

# Testy $\chi^2$

Wojciech Kotłowski

Statystyka i analiza danych 2019/2020

19.05.2020

## Testy $\chi^2$

- Dotyczą zmiennej/zmiennych **dyskretnych**, ze skończoną liczbą możliwych wartości:
  - Płeć, kolor, uporządkowane kategorie, narodowość, wynik rzutu kostką, itp.
- Nie testują jednego parametru rozkładu, ale **cały rozkład prawdopodobieństwa**.
- Tutaj poznamy dwie wersje: test rozkładu **jednej** zmiennej oraz test rozkładu **dwóch zmiennych**.

## Test dla jednej zmiennej

Dyskretna zmienna  $X$  przyjmująca jedną z wartości  $\{x_1, \dots, x_k\}$ .

- **Układ hipotez:**

---

$H_0$  : Zmienna  $X$  ma rozkład  $P$

$H_1$  : Zmienna  $X$  ma rozkład różny od  $P$

---

- **Tabela wartości obserwowanych (*observed*):**

$x_1$	$x_2$	$x_3$	$\dots$	$x_k$	$\Sigma$
$o_1$	$o_2$	$o_3$	$\dots$	$o_k$	$n$

- **Tabela wartości oczekiwanych (*expected*) z  $H_0$ :**

$x_1$	$x_2$	$x_3$	$\dots$	$x_k$	$\Sigma$
$e_1$	$e_2$	$e_3$	$\dots$	$e_k$	$n$

$$e_i = P(X = x_i) \cdot n$$

- **Statystyka testowa:**

$$\chi^2 = \sum_{i=1}^k \frac{(o_i - e_i)^2}{e_i} \sim \chi^2(k - 1)$$

Jeśli  $\chi^2 > \chi_{kr}^2$ , odrzucamy  $H_0$ .

## Test dla jednej zmiennej – przykład

$X$  – wynik rzutu kostką. Testujemy, czy kostka jest uczciwa

## Test dla jednej zmiennej – przykład

$X$  – wynik rzutu kostką. Testujemy, czy kostka jest uczciwa

- **Układ hipotez:**

---

$H_0$  : Zmienna  $X$  ma rozkład jednostajny na  $\{1, 2, 3, 4, 5, 6\}$

$H_1$  : Zmienna  $X$  nie ma rozkładu jednostajnego

---

## Test dla jednej zmiennej – przykład

$X$  – wynik rzutu kostką. Testujemy, czy kostka jest uczciwa

- **Układ hipotez:**

---

$H_0$  : Zmienna  $X$  ma rozkład jednostajny na  $\{1, 2, 3, 4, 5, 6\}$

$H_1$  : Zmienna  $X$  nie ma rozkładu jednostajnego

---

- **Tabela wartości obserwowanych** przy  $n = 30$  rzutach:

$X = 1$	$X = 2$	$X = 3$	$X = 4$	$X = 5$	$X = 6$	$\Sigma$
4	6	3	6	8	3	30

## Test dla jednej zmiennej – przykład

$X$  – wynik rzutu kostką. Testujemy, czy kostka jest uczciwa

- **Układ hipotez:**

---

$H_0$  : Zmienna  $X$  ma rozkład jednostajny na  $\{1, 2, 3, 4, 5, 6\}$

$H_1$  : Zmienna  $X$  nie ma rozkładu jednostajnego

---

- **Tabela wartości obserwowanych** przy  $n = 30$  rzutach:

$X = 1$	$X = 2$	$X = 3$	$X = 4$	$X = 5$	$X = 6$	$\Sigma$
4	6	3	6	8	3	30

- **Tabela wartości oczekiwanych z  $H_0$ :**

$X = 1$	$X = 2$	$X = 3$	$X = 4$	$X = 5$	$X = 6$	$\Sigma$
						30

## Test dla jednej zmiennej – przykład

$X$  – wynik rzutu kostką. Testujemy, czy kostka jest uczciwa

- **Układ hipotez:**

---

$H_0$  : Zmienna  $X$  ma rozkład jednostajny na  $\{1, 2, 3, 4, 5, 6\}$

$H_1$  : Zmienna  $X$  nie ma rozkładu jednostajnego

---

- **Tabela wartości obserwowanych** przy  $n = 30$  rzutach:

