

**Eksploracja  
Danych**

(c) Marcin  
Sydow

Wstęp

Tabela

Zmienne

Statystyki  
Opisowe

Wskaźniki  
rozproszenia

Rozkład  
zmiennej

Rozkład  
normalny

Zależności  
zmiennych

Momenty  
centralne

Podsumowanie

# Eksploracja Danych

## Dane

(c) Marcin Sydow

# Zawartość wykładu

Eksploracja  
Danych

(c) Marcin  
Sydow

Wstęp

Tabela

Zmienne

Statystyki  
Opisowe

Wskaźniki  
rozproszenia

Rozkład  
zmiennnej

Rozkład  
normalny

Zależności  
zmiennych

Momenty  
centralne

Podsumowanie

- źródła i formaty danych
- tabela jako podstawowy format danych
- zmienne
- wskaźniki położenia
- wskaźniki rozproszenia
- rozkłady zmiennych
- podsumowanie

# Źródła danych

Eksploatacja  
Danych

(c) Marcin  
Sydow

Wstęp

Tabela

Zmienne

Statystyki  
Opisowe

Wskaźniki  
rozproszenia

Rozkład  
zmiennej

Rozkład  
normalny

Zależności  
zmiennych

Momenty  
centralne

Podsumowanie

- bazy danych
- hurtownie danych
- dokumenty i kolekcje dokumentów
  - tekst
  - html
  - xml
  - xls
- sieć WWW
- serwery WWW (np. logi)
- portale społecznościowe
  - fb
  - twitter, etc.
- pliki graficzne i multimedialne
- czujniki
- sieci sensorów
- usługi sieciowe, etc.

# Formaty danych

Ekploracja  
Danych

(c) Marcin  
Sydow

Wstęp

Tabela

Zmienne

Statystyki  
Opisowe

Wskaźniki  
rozproszenia  
Rozkład  
zmiennnej

Rozkład  
normalny

Zależności  
zmiennych

Momenty  
centralne

Podsumowanie

Możliwych jest bardzo wiele formatów danych, np:

- csv
- text
- xls
- xml
- pdf
- binarne
- grafowe
- arff

# Jakość danych źródłowych

Eksploatacja  
Danych

(c) Marcin  
Sydow

Wstęp

Tab e la

Zmienne

Statystyki  
Opisowe

Wskaźniki  
rozproszenia

Rozkład  
zmiennej

Rozkład  
normalny

Zależności  
zmiennych

Momenty  
centralne

Podsumowanie

Rola jakości danych wejściowych (reguła “GIGO”).

Jeśli dane będą nieodpowiednie (złej jakości, niereprezentatywne, etc.) to niezależnie od poprawnych procedur i modeli wyniki też będą niskiej jakości.

Jako pierwszy etap pracy z danymi niezbędna jest więc faza oceny przydatności/jakości danych.

# Ocena przydatności danych

Eksploatacja  
Danych

(c) Marcin  
Sydow

Wstęp

Tabela

Zmienne

Statystyki  
Opisowe

Wskaźniki  
rozproszenia  
Rozkład  
zmiennnej

Rozkład  
normalny

Zależności  
zmiennych

Momenty  
centralne

Podsumowanie

Na tym etapie nie należy jeszcze modyfikować danych źródłowych.

Należy sprawdzić m.in.:

- jakie informacje można odczytać z danych źródłowych?
- czy na podstawie tych danych można uzyskać odpowiedzi na pytania postawione w ramach eksperymentu?
- jakie potencjalne problemy mogą wystąpić podczas próby realizacji eksperymentu

# Tabela danych i inne formaty

Ekploracja  
Danych

(c) Marcin  
Sydow

Wstęp

**Tabela**

Zmienne

Statystyki  
Opisowe

Wskaźniki  
rozproszenia  
Rozkład  
zmiennej

Rozkład  
normalny

Zależności  
zmiennych

Momenty  
centralne

Podsumowanie

Podstawowym formatem danych w eksploracji danych jest format **tabeli**

wiersze: obserwacje (przypadki)

kolumny: zmienne (atrybuty)

