

WSTĘPNA ANALIZA DANYCH

KIEDY PO RAZ PIERWSZY SPOTYKAMY SIĘ Z NOWYM ZESTAWEM DANYCH, NASZYM ZADANIEM JEST OPIS PODSTAWOWYCH ICH CECH. GŁÓWNE CECHY DANYCH MÓWIĄ NAM O ZASADNICZYCH WŁASNOŚCIACH ZJAWISK LUB EKSPERYMENTU, KTÓRY BADAMY. PONADTO PRAWIE ZAWSZE POTRZEBNY JEST NAM SYNTETYCZNY OPIS DANYCH: BARDZO TRUDNO JEST NA PRZYKŁAD ANALIZOWAĆ "SUROWE" WYNIKI SPISU POWSZECHNEGO W POLSCE. KONIECZNE JEST DOKONANIE ODPOWIEDNIEGO ICH PRZEKSZTAŁCENIA I UPROSZCZENIA UMOŻLIWIĄJĄCEGO ANALIZĘ. PRZED WSZYSTKIM MUSIMY JEDNAK USTALIĆ, JAKI JEST TYP DANYCH. JEŚLI MAMY DO CZYNNIENIA Z LICZBAMI ODPOWIADAJĄCYMI WARTOŚCIOM MIERZONYCH WIELKOŚCI, JAK NA PRZYKŁAD W PRZYPADKU POMIARU TEMPERATURY PRZY GRUNCIE O GODZINIE ÓSMEJ RANO NA ŚNIEŻCE W KOLEJNYCH DNIACH LISTOPADA, TO MÓWIMY WTEDY O DANYCH ILOŚCIOWYCH. W PRZYPADKU, GDY REJESTRUJEMY CECHĘ JAKOŚCIOWĄ OBIEKTÓW, NA PRZYKŁAD PŁEĆ LUB TYP SCHORZENIA PACJENTÓW, MÓWIMY O DANYCH JAKOŚCIOWYCH. OCZYWIŚCIE, JEŚLI DLA JEDNEGO OBIEKTU DOKONUJEMY KILKU POMIARÓW, TO CZĘŚĆ Z NICH MOŻE BYĆ TYPU ILOŚCIOWEGO, A CZĘŚĆ JAKOŚCIOWEGO. MOŻEMY REJESTROWAĆ JEDNOCZEŚNIE WIEK PACJENTA (CECHA ILOŚCIOWA) I TO, CZY MA ON LUB NIE PROBLEMY ZE SNEM (CECHA JAKOŚCIOWA). OKREŚLENIE TYPU DANYCH JEST NIEZBĘDNE PRZED PRZYSTĄPIENIEM DO ICH WSTĘPNEJ ANALIZY.

GRAFICZNE PRZEDSTAWIENIE DANYCH

WYKRES ZAWIERA ZNACZNIE WIĘCEJ INFORMACJI NIŻ JEDEN, A NAWET KILKA WSKAŹNIKÓW LICZBOWYCH OBLICZONYCH NA PODSTAWIE DANYCH. CZĘSTO JEST TAK, ŻE WARTOŚĆ PEWNEGO WSKAŹNIKA ODPOWIADA DWÓM ZUPEŁNIE ROŻNYM WYKRESOM I DLATEGO OPIERANIE SIĘ WYŁĄCZNIE NA WARTOŚCI TEGO WSKAŹNIKA MOŻE BYĆ MYLĄCE. ZARAZEM, WYKRES TEŻ JEST PEWNĄ REDUKCJĄ INFORMACJI W STOSUNKU DO ORYGINALNYCH DANYCH, ALE JEST TO REDUKCJA BEZ PORÓWNIANIA MNIEJ DRASTYCZNA.

<http://en.wikipedia.org/wiki/Chart>

HISTOGRAMY

<http://en.wikipedia.org/wiki/Histogram>

OKREŚLANIE SZEROKOŚCI PRZEDZIAŁÓW (BINÓW, SŁUPKÓW) HISTOGRAMU

$$h = 2,64 \cdot \frac{IQR(x)}{n^{\frac{1}{3}}}$$

WYBÓR POCZĄTKU PIERWSZEGO PRZEDZIAŁU – NAJLEPSZĄ METODĄ JEST TAKIE DOBRANIE PIERWSZEGO PRZEDZIAŁU, ABY NAJMNIEJSZA WARTOŚĆ WYSTĘPUJĄCA W ZBIORZE BYŁA ŚRODKIEM PIERWSZEGO PRZEDZIAŁU.

PODSTAWOWE POJĘCIA STATYSTYKI

MODA (DOMINANTA) – WARTOŚĆ O NAJWIĘKSZYM PRAWDOPODOBIENSTWIE WYSTĄPIENIA. JEST TO WARTOŚĆ, KTÓRA WYSTĘPUJE NAJCZĘŚCIEJ W ZBIORZE DANYCH.

MEDIANA – WARTOŚĆ ŚRODKOWA (NIE ŚREDNIA!!) W ZBIORZE DANYCH. DRUGI KWARTYL.

KWARTYLE – WARTOŚCI ŚRODKOWE DLA 4 PRZEDZIAŁÓW W ZBIORZE.

