essence**mediacom**

# Wprowadzenie do Tidyverse

02
03
23

Część II

# Co dzisiaj

Operacje na zmiennych

Operacje na wierszach

Pivoty danych

Praca z tekstem

# Operacje na zmiennych

# Biblioteki w ramach Tidyverse'a

Import danych

Uporządkowanie danych

**Przetwarzanie danych**

Programowanie

Wizualizacja

# Przetwarzanie danych z pomocą biblioteki dplyr

Uporządkowany zbiór czasowników

| Czasownik | Opis |
|---|---|
| group_by() | grupuje na podstawie zadanych zmiennych |
| summarise() | na podstawie group_by tworzy agregaty zmiennych |
| mutate() | dodaje nową zmienną przy pomocy dodatkowych funkcji |
| select() | wybiera zmienną na podstawie jej nazwy |
| rename() | zmienia nazwy kolumn |
| filter() | filtruje zmienne na podstawie ich wartości |
| arrange() | zmienia kolejność wierszy |

em

# Operacje na wierszach

# Biblioteki w ramach Tidyverse'a

Import danych

Uporządkowanie danych

**Przetwarzanie danych**

Programowanie

Wizualizacja

# Przetwarzanie danych z pomocą biblioteki dplyr

Uporządkowany zbiór czasowników

| Czasownik | Opis |
|---|---|
| group_by() | grupuje na podstawie zadanych zmiennych |
| summarise() | na podstawie group_by tworzy agregaty zmiennych |
| mutate() | dodaje nową zmienną przy pomocy dodatkowych funkcji |
| select() | wybiera zmienną na podstawie jej nazwy |
| rename() | zmienia nazwy kolumn |
| filter() | filtruje zmienne na podstawie ich wartości |
| arrange() | zmienia kolejność wierszy |

# Pivotowanie danych

# Biblioteki w ramach Tidyverse'a

**Import danych**

**Uporządkowanie danych**

**Przetwarzanie danych**

**Programowanie**

**Wizualizacja**

# Tworzenie uporządkowanego zbioru danych

Funkcje w ramach biblioteki tidyr

| Funkcja | Opis |
|---|---|
| separate() | Oddziela / dzieli pojedynczą kolumnę na wiele kolumn |
| unite() | W przeciwieństwie do separate() - łączy dwie lub więcej kolumn w jedną |
| pivot_wider() | Funkcja „rozszerza'' wiele kolumn ze zbioru danych i konwertuje je na pary klucz-wartość — POZIOMO |
| pivot_longer() | Funkcja zajmuje dwie kolumny i „wydłuża'' je w kilka kolumn - PIONOWO |

# Dane długie (long) a dane szerokie (wide)

| | Date | SE_TEMP_MAX | SE_TEMP_AVG | SE_TEMP_MIN |
|---|---|---|---|---|
| 1 | 2019-12-30 | 2.2857143 | 0.4285714 | -1.1428571 |
| 2 | 2020-01-06 | 10.1428571 | 6.0000000 | 2.2857143 |
| 3 | 2020-01-13 | 16.2857143 | 10.7142857 | 5.0000000 |
| 4 | 2020-01-20 | 18.4285714 | 12.2857143 | 6.0000000 |
| 5 | 2020-01-27 | 15.8571429 | 11.8571429 | 8.4285714 |
| 6 | 2020-02-03 | 22.5714286 | 17.7142857 | 13.0000000 |
| 7 | 2020-02-10 | 22.8571429 | 16.7142857 | 11.1428571 |
| 8 | 2020-02-17 | 16.8571429 | 13.0000000 | 9.2857143 |
| 9 | 2020-02-24 | 19.0000000 | 15.2857143 | 11.2857143 |
| 10 | 2020-03-02 | 22.4285714 | 18.4285714 | 13.8571429 |
| 11 | 2020-03-09 | 22.4285714 | 17.8571429 | 13.1428571 |
| 12 | 2020-03-16 | 27.0000000 | 21.2857143 | 15.8571429 |
| 13 | 2020-03-23 | 20.8571429 | 16.7142857 | 12.7142857 |
| 14 | 2020-03-30 | 25.7142857 | 19.2857143 | 13.2857143 |
| 15 | 2020-04-06 | 23.1428571 | 18.5714286 | 13.8571429 |
| 16 | 2020-04-13 | 23.8571429 | 18.4285714 | 12.7142857 |
| 17 | 2020-04-20 | 25.4285714 | 20.0000000 | 14.2857143 |
| 18 | 2020-04-27 | 28.4285714 | 22.4285714 | 16.4285714 |
| 19 | 2020-05-04 | 29.2857143 | 23.5714286 | 17.5714286 |
| 20 | 2020-05-11 | 23.7142857 | 18.1428571 | 12.2857143 |