(przykłady)

(R: format taki nazywany jest “tidy”)

# Zmienne

- jakościowe (kategoryczne)
  - porządkowe (np. rozmiar koszulki)
  - regularne (np. kolor auta)
- ilościowe (numeryczne)
  - ciągłe (np. ilość wydobytego gazu ziemnego)
  - dyskretne (np. liczba wyprodukowanych aut)

(przykłady)

Stałe mogą być usunięte (ale należy się upewnić, że są faktycznie stałe w populacji a nie tylko w próbie)

Wartości niepowtarzalne (o ile nie są kluczem), też mogą być usunięte.

Specjalny przypadek: zmienne monotoniczne (ale mogą np. oznaczać czas).

Ekploracja  
Danych

(c) Marcin  
Sydow

Wstęp

Tabela

Zmienne

Statystyki  
Opisowe

Wskaźniki  
rozproszenia

Rozkład  
zmiennej

Rozkład  
normalny

Zależności  
zmiennych

Momenty  
centralne

Podsumowanie



# Zmienne anachroniczne

Eksploatacja  
Danych

(c) Marcin  
Sydow

Wstęo

Tabela

Zmienne

Statystyki  
Opisowe

Wskaźniki  
rozproszenia

Rozkład  
zmiennej

Rozkład  
normalny

Zależności  
zmiennych

Momenty  
centralne

Podsumowanie

To takie, które zostały uzupełnione w późniejszym czasie, gdy znano już inne wydarzenia.

Pozornie wyglądają na bardzo wartościowe (silna korelacja ze zmienną przewidywaną), ale faktycznie są bezwartościowe.

(przykład)

# Brakujące dane

Eksploracja  
Danych

(c) Marcin  
Sydow

Wstęp

Tabela

Zmienne

Statystyki  
Opisowe

Wskaźniki  
rozproszenia

Rozkład  
zmiennnej

Rozkład  
normalny

Zależności  
zmiennych

Momenty  
centralne

Podsumowanie

Należy odróżnić dwie sytuacje:

- nieznane wartości (brak pomiaru/informacji)
  - nieistniejące wartości (wynika ze stanu faktycznego)
- 
- czy występują w danych wartości NULL?
  - czy zostały zidentyfikowane i zastąpione (np. flagą “missing”)
  - czy można uzupełnić braki?
  - czy brak niesie wartościową informację?

# Reprezentatywność danych

Ekploracja  
Danych

(c) Marcin  
Sydow

Wstęp

Tabela

Zmienne

Statystyki  
Opisowe

Wskaźniki  
rozproszenia

Rozkład  
zmiennnej

Rozkład  
normalny

Zależności  
zmiennych

Momenty  
centralne

Podsumowanie

Generalnie im większa próbka, tym lepsza reprezentacja populacji, pod warunkiem, że próbe nie jest obciążona tzn. nie została wybrana przy pewnych specyficznych założeniach, które mogą ją warunkować.

Np. ankieta wykonana tylko przez telefony stacjonarne dotyczy specyficznej podpopulacji ludzi mających takie telefony.

Błędy pomiaru:

- błąd systematyczny (obciążenie)
- błąd przypadkowy (szum)

# Statystyki opisowe

Ekploracja  
Danych

(c) Marcin  
Sydow

Wstęp

Tabela

Zmienne

Statystyki  
Opisowe

Wskaźniki  
rozproszenia

Rozkład  
zmiennej

Rozkład  
normalny

Zależności  
zmiennych

Momenty  
centralne

Podsumowanie

Zmienną można podsumować pojedynczą statystyką. Istnieją różne rodzaje statystyk:

