

# Statystyka i Analiza Danych

## W2: Zmienna losowa

dr hab. inż. Katarzyna Filipiak, prof. PP

Instytut Matematyki  
Politechnika Poznańska

2023/2024

# Zmienna losowa

**Zmienną losową**  $X$  nazywamy funkcję  $X = X(\omega)$  określoną na przestrzeni zdarzeń elementarnych  $\Omega$  o wartościach w zbiorze liczb rzeczywistych  $\mathbb{R}$

Niech zmienna losowa **dyskretna** (skokowa)  $X$  przyjmuje wartości  $x_1, x_2, \dots$  odpowiednio z prawdopodobieństwami  $p_1, p_2, \dots$ ;  $\sum_{i=1}^{\infty} p_i = 1$ .

**Zmienna losowa ciągła** – zmienna przyjmująca wszystkie wartości z pewnego przedziału liczbowego.

# Zmienna losowa dyskretna

# Rozkład prawdopodobieństwa

**Rozkładem prawdopodobieństwa** zmiennej losowej dyskretnej  $X$ , nazywamy funkcję przyporządkowującą wartościom zmiennej  $x_i$  ( $i = 1, 2, \dots$ ) prawdopodobieństwa ich przyjęcia:

$$P(X = x_i) = P(\{\omega : X(\omega) = x_i\}) = p_i.$$

# Przykład

W pewnym eksperymencie wykorzystano **trzy** automatyczne aparaty fotograficzne w celu dokumentowania jego przebiegu. W danych warunkach prawdopodobieństwo wykonania poprawnej fotografii dla każdego aparatu jest takie samo i wynosi  $p = 0,6$ . Oblicz prawdopodobieństwo:

- (a) nieudokumentowania eksperymentu;
- (b) zarejestrowania eksperymentu przez co najmniej dwa aparaty.

rozkład prawdopodobieństwa (otrzymany za pomocą drzewa probabilistycznego):

$x_i$	0	1	2	3
$p_i$				

# Przykład - rozkład dwumianowy $\text{bin}(n, p)$

- dla pojedynczego aparatu możliwe są tylko dwa zdarzenia: zrobienie zdjęcia ('sukces') lub awaria ('porażka') – na drzewie probabilistycznym na każdym poziomie rysujemy tylko dwie gałęzie
- prawdopodobieństwo, że aparat zadziała (nastąpi awaria), nie zależy od zadziałania (awarii) pozostałych aparatów – na drzewie probabilistycznym na każdym poziomie prawdopodobieństwa, które wpisujemy na gałęziach, są takie same dla każdego sukcesu (0,6) i dla każdej porażki (0,4), odpowiednio
- najmniejszą wartością, jaką przyjmuje zmienna losowa  $X$  zliczająca aparaty, które zadziałały, jest 0, natomiast największą jest 3 - liczba wszystkich aparatów ( $X = 0, 1, 2, 3$ )

Wniosek:  $X \sim \text{bin}(n, p) \Rightarrow X \sim \text{bin}(3, 0,6)$

$$\begin{aligned} P(X = 0) &= \text{dbinom}(0, 3, 0.6) = \\ P(X \geq 2) &= P(X = 2) + P(X = 3) \\ &= \end{aligned}$$

# Dystrybuanta

Dystrybuantą zmiennej losowej  $X$  nazywamy funkcję  $F(X)$  określoną na zbiorze liczb rzeczywistych taką, że

$$F(x) = P(X \leq x).$$

Dystrybuanta zmiennej losowej dyskretnej:  $F(x) = \sum_{x_i \leq x} p_i$     `pname(x, param)`

Dla dowolnych liczb rzeczywistych  $a$  i  $b$  zachodzi:

$$P(X \leq a) = F(a)$$

$$P(X > a) = 1 - F(a)$$

$$P(a < X \leq b) = F(b) - F(a)$$

# Przykład

W pewnym eksperymencie wykorzystano **trzy** automatyczne aparaty fotograficzne w celu dokumentowania jego przebiegu. W danych warunkach prawdopodobieństwo wykonania poprawnej fotografii dla każdego aparatu jest takie samo i wynosi  $p = 0,6$ . Oblicz prawdopodobieństwo:

- (a) nieudokumentowania eksperymentu;
- (b) zarejestrowania eksperymentu przez co najmniej dwa aparaty.

$$X \sim \text{bin}(3, 0,6)$$

$$P(X = 0) = \text{dbinom}(0, 3, 0.6) =$$

$$P(X \geq 2) = P(X > 1) = 1 - F(1) = 1 - \text{pbinom}(1, 3, 0.6) =$$



# Wartość oczekiwana

Niech zmienna losowa dyskretna  $X$  przyjmuje wartości  $x_1, x_2, \dots$  odpowiednio z prawdopodobieństwami  $p_1, p_2, \dots$

## Definicja

**Wartością oczekiwaną** zmiennej losowej  $X$  nazywamy liczbę oznaczoną symbolem  $E(X)$  i określoną wzorem

$$E(X) = \sum_{i=1}^{\infty} x_i p_i$$

o ile nieskończony szereg  $\sum_{i=1}^{\infty} x_i p_i$  jest zbieżny.

# Własności $E(X)$

Niech  $a, b \in \mathbb{R}$ .

$$(1) \quad E(a) = a$$

$$(2) \quad E(bX) = b \cdot E(X)$$

$$(3) \quad E(X + a) = E(X) + a$$

(4) Niech  $g : \mathcal{D} \rightarrow \mathbb{R}$ , gdzie  $x_i \in \mathcal{D}$ ,  $i = 1, 2, \dots$

$$\text{Wówczas:} \quad E[g(X)] = \sum_{i=1}^{\infty} g(x_i) \cdot p_i$$

# Przykład

W pewnym eksperymencie wykorzystano trzy automatyczne aparaty fotograficzne w celu dokumentowania jego przebiegu. W danych warunkach prawdopodobieństwo wykonania poprawnej fotografii dla każdego aparatu jest takie samo i wynosi  $p = 0,6$ . Zdjęć zrobionych przez ile aparatów można się spodziewać (ile średnio aparatów udokumentuje przebieg eksperymentu)?

$x_i$	0	1	2	3	
$p_i$	0,064	0,288	0,432	0,216	1

$$E(X) =$$

# Wariancja

Niech  $\mu = E(X)$ .

## Definicja

**Wariancję** zmiennej losowej  $X$  nazywamy liczbę oznaczoną symbolem  $D^2(X)$  i określoną wzorem

$$D^2(X) = E[(X - \mu)^2] = E(X^2) - \mu^2.$$

Jeżeli  $X$  – zmienna losowa dyskretna:

$$D^2(X) = \sum_{i=1}^{\infty} x_i^2 \cdot p_i - \mu^2$$

# Własności $D^2(X)$

Niech  $a, b \in \mathbb{R}$ .

$$D^2(X) \geq 0 \quad (!ZAWSZE!)$$

$$D^2(X + a) = D^2(X)$$

$$D^2(bX) = b^2 \cdot D^2(X)$$

$$D^2(a + bX) = b^2 \cdot D^2(X)$$

## Definicja

**Odchyleniem standardowym** zmiennej losowej  $X$  nazywamy liczbę oznaczoną symbolem  $D(X)$  i określoną wzorem

$$D(X) = \sqrt{D^2(X)}.$$

# Przykład

W pewnym eksperymencie wykorzystano trzy automatyczne aparaty fotograficzne w celu dokumentowania jego przebiegu. W danych warunkach prawdopodobieństwo wykonania poprawnej fotografii dla każdego aparatu jest takie samo i wynosi  $p = 0,6$ . Oblicz odchylenie standardowe liczby aparatów dokumentujących przebieg eksperymentu.

$x_i$	0	1	2	3	
$p_i$	0,064	0,288	0,432	0,216	1

$$D^2(X) =$$

# Inne rozkłady dyskretne

- rozkład równomierny
- rozkład zero-jedynkowy (Bernoulliego)
- rozkład Poissona
- rozkład geometryczny
- rozkład Pascala
- rozkład hipergeometryczny
- ...

