

ŁUKASZ ŻYŁA

**DZIENNIKARSTWO
DANYCH
I DATA
STORYTELLING**

Wszelkie prawa zastrzeżone. Nieautoryzowane rozpowszechnianie całości lub fragmentu niniejszej publikacji w jakiegokolwiek postaci jest zabronione. Wykonywanie kopii metodą kserograficzną, fotograficzną, a także kopiowanie książki na nośniku filmowym, magnetycznym lub innym powoduje naruszenie praw autorskich niniejszej publikacji.

Wszystkie znaki występujące w tekście są zastrzeżonymi znakami firmowymi bądź towarowymi ich właścicieli.

Autor oraz wydawca dołożyli wszelkich starań, by zawarte w tej książce informacje były kompletne i rzetelne. Nie biorą jednak żadnej odpowiedzialności ani za ich wykorzystanie, ani za związane z tym ewentualne naruszenie praw patentowych lub autorskich. Autor oraz wydawca nie ponoszą również żadnej odpowiedzialności za ewentualne szkody wynikłe z wykorzystania informacji zawartych w książce.

Redaktor prowadzący: Barbara Gancarz-Wójcicka

Helion S.A.

ul. Kościuszki 1c, 44-100 Gliwice

tel. 32 231 22 19, 32 230 98 63

e-mail: onepress@onepress.pl

WWW: <http://onepress.pl> (księgarnia internetowa, katalog książek)

Drogi Czytelniku!

Jeżeli chcesz ocenić tę książkę, zajrzyj pod adres

<http://onepress.pl/user/opinie?dzidan>

Możesz tam wpisać swoje uwagi, spostrzeżenia, recenzję.

ISBN: 978-83-283-8312-8

Copyright © Łukasz Żyła 2022

Printed in Poland.

- [Kup książkę](#)
- [Poleć książkę](#)
- [Oceń książkę](#)

- [Księgarnia internetowa](#)
- [Lubię to! » Nasza społeczność](#)

Spis treści

Wstęp	7
--------------------	----------

Rozdział 1. Dane i dziennikarstwo — wprowadzenie	10
---	-----------

1.1. Czym jest dziennikarstwo danych	10
1.2. Dziennikarstwo danych jako metoda	11
1.3. Dane	14
1.4. Dziennikarstwo a świat cyfrowy	15
1.5. Otwarte dane	16
1.6. Warsztat dziennikarza pracującego z danymi	17
1.7. Dokąd zmierza dziennikarstwo cyfrowe?	18

Rozdział 2. Dziennikarstwo danych — przykłady	20
--	-----------

2.1. Odkrywamy	20
2.2. Pokazujemy	25
2.3. Angażujemy	27

Rozdział 3. Internet, czyli dziennikarstwo online	29
--	-----------

3.1. Internet	29
3.2. Jak działają strony internetowe	30
3.3. Język HTML i CSS	32

Rozdział 4. Dane — formaty, źródła, zdobywanie	35
---	-----------

4.1. Przetwarzanie danych	35
4.2. Źródła informacji	41
4.3. Dane na wniosek	42
4.3.1. Bądź przygotowany na opóźnienie	43
4.3.2. Kiedy informacja jest publiczna?	44
4.3.3. Kto ma obowiązek udzielenia informacji?	44

4.3.4. Wnioskujej o dane przetwarzalne	45
4.3.5. Precyzja i cierpliwość	46
4.3.6. Pytaj w różnych miejscach i testuj	47
4.3.7. Odmowa i co dalej?	47
4.4. A może zautomatyzować?	48
4.4.1. Wykorzystanie Google Sheets	51
4.4.2. Gdzie przechowywać dane?	52
4.4.3. Jakość i źródła ogólnodostępne	53

Rozdział 5. Czyszczenie i agregacja danych 55

5.1. Pozbądź się śmieci	55
5.1.1. Sortowanie	57
5.1.2. Wyszukiwanie fasetowe	57
5.1.3. Wykrywanie duplikatów	59
5.1.4. Zastosowanie filtru tekstowego	59
5.1.5. Transformacje komórek	59
5.1.6. Usuwanie niewłaściwych danych	61
5.1.7. Grupowanie danych, czyli klastrowanie	61
5.2. Szukanie wzorców, czyli analiza	63
5.2.1. Odrobina statystyki	63
5.2.2. Powiązania i korelacje	66
5.3. Metody analizy w praktyce	67
5.3.1. Przygotowanie tabeli w Google Sheets	67
5.3.2. Analizy bazodanowe	73

Rozdział 6. Jak tworzyć zrozumiałe i przyciągające uwagę wizualizacje 78

6.1. Po co i jak wizualizować?	80
6.2. Jak wygląda dobra wizualizacja?	87
6.2.1. Bądź uczciwy	90
6.2.2. Dobierz odpowiedni wykres	90
6.2.3. Wykorzystaj właściwe skale	90
6.2.4. Zastosuj odpowiednie kolory	91
6.2.5. Spraw, aby wizualizacja była czytelna	91
6.3. Najczęstsze błędy	92
6.3.1. Wykres 3D	92
6.3.2. Brak odwzorowania danych	92
6.3.3. Zbyt duże opisy i legenda	93
6.3.4. Brak różnic	93
6.3.5. Zmiana skali	93

6.4. Rodzaje wizualizacji	94
6.4.1. Wizualizacja proporcji	94
6.4.2. Porównanie kilku wartości	98
6.4.3. Śledzenie zmian w czasie	100
6.4.4. Obserwacja zależności między danymi	102
6.4.5. Zależności przestrzenne	107
6.4.6. Rozkład cechy statystycznej	107

Rozdział 7. Data storytelling 111

7.1. Składowe	112
7.2. Wzorce, czyli szkielet data story	115
7.2.1. Wzorce argumentacji	116
7.2.2. Wzorce przepływu	117
7.2.3. Wzorce kadrowania narracji	117
7.2.4. Wzorce empatii i emocji	119
7.2.5. Wzorce zaangażowania	121
7.3. Elementy GUI	123
7.4. Data storytelling w praktyce	124