- statystyki pozycyjne (wskaźniki położenia)
- statystyki rozproszenia (wskaźniki rozproszenia)
- statystyki symetrii rozkładu (np. skośność)
- statystyki spłaszczenia rozkładu (np. kurtoza)

# Statystyki pozycyjne

Ekploracja  
Danych

(c) Marcin  
Sydow

Wstęp

Tabela

Zmienne

Statystyki  
Opisowe

Wskaźniki  
rozproszenia

Rozkład  
zmiennnej

Rozkład  
normalny

Zależności  
zmiennych

Momenty  
centralne

Podsumowanie

Statystyki, które podsumowują bezwzględne położenie wartości:

- moda (dominanta): najczęstsza wartość danej zmiennej
- średnia
- mediana (wartość środkowa)
- kwartyle (są trzy: mediana jest drugim kwartyłem)

# Wartość średnia w próbie (ang. mean)

Eksploracja  
Danych

(c) Marcin  
Sydow

Wstęp

Tabela

Zmienne

Statystyki  
Opisowe

Wskaźniki  
rozproszenia

Rozkład  
zmiennnej

Rozkład  
normalny

Zależności  
zmiennych

Momenty  
centralne

Podsumowanie

Zakładamy, że próba jest wektorem  $x = (x_1, \dots, x_n)$ :

$$\text{mean}(x) = \frac{1}{n} \sum_{i=1}^n x_i$$

(jaka jest złożoność czasowa obliczenia średniej?)

(R: `mean(x)`)

# Mediana (ang. median)

Mediana wyznacza wartość środkową w ciągu uporządkowanych wartości danej zmiennej. W przypadku parzystej liczności jest to średnia z 2 środkowych wartości.

Założmy, że próba została posortowana niemalejąco  
( $x_{(1)}, \dots, x_{(n)}$ ):

$$\text{median}(x) = x_{(\frac{n+1}{2})}$$

(n nieparzyste)

$$\text{median}(x) = \frac{(x_{(\frac{n}{2})} + x_{(\frac{n}{2}+1)})}{2}$$

(n parzyste)

(jaka jest złożoność czasowa obliczenia mediany?)

(R:  $\text{median}(x)$ )

# Średnia vs mediana

Ekploracja  
Danych

(c) Marcin  
Sydow

Wstęp

Tabela

Zmienne

Statystyki  
Opisowe

Wskaźniki  
rozproszenia  
Rozkład  
zmiennnej

Rozkład  
normalny

Zależności  
zmiennych

Momenty  
centralne

Podsumowanie

Średnia jest podstawową statystyką pozycyjną i łatwą w policzeniu, ale jest mniej niż mediana odporna na wartości odstające.

Przykład: gdy  $x$  oznacza pensję w danej firmie, a nieliczna grupa zarabia dużo więcej niż inni (np. prezesi, rektorzy, etc.) wartość średnia może dawać nieco zmylający obraz zarobków. Mediana jest bardziej odporna na nieliczne wartości odstające.

(w przypadku rozkładów nieunimodalnych zarówno wartość średnia jak i mediana mogą być mało przydatne - lepiej wtedy pogrupować dane i odrębnie podsumować oddzielne grupy)



# Średnia ucinana (ang. trimmed mean) i średnia winsorowska

Warianty średniej bardziej odporne na wartości odstające.

Średnia ucinana: (pomija  $k$  wartości najmniejszych i  $k$  największych)

$$\text{mean}_t(x) = \frac{1}{n - 2k} \sum_{i=k+1}^{n-k} x_{(i)}$$

(R: `mean(x, trim= $\alpha$ )`, gdzie  $\alpha$  to ułamek wartości do pominięcia)

Średnia winsorowska: (zastępuje te wartości najbliższym sąsiadem co do wartości)

$$\text{winsor}(x) = \frac{1}{n} [(k + 1)x_{(k+1)} + \sum_{i=k+2}^{n-k-1} x_{(i)} + (k + 1)x_{(n-k)}]$$