| | Date | type | temperature |
|---|---|---|---|
| 1 | 2019-12-30 | MAX | 2.2857143 |
| 2 | 2019-12-30 | AVG | 0.4285714 |
| 3 | 2019-12-30 | MIN | -1.1428571 |
| 4 | 2020-01-06 | MAX | 10.1428571 |
| 5 | 2020-01-06 | AVG | 6.0000000 |
| 6 | 2020-01-06 | MIN | 2.2857143 |
| 7 | 2020-01-13 | MAX | 16.2857143 |
| 8 | 2020-01-13 | AVG | 10.7142857 |
| 9 | 2020-01-13 | MIN | 5.0000000 |
| 10 | 2020-01-20 | MAX | 18.4285714 |
| 11 | 2020-01-20 | AVG | 12.2857143 |
| 12 | 2020-01-20 | MIN | 6.0000000 |
| 13 | 2020-01-27 | MAX | 15.8571429 |
| 14 | 2020-01-27 | AVG | 11.8571429 |
| 15 | 2020-01-27 | MIN | 8.4285714 |
| 16 | 2020-02-03 | MAX | 22.5714286 |
| 17 | 2020-02-03 | AVG | 17.7142857 |
| 18 | 2020-02-03 | MIN | 13.0000000 |
| 19 | 2020-02-10 | MAX | 22.8571429 |
| 20 | 2020-02-10 | AVG | 16.7142857 |

Do modelowania potrzeba danych **szerokich**

Ale dane **długie** są łatwiejsze do obróbki, wygodniejsze do trzymania w bazach danych oraz wymagane przez niektóre programy.

# Praca z tekstem

# Biblioteki w ramach Tidyverse'a

**Import danych**

**Uporządkowanie danych**

**Przetwarzanie danych**

**Programowanie**

**Wizualizacja**

# Praca z tekstem

Funkcje w ramach pakietu stringr

| Funkcja | Opis |
|---------|------|
| str_length() | Długość wyrażenia tekstowego |
| str_c() | Sklejenie kilku wyrażeń tekstowych |
| str_replace()<br>str_replace_all() | Podmiana fragmentu tekstu |
| str_detect() | Sprawdzenie występowania w jednym tekście innego tekstu |

# Wyrażenia regularne – wstęp do wstępu

Używanie wyrażeń regularnych

| Wyrażenie | Co to jest |
| --- | --- |
| **a** | Litera ☺ |
| **\ \** | \ |
| **\ "** | " |
| **[:upper:]** | Wielkie litery |
| **[:punct:]** | Znaki przestankowe |
| … i inne | |

# StringR cheat sheet

**essencemediacom**
business science

Jarek Dejneka, Managing Partner
jaroslaw.dejneka@essencemediacom.com

Bartek Kowalski, Business Science Director
bartosz.kowalski@essencemediacom.com

@ mbs@mediacom.com

@MBSWarsaw

@MBSWarsaw

business-science.pl

*„An investment in knowledge always pays the best interest"*

Benjamin Franklin