

Wstęp do eksploracji danych

Matematyka i Analiza Danych, II rok



Anna Kozak

Data Scientist | Data Visualization | Responsible Machine Learning

Data Scientist (5 years)

Research and teaching assistant - MI2 Lab (3 years)

Associate Teacher - WUT (1 year)

Warsaw RUG Meetup - Spotkania Entuzjastów R

R-Ladies Warsaw

Koło Naukowe Data Science PW

 MS Teams
anna.kozak@mini.pw.edu.pl



PhD student | Computer Science
AutoML | Meta-learning | RML
Research & Business Experience

Data Scientist (5 years)
Research assistant - MI2 Lab (3 years)
Coordinator of Case Studies 2020/2021



MS Teams
katarzyna.woznica.dokt@pw.edu.pl

Wstęp do eksploracji danych składa się z:

- wykładu
- zajęć laboratoryjnych

Wykłady - wtorki 12:15

Laboratoria - wtorki 14:15, 16:15

Konsultacje - prośba o kontakt

Wykład

Na wykładzie będą przedstawione teoretyczne aspekty pracy z danymi, jak i praktyczne.

Materialy

1. <https://github.com/MI2-Education/2022L-ExploratoryDataAnalysis>
2. Biecek Przemysław, Odkrywać! Ujawniać! Objasniać! Zbiór esejów o sztuce prezentowania danych, <http://www.biecek.pl/Eseje/index.html>
3. Rosling Hans, Factfulness
4. Tufte E.R., The Visual Display of Quantitative Information
5. Biecek Przemysław, Przewodnik po pakiecie R

Ocena

- prace domowe (40 pkt)
- projekty (44 pkt)
- wejściówki (6 pkt)

Suma 90 pkt

Ocena	3	3.5	4	4.5	5
Punkty	(50, 60]	(60, 70]	(70, 80]	(80, 90]	(90, ∞)

Z każdego projektu należy uzyskać ponad 50% możliwych punktów.

Laboratorium

- praca w R i Python (głównie R)
- powtórzenie operacji na danych (R: dplyr, tidyr, forcats; Python: numpy, pandas)
- wstęp do narzędzi pozwalających na estetyczne prezentowanie danych (R: ggplot2 + rozszerzenia, Python: matplotlib, seaborn)
- przygotowywanie interaktywnych wizualizacji i raportów (R: Rmarkdown, plotly, shiny)
- różne sposoby oceny zmiennych, danych, wizualizacji

Cel zajęć projektowych

- wykorzystanie i utrwalenie zdobytej wiedzy z wykładu oraz laboratoriów
- praktyczna praca z danymi
- ćwiczenie sposobu prezentacji wyników

Zasady

- 2 projekty w ciągu semestru
- zespoły 3 osobowe, różne podczas 1 i 2 projektu
- projekt trwa 5-6 tygodni
- 24 pkt i 20 pkt (w tym 5 pkt za pracę na zajęciach projektowych)

Projekt 1

Zadanie: Przygotowanie plakatu na zadany temat.

Rezultat: Plakat w formacie A2 w wersji pdf.

Zajęcia:

- wspólne dyskusje
- prezentacje kolejnych etapów

**CENTRUM
NAUKI
KOPERNIK**

Za tydzień

- podział na zespoły 3 osobowe w obrębie grup laboratoryjnych

Ocena

Za projekt można otrzymać od 0 do 24 punktów, z czego:

- 5p (1 x 1p, 2 x 2p) uzyskuje się za przedstawienie postępu prac w danym tygodniu
- 5p uzyskuje się za przygotowanie estetycznych wykresów (dwa lub więcej)
- 5p uzyskuje się, jeżeli przygotowane wykresy mają wszystkie niezbędne elementy do poprawnego odczytania danych (tytuł, podtytuł, adnotacje na osiach, legenda, jednostki, opis jak czytać wykres)
- 5p uzyskuje się za estetykę i pomysłowość aranżacji wykresów i opisów w jedną całość
- 4p uzyskuje się za prezentację projektu

Projekt 2

Zadanie: Przygotowanie interaktywnego raportu .

Temat: Spojrzenie na klimat i środowisko.

Rezultat: Raport lub aplikacja Shiny.

Zajęcia:

- wspólne dyskusje
- prezentacje kolejnych etapów

Ocena

Za projekt można otrzymać od 0 do 20 punktów, z czego:

- 5p (1 x 1p, 2 x 2p) uzyskuje się za przedstawienie postępu prac w danym tygodniu
- 7p uzyskuje się za postawienie pytań badawczych oraz znalezienie adekwatnych danych
- 6p uzyskuje się za jakość wizualizacji, tabel, opisów, element interaktywności (wykresy mają wszystkie niezbędne elementy do poprawnego odczytania danych (tytuł, podtytuł, adnotacje na osiach, legenda, jednostki, opis jak czytać wykres))
- 2p uzyskuje się prezentację projektu

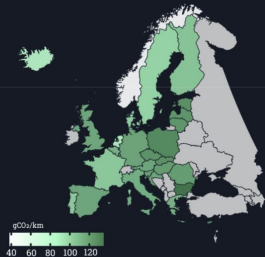
Plakaty, które zmieniają spojrzenie na klimat i środowisko



BRIGHT RIDE

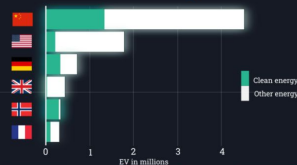
THE LIGHT AT THE END OF THE TUNNEL FOR OUR ENVIRONMENT

Average gCO₂/km emissions by cars registered in 2020



Why does Norway occupy the first position in terms of emitting CO₂ by newly registered cars in Europe with 38.2 gCO₂/km emissions? Not only are these electric cars relatively cheap due to lowering taxes in EV but drivers also benefit from reduced road and ferry tolls as well as discounted parking. In contrast to Norwegian pro-ecological attitude *Poland* is among the countries of the EU with the smallest number of charging stations for electric vehicles per 100 km of roads that leads to having the highest average emission of carbon dioxide per km which indicates 134.4 gCO₂/km.