# Rozkłady dyskretne w R

*name* = nazwa rozkładu

*param* = parametry rozkładu

Gęstość:	<b>d</b> (density)	+	<i>name</i>	=	<b>d</b> <i>name</i> ( <i>x</i> , <i>param</i> )
Dystrybuanta:	<b>p</b> (probability)	+	<i>name</i>	=	<b>p</b> <i>name</i> ( <i>x</i> , <i>param</i> )
Kwantyl:	<b>q</b> (quantile)	+	<i>name</i>	=	<b>q</b> <i>name</i> ( $\alpha$ , <i>param</i> )
Losowa obserwacja:	<b>r</b> (random)	+	<i>name</i>	=	<b>r</b> <i>name</i> ( <i>N</i> , <i>param</i> )

Przykładowe nazwy rozkładów dyskretnych:

- dwumianowy: **binom**
- Poissona: **pois**

Histogram rozkładu dyskretnego (wykres liniowy):

```
plot(x, dname(x, param), type = "h")
```



# Zmienna losowa ciągła

# Funkcja gęstości prawdopodobieństwa

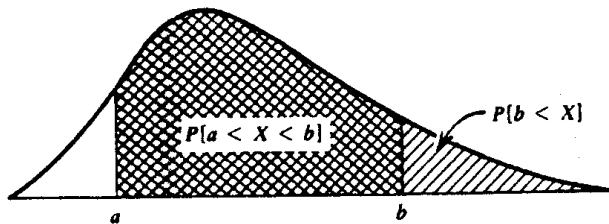
**Funkcja gęstości prawdopodobieństwa**  $f(x)$  – funkcja opisująca rozkład prawdopodobieństwa zmiennej losowej ciągłej posiadająca następujące cechy:

- a)  $f : \mathbb{R} \longrightarrow \mathbb{R}^+ \cup \{0\}$ ,
- b)  $P(a \leq X \leq b) = \text{pole pod krzywą } f(x) \text{ między } a \text{ i } b.$

Własności:

- $f(x) \geq 0$
- $P(a \leq X \leq b) = \int_a^b f(x) dx$
- $\int_{-\infty}^{\infty} f(x) dx = 1$
- $P(X > b) = 1 - P(X \leq b)$

# Funkcja gęstości



$$\begin{aligned}P(a < X \leq b) &= P(a < X < b) \\&= P(a \leq X \leq b) \\&= P(a \leq X < b)\end{aligned}$$

$$P(X = a) = 0$$

# Przykład

Czujnik śledzący stację wymaga dużej liczby wysokiej jakości taśm magnetycznych. Na taśmie magnetycznej mogą pojawić się rysy. Niech zmienna losowa  $X$  oznacza odległość (w cm) między kolejnymi rysami na powierzchni taśmy, a jej rozkład opisany jest funkcją gęstości

$$f(x) = \begin{cases} 0,01e^{-0,01x} & \text{dla } x \geq 0, \\ 0 & \text{dla } x < 0. \end{cases}$$

Założmy, że została znaleziona pierwsza rysa na taśmie. Oblicz prawdopodobieństwo, że kolejna zostanie znaleziona na kolejnych 50 cm taśmy.

$X$  - zmienna losowa oznaczająca dystans między rysami

$$\begin{aligned} P(X \leq 50) &= \int_0^{50} 0,01e^{-0,01x} dx = \text{integrate}(f, 0, 50) \\ f &= \text{function}(x)\{0.01 * \exp(-0.01 * x)\} \\ &= 0,3934693 \end{aligned}$$

# Dystrybuanta

Dystrybuantą zmiennej losowej  $X$  nazywamy funkcję  $F(X)$  określoną na zbiorze liczb rzeczywistych taką, że

$$F(x) = P(X \leq x).$$

Dystrybuanta zmiennej losowej ciągłej:  $F(x) = \int_{-\infty}^x f(t) dt$   $\text{pname}(x, \text{param})$