$X = 1$	$X = 2$	$X = 3$	$X = 4$	$X = 5$	$X = 6$	$\Sigma$
4	6	3	6	8	3	30

- **Tabela wartości oczekiwanych z  $H_0$ :**

$X = 1$	$X = 2$	$X = 3$	$X = 4$	$X = 5$	$X = 6$	$\Sigma$
5	5	5	5	5	5	30



## Test dla jednej zmiennej – przykład

$X$  – wynik rzutu kostką. Testujemy, czy kostka jest uczciwa

- **Układ hipotez:**

---

$H_0$  : Zmienna  $X$  ma rozkład jednostajny na  $\{1, 2, 3, 4, 5, 6\}$

$H_1$  : Zmienna  $X$  nie ma rozkładu jednostajnego

---

- **Tabela wartości obserwowanych** przy  $n = 30$  rzutach:

$X = 1$	$X = 2$	$X = 3$	$X = 4$	$X = 5$	$X = 6$	$\Sigma$
4	6	3	6	8	3	30

- **Tabela wartości oczekiwanych z  $H_0$ :**

$X = 1$	$X = 2$	$X = 3$	$X = 4$	$X = 5$	$X = 6$	$\Sigma$
5	5	5	5	5	5	30

- **Statystyka testowa:**

## Test dla jednej zmiennej – przykład

$X$  – wynik rzutu kostką. Testujemy, czy kostka jest uczciwa

- **Układ hipotez:**

---

$H_0$  : Zmienna  $X$  ma rozkład jednostajny na  $\{1, 2, 3, 4, 5, 6\}$

$H_1$  : Zmienna  $X$  nie ma rozkładu jednostajnego

---

- **Tabela wartości obserwowanych** przy  $n = 30$  rzutach:

$X = 1$	$X = 2$	$X = 3$	$X = 4$	$X = 5$	$X = 6$	$\Sigma$
4	6	3	6	8	3	30

- **Tabela wartości oczekiwanych z  $H_0$ :**

$X = 1$	$X = 2$	$X = 3$	$X = 4$	$X = 5$	$X = 6$	$\Sigma$
5	5	5	5	5	5	30

- **Statystyka testowa:**

$$\begin{aligned}\chi^2 &= \frac{(4-5)^2}{5} + \frac{(6-5)^2}{5} + \frac{(3-5)^2}{5} + \frac{(6-5)^2}{5} + \frac{(8-5)^2}{5} + \frac{(3-5)^2}{5} \\ &= \frac{1}{5} + \frac{1}{5} + \frac{4}{5} + \frac{1}{5} + \frac{9}{5} + \frac{4}{5} = \frac{20}{5} = 4\end{aligned}$$

## Test dla jednej zmiennej – przykład

$X$  – wynik rzutu kostką. Testujemy, czy kostka jest uczciwa

- **Układ hipotez:**

---

$H_0$  : Zmienna  $X$  ma rozkład jednostajny na  $\{1, 2, 3, 4, 5, 6\}$

$H_1$  : Zmienna  $X$  nie ma rozkładu jednostajnego

---

- **Tabela wartości obserwowanych** przy  $n = 30$  rzutach:

$X = 1$	$X = 2$	$X = 3$	$X = 4$	$X = 5$	$X = 6$	$\Sigma$
4	6	3	6	8	3	30

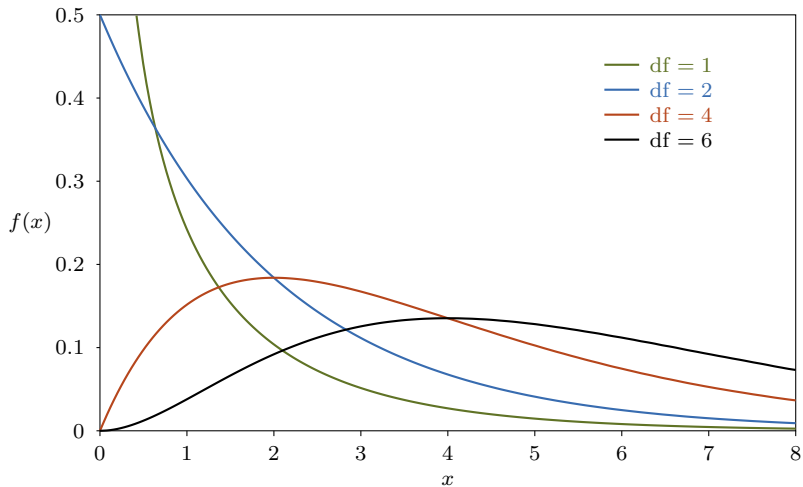
- **Tabela wartości oczekiwanych z  $H_0$ :**

$X = 1$	$X = 2$	$X = 3$	$X = 4$	$X = 5$	$X = 6$	$\Sigma$
5	5	5	5	5	5	30

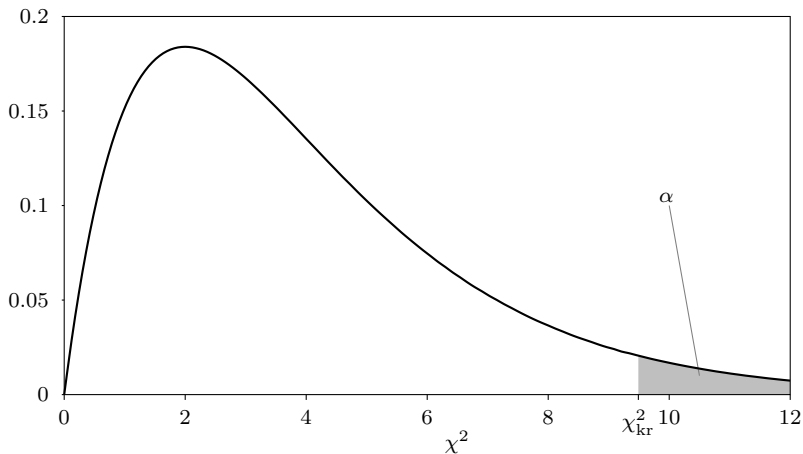
- **Statystyka testowa:**

- Wartość statystyki:  $\chi^2 = 4$
- Stopnie swobody:  $k - 1 = 5$
- Dla  $\alpha = 0.01$ ,  $\chi_{kr}^2 = 15.08$  (z tablic)
- **Wniosek:**  $\chi^2 < \chi_{kr}^2$ , więc brak podstaw do odrzucenia  $H_0$ .

# Rozkład $\chi^2(k)$



## Rozkład $\chi^2(4)$



Obszar krytyczny zawsze z prawej strony:  $C_{kr} = (\chi_{kr}^2, \infty)$ .

## Test dla dwóch zmiennych

$X \in \{x_1, \dots, x_w\}$  i  $Y \in \{y_1, \dots, y_k\}$ .

**Układ hipotez:**

---

$H_0$  : Zmienne  $X$  i  $Y$  są **niezależne**

$H_1$  : Zmienne  $X$  i  $Y$  są **zależne**

---

**Tabela w. obserwowanych**

	$y_1$	$y_2$	$\dots$	$y_k$	$\Sigma$
$x_1$	$o_{1,1}$	$o_{1,2}$	$\dots$	$o_{1,k}$	$W_1$
$x_2$	$o_{2,1}$	$o_{2,2}$	$\dots$	$o_{2,k}$	$W_2$
$\dots$	$\dots$	$\dots$	$\dots$	$\dots$	$\dots$
$x_w$	$o_{w,1}$	$o_{w,2}$	$\dots$	$o_{w,k}$	$W_w$
$\Sigma$	$K_1$	$K_2$	$\dots$	$K_k$	$n$