# Inne przydatne statystyki

Ekploracja  
Danych

(c) Marcin  
Sydow

Wstęp

Tabela

Zmienne

Statystyki  
Opisowe

Wskaźniki  
rozproszenia

Rozkład  
zmiennej

Rozkład  
normalny

Zależności  
zmiennych

Momenty  
centralne

Podsumowanie

- liczba istniejących wartości
- liczba unikatowych wartości
- liczba brakujących wartości

# Wskaźniki rozproszenia

Ekploracja  
Danych

(c) Marcin  
Sydow

Wstęp

Tabela

Zmienne

Statystyki  
Opisowe

**Wskaźniki  
rozproszenia**

Rozkład  
zmiennej

Rozkład  
normalny

Zależności  
zmiennych

Momenty  
centralne

Podsumowanie

Statystyki informujące o tym jak bardzo poszczególne wartości są rozproszone (zmienne):

- rozstęp próby
- wariancja i odchylenie standardowe
- rozstęp ćwiartkowy (IQR: inter-quartile range)

# Wariancja w próbie

Zakładamy, że próba jest wektorem  $x = (x_1, \dots, x_n)$ .

Wariancja:

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \text{mean}(x))^2$$

Jest miarą rozrzutu wartości danej zmiennej wokół wartości średniej.

Pewnym problemem technicznym jest fakt, że jednostką wariancji jest kwadrat jednostki badanej zmiennej.

Odchylenie standardowe:

$$s = \sqrt{s^2}$$

(odchylenie standardowe ma już te same jednostki co badana zmienna)

(R: `var(x)`, `sd(x)`)

Eksploatacja  
Danych

(c) Marcin  
Sydow

Wstęp

Tab e la

Zmienne

Statystyki  
Opisowe

Wskaźniki  
rozproszczenia

Rozkład  
zmiennej

Rozkład  
normalny

Zależności  
zmiennych

Momenty  
centralne

Podsumowanie

# Odchylenie przeciętne

Eksploatacja  
Danych

(c) Marcin  
Sydow

Wstęp

Tabela

Zmienne

Statystyki  
Opisowe

Wskaźniki  
rozproszczenia

Rozkład  
zmiennnej

Rozkład  
normalny

Zależności  
zmiennych

Momenty  
centralne

Podsumowanie

Obecność kwadratów we wzorze na wariancję powiększa wpływ wartości znacznie odbiegających od średniej (i jednocześnie pomniejsza tych będących b.blisko)

Odchylenie przeciętne nie ma takiej wady:

$$d_1 = \frac{1}{n} \sum_{i=1}^n |x_i - \text{mean}(x)|$$

(jednak jest to wyrażenie trudniejsze analitycznie)

# Kwartyle i IQR

Kwartyl może być rozumiany jako uogólnienie mediany.

Oznacza on wartość będącą w odpowiednio w:

- jednej czwartej (dolny kwartył) ( $Q_1$ )
- połowie (mediana)
- trzech czwartych (górny kwartył) ( $Q_3$ )

posortowanych niemalejąco wartości danej zmiennej.

Rozstęp ćwiartkowy (IQR):

$$IQR = Q_3 - Q_1$$

(miara rozproszenia)

Zauważmy, że pomiędzy dolnym i górnym kwartyłem znajduje się dokładnie połowa obserwowanych wartości.

Eksploatacja  
Danych

(c) Marcin  
Sydow

Wstęp

Tab e la

Zmienne

Statystyki  
Opisowe

Wskaźniki  
rozproszenia

Rozkład  
zmiennej

Rozkład  
normalny

Zależności  
zmiennych

Momenty  
centralne

Podsumowanie

# Wykres ramkowy

Jest to graficzna forma podsumowania kilku podstawowych wskaźników położenia i rozproszenia jednocześnie:

- Q1: dolny bok ramki
- mediana: pozioma linia w środku ramki
- Q3: górny bok ramki
- górny “wąs” czyli najwyższa obserwacja od mediany będąca nie dalej niż  $Q_3 + 1,5 \times IQR$
- dolny “wąs” czyli najniższa obserwacje od mediany będąca nie dalej niż  $Q_1 - 1,5 \times IQR$
- wartości odstające: zaznaczane jako kropki poza wąsami.