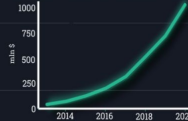
Number of EV that drive on clean energy



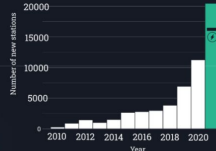
Some say that the answer to the emission problems are electric vehicles. The concept of clean energy makes us believe that we actually save our environment by choosing it over other energy sources. The energy used to charging electric cars mostly comes from conventional energy sources. In fact, only Norway produces the proper amount of EV to serve electric cars in green way.

Due to the new legislation that focuses on limiting CO₂ emissions and technology breakthrough, electric cars are becoming as popularly. More and more motor companies invest in this branch meanwhile limiting their gasoline and diesel car production, but is the market ready for this sudden changes?

Stock value of electric cars



Number of new stations

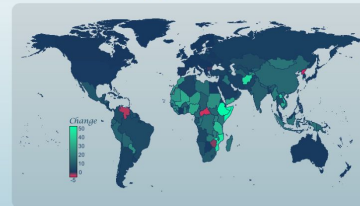


Critics of electric vehicles say that there is a problem with insufficient number of charging stations. But is it still a problem? The chart above presents number of new charging stations in the United States each year. As one may see, that number grows almost exponentially. Due to the fact that data was collected only to September 2020, the number of new charging stations in this year may significantly increase.



Sea the Change

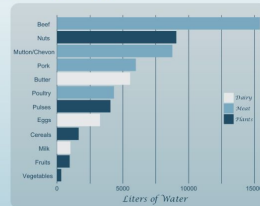
Change In Access To Improved Water since 2000 (in percentage points)



<https://ncesdata.org/explorers/water-and-sanitation>

With just a few exceptions, the access to improved drinking water is on the rise. The improvement is especially visible in Africa and south-east Asia. Despite these changes, the access in those regions is still not on par with the access seen in developed countries where it is almost universal.

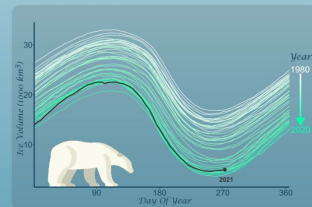
Water Used In Food Production per kilogram of product



https://waterfootprint.org/media/downloads/footprint-46-waterfootprint-access-products-Vol1_3.pdf

Agriculture is the largest water user worldwide, accounting for 70 percent of total freshwater withdrawals, but these amounts can reach as much as 95 percent in some developing countries. It is 3496 litres of water 'eaten' per person, per day.

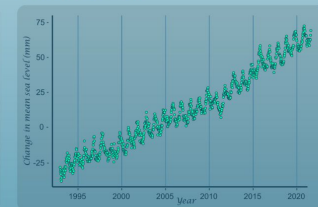
Arctic Ice Volume between 1980 and 2020



<https://psc.apl.usc.edu/research/projects/arctic-sea-ice-volume-anomaly/data/>

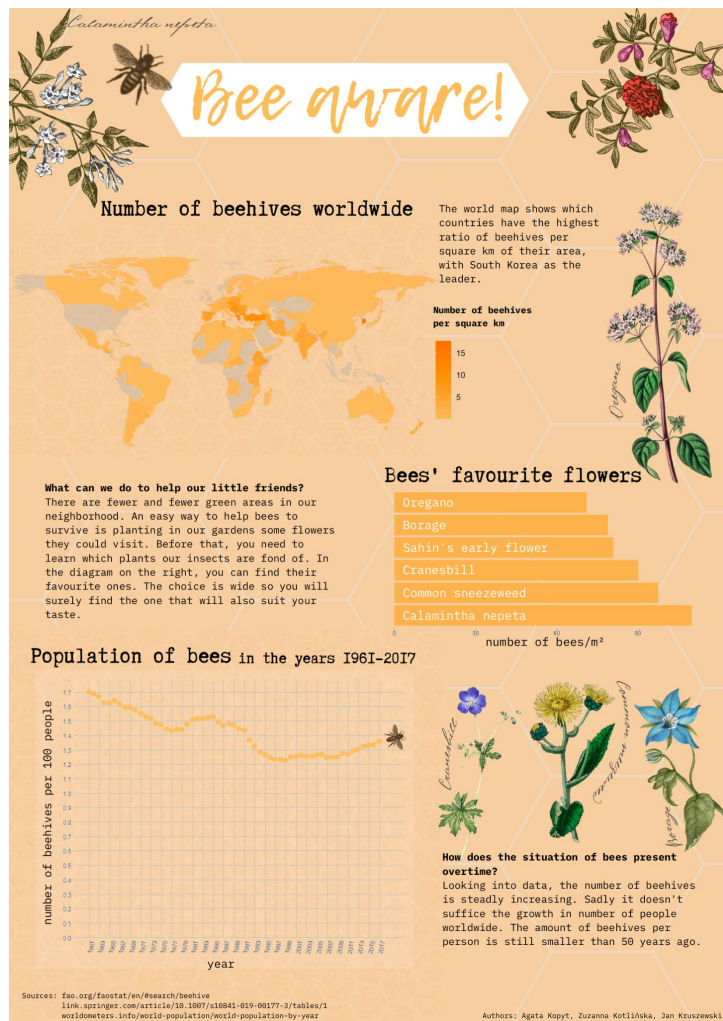
Polar ice caps are melting as global warming causes climate change. We lose Arctic sea ice at a rate of almost 13% per decade. If emissions continue to rise unchecked, the Arctic could be ice-free in the summer by 2040. But what happens in the Arctic does not stay in the Arctic. Sea ice loss has far-reaching effects around the world.