Dla dowolnych liczb rzeczywistych  $a$  i  $b$  zachodzi:

$$P(X \leq a) = F(a)$$

$$P(X > a) = 1 - F(a)$$

$$P(a < X \leq b) = F(b) - F(a)$$

# Przykład - rozkład wykładniczy $\text{Exp}(\lambda)$

Czujnik śledzący stację wymaga dużej liczby wysokiej jakości taśm magnetycznych. Na taśmie magnetycznej mogą pojawić się rysy. Niech zmienna losowa  $X$  oznacza odległość (w cm) między kolejnymi rysami na powierzchni taśmy, a jej rozkład opisany jest funkcją gęstości

$$f(x) = \begin{cases} 0,01e^{-0,01x} & \text{dla } x \geq 0, \\ 0 & \text{dla } x < 0. \end{cases}$$

Założmy, że została znaleziona pierwsza rysa na taśmie. Oblicz prawdopodobieństwo, że kolejna zostanie znaleziona na kolejnych 50 cm taśmy.

$X$  - zmienna losowa oznaczająca dystans między rysami

$$X \sim \text{Exp}(0,01)$$

$$P(X \leq 50) = F(50) = \text{pexp}(50, 0.01) =$$

# Wartość oczekiwana

Niech  $f(x)$  – funkcja gęstości zmiennej losowej ciągłej  $X$ .

## Definicja

**Wartością oczekiwaną** ciągłej zmiennej losowej  $X$  nazywamy

$$E(X) = \int_{-\infty}^{\infty} x \cdot f(x) dx$$

(o ile całka istnieje).

Własności: (1) – (3) jak dla zmiennej losowej dyskretnej

(4) Niech  $g : \mathcal{D} \rightarrow \mathbb{R}$ , gdzie  $x \in \mathcal{D}$ . Wówczas:

$$E[g(X)] = \int_{-\infty}^{\infty} g(x) \cdot f(x) dx.$$

# Przykład

Czujnik śledzący stację wymaga dużej liczby wysokiej jakości taśm magnetycznych. Na taśmie magnetycznej mogą pojawić się rysy. Niech zmienna losowa  $X$  oznacza odległość (w cm) między kolejnymi rysami na powierzchni taśmy, a jej rozkład opisany jest funkcją gęstości

$$f(x) = \begin{cases} 0,01e^{-0,01x} & \text{dla } x \geq 0, \\ 0 & \text{dla } x < 0. \end{cases}$$

Jaka przeciętnie odległość dzieli kolejne rysy na taśmie?

$$E(X) = \int_0^{\infty} x \cdot 0,01e^{-0,01x} dx = \text{integrate}(f, 0, \text{Inf})$$

$$f = \text{function}(x)\{x * 0.01 * \exp(-0.01 * x)\}$$

$$= 100$$

$$X \sim \text{Exp}(0,01) \quad \Rightarrow \quad E(X) = \frac{1}{\lambda} = \frac{1}{0,01} = 100$$



# Wariancja

Niech  $\mu = E(X)$ .

## Definicja

**Wariancję** zmiennej losowej  $X$  nazywamy liczbę oznaczoną symbolem  $D^2(X)$  i określoną wzorem

$$D^2(X) = E[(X - \mu)^2] = E(X^2) - \mu^2.$$

Jeżeli  $X$  – zmienna losowa ciągła:

$$D^2(X) = \int_{-\infty}^{\infty} x^2 \cdot f(x) dx - \mu^2$$

Własności: jak dla zmiennej losowej dyskretnej

# Rozkład normalny

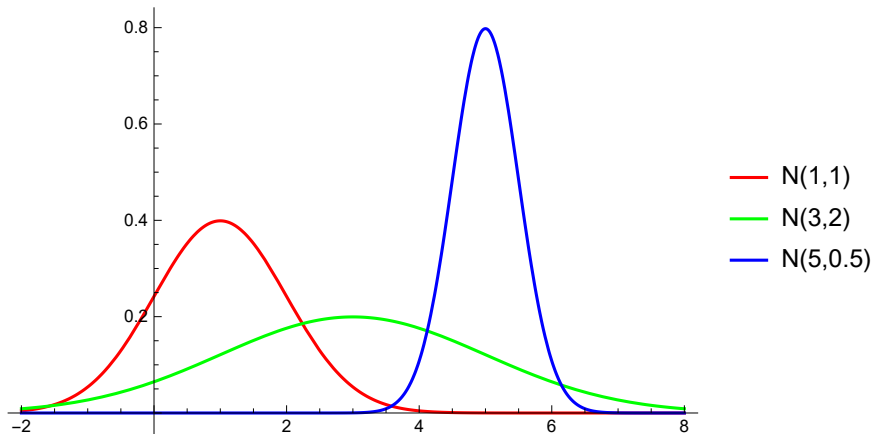
Funkcja gęstości ( $\mu \in \mathbb{R}$ ,  $\sigma \in \mathbb{R}^+$ ):

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

Charakterystyki liczbowe:

$$E(X) = \mu, \quad D^2(X) = \sigma^2$$

# Rozkład normalny



# Standardowy rozkład normalny $N(0, 1)$

$\phi(x)$  – dystrybuanta  $N(0, 1)$   $\diamond$

$$X \sim N(\mu, \sigma) \Rightarrow Z = \frac{X - \mu}{\sigma} \sim N(0, 1)$$

$$P(X \leq b) = P\left(\frac{X - \mu}{\sigma} \leq \frac{b - \mu}{\sigma}\right) = P\left(Z \leq \frac{b - \mu}{\sigma}\right) = \phi\left(\frac{b - \mu}{\sigma}\right)$$

$$P(X \leq b) = F(b) = \text{pnorm}(b, \mu, \sigma)$$

$$P(a \leq X \leq b) = P\left(\frac{a - \mu}{\sigma} \leq Z \leq \frac{b - \mu}{\sigma}\right) = \phi\left(\frac{b - \mu}{\sigma}\right) - \phi\left(\frac{a - \mu}{\sigma}\right)$$

$$P(a \leq X \leq b) = F(b) - F(a) = \text{pnorm}(b, \mu, \sigma) - \text{pnorm}(a, \mu, \sigma)$$

# Przykład

Niech  $X$  (w calach) będzie średnicą łożysk kulkowych produkowanych w pewnym zakładzie. Wiedząc, że  $X$  podlega rozkładowi normalnemu z wartością oczekiwaną 1 cal oraz odchyleniem standardowym 0,001 cala, tzn.  $X \sim N(1, 0,001)$ , oblicz prawdopodobieństwo, że średnica łożyska

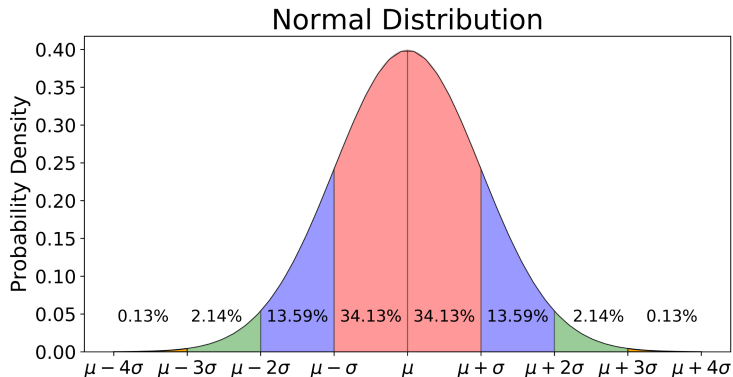
- (a) nie przekracza 1,0015 cala;
- (b) przekracza 0,9995 cala;
- (c) znajduje się w przedziale od 0,9998 do 1,0004 cala.