(R: `boxplot(x)`)

# Rozkład zmiennej (dyskretnej) i Histogram

Eksploatacja  
Danych

(c) Marcin  
Sydow

Wstęp

Tabela

Zmienne

Statystyki  
Opisowe

Wskaźniki  
rozproszenia

**Rozkład  
zmiennej**

Rozkład  
normalny

Zależności  
zmiennych

Momenty  
centralne

Podsumowanie

Jak często zdarza się dana wartość.

Histogram: Graficzne narzędzie do wizualizacji rozkładu.

(R: `hist(x)`)

Uwaga: bardzo ważnym parametrem histogramu jest szerokość koszyka.

(przykład)

Rozkład:

- jednomodalny (istnieje jedno maksimum)
- bimodalny (dwa maksima)
- wielo-modalny (wiele maksimumów)



# Gęstość rozkładu zmiennej ciągłej

Ekploracja  
Danych

(c) Marcin  
Sydow

Wstęp

Tabela

Zmienne

Statystyki  
Opisowe

Wskaźniki  
rozproszenia

**Rozkład  
zmiennej**

Rozkład  
normalny

Zależności  
zmiennych

Momenty  
centralne

Podsumowanie

Jest to ciągła funkcja nieujemna, która determinuje prawdopodobieństwo, że zmienna ta przybiera daną wartość.

- pole powierzchni pod wykresem gęstości wynosi 1.
- prawdopodobieństwo, że zmienna ciągła ma wartość pomiędzy  $a$  i  $b$  równa się polu powierzchni pod wykresem (całce  $z$ ) gęstości pomiędzy wartościami  $a$  i  $b$ .

Histogram można uważać za pewne przybliżenie wykresu funkcji gęstości.

# Dystrybuanta rozkładu zmiennej ciągłej

Eksploracja  
Danych

(c) Marcin  
Sydow

Wstęp

Tabela

Zmienne

Statystyki  
Opisowe

Wskaźniki  
rozproszenia

**Rozkład  
zmiennej**

Rozkład  
normalny

Zależności  
zmiennych

Momenty  
centralne

Podsumowanie

Jest to funkcja  $\Phi(a)$  informująca jakie jest prawdopodobieństwo, że zmienna o danym rozkładzie przyjmie wartości niewiększe niż  $a$  (pole pod krzywą gęstości od  $-\infty$  do  $a$ ).

# Rozkład Normalny (Gaussa)

Funkcja gęstości rozkładu normalnego  $N(\mu, \sigma)$  dana jest wzorem:

$$\phi_{\mu, \sigma} = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

( $\mu, \sigma$  to parametry będące odpowiednio wartością średnią rozkładu i jego odchyleniem standardowym)

Rozkład  $N(0, 1)$  nazywamy standardowym rozkładem normalnym.

Rozkład normalny ma pewne bardzo specjalne własności powodujące, że jego przybliżenia są często spotykane (dzięki centralnemu twierdzeniu granicznemu).