Rozdział 8. Jak stworzyć wizualizacje 136

8.1. Na początek coś prostego — Datawrapper	137
8.2. Flourish	139
8.3. Sieci społecznościowe	141
8.4. Coś bardziej zaawansowanego, czyli Tableau	143
8.5. Ekosystem, czyli Google Charts i Google Maps	149
8.6. Infogram	151
8.7. Łączenie wizualizacji	153

Rozdział 9. Wstęp do programowania 159

9.1. Zmienne	161
9.2. Obiekty	162
9.3. Tablice	162
9.4. Tablice i obiekty	164
9.5. Funkcje	165
9.6. Instrukcje warunkowe	166

Rozdział 10. Wizualizacja zaprogramowana	169
10.1. Chart.js	170
10.2. Highcharts i AmCharts	174
10.3. p5.js	175
10.4. ECharts	180
10.5. Vega-Lite	182
Rozdział 11. Data story w sieci	189
Dodatek A. Wywiad z Dominikiem Uhligiem, redaktorem naczelnym serwisu BIQdata „Gazety Wyborczej”	201
Zakończenie	205
Bibliografia	206

Rozdział 5.

Czyszczenie i agregacja danych

Kiedy przymierzamy się do stworzenia materiału dziennikarskiego, najpierw, rzecz jasna, szukamy informacji. Czasami poszukujemy jej w ciemno, czasami kierując się intuicją, a innym razem po prostu logicznie wyciągając wnioski, potrafimy je z pewnym prawdopodobieństwem przewidzieć.

W przypadku dziennikarstwa danych zamiast informacji poszukujemy danych. Nie pozostajemy przy samej informacji; zamiast tego sprawdzamy, na podstawie jakich danych została utworzona (o tym, że informacje składają się z danych, pisałem w rozdziale 4.). Natomiast poszukiwanie źródeł informacji, czyli danych, nie różni się niczym innym od poszukiwania informacji, tzn. tak samo często kierujemy się intuicją w poszukiwaniach, co również zostało opisane. W przypadku tradycyjnego dziennikarstwa to poszukiwanie zajmuje najwięcej czasu, potem czeka nas już tylko przygotowywanie materiału (tekstu, wideo, reportażu radiowego) i oddanie do redakcji.

Przygotowanie i oddanie do składu zajmuje zdecydowanie mniej czasu aniżeli wyszukiwanie informacji. Odmienne jest z materiałem dziennikarstwa danych. Po wyszukaniu danych konieczna jest żmudna praca przy ich przygotowywaniu i analizie. Zdarza się często, że jakość danych, które pozyskaliśmy (zwłaszcza danych publicznych), pozostawia wiele do życzenia. Trzeba sobie zdawać sprawę, że instytucje (zwłaszcza publiczne), które gromadzą dane, nie robią tego, aby można było na ich podstawie prowadzić jakies analizy. Dane gromadzone są często do celów administracyjnych czy biurokratycznych. Dlatego urzędnicy niespecjalnie przywiązują wagę do jakiegokolwiek standaryzacji prezentacji danych. Mając jednak dane w postaci tabelarycznej w takich formatach jak `.csv` czy `.xls`, możemy je odpowiednio wyczyścić, sformatować, a następnie posortować i przefiltrować.

5.1. Pozbądź się śmieci

Na początek zajmiesz się czyszczeniem danych. Aby odnaleźć jakieś wzorce w danych, wywnioskować jasne i klarowne informacje, musisz te dane wyczyścić. Co to oznacza? Musisz zrobić wszystko, żeby dane stały się zrozumiałe i przejrzyste w formie tabeli. Tabele nie mogą zawierać pustych miejsc (np. spacji czy pustych komórek), wszystkie litery powinny być ustandaryzowane, w kolumnach mają się znajdować dane tylko jednego typu. Nie da się przecież pracować na tabeli, jeśli w kolumnie *X* umieszczone są numery i nazwy miast — takiej kolumny nie da się np. zsumować.

Na tym etapie pracy zachęcam Cię do korzystania z różnych narzędzi i do raczej intuicyjnego podejścia. Jeżeli dysponujesz małą tabelą w formacie *.xls*, w której znajdują się połączone kolumny, w dodatkowych wierszach jakieś opisy itp., to może wystarczy po prostu wyczyścić ją ręcznie w Google Sheets? Lecz jeśli wierszy masz ponad 100, lepiej skorzystać z automatycznych rozwiązań. Jednym z najbardziej popularnych narzędzi do sporych zbiorów jest Open Refine. Ta aplikacja (do 2012 r. funkcjonowała jako Google Refine) to narzędzie ułatwiające porządkowanie i przetwarzanie danych, umożliwiające import danych z wielu formatów: *.tsv*, *.csv*, **.sv*, *.xls* i *.xlsx*, *.json*, *.xml*, *.rdf*, *.xml*. Open Refine ułatwia sortowanie i dzielenie danych, a także ich przetwarzanie z wykorzystaniem tzw. wyrażeń regularnych. Swoją pracę można w łatwy sposób eksportować do popularnych formatów (*.csv*, tabela *.html* czy pliki Excela i *.odt*). Z Open Refine można korzystać w systemach linuksowych, macOS oraz w Windowsach. Narzędzie dystrybuowane jest na licencji *open source*.

Aplikację pobierzesz ze strony openrefine.org. Żeby ją zainstalować, postępuj zgodnie z instrukcjami. Aplikacji używa się w przeglądarce. Aby zacząć pracę z plikiem z danymi, kliknij *Create project*, a następnie skorzystaj z możliwości narzędzia i wybierz sposób pobrania danych. Oprócz tradycyjnego pobrania pliku z komputera możesz także umieścić link do adresu URL czy bezpośrednio połączyć program z arkuszem Google'a. W aplikacji pojawi się tabela, a pod nią dodatkowe możliwości jej przetworzenia (rysunek 5.1).