$$P(X < 1,0015) = F(1,0015) = \text{pnorm}(1.0015, 1, 0.001) =$$

$$P(X > 0,9995) =$$

$$P(0,9998 < X < 1,0004) =$$

# Reguła trzech sigm



# Reguła trzech sigm

Niech  $X \sim N(\mu, \sigma)$ . Wówczas

$$P(\mu - \sigma < X < \mu + \sigma) = 0.683$$

$$P(\mu - 2\sigma < X < \mu + 2\sigma) = 0.954$$

$$P(\mu - 3\sigma < X < \mu + 3\sigma) = 0.997$$

# Przybliżenie rozkładu dwumianowego

Niech  $X \sim \text{bin}(n, p)$ , gdzie  $n$  jest duże i  $p \in (0, 1)$ . Wiadomo, że

$$E(X) = n \cdot p, \quad D^2(X) = n \cdot p \cdot q.$$

Wówczas

$$X \underset{\text{app}}{\sim} N(n \cdot p, \sqrt{n \cdot p \cdot q})$$



# Przykład

Długoterminowe obszerne badania przeprowadzone w USA kilka lat temu wykazały, że 30% populacji dorosłych osób regularnie pije alkohol. Jeśli wyniki te są prawdziwe, to jakie jest prawdopodobieństwo, że w losowej próbie 1000 dorosłych liczba osób regularnie pijących alkohol jest nie większa niż 280?

$A$  - zmienna losowa zliczająca osoby regularnie pijące alkohol

$$A \sim \text{bin}(1000, 0,3)$$

$$A \underset{\text{app}}{\sim} N(1000 \cdot 0,3, \sqrt{1000 \cdot 0,3 \cdot 0,7})$$

$$P(A \leq 280) = F(280) =$$

$$P(A \leq 280) = F(280) \approx$$

# Inne rozkłady ciągłe

- rozkład jednostajny (prostokątny)
- rozkład trójkątny
- rozkład t-Studenta
- rozkład Beta
- rozkład Gamma
- rozkład chi-kwadrat
- ...

# Rozkłady ciągłe w R

*name* = nazwa rozkładu

*param* = parametry rozkładu

Gęstość: **d** (density) + *name* = **dname**(*x*, *param*)

Dystrybuanta: **p** (probability) + *name* = **pname**(*x*, *param*)

Kwantyl: **q** (quantile) + *name* = **qname**( $\alpha$ , *param*)

Losowa obserwacja: **r** (random) + *name* = **rname**(*N*, *param*)

Przykładowe nazwy rozkładów ciągłych:

wykładniczy: **exp**                      *t*-Studenta: **t**

normalny: **norm**                      chi-kwadrat: **chisq**

*F*-Snedecora: **f**

Wykres rozkładu ciągłego: **curve**(**dname**(*x*, *param*))

# Najważniejsze informacje

Rozkłady dyskretne:

- dwumianowy  $\text{bin}(n, p)$ :
  - prawdopodobieństwo:  $\text{dbinom}(\text{punkt}, n, p)$
  - wykres:  $\text{plot}(x, \text{dbinom}(x, n, p), \text{type} = "h")$
  - wartość oczekiwana i odch. std.:  $E(X) = n \cdot p$ ,  $SD(X) = \sqrt{n \cdot p \cdot (1 - p)}$

Rozkłady ciągłe:

- wykładniczy  $\text{EXP}(\lambda)$ :
  - prawdopodobieństwo (czerwony wzór):  $\text{pexp}(\text{punkt}, \lambda)$
  - wykres:  $\text{curve}(\text{dexp}(x, \lambda))$
  - wartość oczekiwana i odch. std.:  $E(X) = 1/\lambda$ ,  $SD(X) = 1/\lambda$
- normalny  $N(\mu, \sigma)$ :
  - prawdopodobieństwo (czerwony wzór):  $\text{pnorm}(\text{punkt}, \mu, \sigma)$
  - wykres:  $\text{curve}(\text{dnorm}(x, \mu, \sigma))$
  - wartość oczekiwana i odch. std.:  $E(X) = \mu$ ,  $SD(X) = \sigma$