Global Mean Sea Level between 1992 and 2021

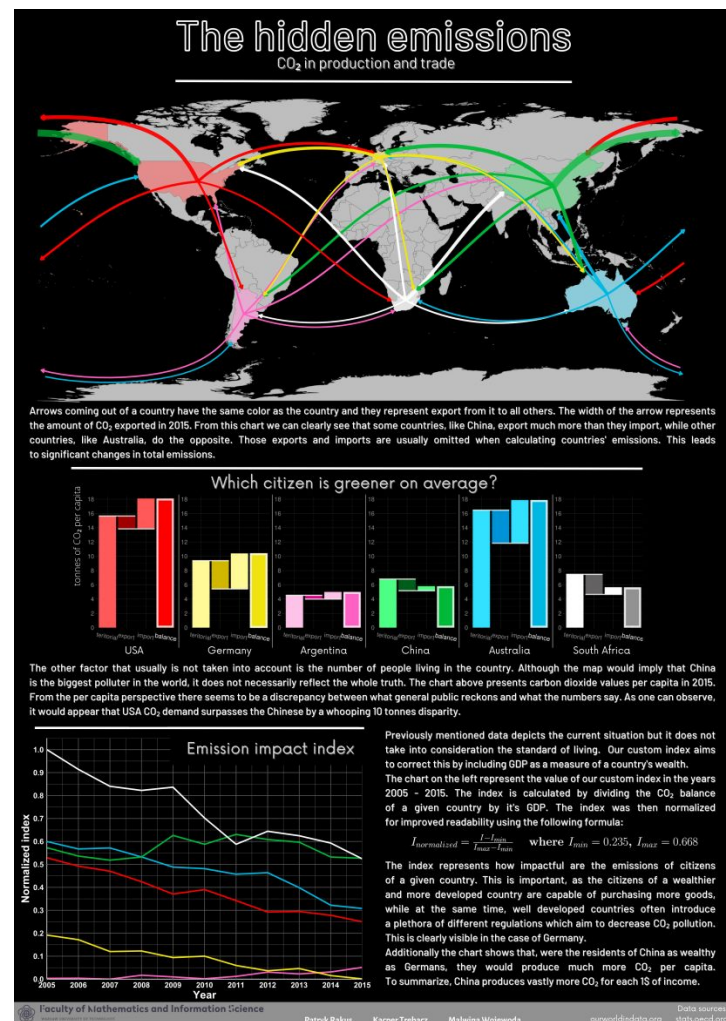


<https://sealevel.colorado.edu/index.php/data/2021/m2-0>

The rising water level is mostly due to a combination of melt water from glaciers and ice sheets as well as thermal expansion of seawater as it warms. In urban settings along coastlines around the world, rising seas threaten infrastructure necessary for local jobs and regional industries.



Source: doi.org/10.1007/978-1-4939-9841-8-1007
www.worldometers.info/world-population/world-population-by-year



ECOLOGICAL ENERGY AND LIFE OF EUROPEANS

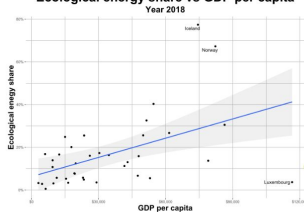
Due to the global warming, being ecological is becoming more and more important. Countries are constantly aiming to reduce carbon emissions by utilizing renewable sources of energy. How does it affect their residents?

To answer that, we'll be taking a closer look at European countries.

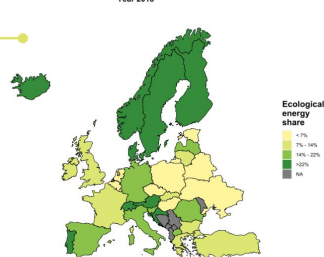
HOW ECOLOGICAL IS EUROPE?

There are some countries that care a lot about being ecological, while others have a very different approach. In Europe, we've got both ends of this spectrum, starting with Iceland, which gets around 77% of its energy from renewable sources, and ending with Belarus which uses such sources to obtain less than a percent of its energy.

Ecological energy share vs GDP per capita



Ecological energy share in European countries

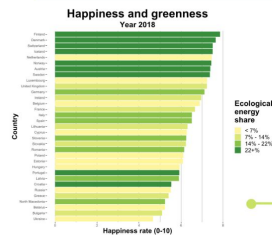


RICHER = MORE ECOLOGICAL?

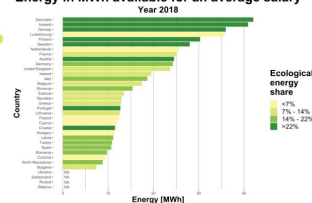
We can see the main trend in data. The richer the country is, the more ecological energy it uses. Iceland and Norway are perfect examples of this trend. They both have relatively high GDP as well as ecological energy share. Luxembourg, however, is an exception. It has only about 4 percent of energy that comes from renewables sources, despite having the highest GDP per capita in Europe.

HOW MUCH ENERGY CAN EUROPEANS AFFORD?

In 2018 the price of energy was relatively the cheapest in Denmark which uses a fair amount of renewable energy sources. Does the method of producing energy impact its prices? We can see that the top of the plot corresponding to the best prices is filled mostly with countries that gain energy from more eco-friendly sources. There are some exceptions but overall, it seems like the eco-friendlier your country is, the more energy you can afford.



Energy in MWh available for an average salary



WHAT ABOUT HAPPINESS?

In 2018 seven out of nine top ecological European countries were also leading in happiness for the region. The farther down the happiness rate, the less of a distinction can be made considering the greenness percentage. At all in all, it is probably the rich that can afford to be both happy and ecological.

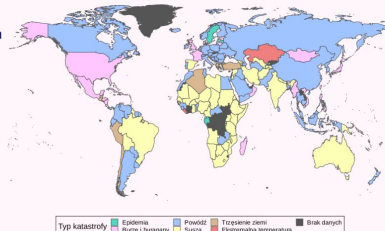
Katastrofy naturalne na mapach świata

Natura objawia swoją potęgę w wielu różnych zjawiskach. Niektóre z nich pokazują jej piękno, inne budzą grozę. My zainteresowaliśmy się katastrofami naturalnymi, które mogą być niezwykle niszczycielskie i śmiertelne. Postanowiliśmy zbadać, które kraje są najbardziej dotknięte przez katastrofy naturalne i które kataklizmy występują w konkretnych częściach świata.