**Tabela w. oczekiwanych**

	$y_1$	$y_2$	$\dots$	$y_k$	$\Sigma$
$x_1$	$e_{1,1}$	$e_{1,2}$	$\dots$	$e_{1,k}$	$W_1$
$x_2$	$e_{2,1}$	$e_{2,2}$	$\dots$	$e_{2,k}$	$W_2$
$\dots$	$\dots$	$\dots$	$\dots$	$\dots$	$\dots$
$x_w$	$e_{w,1}$	$e_{w,2}$	$\dots$	$e_{w,k}$	$W_w$
$\Sigma$	$K_1$	$K_2$	$\dots$	$K_k$	$n$

Wartości oczekiwane:  $e_{ij} = \frac{W_i K_j}{n}$  (  $\frac{\text{suma wiersza} \times \text{suma kolumny}}{\text{podsumowanie tabeli}}$  )

**Statystyka testowa:**

$$\chi^2 = \sum_{i=1}^w \sum_{j=1}^k \frac{(o_{ij} - e_{ij})^2}{e_{ij}} \sim \chi^2((w-1) \cdot (k-1))$$

## Wartości oczekiwane

Skąd wzór  $e_{ij} = \frac{W_i K_j}{n}$ ?

## Wartości oczekiwane

$$\text{Skąd wzór } e_{ij} = \frac{W_i K_j}{n} ?$$

Spodziewamy się wystąpienia:

$$n \cdot P(X = x_i, Y = y_j)$$

obserwacji dla których  $X = x_i$  i  $Y = y_j$ .



## Wartości oczekiwane

$$\text{Skąd wzór } e_{ij} = \frac{W_i K_j}{n} ?$$

Spodziewamy się wystąpienia:

$$n \cdot P(X = x_i, Y = y_j)$$

obserwacji dla których  $X = x_i$  i  $Y = y_j$ .

Przy założeniu  $H_0$  zmienne są **niezależne**, a więc:

$$\begin{aligned} P(X = x_i, Y = y_j) &= P(X = x_i) \cdot P(Y = y_j) \\ &= \frac{W_i}{n} \cdot \frac{K_j}{n} \end{aligned}$$

Pomnożenie przez  $n$  to właśnie ten wzór.

## Test dla dwóch zmiennych – przykład

W USA przeprowadzono sondaż opinii na 1000 losowo wybranych osób. Sprawdź, czy istnieje zależność między płcią odpytanych osób a ich preferencjami politycznymi.

	republican	democrat	independent	$\Sigma$
male	200	150	50	400
female	250	300	50	600
$\Sigma$	450	450	100	1000

## Test dla dwóch zmiennych – przykład

$X$  – płeć,  $Y$  – preferencje wyborcze.

**Układ hipotez:**

## Test dla dwóch zmiennych – przykład

$X$  – płeć,  $Y$  – preferencje wyborcze.

### Układ hipotez:

---

$H_0$  : Brak zależności między płcią a pref. wyborczymi

$H_1$  : Istnieje zależność

---

## Test dla dwóch zmiennych – przykład

$X$  – płeć,  $Y$  – preferencje wyborcze.

### Układ hipotez:

---

$H_0$  : Brak zależności między płcią a pref. wyborczymi

$H_1$  : Istnieje zależność

---

### Tabela w. obserwowanych

	rep	dem	ind	$\Sigma$
M	200	150	50	400
F	250	300	50	600
$\Sigma$	450	450	100	1000

### Tabela w. oczekiwanych

	rep	dem	ind	$\Sigma$
M				400
F				600
$\Sigma$	450	450	100	1000

## Test dla dwóch zmiennych – przykład

$X$  – płeć,  $Y$  – preferencje wyborcze.

### Układ hipotez:

---

$H_0$  : Brak zależności między płcią a pref. wyborczymi

$H_1$  : Istnieje zależność

---

### Tabela w. obserwowanych

	rep	dem	ind	$\Sigma$
M	200	150	50	400
F	250	300	50	600
$\Sigma$	450	450	100	1000

### Tabela w. oczekiwanych

	rep	dem	ind	$\Sigma$
M	180			400
F				600
$\Sigma$	450	450	100	1000

## Test dla dwóch zmiennych – przykład

$X$  – płeć,  $Y$  – preferencje wyborcze.

