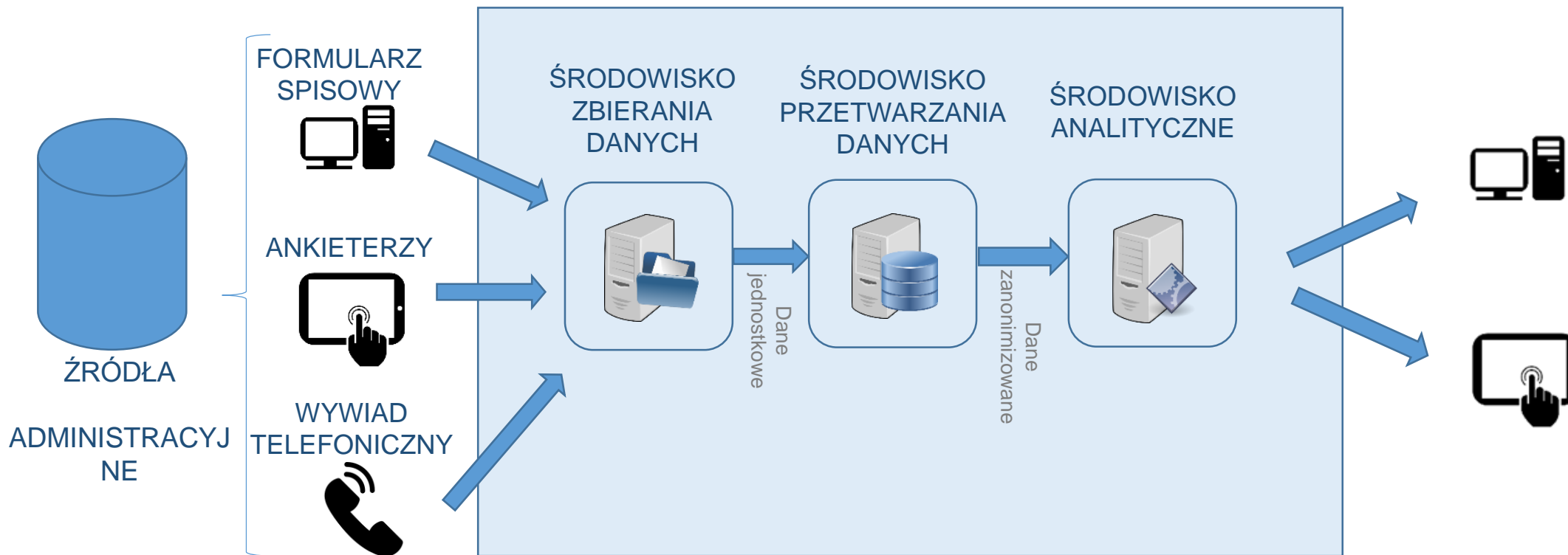
Dyrektor Departamentu Systemów Teleinformatycznych,
Geostatystyki i Spisów

Elektronizacja spisów - odejście od papieru

Spis powszechny bez papieru - elektroniczne kanały zbierania danych w spisach



Architektura - przepływ danych



Porównanie spisów 2011 i 2021

2011



18 tys. rachmistrzów
0 formularzy
0 ton papieru



2 rodzaje formularzy spisowych:
- badanie pełne - 15 pytań
- badanie reprezentacyjne - 100 pytań



Koszt spisu –
395 mln zł (83,1 USD)



Wykorzystano dane
z 27 rejestrów

2021



16 tys. rachmistrzów
0 formularzy
0 ton papieru



Jeden formularz spisowy (58 pytań):
- kwestionariusz osobowy – 25 pytań
- kwestionariusz mieszkania – 12 pytań



Koszt spisu –
Plan 386 mln zł (81,2 USD)
Wykonanie planu 273 mln zł (57,5 USD)



Wykorzystano dane
z 35 rejestrów

Wykorzystanie danych administracyjnych w spisach

1. Źródło danych do budowy wykazu osób do spisu.
2. Źródło danych do budowy wykazu adresowego.
3. Źródło danych do kontroli informacji formularza spisowego.
4. Z danych administracyjnych pochodzą również informacje pozwalające na identyfikację i autentykację osób oraz optymalne zarządzanie procesem spisowym i pracą rachmistrzów w terenie.

Aplikacja formularzowa - logowanie

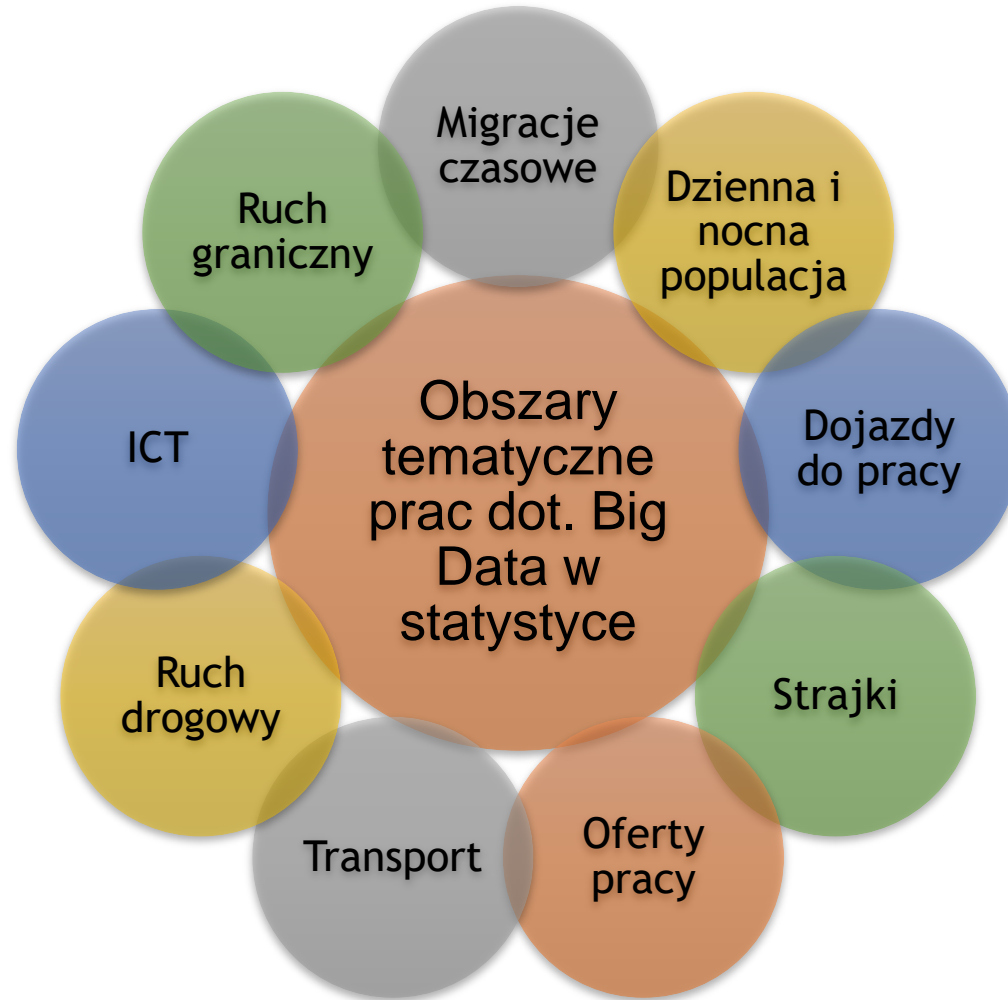
Logowanie do aplikacji formularzowej w celu dokonania samospisu internetowego było możliwe na 3 sposoby:

- **Krajowy Węzeł Identyfikacji Elektronicznej** – logowanie za pomocą środków identyfikacji elektronicznej (Profil Zaufany, bankowość internetowa);
- **numer PESEL (PIN) i nazwisko rodowe matki**, które wymagało od respondenta natychmiastowego zdefiniowania indywidualnego hasła dostępu;
- **adres email oraz indywidualne hasło dostępu** – metoda ta była przeznaczona dla cudzoziemców nieposiadających numeru PESEL.

Wykorzystanie danych administracyjnych w spisach (cd.)

5. Obliczenie wyników spisu:
 - jako podstawowe źródło danych;
 - jako pomocnicze źródło danych dla części informacji.
6. Udostępnienie danych wynikowych na bardzo niskich poziomach agregacji terytorialnej (w tym na poziomie siatki 1 km - grid).
7. W przypadku niepełnego pokrycia badanej populacji zebranymi informacjami – wsparcie imputacja danych.

Prace nad łączeniem danych administracyjnych ze źródłami Big Data - obszary tematyczne



Wyzwania



Spisy powszechne

**Cyfryzacja rejestrów publicznych podstawą do tworzenia
wykazu
adresowo-mieszkaniowo-osobowego**

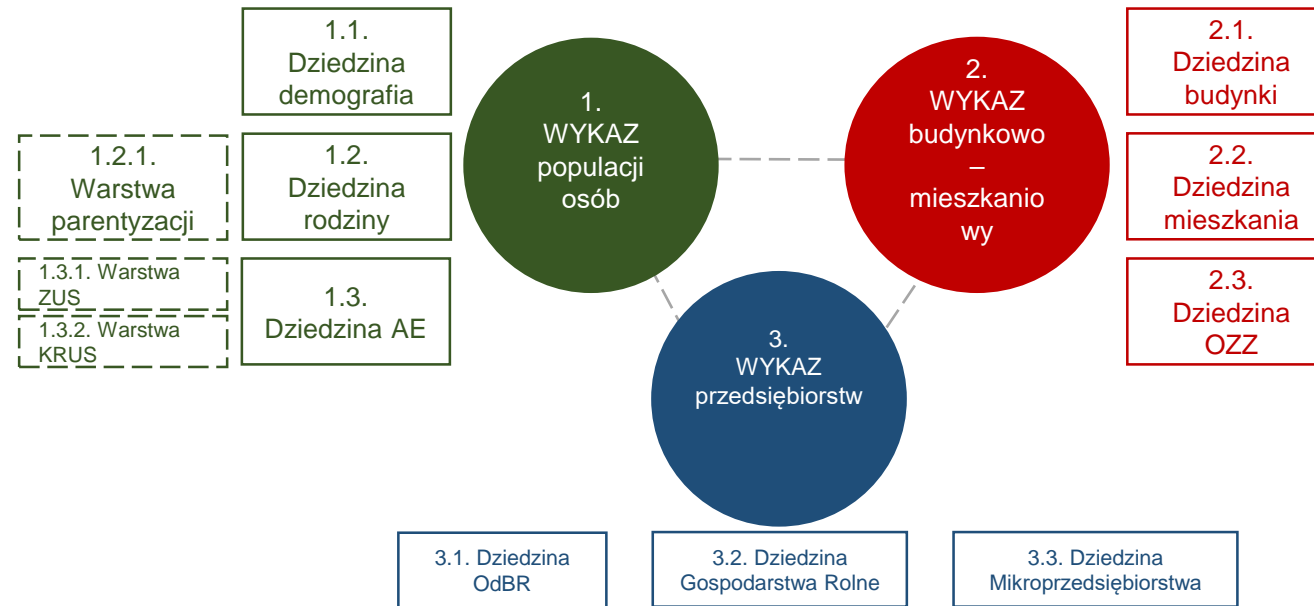
Statystyczna Baza Spisowa

Statystyczna Baza Spisowa - system zmiennych statystycznych powstałych głównie na bazie rejestrów administracyjnych w wyniku wykorzystania jednego źródła lub połączenia kilku źródeł.