The screenshot shows the OpenRefine web interface. At the top, there's a navigation bar with 'Open Project', 'Import Project', and 'Language Settings'. Below that is a table with columns: Record ID, Object Title, Registration Number, Description, Marks, Production Date, Provenance (Production), and Proven. The table contains six rows of data related to space exploration artifacts. Below the table, there are sections for 'Parse data as' (listing various file formats like CSV, TSV, etc.), 'Character encoding' (with a dropdown menu), and a 'Parse options' section with various checkboxes for how to handle data (e.g., 'Ignore first line(s)', 'Parse next line(s) as column headers', etc.).

Record ID	Object Title	Registration Number	Description	Marks	Production Date	Provenance (Production)	Proven
1	Rocket motor on loan from Roswell Museum and Art Center, USA	L2106-31	Rocket motor, liquid fuelled combustion chamber, steel / aluminum wrapped tubes with insulation, included in flight of Dec 28 1952, Robert H Goddard, USA, 1952 (Roswell Acc No: 1956-28-12)				
2	Fragment of moon rock on loan from National Aeronautics and Space Administration (NASA), USA	L41151	Fragment of lunar sample (moon rock), NASA No.61018, 116 (D11) and box, weight of rock 89 grams, collected by Apollo 16 cosmonaut Charles Duke on the east rim of Plum Crater 30 meters west north of the landing site, 1972. This fragment is 3.9 billion years old. This is older than most of the surface rock on earth. The rock is breccia, fragments held together by a cementing substance.				
3	F-1 rocket engine on loan from National Air and Space Museum (NASM), USA	L2040-3	F-1 rocket engine (NASM 6311 Cat:1978-0199)				
4	Shuttle garments worn by Dr P Scully-Power on loan from National Aeronautics and Space Administration (NASA), USA, 1984	L2066-13	Protective clothing, space shuttle "in-flight Covered Demand" PMS10186 10055-01, 87x, 10x, with name tag PMS 8710788-01, cotton, used by Paul Scully-Power / made for NASA, USA, 1984 (NASA: A37 28184278)		1984	Maker: National Aeronautics and Space Administration (NASA); Florida, USA, 1984)	User: St Paul, Pa 1984)
5	Replica Manned Maneuvering Unit (MMU) on loan from Lockheed Martin Space Systems Company, USA, 1984	L1870	Model (full size replica), Manned Maneuvering Unit (MMU), metal, Martin Marietta Aeronautics Group, USA, [1984]		1984	Designer: Martin Marietta Aeronautics Group, Maryland, USA; 1984)	
6	Mural, Zheng He and Columbus by Guan Wei, 2008 - 2007	20105/1	Mural painting, "Zheng He and Columbus", acrylic on medium density fibre board (MDF), painted by Guan Wei for the exhibition, "Other histories: Guan Wei's table for a contemporary world", Powerhouse Museum, Sydney, New South Wales, Australia, 2008-2007. A mural depicting a section of the globe and constellation against a black background. The mural was painted in two parts on medium density fibre (MDF) board and needs to be displayed together. Two expeditionary ships are at the centre, the small boat symbolising the ship of Christopher Columbus (1481-1506) and the larger boat symbolising the ship of Zheng He (1371-142). A hand and red arrow point to the left. The painting includes maps, human figures and a mysterious animal. This particular mural painting was displayed at the beginning of the exhibition as an introduction and shows discoverer stories of		2008 - 2007	unknown; 2008 - 2007)	

Rysunek 5.1. Open Refine to rozbudowany program do czyszczenia i agregacji danych. Na początku z pewnością wymaga sporej uwagi i skupienia.

Źródło: materiały własne

Program domyślnie ustawi pierwszy wiersz jako nagłówek. Możesz także ustawić sposób parsowania Twojego pliku, np. w przypadku *.tsv* separacja komórek powinna być zastosowana za pomocą tabulatorów, a w przypadku *.csv* za pomocą przecinków. Program zazwyczaj wszystkie te informacje wykrywa i automatycznie Ci w tym pomaga, poprzez dolny panel możesz jednak dokonać manualnie kilku zmian w sposobie pobierania Twojej tabeli.

Gdy uda Ci się już zaimportować całą tabelę, możesz sobie poklikać różne miejsca, aby bliżej zapoznać się z interfejsem. Do manipulowania danymi służą funkcje dostępne poprzez kliknięcie ikonki trójkąta przy nagłówkach kolumn. Daną kolumnę możesz np. zwinąć, klikając *view* i dalej *Collapse this column*; jeżeli chcesz, aby kolumna z powrotem się pojawiła, wystarczy kliknąć puste miejsce po niej. Kolumny możesz też przesuwac w prawo lub w lewo, aby wygodniej Ci się na nich pracowało. Możesz także zmieniać nazwy kolumn przy użyciu funkcji *rename*. Na bocznym panelu z lewej strony możesz prześledzić historię swoich zmian i powrócić do wcześniejszych ustawień.

Po zakończeniu pracy możesz wyeksportować projekt do różnych formatów: *.csv*, *.tsv*, *.xls*, *.xlsx* czy tabeli w dokumencie *.html*. Jeżeli pracujesz z bardzo dużymi zbiorami danych, możesz zwiększyć alokację zużywania pamięci dla swojej aplikacji. Możesz to zrobić zgodnie z instrukcją znajdującą się na stronie kodu źródłowego: <https://github.com/OpenRefine/OpenRefine/wiki/FAQ%3A-Allocate-More-Memory>.

W celu pokazania możliwości programu wykorzystam zbiór danych ze strony dane.gov.pl, zawierający listę beneficjentów (projektów) Funduszy Europejskich w latach 2014 – 2020. Plik bez problemu wyszukasz na stronie (ponieważ instytucje publiczne często zmieniają adresy URL swych zasobów, nie będę podawał aktualnego linku).