Najbardziej dotkliwe katastrofy naturalne w poszczególnych krajach

Na przestrzeni lat 1971-2008

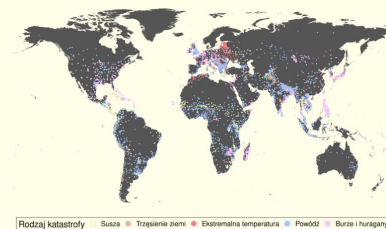
Która z katastrof naturalnych dotyka najczęściej obywateli danego państwa? Okazuje się, że wyniki są zgodne z naszymi przewidywaniami. W większości krajów europejskich najwięcej osób nękanych jest przez powodzie, natomiast w Afryce dominują susze. W USA są to burze i huragany, a w Chile trzęsienia ziemi.



Występowanie katastrof naturalnych

Na przestrzeni lat 2010-2018

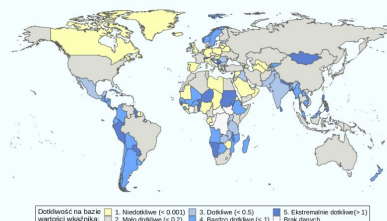
Gdzie występują konkretne katastrofy? Po naniesieniu pojedynczych wystąpień na mapę widać, że w dużej mierze pokrywa się to z oczekiwaniami które moglibyśmy mieć po przeanalizowaniu poprzedniej mapki. Niestety w niektórych częściach świata zbieranych jest znacznie mniej danych niż w innych. Dlatego na przykład w Europie jest znacznie więcej zaznaczonych wydarzeń niż w Brazylii czy Rosji.



Dotkliwość katastrof naturalnych

Na przestrzeni lat 2010-2018

W którym kraju katastrofy naturalne są najbardziej dotkliwe dla człowieka? Aby to sprawdzić, dla każdego z nich policzyliśmy nasz własny wskaźnik. Jest to liczba osób dotkniętych w wyniku kataklizmów przez populację danego kraju, uśredniona na przestrzeń lat. Następnie pogrupowaliśmy ten wskaźnik na 5 stopni dotkliwości i umieściliśmy na mapie. Okazuje się, że w najgroźniejszych krajach, takich jak Filipiny, Niger czy Peru, wskaźnik osiąga wartości powyżej 1.



Faculty of Mathematics and Information Science

AUTHORS:

Katrzewska Julia
Piórczyński Miłosz
Sokołowski Jędrzej

SOURCES:

<https://data.worldbank.org>
<https://happyplanetindex.org>
<https://euripa.europa.eu/euripa-stat>
<https://github.com/bowdlenenergy/data>

Współfinansowane przez Ministerstwo Edukacji i Nauki
Ministerstwo Infrastruktury i Budownictwa
Ministerstwo Gospodarki i Pracy

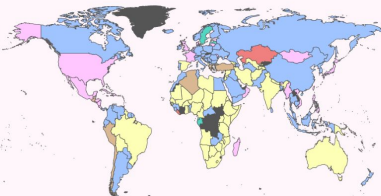
Jakub Piskow
Krzysztof Woźniak
Łukasz Tomaszewski

Katastrofy naturalne na mapach świata

Natura objawia swoją potęgę w wielu różnych zjawiskach. Niektóre z nich pokazują jej piękno, inne budzą grozę. My zainteresowaliśmy się katastrofami naturalnymi, które mogą być niezwykle niszczycielskie i śmiertelne. Postanowiliśmy zbadać, które kraje są najbardziej dotknięte przez katastrofy naturalne i które kataklizmy występują w konkretnych częściach świata.

Najbardziej dotknięte katastrofy naturalne w poszczególnych krajach Na przestrzeni lat 1971-2008

Która z katastrof naturalnych dotyka najczęściej obywateli danego państwa? Okazuje się, że wyniki są zgodne z naszymi przewidywaniami. W większości krajów europejskich najczęściej osób nękanych jest przez powódzie, natomiast w Afryce dominują susze. W USA są to burze i huragany, a w Chile trzęsienia ziemi.



Typ katastrofy: Epidemia, Susza, Powódź, Trzęsienie ziemi, Brak danych

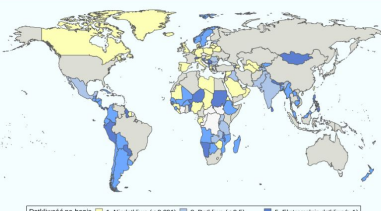
Występowanie katastrof naturalnych Na przestrzeni lat 2010-2018

Gdzie występują konkretne katastrofy? Po naniesieniu pojedynczych wystąpień na mapę widać, że w dużej mierze pokrywa się to z oczekiwaniami, które moglibyśmy mieć po przeanalizowaniu poprzedniej mapki. Niestety w niektórych częściach świata zbieranych jest znacznie mniej danych niż w innych. Dlatego na przykład w Europie jest znacznie więcej zaznaczonych wydarzeń niż w Brazylii czy Rosji.

Podział katastrof: Susza, Trzęsienie ziemi, Powódź, Burze i huragany

Dotkliwość katastrof naturalnych Na przestrzeni lat 2010-2018

W którym kraju katastrofy naturalne są najbardziej dotkliwe dla człowieka? Aby to sprawdzić, dla każdego z nich policzyliśmy nasz własny wskaźnik. Jest to liczba osób dotkniętych w wyniku kataklizmów przez populację danego kraju, uśredniona na przestrzeni lat. Następnie pogrupowaliśmy ten wskaźnik na 5 stopni dotkliwości i umieściliśmy na mapie. Okazuje się, że w najgroźniejszych krajach, takich jak Filipiny, Niger czy Peru, wskaźnik osiąga wartości powyżej 1.



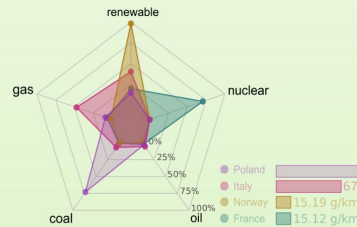
Dotkliwość na bazie wartości wskaźnika: 1. Niekłkliwe (< 0.001), 2. Mało dotkliwe (< 0.2), 3. Dotkliwe (< 0.5), 4. Bardzo dotkliwe (< 1), 5. Ekstremalnie dotkliwe (> 1), Brak danych

Źródło:
https://ourworldindata.org/natural-disasters
https://ourworldindata.org/cheatsheet/energy-consumption-electric-car
https://ec.europa.eu/eurostat/web/transport/data/database

Jakub Piwko
Krzysztof Wodnicki
Łukasz Tomaszewski

Are Electric cars that Eco?