System gromadzący w swoich zasobach dane m.in. w zakresie informacyjnym przewidzianym w spisie 2021 oraz w spisach rocznych tj. po roku 2021.

Kompleksowe narzędzie umożliwiające generowanie wynikowych informacji statystycznych w sposób zautomatyzowany i dynamiczny na podstawie cyfrowych danych wejściowych oraz algorytmów opisujących metodę przetwarzania tych danych.

Statystyczna Baza Spisowa



Statystyczna Baza Spisowa

Zakłada istnienie 2 typów rejestrów:

- systemy referencyjne („rdzenie”) - cechują się najwyższą jakością danych, kompletnością i aktualnością oraz wysokim pokryciem podmiotowym;
- systemy zasilające - pozostałe systemy, rejestry, bazy danych, wykazy, ewidencje.

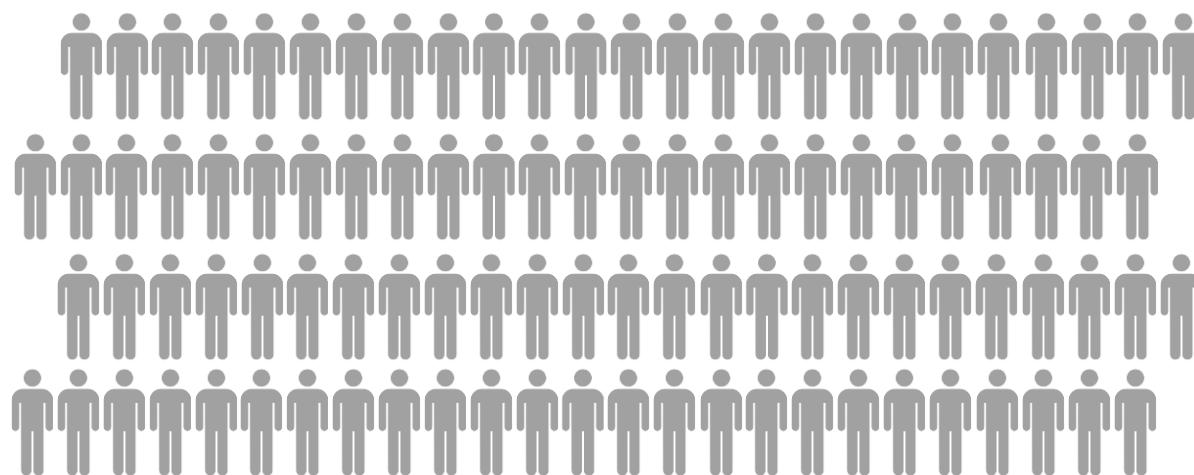
Rejestry wykorzystane do budowy wykazu

- Budowa populacji – 6 rejestrów (PESEL, KEP, NFZ, ZUS i KRUS, Energetyka);
- Budowa wykazu adresowo-mieszkaniowego – 47 rejestrów.

Wykaz populacji osób

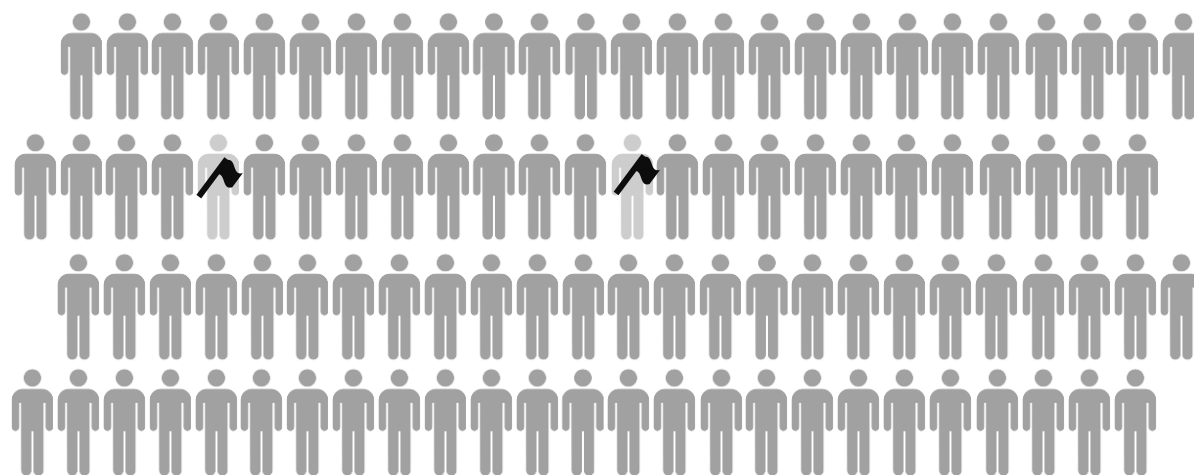
1. Określenie liczby osób aktywnych w Polsce.
2. W SBS przechowywane są informacje o wszystkich osobach, które oznaczane są odpowiednimi wyróżnikami (tzw. flagi).
3. Określenie, dla osób przebywających w Polsce, faktycznego adresu zamieszkania, który posłuży do tworzenia bilansów ludności.