5.1.1. Sortowanie

Dane z tabeli możesz posortować, co przyda Ci się do późniejszej analizy. Do tego wykorzystasz funkcję *sort*, która pojawi się po kliknięciu ikony trójkąta przy nazwie kolumny. Daną kolumnę możesz posortować ze względu na tekst (od A do Z lub na odwrót), liczby (od najmniejszych do największych lub na odwrót), daty (od najwcześniejszych lub najpóźniejszych) i format *boolean*, czyli dane, które mają wartość prawdy lub fałszu. Metodą „przeciągnij i upuść” (ang. *drag & drop*) możesz ustalić, w jakiej kolejności mają się pojawić np. puste komórki lub te z błędem.

5.1.2. Wyszukiwanie fasetowe

Open Refine daje Ci możliwość analizowania każdej kolumny np. pod względem powtarzanych cech w wierszach. Cechy informacji, które powtarzają się w komórkach, to fasety i na ich podstawie w Open Refine tworzy się filtry. Proces ten nazywa się wyszukiwaniem fasetowym. Możliwość wyszukiwania fasetowego to przede wszystkim sposobność spojrzenia na dane z różnych perspektyw. Metoda ta nie zmienia wartości komórek. Przetestuj ją na kolumnie *Nazwa beneficjenta*.

Wybierz *facet*, a następnie *Text facet*. W lewym panelu pojawia się efekt funkcji. Zmień *Sort by* na *count*. Jak widzisz, na pierwszym miejscu jest Łódź, która otrzymała fundusze na 121 projektów. Widzisz także 21497 *choices*, co oznacza, że łącznie 21 497 podmiotów uzyskało dofinansowanie na swoje projekty.

Funkcję tę można zastosować także do wartości liczbowych. Aby jednak tego dokonać, musisz się zapoznać z inną możliwością Open Refine. Zwróć uwagę na kolumnę *Poziom unijnego dofinansowania w procentach*. Wartości zapisane są z przecinkiem, co oznacza, że nie są one liczbami. Cyfry dziesiętne w programie Open Refine oddzielane są kropką. W tym przypadku są to przecinki, zatem wartości w tej kolumnie to tekst. Aby wykorzystać funkcję do wyszukiwania fasetowego na liczbach, trzeba wartości z kolumny przekonwertować na wartość liczbową.

W programie masz do czynienia z czterema typami danych: tekst (*string*), liczba (*integer*), *boolean* (wartość prawda lub fałsz), *data*. Każda komórka zawierająca jakąś wartość musi być jednym z tych czterech typów danych. W celu analizy danych możesz dokonać transformacji jednego typu na inny.

Aby przekonwertować wartość tekstową na liczbową, musisz najpierw zamienić przecinek na kropkę, a następnie dokonać transformacji. Do tego możesz użyć zaawansowanej możliwości tworzenia tzw. wyrażeń regularnych za pomocą Google Refine Expression Language, czyli GREL. GREL działa tak samo jak tworzenie formuł w Excelu. Kliknij więc ikonę trójkąta kolumny, potem *Edit cells* i wybierz *Transform*. W ramce, gdzie zapisana jest nazwa *value*, musisz wpisać następującą regułę: `value.replace(,,'.').toNumber()`.

`value.replace(,,'.')` oznacza, że w wartości komórki (*value*) odnajdujesz przecinek (pierwszy znak w apostrofie — ','), który zastępujesz kropką (drugi znak w apostrofie '.'). Przecinek i kropka w apostrofach oddzielone są przecinkiem.

`toNumber()` — program dokonuje transformacji wartości (*value*) na typ liczbowy (*number*).

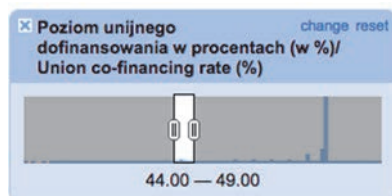
Kliknij OK. Jak widać, wartości zmieniły kolor na zielony. Oznacza to, że są teraz typem liczbowym.

Możesz teraz skorzystać z funkcji fasetowej. Kliknij *facet*, a następnie *numeric facet*. W lewym panelu pojawiła się ramka z grafiką (rysunek 5.2).

Rysunek 5.2.

Suwak regulujący zakres liczbowy.

Źródło:
materiały własne



Suwakiem możesz np. ustawić obszar 40% dotacji. Możesz teraz dodać kolejny filtr, czyli *facet*, z nazwami beneficjentów. Kliknij przy kolumnie z nazwami beneficjentów *facet* i *text facet*. W ten sposób w lewym panelu pojawi się lista podmiotów, które otrzymały dotacje na poziomie 40% wartości projektu.

5.1.3. Wykrywanie duplikatów

W zbiorach danych często będziesz miał do czynienia z powtarzającymi się rekordami. Duplikaty to rekordy, które pojawiają się minimum podwójnie. W tym fragmencie posłużysz się przykładem z zagranicznej bazy danych, a konkretnie z The Museum of Applied Arts and Sciences z Sydney, którą można znaleźć w serwisie data.nsw.gov.au (wszelkiego rodzaju zbiory muzealne to dobre miejsca, na których można szlifować umiejętności czyszczenia baz). Obecnie muzeum umożliwia dostęp do swych zbiorów poprzez API. Ty jednak skorzystasz z pliku *.tsv*, który wcześniej został pobrany z API jako plik *.json*.