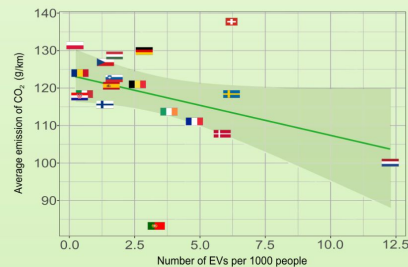
Study on in which countries electric vehicles really reduce CO₂ emissions



Average CO₂ emissions of Renault Zoe by country

To check how different methods of energy production impact emissions from electric cars we chose 4 countries with different primary way of producing energy and calculated how many grams of CO₂ are produced by driving 1km.

Poland: 156.55 g/km
Italy: 67.72 g/km
Norway: 15.19 g/km
France: 15.12 g/km



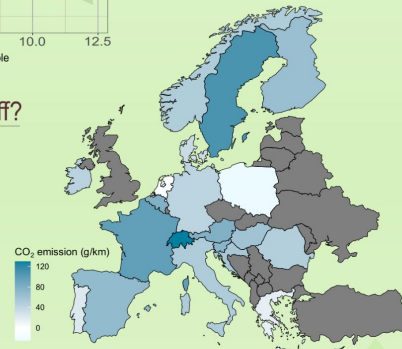
Average CO₂ emissions of passenger cars by country

To further examine beneficiality of electric vehicles we compared average CO₂ emission of all types of cars and number of electric vehicles in each country. It's easy to observe that generally the more electric cars there are in a country, the less CO₂ is produced by all vehicles on average.

Does an electric car pay off?

The difference in CO₂ production of electric and internal combustion cars

In previous visualizations we showed that in some countries it's more beneficial to own an electric vehicle than in others. We calculated the difference of CO₂ production of an electric and combustion cars in each country. The bigger the difference, the bigger is the positive influence on environment of buying an electric car.



CO₂ emission (g/km)
120
80
40
0

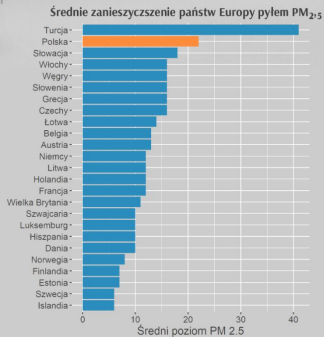


Authors:
Filip Szympliński
Michał Tomczyk
Piotr Wilczyński

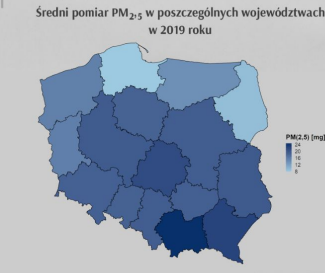
Data sources:
<https://ourworldindata.org>
<https://ev-database.org/cheatsheet/energy-consumption-electric-car>
<https://ec.europa.eu/eurostat/web/transport/data/database>

ZANIECZYSZCZENIE POWIETRZA W POLSCE

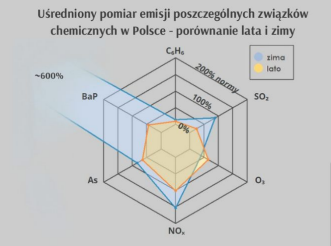
Polskie powietrze jest jednym z najbardziej zanieczyszczonych w EU, a mimo tego, nie są prowadzone adekwatne działania zapobiegające temu problemowi. Największy problem występuje z wielopierścieniowymi węglowodorami aromatycznymi (WVA) (m.in. z BaP), pyłem zawieszonym $PM_{2.5}$ oraz drobniejszą frakcją pyłu $PM_{2.5}$. Ich normy w licznych polskich miejscowościach są przekraczane kilku a nawet kilkunastokrotnie.



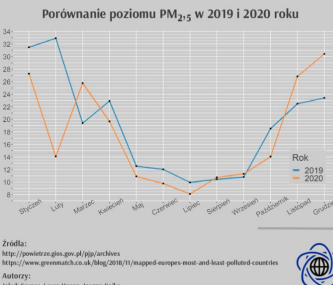
PYŁ ZAWIESZONY $PM_{2.5}$
Nasz niepokój może wzburzać poziom $PM_{2.5}$. Ze względu na mały rozmiar, dociera do dróg oddechowych, skąd przenika do krwiobiegu. Jego stężenie roczne nie powinno być wyższe niż 10 $\mu g/m^3$, a według mapy zanieczyszczenia $PM_{2.5}$ Polska (szczególnie południowa) boryka się z jego wysokimi średniorocznymi stężeniami.



WVA I INNE SKŁADNIKI SMOGU
Najlepiej przebadanym wielopierścieniowym węglowodorem aromatycznym jest benzo(a)piren (BaP). W 1987 roku Międzynarodowa Agencja Badań nad Rakiem (IARC) uznała BaP za główny ludzki kancerogen. W naszych badaniach porównaliśmy utrudnione pomiary składników smogu z ich normami prawnymi. Jak widać, tak toksyczna substancja jak BaP mocno jest przekraczana. Tak samo poziomy SO_2 i tlenek azotu (NO_x) mogą wzburzać nasz niepokój.



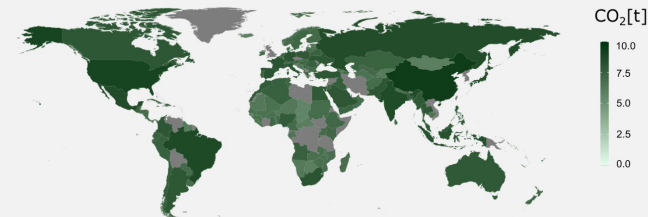
WPŁYW LOCKDOWNU NA STAN POWIETRZA
Czy siedząc w domu poprawiliśmy jakość naszego powietrza? Celem ostatniego badania było porównanie wpływu lockdownu na poziom wytwarzanego zanieczyszczenia. Jak się okazuje, wiele się nie poprawiło. Być może, w Polsce to nie komunikacja samochodami jest głównym źródłem tego problemu, a jest nim ogrzewanie gospodarstw domowych.
Stan powietrza w naszym kraju jest alarmujący - brakuje regulacji ku tej kwestii: ogrzewanie piecami, kotłami czy kominkami nie powinno być na porządku dziennym. Musimy zacząć zwracać uwagę na to, czym palimy i co palimy.