Pula numerów PESEL



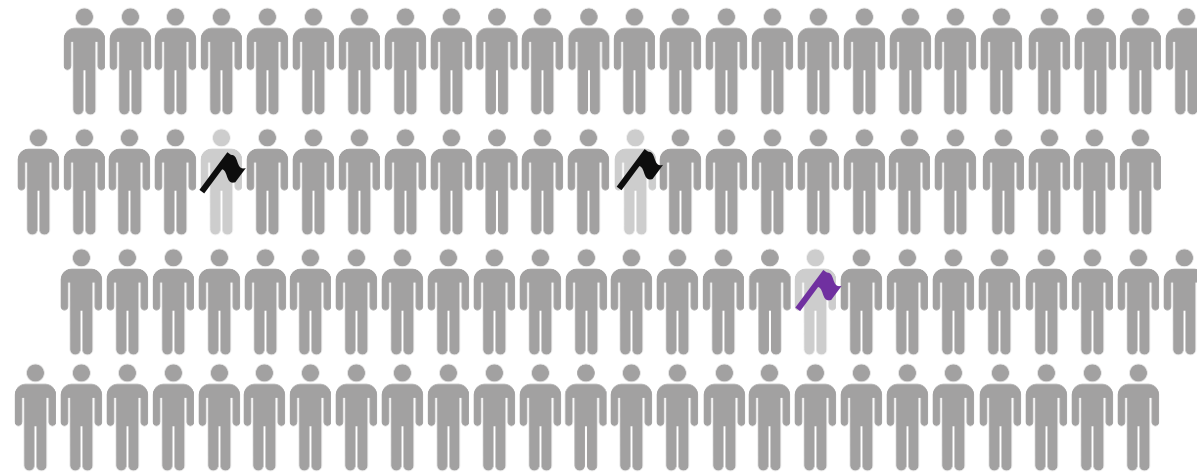
40 829 763

Weryfikacja osób zmarłych tzw. twardy zgon



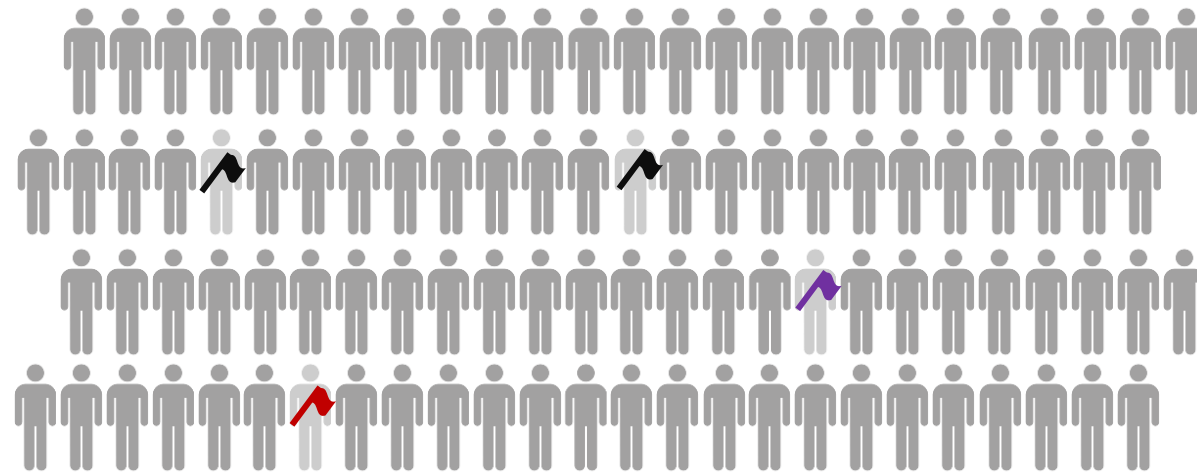
40 522 011

Weryfikacja osób 90+ tzw. miękki zgon



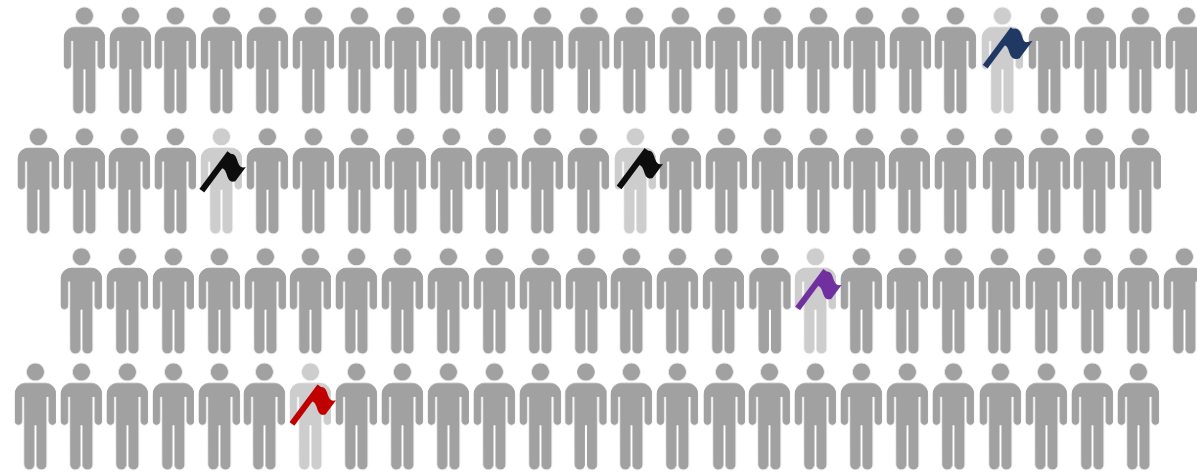
40 396 073

Weryfikacja osób, które wyjechały z Polski



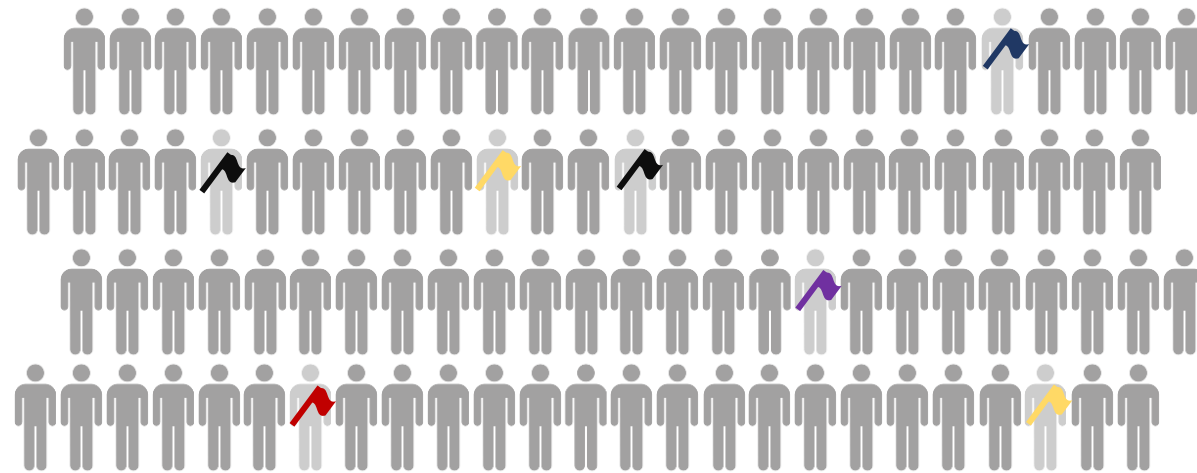
40 324 811

Weryfikacja osób zamieszkałych i pobierających świadczenia poza RP



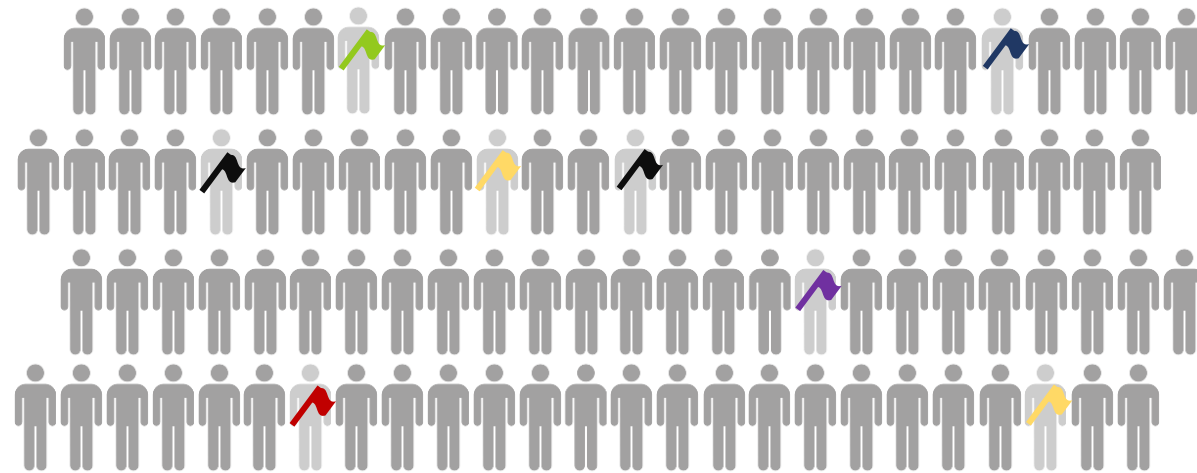
40 309 512

Weryfikacja o osoby występujące tylko w jednym rejestrze



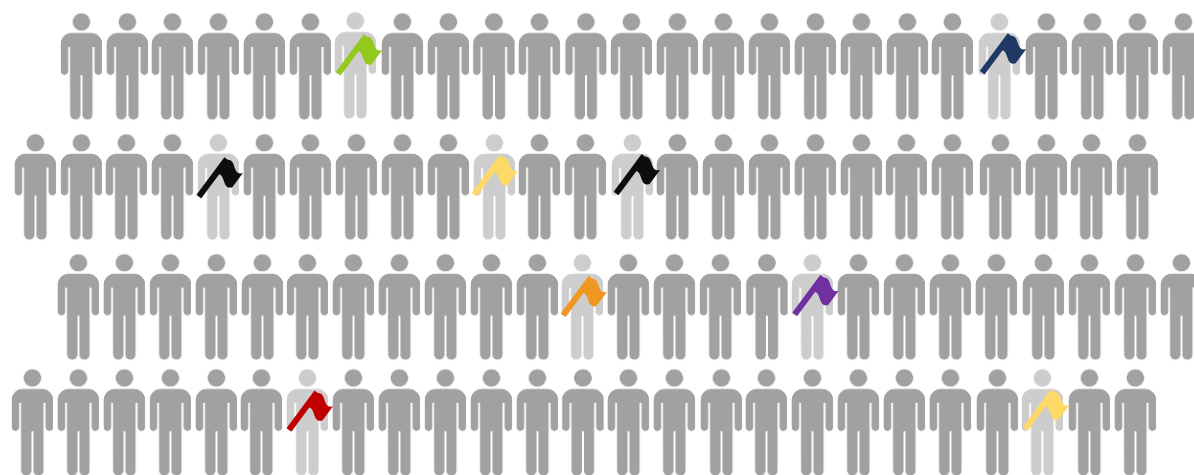
38 566 549

Weryfikacja o osoby posiadające adres zamieszkania poza granicami Polski



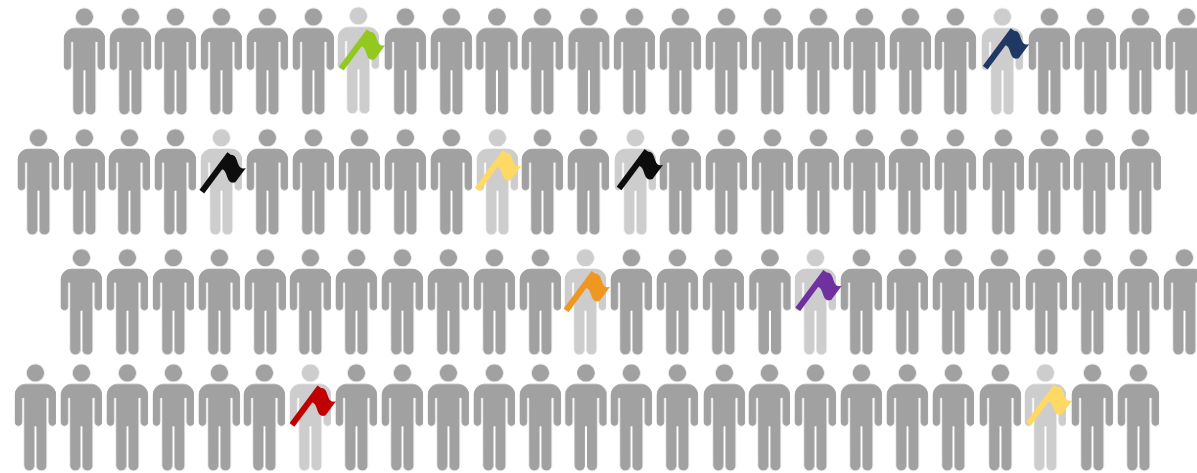
38 400 961

Weryfikacja o dzieci, które nie realizują obowiązku szkolnego w Polsce



38 389 320

Populacja celu



38 389 320

Wykaz populacji osób

Liczba numerów PESEL w populacji celu według pochodzenia ze zbiorów

Lp.	Pochodzenie	Liczba nr PESEL
1	PESEL	37 014 453
2	KEP	1 313 436
3	ZUS	59 534
4	KRUS	1 451
5	NFZ	446
6	ARiMR	0
7	WYMELD	0
	ogółem	38 389 320