Plik umieść w programie w podany wcześniej sposób. Aby sprawdzić liczbę zduplikowanych rekordów, posłuż się w tej bazie polem *Registration Number*. Następnie wybierz *Facet/Customized facets/Duplicates facet*. Jak widzisz, w bocznym panelu z wartością *true* znajduje się 667 rekordów. Oznacza to, że 667 rekordów bazy danych jest zduplikowanych. Po kliknięciu *true* w kolumnie *Registration Number* zauważysz na górze puste pola. Oznacza to, że program traktuje puste pola jako zduplikowane. Musisz je wyodrębnić. Kliknij więc *Registration Number/Facet/Customized facets/Facet by blank*. Jeśli zaznaczysz *false*, przekonasz się, że teraz rekordów zduplikowanych jest 37. Aby przyjrzeć się tym rekordom, wybierz dodatkowo *text facet*. Posortuj teraz te rekordy liczbą (*count*). Zauważ, że na pierwszym miejscu znalazł się rekord, który ma nawet 4 duplikaty. Następnie spróbuj nieco inną metodą sprawdzić duplikaty w kolumnie *Record ID*. Jako że fasetowe wyszukiwanie duplikatów nie przeszukuje wartości liczbowych (*integers*), zacznij od posortowania tej kolumny. Wybierz *Sort by* na pasku nad nagłówkami kolumn i dalej *By Rekord ID*, zaznacz *numbers* i wybierz *smallest first*. Podczas sortowania zaznacz też *Reorder rows permanently*, aby zmiany w sortowaniu zostały zapisane na trwałe. Niebawem dowiesz się, jak usunąć duplikaty.

5.1.4. Zastosowanie filtru tekstowego

Spróbuj teraz znaleźć w zbiorach muzeum jakieś rekordy, które mogą dotyczyć naszego kraju. Sprawdź tytuły prac. Wybierz *Object Title*, a następnie *Text filter*. W bocznym panelu wpisz Poland. Brawo, znalazłeś 7 rekordów! Filtrem tekstowym możesz teraz sprawdzić kolumnę kategorii zbioru dzieł. Po wybraniu kolumn *Categories* i *Text filter* wprowadź |. W ten sposób sprawdzisz rekordy, które dotyczą więcej niż jednej kategorii. Jak widzisz, większość rekordów dotyczy więcej niż jednej kategorii. Wpisz teraz podwójny podziałnik ||. Okazuje się, że jest 8 rekordów z błędnie wpisanymi podziałkami kategorii. Dzięki filtrowaniu tekstowemu nie tylko możesz wyszukać różnego rodzaju dane, ale także wykryć nieprawidłowości.

5.1.5. Transformacje komórek

Oprócz filtrowania i analizy zbioru danych Open Refine pozwala Ci także na transformacje i edycję komórek. W programie umożliwia to zakładka *Edit cells*. Czasami musisz np. „oczyścić” wartości w komórkach, połączyć je czy usunąć niepotrzebne spacje.

Transformacje możesz zacząć od usuwania białych znaków (np. niepotrzebnych spacji znajdujących się na początku lub na końcu wartości w komórkach). Sprawdź więc tym razem w bazie dotacji unijnych kolumnę z tytułami projektów. Aby wykryć białe znaki, wybierz *Tytuł projektu/Edit cells/Common transforms/Trim leading and trailing whitespace*. Jak widzisz na rysunku 5.3, program usunął białe znaki w 1217 komórkach.

The screenshot shows the OpenRefine interface with a table of 46430 rows. A filter is applied to the 'Tytuł projektu' column, showing a list of filter rules on the left. The first rule is selected: 'Test transform on 0 cells in column Tytuł projektu (value: <string>)'. The table columns include: Tytuł projektu, Numer umiarów, Nazwa beneficjenta, Fundacja Fund, Program Projekt, Projekt Projekt, Opcjonalne Nazwa, Publikowane w, Market projekt, Wykalkulacja, Market udzieleny, and Płatność udzieleny. The 'Tytuł projektu' column contains various project titles, some of which are highlighted in blue, indicating they have been processed by the filter.

Rysunek 5.3. Program umożliwia szybkie usuwanie zbędnych komórek.

Źródło: materiały własne

Należy dodać, że tego typu operacje przeprowadza się tylko na typach tekstowych, czyli *string*, a nie numerycznych, czyli *integer*. Możesz także skorzystać z metody *Collapse consecutive whitespace*, która usuwa białe znaki ze środka stringów.

Jeżeli Twoja baza oparta jest na danych pochodzących ze strony internetowej, może się zdarzyć, że komórki będą zawierać pozostałości kodu HTML-a. Na przykład litera „s” w kodzie HTML-a zapisywana jest jako `"`. Jeżeli znajdziesz w swojej bazie jakieś komórki, w których treść zaczyna się od `&` i kończy się średnikiem (;), możesz zmienić sposób wyświetlania takich znaków na czytelny. Zrobisz to, klikając w danej kolumnie *Edit cells/Common transforms/Unescape HTML entities*. W bazie dotacji unijnych możesz sprawdzić w ten sposób kolumnę *Skrócony opis*. Okazuje się, że program zamienił 466 komórek, wykrywając m.in. takie znaki jak cudzysłów, który wcześniej był zapisany w formie kodu HTML-a jako `”`.

Kolejną możliwością, jaką dają transformacje, jest zmiana wielkości liter w komórkach tekstowych (np. z małych na duże i odwrotnie). Aby zmienić tekst komórki na duże litery, kliknij daną kolumnę, a następnie *Edit cells/Common transforms/To uppercase*. Wykorzystanie odpowiednich metod zależy od tego, w jaki sposób chcesz przekształcić teksty w poszczególnych kolumnach.

Więcej możliwości edycji komórek daje z pewnością użycie wyrażeń regularnych. Jeżeli zaczynasz pracować z Open Refine, warto, byś zapoznał się z możliwościami, jakie daje wykorzystanie właśnie wyrażeń regularnych. Więcej informacji na temat General Refine Expression Language znajdziesz pod linkiem <https://github.com/OpenRefine/OpenRefine/wiki/General-Refine-Expression-Language>.

5.1.6. Usuwanie niewłaściwych danych

Wiesz już, jak filtrować dane i jak wykorzystać wyszukiwanie fasetowe. Wykrycie duplikatów i powtarzających się elementów w treści komórek to dopiero część Twojej pracy. W dalszej kolejności powinieneś pozbyć się tych niepotrzebnych komórek. I tutaj zaczyna się to, co nazywamy czyszczeniem danych, czyli usunięcie konkretnych komórek z bazy danych. Pamiętaj, że usuwanie komórek powinno się odbywać po uprzednim wykorzystaniu metod *facet* i filtrów, w przeciwnym razie narazisz się na usunięcie wszystkich komórek.