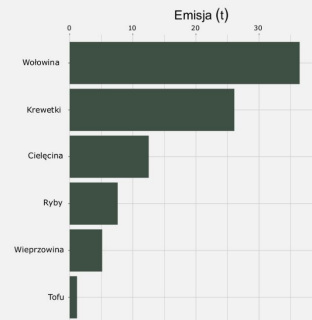
Produkcja żywności a ekologia

Dawid Pludowski
Antoni Zajko
Grzegorz Kiersnowski
Źródła: FAOSTAT, Kaggle, OWID

Emisja CO_2 per capita w roku 2013 w wyniku produkcji żywności

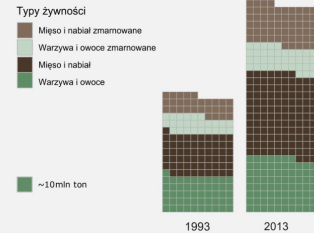


Emisja CO_2 w wyniku produkcji żywności w przeliczeniu na 1000 kcal

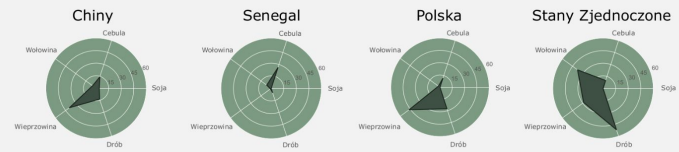


Badania sugerują, że 20% dwutlenku węgla, który wytwarzamy, pochodzi z produkcji żywności. Marnowanie jedzenia również przyczynia się do transmisji nadmierowego dwutlenku węgla do atmosfery. Najbardziej odpowiedzialne za ten stan rzeczy są kraje rozwinięte.

Produkcja i marnowanie żywności w latach 1993 i 2013



Spożycie wybranych produktów w 2013 roku w kilogramach per capita



Pytania?

Eksploracja danych

Dane

Mogą być generowane przez:

- ?

Dane

Mogą być generowane przez:

- banki,
- ubezpieczenia,
- portale społecznościowe,
- firmy telekomunikacyjne,
- szpitale,
- dane eksperymentalne,
- tekst,
- mapy,
- sklepy internetowe,
- ...

Eksploracja danych - czym jest?

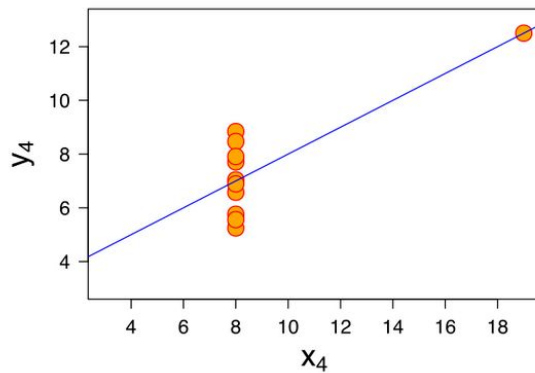
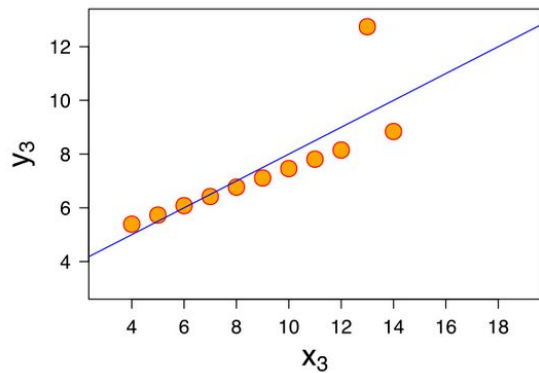
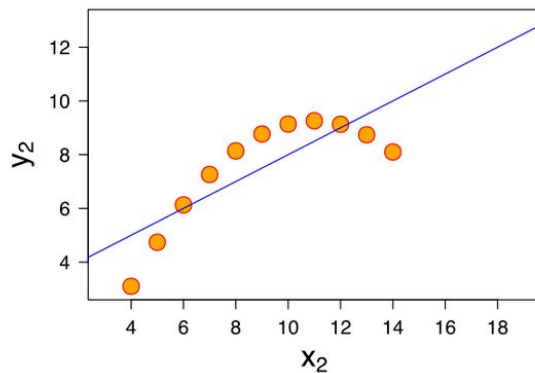
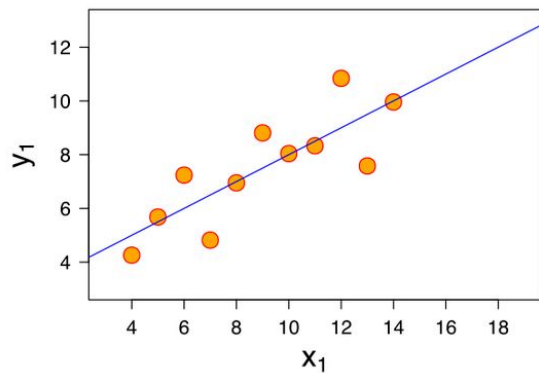
“proces odkrywania nietrywialnych, dotychczas nieznanych, potencjalnie użytecznych reguł, zależności, trendów”