(R: `dnorm`, example: `plot(dnorm(seq(-5, 5, 0.01)))`)

(R: `rnorm`, `qnorm`, `pnorm`: generator losowy, kwantyl, dystrybuanta)

# Kwantyle

Eksploracja  
Danych

(c) Marcin  
Sydow

Wstęp

Tabela

Zmienne

Statystyki  
Opisowe

Wskaźniki  
rozproszenia

Rozkład  
zmiennej

**Rozkład  
normalny**

Zależności  
zmiennych

Momenty  
centralne

Podsumowanie

Kwantyl to uogólnienie kwartyła na przypadek ciągły.  
Kwantyl rzędu  $\alpha \in \{0, 1\}$ .  
(R: `quantile(x)`)

# Przykładowe znane rozkłady

Ekploracja  
Danych

(c) Marcin  
Sydow

Wstęp

Tabela

Zmienne

Statystyki  
Opisowe

Wskaźniki  
rozproszenia

Rozkład  
zmiennej

Rozkład  
normalny

Zależności  
zmiennych

Momenty  
centralne

Podsumowanie

## Rozkłady zmiennej:

- jednostajny  $U(0, 1)$  (każda wartość na przedziale  $[0, 1]$  jest jednakowo prawdopodobna)
- dwupunktowy (są tylko 2 możliwe wartości z prawdopodobieństwem  $p$  i  $1-p$ )
- dwumianowy (liczba "1" w  $n$  próbach dwupunktowego)
- Poissona ( $p(k; \lambda) = \frac{e^{-\lambda} \lambda^k}{k!}$ ) (R: `plot(dpois(seq(0,10),2))`)

# Standaryzacja zmiennej

Ekploracja  
Danych

(c) Marcin  
Sydow

Wstęp

Tabela

Zmienne

Statystyki  
Opisowe

Wskaźniki  
rozproszenia

Rozkład  
zmiennej

Rozkład  
normalny

Zależności  
zmiennych

Momenty  
centralne

Podsumowanie

Przekształcenie liniowe zmiennej  $x$  tak, aby miała wartość oczekiwaną 0 i odchylenie standardowe 1:

$$\frac{x - \mathit{mean}(x)}{\mathit{sd}(x)}$$

(R: `sx <- (x-mean(x))/sd(x)`)  
sprawdzenie: `mean(sx), sd(sx)`)

# Współczynnik korelacji liniowej (Pearsona) dwóch zmiennych

Jest to jedna z miar współzależności dwóch zmiennych  $x$  i  $y$ :

$$r = \frac{1}{n-1} \sum_{i=1}^n \left( \frac{x_i - \text{mean}(x)}{sd(x)} \right) \left( \frac{y_i - \text{mean}(y)}{sd(y)} \right)$$

Przyjmuje wartości od -1 (pełna negatywna zależność liniowa) przez 0 (brak zależności) do 1 (pełna pozytywna zależność liniowa)

(nieustandaryzowany wariant nazywany jest kowariancją)

(R: `cor(x,y)`)

Istnieją też inne warianty tej miary (nieliniowe) np: kendalla, spearmana.

(R: `cor(x, y, method = "kendall")`, `cor(x,y,method = "spearman")`)

# Korelacje między zmiennymi

Ekploracja  
Danych

(c) Marcin  
Sydow

Wstęp

Tabela

Zmienne

Statystyki  
Opisowe

Wskaźniki  
rozproszenia

Rozkład  
zmiennnej

Rozkład  
normalny

Zależności  
zmiennych

Momenty  
centralne

Podsumowanie

Korelacja nie oznacza przyczynowości

Przykłady przypadkowych korelacji:

<http://www.tylervigen.com/spurious-correlations>



# Momenty centralne

Ekploracja  
Danych

(c) Marcin  
Sydow

Wstęp

Tabela

Zmienne

Statystyki  
Opisowe

Wskaźniki  
rozproszenia

Rozkład  
zmiennnej

Rozkład  
normalny

Zależności  
zmiennych

Momenty  
centralne

Podsumowanie

Moment centralny rzędu  $k$  zmiennej  $x$ :

$$\mu_k = \frac{1}{n} \sum_{i=1}^n (x_i - \text{mean}(x))^k$$

(zauważmy, że np. dla  $k=2$ ,  $\mu_k$  to wariancja<sup>1</sup>)

Momenty wyższych rzędów używane są do innych miar rozkładu.