Wróć teraz do wcześniejszej bazy muzeum z Sydney i wykonaj te same kroki. W kolumnie *Registration Number* wybierz *Facet/Customized facets/Duplicates facet*. W bocznym panelu z wartością *true* znajduje się 667 rekordów. Jak już wiesz, wśród tych rekordów są także puste pola, których jednak nie chcesz się pozbyć. Musisz je więc wyodrębnić. W tym celu kliknij *Registration Number/Facet/Customized facets/Facet by blank*. Gdy zaznaczysz *false*, przekonasz się, że teraz rekordów zduplikowanych jest 37. Nadszedł czas, aby je usunąć. W nagłówku pierwszej kolumny *All* kliknij *Edit rows/Remove all matching rows*. Jak widzisz w komunikacie, właśnie usunąłeś 37 wierszy.

5.1.7. Grupowanie danych, czyli klastrowanie

W ostatnim przykładzie zajmiesz się problemem, który często pojawia się w bazach danych, zwłaszcza gdy masz do czynienia z danymi publicznymi. Tego typu bazy zazwyczaj są tworzone przez wiele osób, dlatego często można napotkać różne nazwy dotyczące tych samych elementów. W kategoriach wydatków możesz natknąć się np. na: *Wydatki na usługi zdrowotne | zdrowie*, *Wydatki na usługi zdrowotne | Zdrowie* i *Wydatki na zdrowie*. Jak można się domyślić, wszystkie dane objęte tymi kategoriami dotyczą tego samego rodzaju wydatku. Problem ten rozwiązuje tzw. klastrowanie. To nic innego jak grupowanie obiektów o podobnych właściwościach. Open Refine daje Ci możliwość automatycznego poradzenia sobie z problemem grupowania.

Wykorzystasz teraz bazę muzeum z Sydney. Aby sprawdzić, czy podobna sytuacja występuje w tej bazie, sprawdź kolumnę kategorii dzieł. Do klastrowania kolumny użyj metody *cluster*. Znajdziesz ją w *Edit cells/Cluster and edit*. Zostanie wyświetlony panel z pogrupowanymi kategoriami. Program Open Refine stosuje algorytmy skonstruowane z wielu metod pozwalających wykryć podobieństwa w odpowiednich kategoriach.

Przeanalizuj ramkę z rysunku 5.4. *Cluster Size* oznacza liczbę różnych oznaczeń tych samych kategorii, jaką znalazł program. *Row Count* pokazuje, ile znaleziono wierszy z wszystkimi nazwami kategorii. W *Values in Cluster* wypisane są wykryte przez program rodzaje kategorii

(w naszym przypadku program wykrył po dwa rodzaje oznaczeń tej samej kategorii dla wszystkich kategorii). W polu *Merge?* możesz zaznaczyć *checkbox* — wówczas nazwy kategorii zmienia się w treść znajdującą się w *New Cell Value*. Następnie wybierz *Merge Selected & Re-Cluster* (jeżeli nie chcesz zamykać okna i np. przyjrzeć się jeszcze raz możliwości klastrowania) lub *Merge & Close* (jeżeli chcesz dokonać klastrowania i zamknąć okno). Przy grupowaniu podobnych elementów powinieneś działać ostrożnie, aby się nie okazało, że połączyłeś dwa zupełnie różne obiekty.

Cluster & Edit column "Categories"

This feature helps you find groups of different cell values that might be alternative representations of the same thing. For example, the two strings "New York" and "new york" are very likely to refer to the same concept and just have capitalization differences, and "Gödel" and "Godel" probably refer to the same person. Find out more ...

Method: *key collision* Keying Function: *fingerprint* 18 clusters found

Cluster Size	Row Count	Values in Cluster	Merge?	New Cell Value
2	17	<ul style="list-style-type: none"> Audio equipment Audio and Visual Equipment (15 rows) Audio and visual equipment Audio and Visual Equipment (2 rows) 	<input type="checkbox"/>	Audio equipment Audio and Vis
2	215	<ul style="list-style-type: none"> Botanical specimens Botanical Specimens (211 rows) Botanical specimens Botanical specimens (4 rows) 	<input type="checkbox"/>	Botanical specimens Botanical
2	3	<ul style="list-style-type: none"> Micrometers Measuring Instruments (2 rows) Micrometers Measuring instruments (1 rows) 	<input type="checkbox"/>	Micrometers Measuring Instrum
2	2	<ul style="list-style-type: none"> Astronomical equipment Testing equipment Scientific Instruments (1 rows) Astronomical equipment Testing equipment Scientific instruments (1 rows) 	<input type="checkbox"/>	Astronomical equipment Testin
2	2	<ul style="list-style-type: none"> Traffic control equipment Transport-Land Electronics (1 rows) Traffic control equipment Transport-Land Electronics (1 rows) 	<input type="checkbox"/>	Traffic control equipment Trans

Rows in Cluster: 0 — 590

Average Length of Choices: 18 — 63

Length Variance of Choices: 0 — 5.5

Select All Unselect All Export Clusters Merge Selected & Re-Cluster Merge Selected & Close Close

Rysunek 5.4. Grupowanie umożliwia wydobycie danych, które zostały np. niewłaściwie opisane.

Źródło: materiały własne

Open Refine z pewnością daje duże możliwości czyszczenia danych, a jego wielka zaleta to fakt, że jest darmowy. Program można rozwijać, dodając do niego kolejne rozszerzenia, lub można (w przypadku wyrażeń regularnych) zastosować język programowania.

Część tych możliwości z pewnością znajdziesz też w innych programach. Przy mniejszych zbiorach danych wystarczy Ci Google Sheets lub Excel (ewentualnie z manualnym czyszczeniem). Tak czy inaczej pokazane wyżej elementy występują niemal w każdej tabeli, tzn. na pewno znajdziesz niepotrzebne puste miejsca, powtarzające się wiersze, dublujące się kategorie itd. Czasami jednak trzeba będzie połączyć dwie tabelę dotyczące tego samego zagadnienia, a sformatowane w zupełnie inny sposób, i tu już Open Refine okaże się nieoceniony.