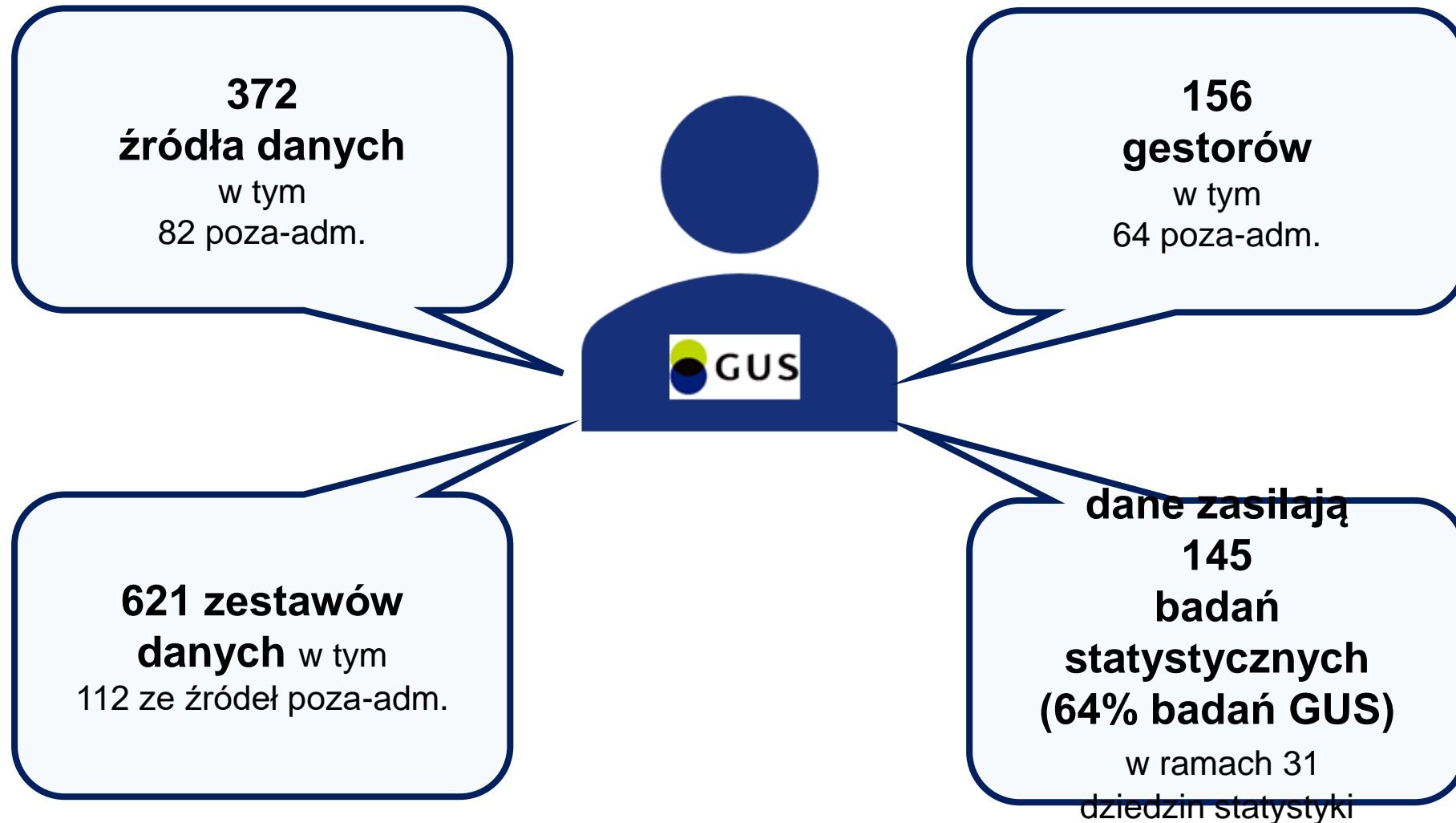
Administracyjne i poza-administracyjne źródła danych w eStatystyce

Wykorzystanie przez GUS administracyjnych i poza-administracyjnych źródeł danych



Na etapie projektowania lub zmiany prowadzonego badania GUS w pierwszej kolejności wykorzystuje dane z rejestrów urzędowych i systemów informacyjnych administracji publicznej.

Wykorzystanie przez GUS administracyjnych i poza-administracyjnych źródeł danych według pbssp2024



Korzyści

większa elastyczność i szybsza reakcja statystyki na potrzeby informacyjne użytkowników
- zapewnienie informacji adekwatnych do potrzeb

redukcja obciążeń administracyjnych

poprawa efektywności (obniżenie kosztów) tworzenia statystyk

możliwość rezygnacji z tradycyjnego zbierania danych na formularzach statystycznych

zwiększenie dokładności danych wynikowych i oszacowań

zwiększenie aktualności danych wynikowych oraz skrócenie czasu przetwarzania danych

uzyskanie możliwości łączenia zbiorów danych administracyjnych i statystycznych

Internetowe źródła danych, webscraping, analiza tekstu, uczenie maszynowe w eStatystyce - na przykładzie badania dóbr konsumpcyjnych

Internetowe źródła danych w CPI (inflacja)

- Wykorzystanie statystyki wieloźródłowej do pomiaru wskaźnika cen towarów i usług konsumpcyjnych CPI czyli popularnej miary inflacji;
- Wyzwania wobec nowych metod i sposobów pracy;
- Rozwój statystyki eksperymentalnej;
- Stosowanie zaawansowanych metod i technik uczenia maszynowego przy klasyfikowaniu produktów i rozpoznawaniu tekstu przy skanowaniu i skrapowaniu danych.

Internetowe źródła danych w CPI (inflacja)

- Ogromny, dynamiczny, różnorodny wolumen nieustrukturyzowanych danych skanowanych i skrapowanych;
- Dostęp do danych w czasie rzeczywistym – economy based on data versus economy driven by data;
- Kwestie interoperacyjności i dostosowania internetowych źródeł danych do celów statystycznych;
- Ocena jakości danych na wejściu i danych na wyjściu, tajemnica statystyczna, tajemnica biznesowa, handlowa, zarządzanie ryzykiem, zasada partnerskiej współpracy sektora publicznego i prywatnego oparta na etyce i fair play.

Dane są wszędzie!

BigData i IoT w produkcji eStatystyk transportu drogowego i morskiego

Wsparcie dla kształtowania i monitorowania polityki transportowej kraju

System TranStat

- Inteligentny system produkcji statystyk transportu drogowego i morskiego z wykorzystaniem wielkich wolumenów danych na rzecz kształtowania polityki transportowej kraju.

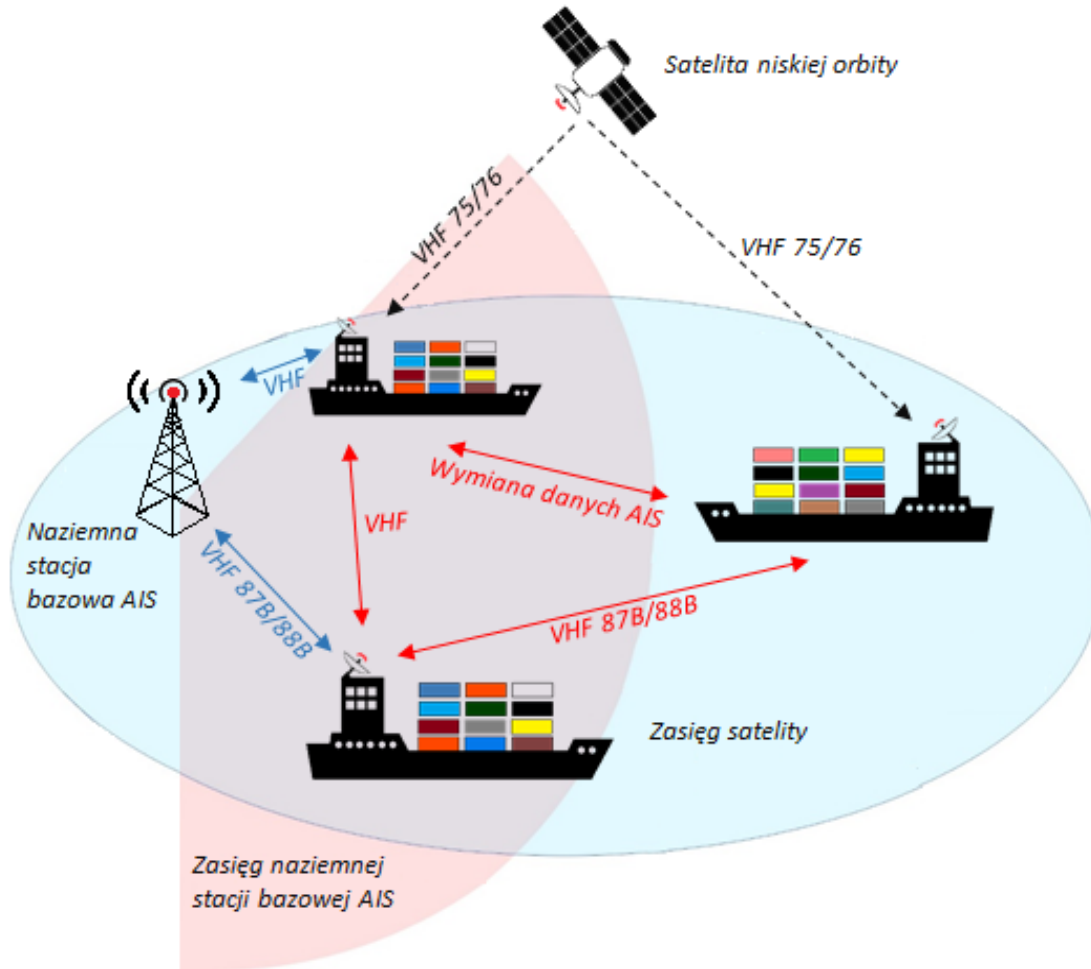
- **Główny cel:**

unowocześnienie systemu produkcji statystyk dot. transportu:

- drogowego;
- morskiego;

poprzez wykorzystanie wielkich zbiorów danych, wprowadzenie nowych produktów oraz stworzenie warunków do prowadzenia analiz funkcjonowania oraz zarządzania systemami transportowymi do nowoczesnego planowania i kształtowania polityki transportowej kraju.

System TranStat - źródła danych - AIS



Źródło: <https://ws.stat.gov.pl/Article/2023/12/001-024>

Dane dynamiczne

- Informacja o ruchu statków
- Automatyczna transmisja danych
- W zależności od prędkości częstotliwość transmisji: 2-10 sek.
- 3 – 6 min. gdy statek jest na kotwicy

- Numer MMSI (Maritime Mobile Service Identity) - dane identyfikacyjne statku;
- Długość geograficzna;
- Szerokość geograficzna;
- Wskazanie klasy dokładności;
- Prędkość nad dnem;
- Aktualne wartości kąta drogi nad dnem;
- Prędkość kątowna zwrotu;
- Status nawigacyjny statku;
- Czas uniwersalny (UTC).