---

<sup>1</sup>Uwaga: w przypadku nieobciążonej estymacji momentów z próbki wzory są nieco zmodyfikowane, ale idea jest b.podobna

# Skośność

Ekploracja  
Danych

(c) Marcin  
Sydow

Wstęp

Tabela

Zmienne

Statystyki  
Opisowe

Wskaźniki  
rozproszenia

Rozkład  
zmiennnej

Rozkład  
normalny

Zależności  
zmiennych

Momenty  
centralne

Podsumowanie

Jest to miara symetrii rozkładu.

$$A = \frac{\mu_3}{\sigma^3}$$

- rozkład symetryczny (np. Gaussa)
- rozkład lewoskośny (ujemna skośność): dłuższy lewy ogon
- rozkład prawoskośny (dodatnia skośność): dłuższy prawy ogon

Uwaga: nieobciążony estymator skośności z próbki ma nieco bardziej skomplikowany wzór, ale podobna idea.

Miara zagęszczenia rozkładu wokół wartości centralnej.

$$Kurt(x) = \frac{\mu_4}{\sigma^4}$$

- wartość zero (np. rozkład Gaussa)
- wartość ujemna (platykurtoza): spłaszczony
- wartość dodatnia (leptokurtoza): wyostrzony

Uwaga: nieobciążony estymator kurtozy z próbki ma nieco bardziej skomplikowany wzór, ale podobna idea.

# Miara zróżnicowania dyskretnej zmiennej: Entropia informacji

Eksploatacja  
Danych

(c) Marcin  
Sydow

Wstęp

Tab e la

Zmienne

Statystyki  
Opisowe

Wskaźniki  
rozproszczenia

Rozkład  
zmiennej

Rozkład  
normalny

Zależności  
zmiennych

Momenty  
centralne

Podsumowanie

Im wyższa entropia tym wyższe zróżnicowanie.

Mając dany rozkład  $P(x = k)$  zmiennej  $x$ , jej entropia dana jest wzorem:

$$H(x) = - \sum_k P(x = k) \log_2(P(x = k))$$

W termodynamice entropia (upraszczając) jest miarą chaosu (im wyższe uporządkowanie tym niższy chaos)

W teorii informacji jest ona miarą wartości informacji.  
(R:  $\text{entropy}(c(0.5, 0.5))$ )

# Podsumowanie

Ekploracja  
Danych

(c) Marcin  
Sydow

Wstęp

Tabela

Zmienne

Statystyki  
Opisowe

Wskaźniki  
rozproszenia

Rozkład  
zmiennnej

Rozkład  
normalny

Zależności  
zmiennych

Momenty  
centralne

Podsumowanie

- źródła i formaty danych
- tabela jako podstawowy format danych
- zmienne
- wskaźniki położenia
- wskaźniki rozproszenia
- rozkłady zmiennych

# Przykładowe pytania/zadania/problemy

- rola fazy oceny przydatności danych. Jakie pytania można w tej fazie postawić?
- opisz format tabelaryczny. Co reprezentują wiersze, co reprezentują kolumny?
- wymień rodzaje zmiennych i dla podanej tabeli dokonaj klasyfikacji zmiennych
- wymienić, podać wzory i umieć obliczyć wskaźniki położenia zmiennej dla niewielkich próbek
- j.w. dla wskaźników rozproszenia
- co to jest rozkład zmiennej
- wzór i podstawowe własności gęstości rozkładu normalnego
- standaryzacja zmiennej
- momenty centralne i ich zastosowanie w podsumowaniach zmiennych

**Ekploracja  
Danych**

(c) Marcin  
Sydow

Wstęp

Tab e la

Zmienne

Statystyki  
Opisowe

Wskaźniki  
rozproszczenia

Rozkład  
zmiennej

Rozkład  
normalny

Zależności  
zmiennych

Momenty  
centralne

**Podsumowanie**

Dziękuję za uwagę.