5.2. Szukanie wzorców, czyli analiza

Jak wspominałem na początku książki, etap pobierania, czyszczenia i agregowania danych zajmuje nawet do 90% czasu tworzenia materiału. Może się okazać, że będziesz obrabiać przez 5 dni jeden zbiór danych, po to, by w godzinę przygotować wykres słupkowy. To jednak esencja pracy dziennikarza danych. To tutaj będziesz odkrywać historie: w liczbach.

Aby je jednak odkryć, musisz wiedzieć, jak szukać i gdzie szukać. W liczbach mogą się kryć różnego rodzaju zależności, których poszukujesz. Jak je znaleźć? Cóż, musisz swoje zbiory przeanalizować, dokonać ich analizy ilościowej. Tymi zagadnieniami zajmuje się statystyka. I to jest moment, przy którym dziennikarzy i dziennikarki ogarnia paniczny strach. Już spieszą z pocieszeniem. To, czego potrzebujesz ze statystyki, aby być dziennikarzem danych, to raptem kilka pojęć i wzorów. W większości przypadków Twoja praca będzie opierać się na tzw. danych zastanych. Oczywiście możesz zaproponować w redakcji przeprowadzenie ankiet, wykorzystać crowdsourcing itp., wówczas jednak proponuję konsultacje z osobą zajmującą się statystyką i socjologią. Przejdźmy zatem do zagadnień, które powinny być częścią warsztatu każdego dziennikarza, czyli podstawowych pojęć ze statystyki.

5.2.1. Odrobina statystyki

Kluczowym pojęciem niewątpliwie jest zmienna. Zanim rozpoczniesz pracę, musisz zobaczyć, z czym masz do czynienia. Co możesz wówczas odkryć? Z pewnością to, jakie zmienne występują w Twoim zbiorze. W tabeli głosowań radnych możesz znaleźć takie zmienne jak płeć, przynależność partyjna, liczba oddanych głosów, lata itd.

Jak można, Twoim zdaniem, pogrupować te zmienne? W dużym uproszczeniu: na te, które da się „zmierzyć”, i na te, których „zmierzyć” się nie da. Zmienna ma jedną kluczową właściwość: jest albo ciągła, albo nieciągła (dyskretna, skokowa). Zmienną ciągłą będzie zatem „liczba oddanych głosów” (tak samo wzrost, kwota dochodu itp.), natomiast zmiennymi nieciągłymi są pozostałe zmienne (także np.: wykształcenie, ocena, stanowisko).

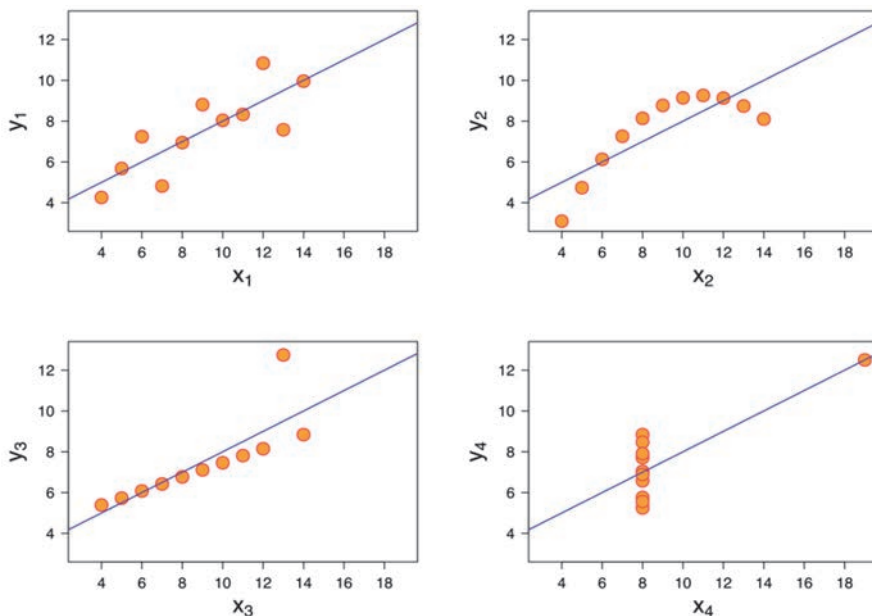
Zmienne z kolei będą Ci narzucać dobór skali, na jakiej będziesz pracować. Ogólnie rzecz biorąc, są cztery skale statystyczne: nominalna, porządkowa, interwałowa i ilorazowa. Często skale te dzieli się na nominalną, porządkową i skalę ilościową:

- **Skala nominalna.** W przypadku tej skali obiekty różnią się od siebie i nie są w żaden sposób ze sobą powiązane. Jest to np. lista miast w Polsce czy lista posłów na Sejm (poseł z jednej partii różni się od posła z innej partii). Nie jesteśmy w stanie oszacować tej różnicy, np. stwierdzić matematycznie, że poseł partii X jest lepszy od posła partii Y.
- **Skala porządkowa.** W tej skali obiekty różnią się od siebie i nie jesteśmy w stanie oszacować tej różnicy (że poseł partii X jest lepszy od posła partii Y), jednak jesteśmy w stanie uszeregować odpowiednie zmienne (np.: minister, wiceminister, podsekretarz stanu).

- **Skale ilościowe.** Te skale dotyczą zmiennych, które różnią się od siebie, można je uszeregować, ale także można je ocenić, a dokładnie ich różnicę (np. ile razy „za” w obecnej kadencji głosowała posłanka X w porównaniu do posłanki Y).

Widząc po raz pierwszy zbiór danych, warto mu się przyjrzeć pod kątem tego, jakiego rodzaju zmienne zawiera. Czy są uporządkowane, czy nie (czy skala nominalna czy porządkowa), a może da się je uporządkować i dokładnie porównać (skale ilościowe)?