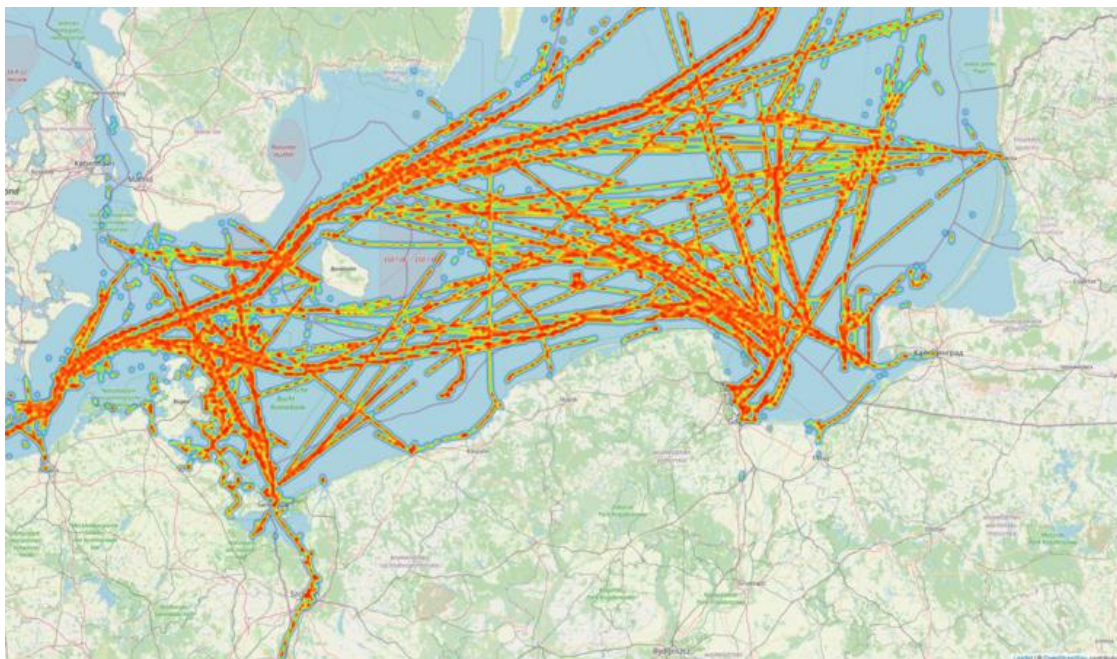
Dane statyczne

- Charakterystyka statku
- Wprowadzane ręcznie
- Co 6 min.

- Numer IMO (International Maritime Organisation) - dane identyfikacyjne statku;
- Numer MMSI (Maritime Mobile Service Identity) - dane identyfikacyjne statku;
- Nazwa statku;
- Sygnał wywoławczy;
- Wymiary statku;
- Typ statku;
- Port docelowy;
- Zanurzenie statku.
- ETA (estimated time of arrival)

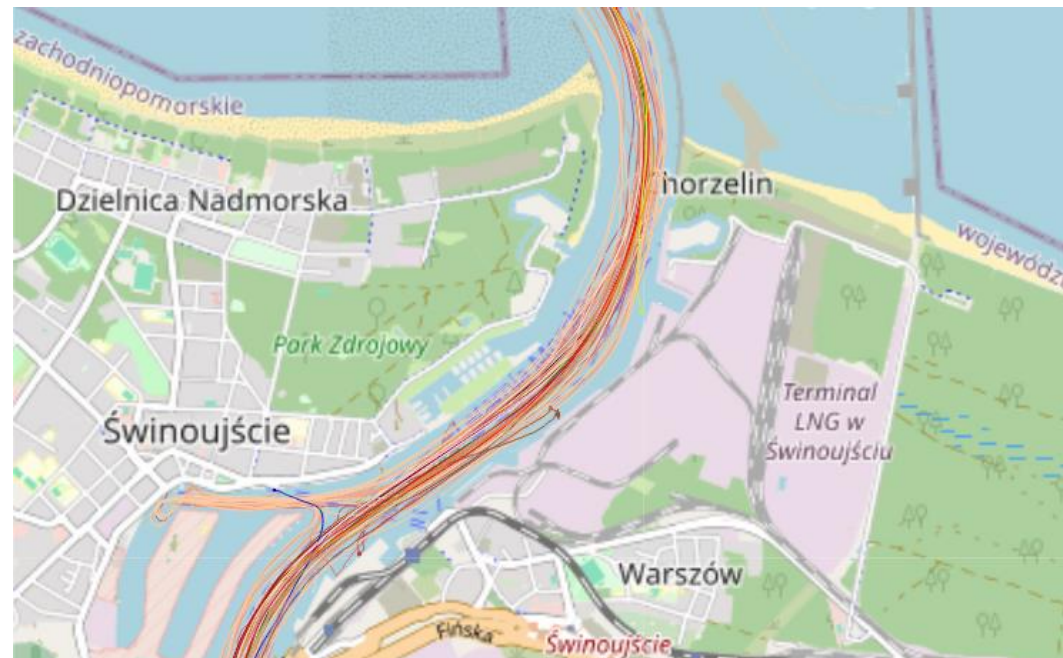
System TranStat - przykładowe wizualizacje

Natężenie ruchu statków



Rys. Natężenie ruchu statków dla Polski (01.06.2024).

Źródło: Opracowane na podst. systemu TranStat



Rys. Natężenie ruchu dla portu w Świnoujściu (01.06.2024).

System TranStat - źródła danych - viaToll/eToll



Źródło: <https://ws.stat.gov.pl/Article/2023/12/001-024>

viaToll/eToll

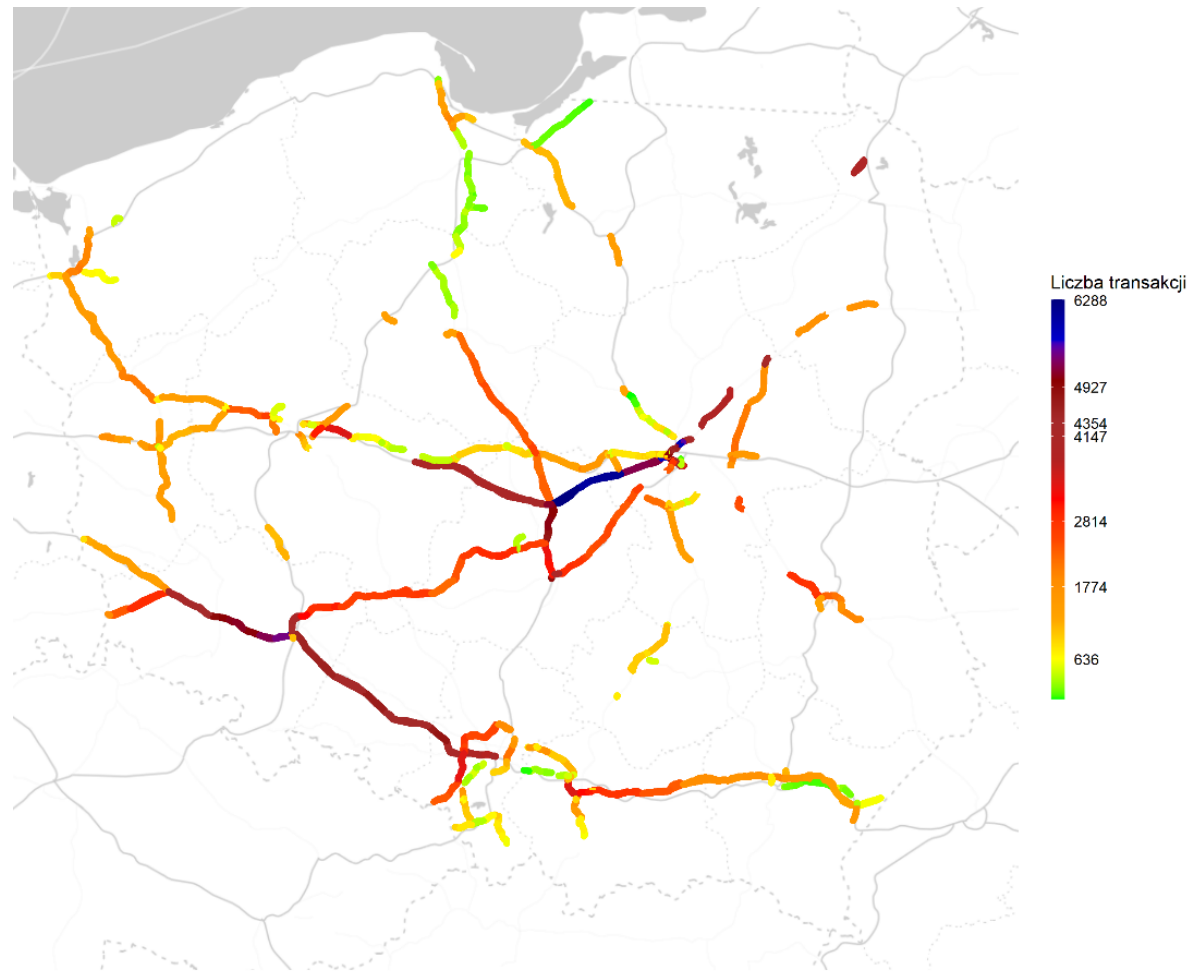
System **viaTOLL** funkcjonował do **30.09.2021 r.** (technologia radiowa).

Na system viaTOLL składały się **951** bramownice oraz urządzenia pokładowe umieszczone w pojazdach.

System **eTOLL** oparty jest o technologię wyznaczania pozycji użytkownika przy zastosowaniu pozycjonowania satelitarnego z wykorzystaniem wirtualnych bramownic.

System TranStat - przykładowe wizualizacje

Natężenie ruchu samochodów



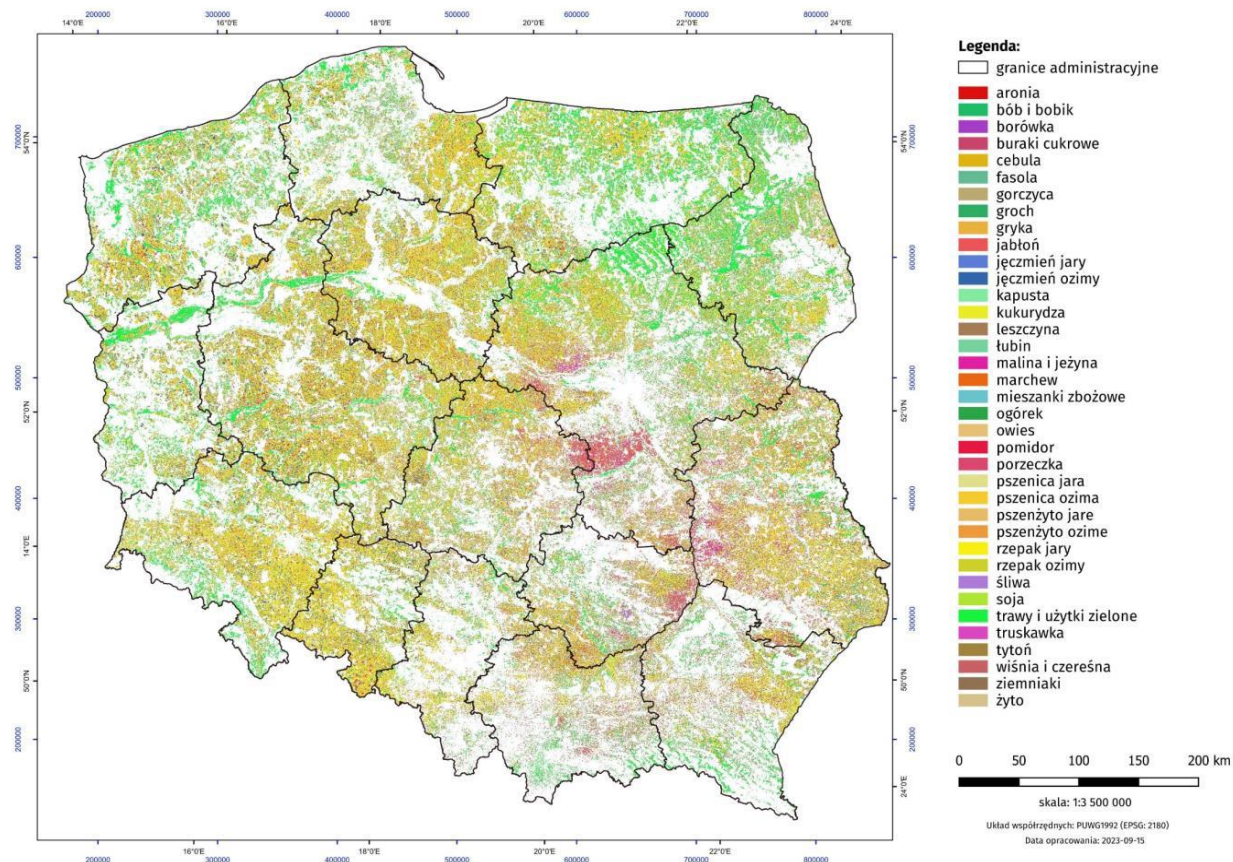
Rys. Natężenie ruchu samochodów (01.06.2024).