Cel: analiza danych w celu lepszego ich zrozumienia

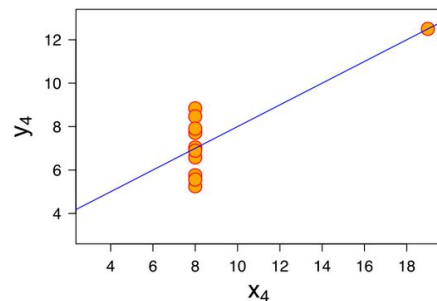
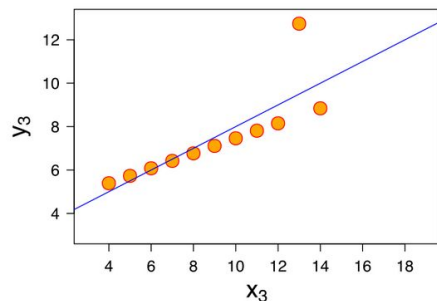
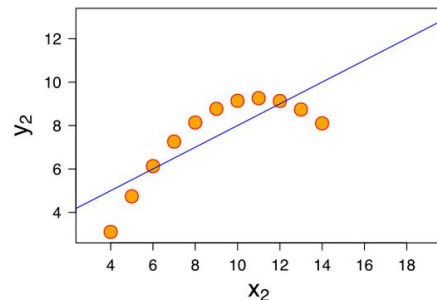
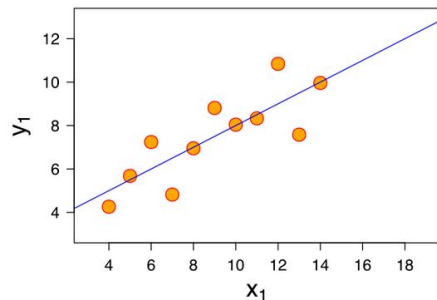
Eksploracja danych - czym jest?

Na eksplorację danych składa się wiele dyscyplin, między innymi:

- bazy danych
- statystyka
- uczenie maszynowe
- wizualizacja danych
- wyszukiwanie informacji



Kwartet
Anscombe'a



Cecha	Wartość
Średnia arytmetyczna zmiennej x	9
Wariancja zmiennej x	11
Średnia arytmetyczna zmiennej y	7.50 (identyczna do dwóch cyfr po przecinku)
Wariancja zmiennej y	4.122 lub 4.127 (identyczna do trzech cyfr po przecinku)
Współczynnik korelacji pomiędzy zmiennymi	0.816 (identyczny do trzech cyfr po przecinku)

Jak rozpoznać rodzaj zmiennej?

"dane liczbowe to nie tylko liczby"

Typy danych

Zmienne jakościowe (nazywane również wyliczeniowymi, czynnikowymi lub kategorycznymi), to zmienne przyjmujące określoną liczbę wartości (najczęściej nie liczbowych). Zmienne te można dalej podzielić na:

- *binarne* (nazywane również dwumianowymi, dychotomicznymi) np. płeć (poziomy: kobieta/mężczyzna),
- *nominalne* (nazywane również zmiennymi jakościowymi nieuporządkowanymi) np. marka samochodu,
- *uporządkowane*, np. wykształcenie (poziomy: podstawowe/średnie/wyższe), ocena z przedmiotu.

Typy danych

Zmienne ilościowe, z których można dodatkowo wyróżnić:

- *zliczenia* (liczba wystąpień pewnego zjawiska, opisywana liczbą całkowitą), np. liczba lat nauki, liczba wypadków,
- *ilorazowe*, czyli zmienne mierzone w skali, w której można dzielić wartości (ilorazy mają sens). Np. długość w metrach (coś jest 2 razy dłuższe, 10 razy krótsze itp.),
- *przedziałowe* (nazywane też interwałowymi), mierzone w skali, w której można odejmować wartości (wyznaczać długość przedziału).

Struktura zbioru danych

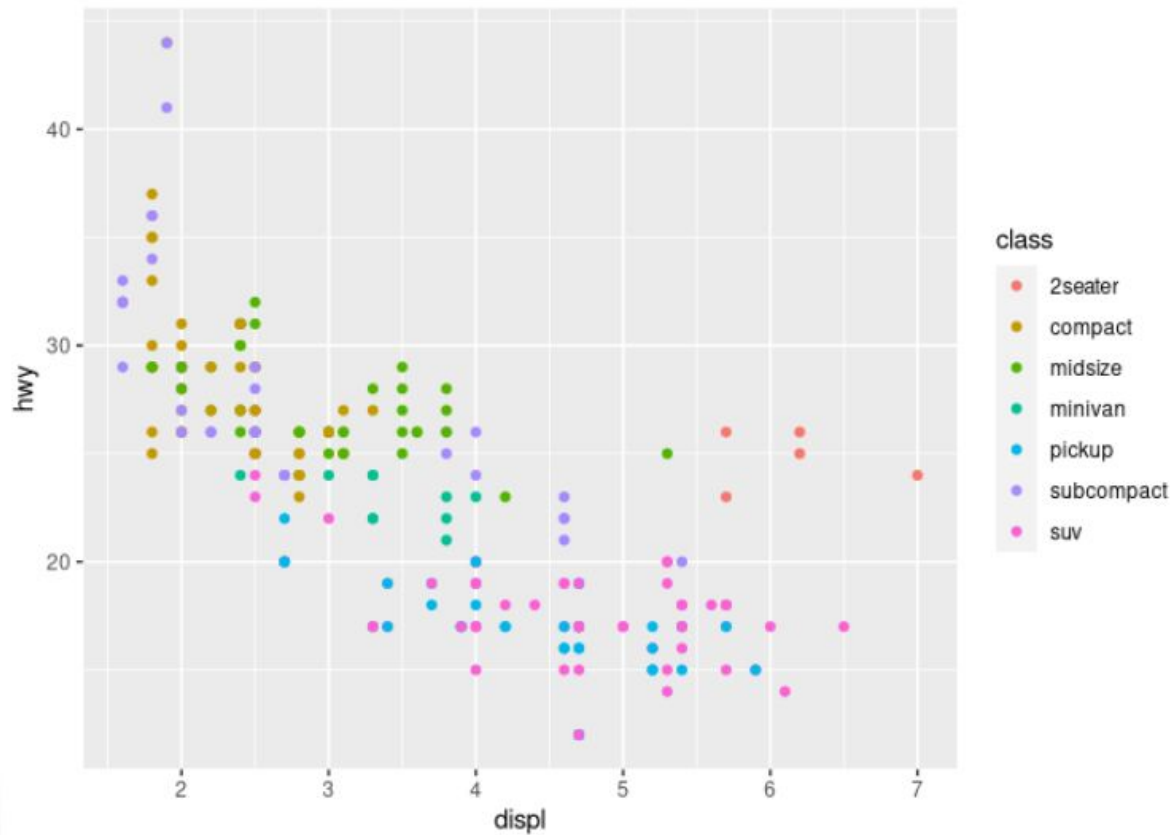
ID	PŁEĆ	ZAWÓD	WZROST	DATA URODZENIA
ID_23	K	INFORMATYK	158	1978-03-12
ID_45	K	PRAWNIK	178	1989-05-29
ID_46	M	MATEMATYK	183	1991-01-19
ID_89	M	INFORMATYK	167	1982-02-20
ID_101	K	LEKARZ	163	1973-02-23