Zrozumienie, z jakimi danymi masz do czynienia, ułatwi Ci dobór odpowiedniego wykresu. W przypadku skal ilościowych na początek należy pochylić się nad wykresami liniowymi (trendy z datami), wszelkiego rodzaju histogramami czy wykresami pudełkowymi. Jeżeli chodzi o skale nominalną i porządkową, najpierw warto sprawdzić wykresy kołowe i słupkowe. Dzięki tej wiedzy i praktyce będziesz w stanie bardzo szybko dobrać odpowiedni wykres do tabeli i sprawdzić pierwsze zależności i informacje, które mogą się za nimi kryć. Na takim etapie, zanim zaczniesz dokonywać analizy, możesz już spróbować zwizualizować swój zbiór na prostym wykresie zgodnie z przyjętą skalą. Możesz także odwrócić ten proces i rozpocząć analizę od szybkich prostych wizualizacji. Eksploracja danych poprzez wizualizację to w końcu odkrywanie czegoś, czego nie zauważyłeś na pierwszy rzut oka, warto więc wykorzystać szybkie eksploracje. Jeden z pierwszych przykładów znajdziesz na rysunku 5.5.



Rysunek 5.5. Francis Anscombe, twórca tych wizualizacji, stworzył cztery grupy danych, które miały te same cechy statystyczne, tj. średnią, wariancję, a nawet współczynnik korelacji. Jednak po zwizualizowaniu każdej z nich osobno okazało się, że znacząco się różnią na wykresie. Eksperyment ten pokazał po raz pierwszy, jak ważna jest wizualna analiza danych.

Źródło: Wikipedia

PROGRAM PARTNERSKI

— GRUPY HELION —

1. ZAREJESTRUJ SIĘ
2. PREZENTUJ KSIĄŻKI
3. ZBIERAJ PROWIZJĘ

Zmień swoją stronę WWW w działający bankomat!

Dowiedz się więcej i dołącz już dzisiaj!

<http://program-partnerski.helion.pl>

GRUPA
Helion 

BEZ DANYCH JESTEŚ JEDYNIĘ KOLEJNĄ OSOBĄ Z OPINIĄ...

Dziennikarstwo danych przeżywa dziś prawdziwy rozkwit. Dzieje się tak dlatego, że nasze życie w dużej mierze przeniosło się do internetu, a internet to... dane. Megabajty, gigabajty, terabajty danych. Misją współczesnego dziennikarza jest przedstawiać je społeczeństwu rzetelnie, a równocześnie atrakcyjnie, czyli w sposób zrozumiały, łatwy do przyswojenia. Nim się jednak owe dane pięknie zestawi, trzeba je znaleźć. Gdzie szukać? Jak je zdobyć? W jaki sposób opowiedzieć dane? Na takie pytania autor odpowiada w tej książce.

Nie przeczytasz w niej o „ładnych wykresach”, bo wbrew pozorom to nie one są esencją dziennikarstwa danych i data storytellingu. Dowiesz się natomiast, gdzie biją źródła potrzebnych Ci informacji, jak je przetwarzać i analizować. Znajdziesz tu także wskazówki, w jaki sposób tworzyć dobre wizualizacje za pomocą prostych aplikacji dostępnych za darmo w internecie i jak kreować angażujące odbiorców *data stories*. Na koniec wejdziesz na wyższy poziom — nauczysz się prezentować dane z wykorzystaniem kodu programistycznego.

Przebiecie się przez gigabajty informacji, przetworzenie ich i stworzenie materiału, który tłumaczy odbiorcy rzeczywistość, jest dziś działaniem obciążonym ogromnym wysiłkiem i jeszcze większym ryzykiem. Dlatego, jeżeli chcemy mieć przynajmniej cień nadziei na dobrze wykonaną pracę, warto sięgnąć po książkę Łukasza Żyły. *Dziennikarstwo danych i data storytelling* to także pozycja dla osób doświadczonych w tym zawodzie. Powód jest oczywisty, technologia zmieniła dziennikarstwo i w pędzie żywiołu, którym ono jest, łatwo popaść w bezpieczną i przez to złudną rutynę, a wtedy jesteśmy o krok od poważnego błędu.

**Bartosz Kurek, były dziennikarz Polsatu,
obecnie menedżer ds. public affairs w Philip Morris**

Co wy tam tak naprawdę robicie? — to częste pytanie, kiedy mówię, że pracuję w dziale danych „Wyborczej”. Teraz, zamiast wchodzić w szczegóły, będę mógł zacząć odpowiedź od słów: „Jest taka książka, warto przeczytać...”, bo Łukasz w bardzo przystępny sposób tłumaczy, czym to się je. Kiedy czytałem tę książkę, wiele razy żałowałem, że czegoś takiego nie było, kiedy ja zaczynałem przygodę z danymi. Dzięki niej widzę, ile jeszcze powinienem się w tej dziedzinie nauczyć.

Dominik Uhlig, szef BIQdata.pl — działu danych „Gazety Wyborczej”

ŁUKASZ ŻYŁA — dziennikarz i specjalista od *data storytelling*, programista, magister prawa (specjalizuje się w prawie do informacji). Prezes Fundacji Media 3.0. Prowadzi zajęcia na kierunku informatyka społeczna w Akademii Górniczo-Hutniczej. Przeszkolił kilkudziesięciu dziennikarzy z takich mediów jak „Gazeta Wyborcza”, Wirtualna Polska, Polska Press. Prowadzi serwis *datablog.pl*. Był członkiem grupy roboczej do spraw otwartości danych przy Ministerstwie Cyfryzacji.

onepress



Księgarnia internetowa:
<http://onepress.pl>



HELION SA
ul. Kościuszki 1c, 44-100 Gliwice
tel.: 32 230 98 63
onepress@onepress.pl

książkiklasybusiness

ebook dostępny na:

ebookpoint

ISBN 978-83-283-8312-8



9 788328 383128

Cena: 49,90 zł