Źródło: Opracowane na podst. systemu TranStat

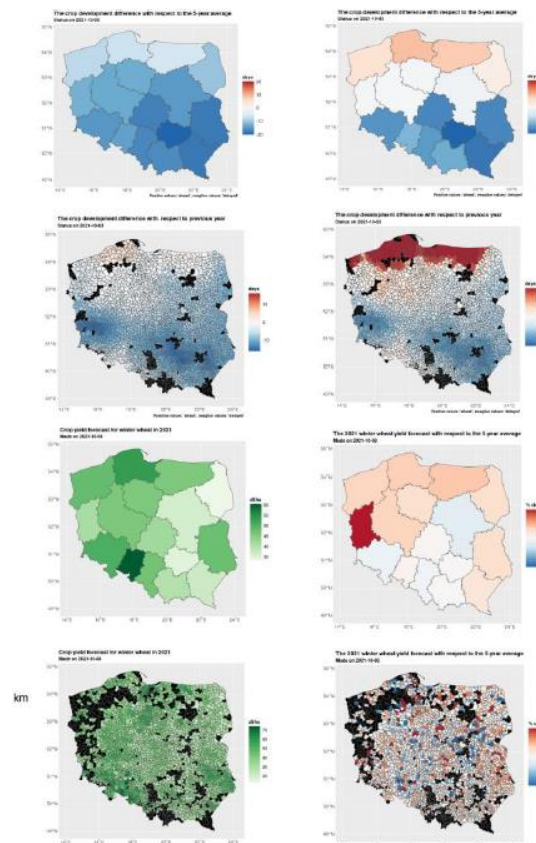
Dane satelitarne i sztuczna inteligencja

wykorzystanie na potrzeby wyznaczenia wskaźników
zrównoważonego rozwoju

Dane satelitarne - wykorzystanie na potrzeby statystyki rolnictwa



Mapa upraw

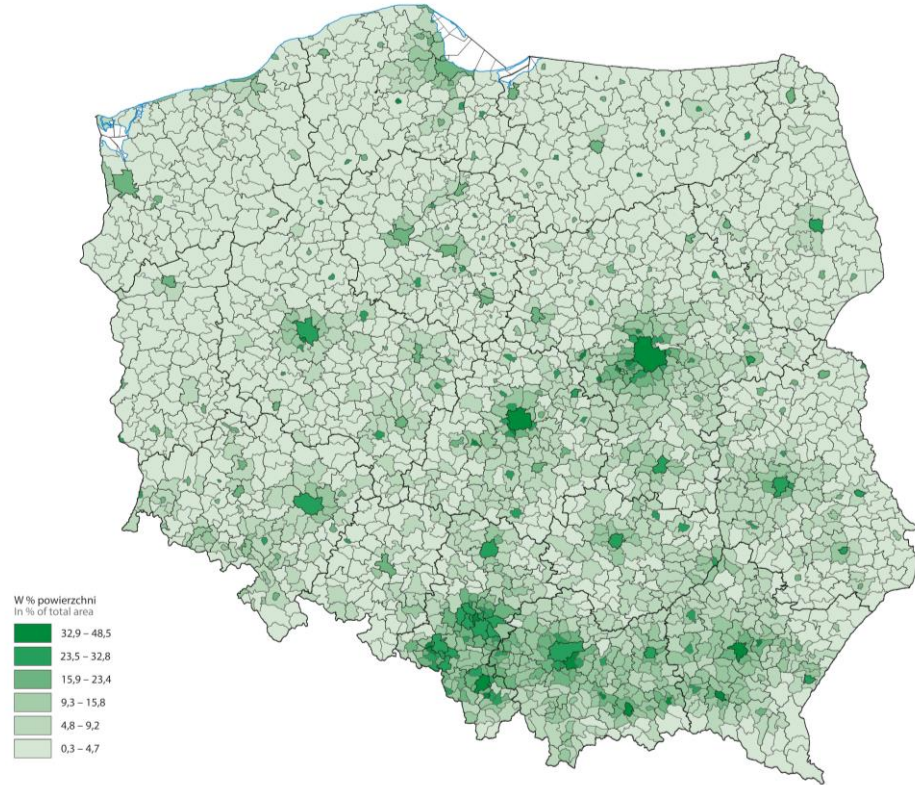


Mapa przesunięcia fazy rozwojowej rzepaku ozimego, pszenicy ozimej i kukurydzy (w dniach) dla województw i gmin w stosunku do średniej z ostatnich 5 lat

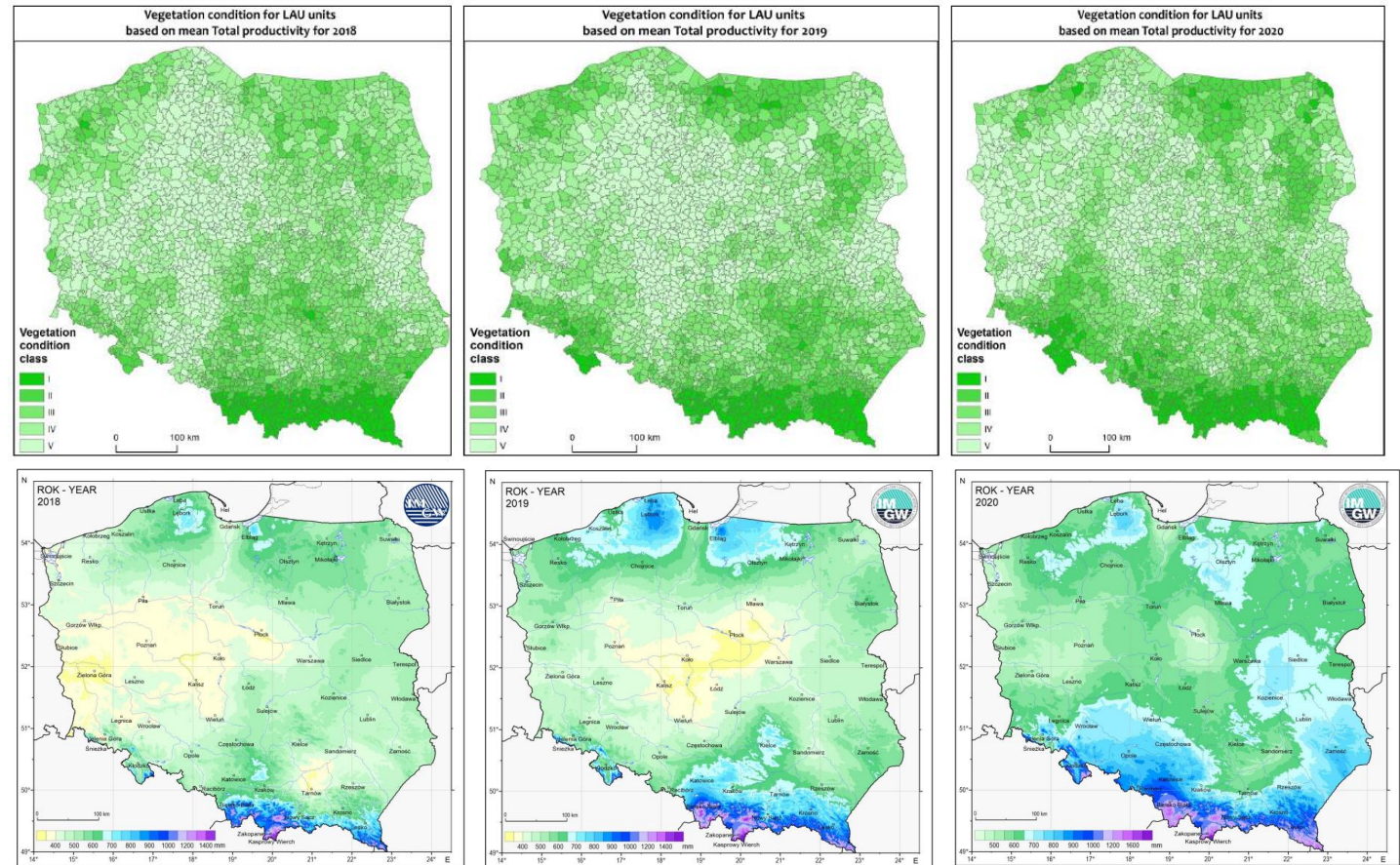
Mapa województw i gmin oraz ich prognoz plonów dla różnic w stosunku do średniej z ostatnich 5 lat.

Wykorzystanie danych satelitarnych w produkcji statystyk na potrzeby planowania przestrzennego

Tereny zieleni w Polsce w 2022 r.



Zróźnicowanie kondycji roślin w latach



Rys. Zróźnicowanie wskaźnika wegetacji i rocznych opadów w latach 2018-20, Źródło: Biuletyn Monitoringu Klimatu Polski <https://klimat.imgw.pl/pl/biuletyn-monitoring> stat.gov.pl

Statystyka rolnictwa przed 2020 r.

GŁÓWNY URZĄD STATYSTYCZNY, al. Niepodległości 208, 00-925 Warszawa www.stat.gov.pl

Nazwa i adres jednostki sprawozdawczej: **R-04** **Sprawozdanie o stanie upraw rolnych według oceny wiosennej w 2015 r.**

Panel sprawozdawczy GUS www.stat.gov.pl

Urząd Statystyczny 10-959 Olsztyn ul. Kościuszki 78/82

Numer identyfikacyjny - REGON: w 2015 r.