### Układ hipotez:

---

$H_0$  : Brak zależności między płcią a pref. wyborczymi

$H_1$  : Istnieje zależność

---

### Tabela w. obserwowanych

	rep	dem	ind	$\Sigma$
M	200	150	50	400
F	250	300	50	600
$\Sigma$	450	450	100	1000

### Tabela w. oczekiwanych

	rep	dem	ind	$\Sigma$
M	180	180		400
F				600
$\Sigma$	450	450	100	1000

## Test dla dwóch zmiennych – przykład

$X$  – płeć,  $Y$  – preferencje wyborcze.

### Układ hipotez:

---

$H_0$  : Brak zależności między płcią a pref. wyborczymi

$H_1$  : Istnieje zależność

---

### Tabela w. obserwowanych

	rep	dem	ind	$\Sigma$
M	200	150	50	400
F	250	300	50	600
$\Sigma$	450	450	100	1000

### Tabela w. oczekiwanych

	rep	dem	ind	$\Sigma$
M	180	180	40	400
F				600
$\Sigma$	450	450	100	1000



## Test dla dwóch zmiennych – przykład

$X$  – płeć,  $Y$  – preferencje wyborcze.

### Układ hipotez:

---

$H_0$  : Brak zależności między płcią a pref. wyborczymi

$H_1$  : Istnieje zależność

---

### Tabela w. obserwowanych

	rep	dem	ind	$\Sigma$
M	200	150	50	400
F	250	300	50	600
$\Sigma$	450	450	100	1000

### Tabela w. oczekiwanych

	rep	dem	ind	$\Sigma$
M	180	180	40	400
F	270	270	60	600
$\Sigma$	450	450	100	1000

## Test dla dwóch zmiennych – przykład

$X$  – płeć,  $Y$  – preferencje wyborcze.

**Układ hipotez:**

---

$H_0$  : Brak zależności między płcią a pref. wyborczymi

$H_1$  : Istnieje zależność

---

**Tabela w. obserwowanych**

	rep	dem	ind	$\Sigma$
M	200	150	50	400
F	250	300	50	600
$\Sigma$	450	450	100	1000

**Tabela w. oczekiwanych**

	rep	dem	ind	$\Sigma$
M	180	180	40	400
F	270	270	60	600
$\Sigma$	450	450	100	1000

**Statystyka testowa:**

$$\begin{aligned}\chi^2 &= \frac{(200-180)^2}{180} + \frac{(150-180)^2}{180} + \frac{(50-40)^2}{40} + \frac{(250-270)^2}{270} + \frac{(300-270)^2}{270} \\ &+ \frac{(50-60)^2}{60} = \frac{20}{9} + 5 + \frac{5}{2} + \frac{40}{27} + \frac{10}{3} + \frac{5}{3} = 16.2\end{aligned}$$

## Test dla dwóch zmiennych – przykład

$X$  – płeć,  $Y$  – preferencje wyborcze.

### Układ hipotez:

---

$H_0$  : Brak zależności między płcią a pref. wyborczymi

$H_1$  : Istnieje zależność

---

### Tabela w. obserwowanych

	rep	dem	ind	$\Sigma$
M	200	150	50	400
F	250	300	50	600
$\Sigma$	450	450	100	1000

### Tabela w. oczekiwanych

	rep	dem	ind	$\Sigma$
M	180	180	40	400
F	270	270	60	600
$\Sigma$	450	450	100	1000

### Statystyka testowa:

- Wartość statystyki:  $\chi^2 = 16.2$
- Stopnie swobody:  $(w - 1)(k - 1) = 1 \cdot 2 = 2$
- Dla  $\alpha = 0.01$ ,  $\chi_{kr}^2 = 9.21$  (z tablic)
- **Wniosek:**  $\chi^2 > \chi_{kr}^2$ , więc odrzucamy  $H_0$ .