Narzędzia do wizualizacji danych

- programistyczne (R, Python, JavaScript)
- programy graficzne (Inkscape)
- programy dedykowane do wizualizacji danych (Tableau)

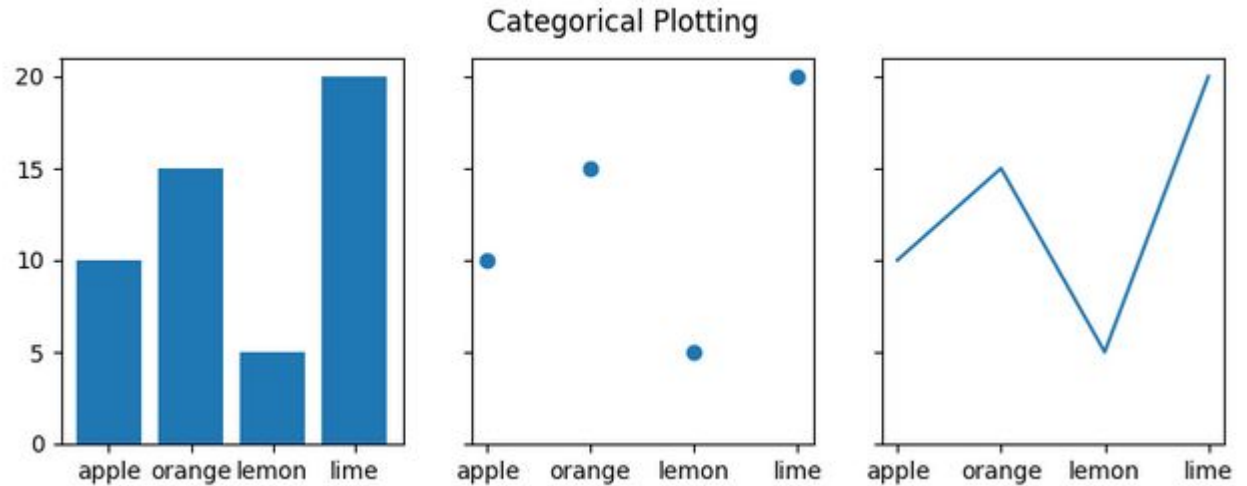
ggplot2 (R)

<https://ggplot2.tidyverse.org>



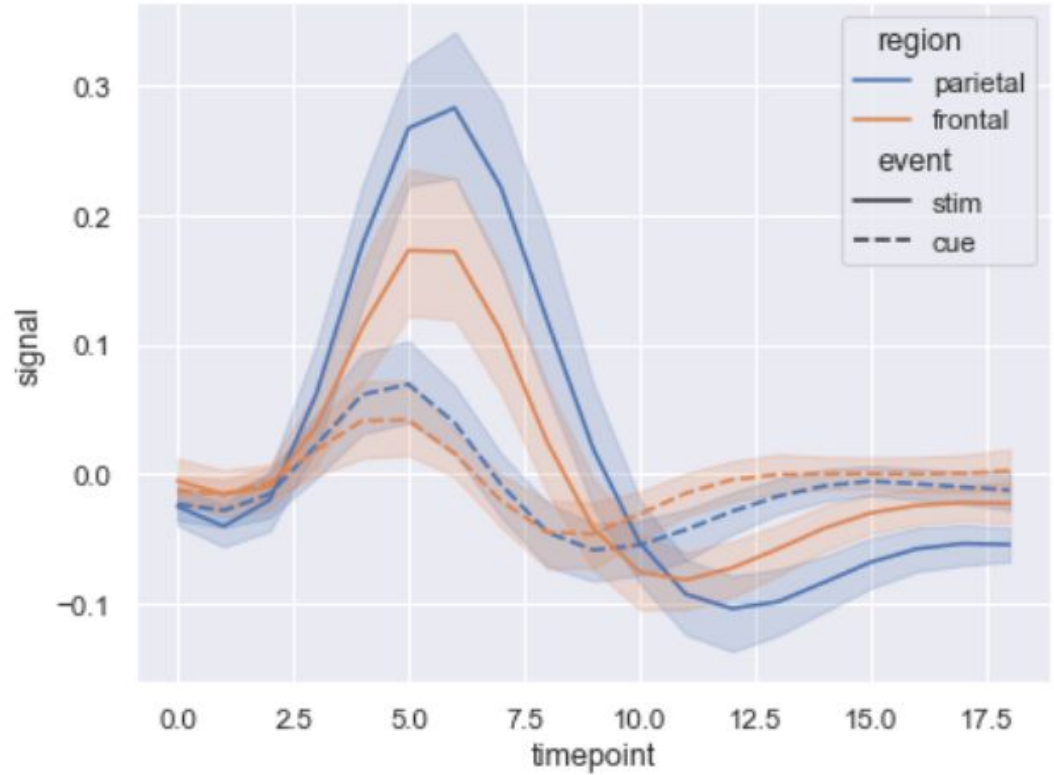
matplotlib (Python)

<https://matplotlib.org/>



seaborn (Python)

<https://seaborn.pydata.org/>



plot.ly

Interaktywne wizualizacje w Javascript z interfejsem w Python i R.

<https://plotly.com/python/line-and-scatter/>

plotly.js: <https://github.com/plotly/plotly.js>

plotly.py: <https://github.com/plotly/plotly.py>

plotly.R: <https://github.com/ropensci/plotly>