Termin przekazania do 15 maja 2016 r. według stanu na dzień 8 maja 2015 r.

Obowiązek przekazywania danych statystycznych wynika z art. 30 pkt 3 ustawy z dnia 29 czerwca 1995 r. o statystyce publicznej (Dz. U. z 2011 r. poz. 591, 4 późn. zm.) oraz rozporządzenia Rady Ministrów z dnia 27 sierpnia 2014 r. w sprawie programu badań statystycznych statystyki publicznej na rok 2015 (Dz. U. poz. 1330).

(e-mail: skrzynka.jednostki.sprawozdawczej@stat.gov.pl - WYPEŁNIĆ WIELKIMI LITERAMI!)

1. Czy gospodarstwo rolne użytkuje obecnie użytki rolne powyżej 1 ha? Tak 1 / Nie 2 → pkt. 2

2. Jaka jest powierzchnia użytków rolnych w gospodarstwie? ha

Dział 1. Powierzchnia zasiewów oraz szkody⁶ powstałe w uprawach ozimych w okresie jesienno-zimowym i wiosennym

Wyszczególnienie	Kod	Powierzchnia			
		ha		w tym ziorana i do ziorania	
		1	2	3	4
RAZEM (od 02 do 10)	01				
Pozostale ozimy	02				
Zyto	03				
Jęczmień ozimy	04				
Paszętno ozimy	05				
Mieszanki siewne ozimy	06				
Rzepak i rzepak ozimy	07				
Koniczyna czerwona ⁶	08				
Pozostałe uprawy rolnic ⁶	09				
Popioły ozime z przenieśnieniem na paszę	10				

⁶ Wymarczenia, wyprania, strąty powodowane, szkody spowodowane przemarzaniem wiosennym. ⁷ Koniczyna czerwona w uprawie z zioraniem. ⁸ Koniczyna czerwona (fazolan, trawy ozime, pastuska polowa itp. bez uprawy siewniczej, upraw trwałych oraz itp.).

Dział 2. Ocena stanu zasiewów w stopach kwalifikacyjnych

Wyszczególnienie	Kod	Powierzchnia ⁶		Przebieg ocena w stopach kwalifikacyjnych	Wyszczególnienie	Kod	Powierzchnia ⁶
		ha	1				
		1	2				1
Pozostale ozimy	01			Mieszanki siewne ozime	09		
	02				10		
Zyto	03			Mieszanki siewne-średniozróżni	11		
	04				12		
Jęczmień ozimy	05			Rzepak i rzepak jary	13		
	06				14		
Paszętno ozime	07			Koniczyna czerwona ⁶	15		
	08				16		

⁶ Powierzchnia aktualnie obsiana (w odniesieniu do zasiewów ozimych po odjęciu siewnicy i do ziorania). ⁷ Koniczyna czerwona z zioraniem.

Liczba badań

12

Liczba wywiadów

304 000

Liczba ankierów

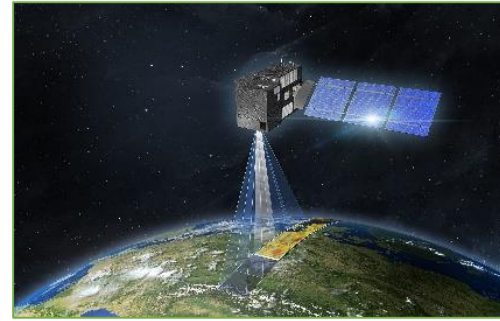
2 500

Poziom agregacji

województwo



Statystyka rolnictwa po 2020 r.



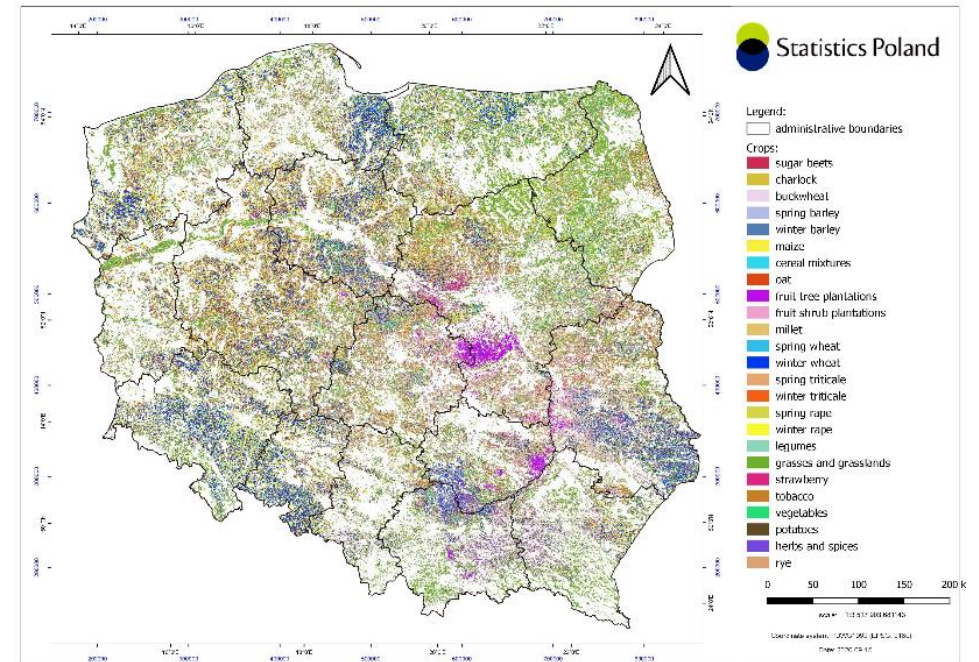
Source: <https://www.eoportal.org/satellite-missions/co2m#eop-quick-facts-section>

Liczba badań **6**

Liczba wywiadów **150 100**

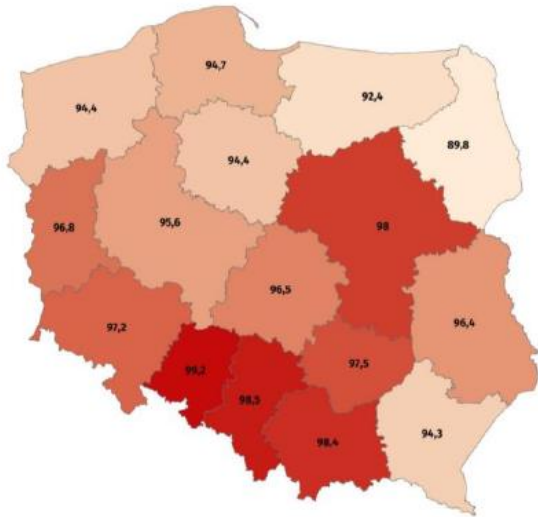
Liczba ankieterów **400**

Poziom agregacji **Województwo, powiat,
gmina, działka**

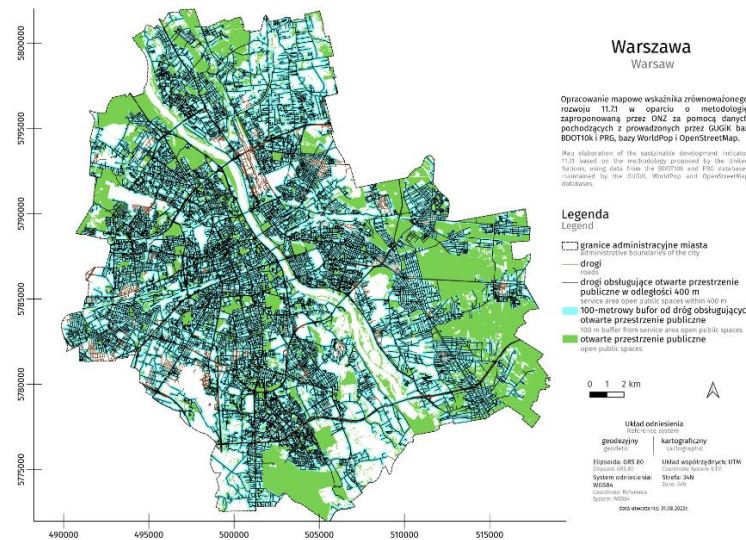


Dane satelitarne - wykorzystanie na potrzeby wyznaczenia wskaźników zrównoważonego rozwoju (SDG)

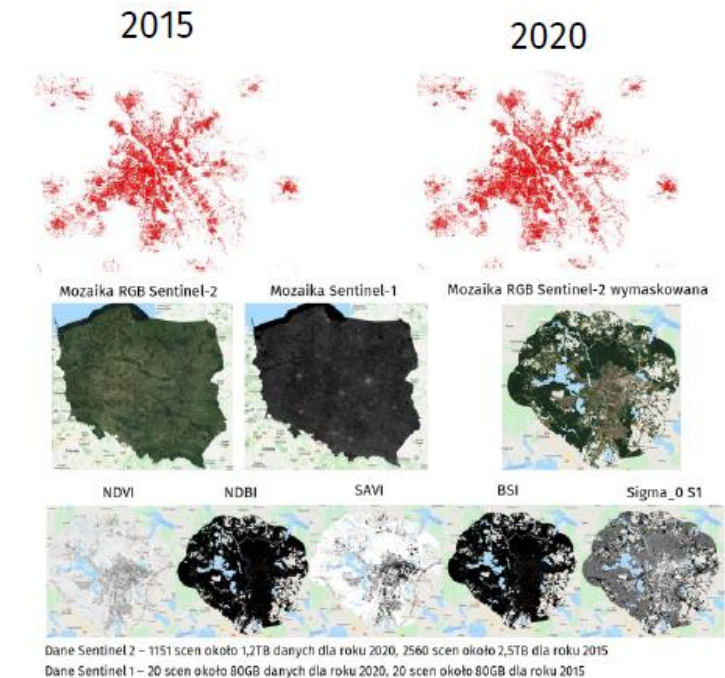
Wskaźnik 9.1.1 – Odsetek ludności wiejskiej zamieszkałej w promieniu 2 km od dostępnej przez cały rok drogi



Wskaźnik 11.7.1 Przeciętny udział terenów stanowiących przestrzeń publiczną dostępną dla wszystkich w powierzchni miasta ogółem



Wskaźnik 11.3.1 - Relacja wskaźnika zużycia gruntów do tempa wzrostu liczby ludności



https://sdg.gov.pl/statistics_exp/

Punkt kompleksowej obsługi na rzecz sztucznej inteligencji - uczenie maszynowe w statystyce publicznej (AIML4OS)

AIML4OS obejmie:

- opracowanie spójnego zestawu produktów – metodologii, wytycznych, środowiska testowego, danych uczących, ram wdrażania i jakości w statystykach oficjalnych w całym ESS przy wdrażaniu rozwiązań opartych na SI;
- utworzenie platformy zapewniającej pojedynczy punkt dostępu dla personelu ESS w celu uzyskania dostępu do ww. produktów;
- zapewnienie wsparcia w zakresie integracji i utrzymania odpowiednich rozwiązań opartych na AI/ML w organizacjach ESS poprzez szkolenia;
- budowanie społeczności wokół rozwiązań open-source rozwijanych i utrzymywanych przez członków ESS.



eStatystyka - otwarte dane statystyczne



API SDP

Składnica Danych Publikacyjnych zapewnia dostęp do pełnego zakresu informacyjnego statystyki publicznej: danych oraz metadanych w spójnym standardzie struktur danych. Dane w API SDP zorganizowane są w dwóch układach, według zmiennych publikacyjnych i według wskaźników.