ROZSTĘP MIĘDZYKWARTYLOWY (IQR) – RÓŻNICA WARTOŚCI POMIĘDZY PIERWSZYM I TRZECIM KWARTYLEM.

ŚREDNIA ARYTMETYCZNA

$$\bar{a} = \frac{\sum_{i=1}^n a_i}{n}$$

ŚREDNIA UCINANA

$$\bar{a} = \frac{\sum_{i=k}^n a_i}{n - k}$$

WARTOŚCI „k” ORAZ „n” USTALAMY ARBITRALNIE, ALE ZAZWYCZAJ SORTUJEMY ZBIÓR WG WARTOŚCI I OBCINAMY WARTOŚCI:

- MIN I MAX
- ODSTAJĄCE OD INNYCH (RÓWNA LICZEBNOŚĆ Z POCZĄTKU I Z KOŃCA)
- 25% PO OBU KOŃCACH

ŚREDNIA WINSOROWSKA

LICZYMY DOKŁADNIE TAK SAMO JAK ŚREDNIĄ ARYTMETYCZNĄ. RÓŻNICA POLEGA NA TYM, IŻ WŚRÓD POSORTOWANYCH ELEMENTÓW WYBIERAMY OKREŚLONĄ LICZBĘ ELEMENTÓW SKRAJNYCH (TYLE SAMO NA POCZĄTKU I KOŃCU) I ZAMIENIAMY ICH WARTOŚCI NA WARTOŚĆ MIN I MAX Z POZOSTAŁYCH ELEMENTÓW.

1. W STU KOLEJNYCH RZUTACH KOSTKĄ OTRZYMANO NASTĘPUJĄCE WYNIKI:

5 2 2 6 3 2 5 3 1 2 5 3 6 2 5 4 4 6 1 6 4 5 5 2 4 6 1 4 4 3 4 2 4 2 4 4 1 1 4 5 3 1 5 6 5 6 1 5 6 2 4
5 5 2 5 4 5 5 1 1 2 2 5 5 2 6 3 5 5 4 1 4 5 5 1 4 3 2 1 2 6 1 2 1 6 5 1 3 6 1 5 6 6 2 2 3 5 5 2 4.

- WYGENERUJ WYKRES ILOŚCIOWY;
- OKREŚL:
 - MODĘ
 - MEDIANĘ
 - KWARTYLE

2. REJESTRUJEMY WIEK 20 PRACOWNIKÓW ZGŁASZAJĄCYCH SIĘ NA OKRESOWE BADANIA W PEWNYM ZAKŁADZIE PRACY. ZAOBSERWOWANE WIELKOŚCI WYNOŚĄ (W LATACH):

36, 41, 33, 34, 38, 26, 33, 36, 30, 48, 39, 31, 38, 37, 22, 31, 25, 32.

- WYZNACZ IQR (ROZSTĘP MIĘDZYKWARTYLOWY);
- OKREŚL PRZEDZIAŁY I POCZĄTEK HISTOGRAMU;
- UTWÓRZ HISTOGRAM;
- OKREŚL:
 - MODĘ
 - MEDIANĘ
 - KWARTYLE
 - ŚREDNIĄ ARYTMETYCZNA

3. DLA PODANEGO ZBIORU DANYCH:

26,40 31,60 29,60 28,20 24,80 26,50 25,85 26,10 26,90 26,05 31,40 28,00 25,55 29,70
26,80 28,80 26,50 28,30 30,50 24,70 25,30 30,20 29,20 28,40 26,90 25,50 26,40 33,00
25,20 26,60 27,50 25,10 24,60 31,80 29,80 27,90 30,20 26,50 31,60 26,60 26,50 27,50
28,40 27,10 30,90 30,30 30,10 28,70 27,60 27,60 28,70 32,90 26,30 26,30 27,40 26,80
24,20 28,70 31,50 26,00 32,60 24,60

- OKREŚL:
 - MODĘ
 - MEDIANĘ
 - KWARTYLE
 - ŚREDNIĄ ARYTMETYCZNA
- WYZNACZ IQR (ROZSTĘP MIĘDZYKWARTYLOWY);
- OKREŚL PRZEDZIAŁY I POCZĄTEK HISTOGRAMU;
- UTWÓRZ HISTOGRAM;
- STWÓRZ WYKRES LINIOWY.

4. DLA ZAŁĄCZONEGO PLIKU Z DANymi:

- WCZYTAJ PLIK .TXT DO EXCELA ROBIĄC ODPOWIEDNIA KONWERSJĘ;
- SFORMATUJ KOMÓRKI TAK, ŻEBY „COŚ” WIDZIEĆ;
- SKOPIUJ DANE DO INNEJ KOLUMNY (WIERSZA) I POSORTUJ;
- POLICZ MIN I MAX, MEDIANĘ, WYZNACZ MODĘ, KWARTYLE, IQR;
- STWÓRZ WYKRES LINIOWY. TERAZ WIDAĆ MIN, MAX I MEDIANĘ.
- POLICZ ŚREDNIĄ
- OSZACOWAĆ LICZBĘ BINÓW I STWORZYĆ WIERSZ Z BINAMI DLA HISTOGRAMOWANIA