API DBW

Dziedziczne Bazy Wiedzy to platforma danych statystycznych, która umożliwia dostęp do bogatego zestawu na bieżąco aktualizowanych informacji z różnych obszarów życia społeczno-gospodarczego i środowiska, prezentowanych w długich szeregach czasowych i maksymalnie pełnych przekrojach klasyfikacyjnych i nomenklaturowych.



API REGON

Rejestr REGON jest bieżąco aktualizowanym zbiorem informacji o podmiotach gospodarki narodowej prowadzonym w systemie informatycznym w postaci centralnej bazy danych.



API TERYT

Rejestr TERYT zawiera bieżące oraz archiwalne informacje o podziale terytorialnym kraju w zakresie: jednostek terytorialnych (województwa, powiaty, gminy), miejscowości i ulic. TERYT zawiera dane identyfikacyjne o ponad 100 tys. miejscowości i ponad 250 tys. ulic.



API BDL

Bank Danych Lokalnych jest największą w Polsce bazą danych o gospodarce, społeczeństwie i środowisku i oferuje tysiące cech statystycznych pogrupowanych tematycznie.



API SDG

Platforma SDG (Sustainable Development Goals) służy do monitorowania celów zrównoważonego rozwoju wyznaczonych przez Agendę 2030. Zawiera dane dla Polski dostępne spośród wskaźników międzynarodowych ONZ oraz umożliwiające pomiar krajowych priorytetów zrównoważonego rozwoju.



API STRATEG

System STRATEG wspiera proces monitorowania rozwoju oraz ocenę efektów działań podejmowanych w ramach polityki spójności. Prezentuje wskaźniki służące do pomiaru celów zapisanych w dokumentach strategicznych i programowych Polski oraz UE.

Open Data Inventory

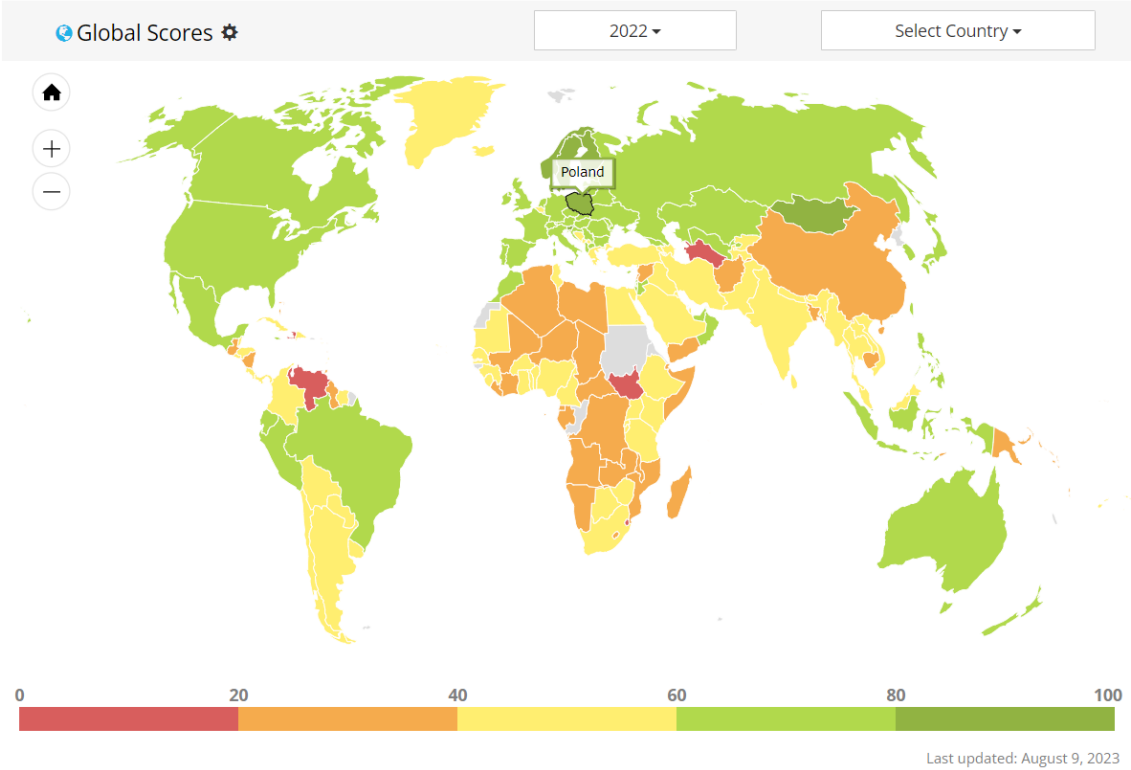
87
ODIN SCORE

Poland

2nd GLOBAL RANK
OUT OF 195

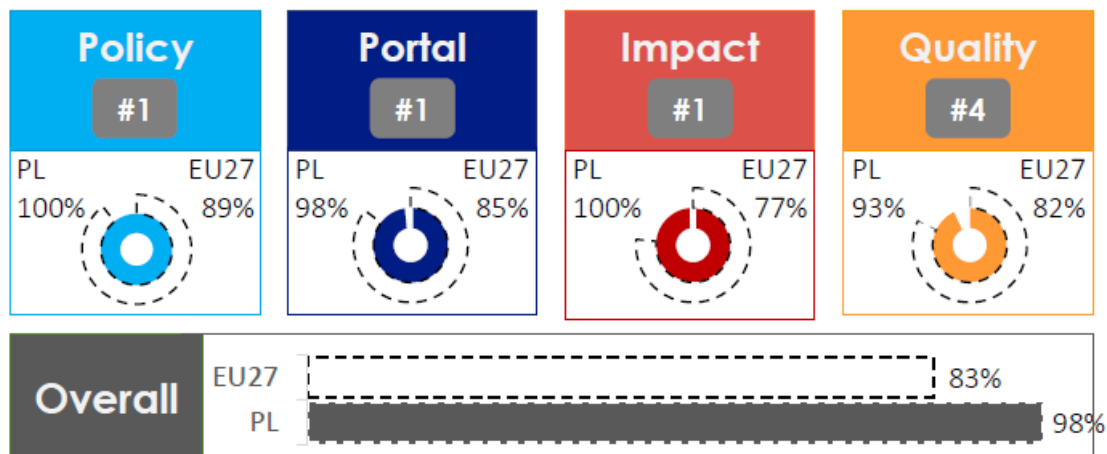
81 COVERAGE SCORE
OUT OF 100

92 OPENNESS SCORE
OUT OF 100

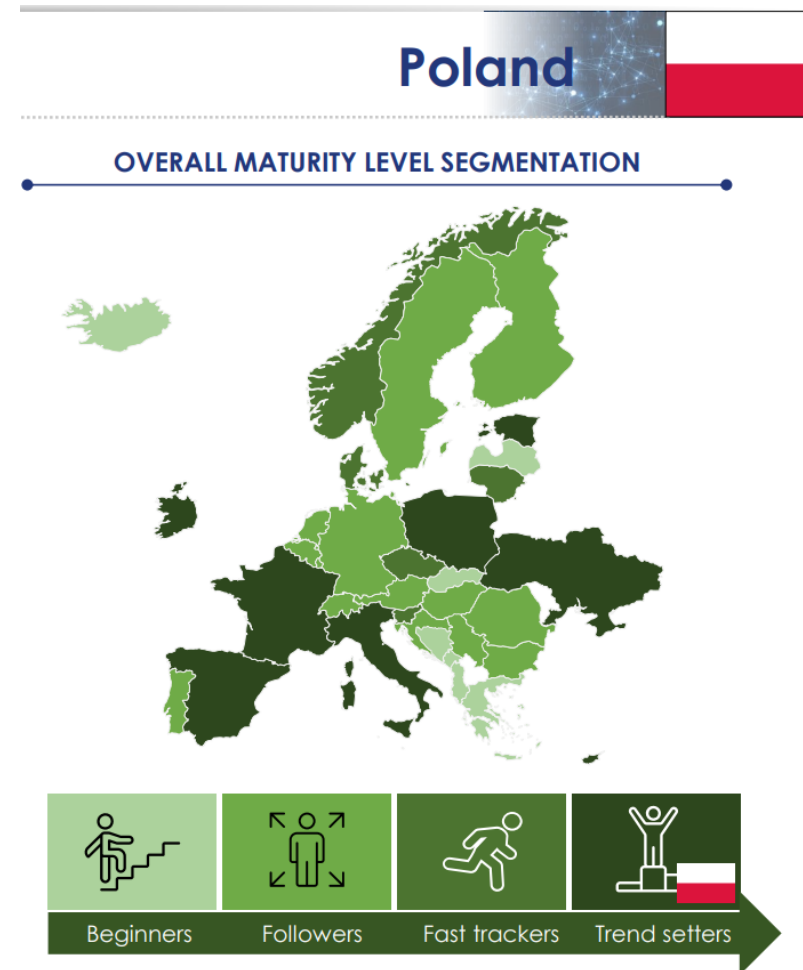


Open Data Maturity

Open data maturity **data.europa.eu**
2023



**2 miejsce, z wynikiem 97,8%,
wśród krajów europejskich,
za Francją**



E-Statystyka

Nowoczesna E-Statystyka to przede wszystkim profesjonalizm statystyków przekładający się na satysfakcję użytkowników wynikającą z wytwarzania przez statystykę oczekiwanej jakości produktów tworzących referencyjne podstawy dla budowy polityk społecznych i gospodarczych poprzez szerokie stosowanie:

- Innowacji;
- Automatyzacji;
- Standaryzacji;
- oraz pełnej elektronizacji procesów wytwórczych.

Dziękuję za uwagę