



# Komputerowa analiza danych doświadczalnych

Wykład podsumowujący  
22.05.2020

dr inż. Łukasz Graczykowski  
[lukasz.graczykowski@pw.edu.pl](mailto:lukasz.graczykowski@pw.edu.pl)

*Semestr letni 2019/2020*



# Parametry rozkładów prawdopodobieństwa

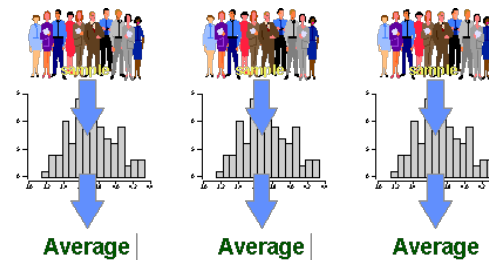
## Generowanie liczb i metody Monte Carlo

## Dopasowanie modelu do danych

## Testy statystyczne

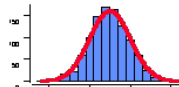
# Statystyczna analiza danych

- Statystyczna analiza danych:
  - traktujemy pomiar jako pewien element zbioru wszystkich możliwych pomiarów (pewnej cechy **populacji** o danym rozkładzie prawdopodobieństwa – najczęściej nieznanym)
  - na podstawie skończonej liczby pomiarów, obserwacji (**próby losowej**, podzbioru populacji), która ma swój rozkład prawdopodobieństwa (znany z pomiarów czy obserwacji), próbujemy dowiedzieć się czegoś (czyli **estymować**) na temat parametrów rozkładu całej populacji
  - innymi słowy, na podstawie próby losowej (pomiarów, obserwacji) stawiamy hipotezy i wyciągamy wnioski dotyczące interesującej nas cechy całej populacji



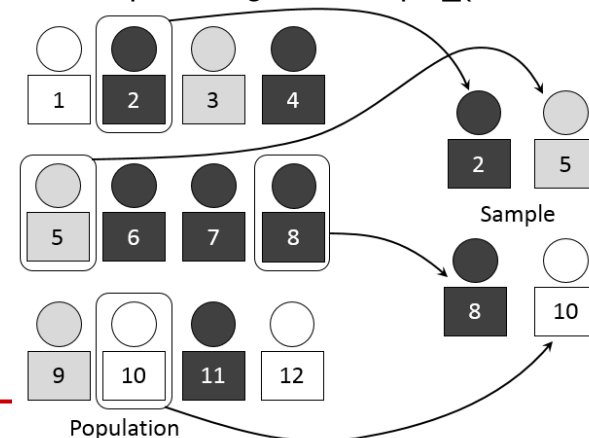
<https://www.proprofs.com/quiz-school/story.php?title=3d-q-sampling-distributions>

The Sampling Distribution...



...is the distribution of a statistic across an infinite number of samples

[https://en.wikipedia.org/wiki/Sample\\_\(statistics\)](https://en.wikipedia.org/wiki/Sample_(statistics))



# Rozkłady 1D - momenty

$$m_l = E(X^l) = \int_{-\infty}^{\infty} x^l f(x) dx \quad \text{momenty zwykłe}$$

$$m_1 = \mu = E(X) = \int_{-\infty}^{\infty} x f(x) dx$$

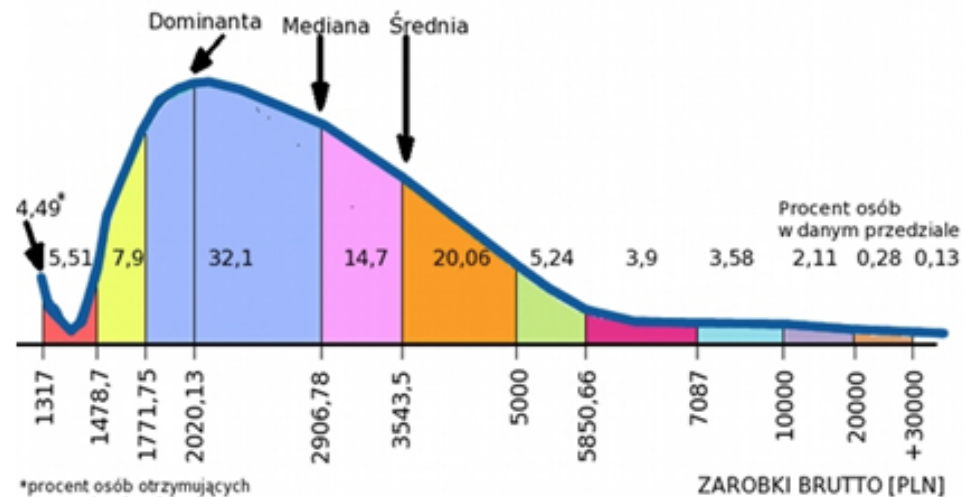
$$\mu_l = E((X - \hat{x})^l) \quad \text{momenty centralne}$$

$$\mu_2 \equiv \sigma^2(X) \equiv E((X - \hat{x})^2) = \int_{-\infty}^{\infty} (x - \hat{x})^2 f(x) dx$$

$$\gamma = \frac{\mu_3}{\sigma^3} \quad K = \frac{\mu_4}{\sigma^4} - 3 \quad \text{wsp. asymetrii, kurtoza}$$

## Procent osób zarabiających dane kwoty brutto

Na podstawie danych GUS za 2010 rok, like-a-geek.jogger.pl

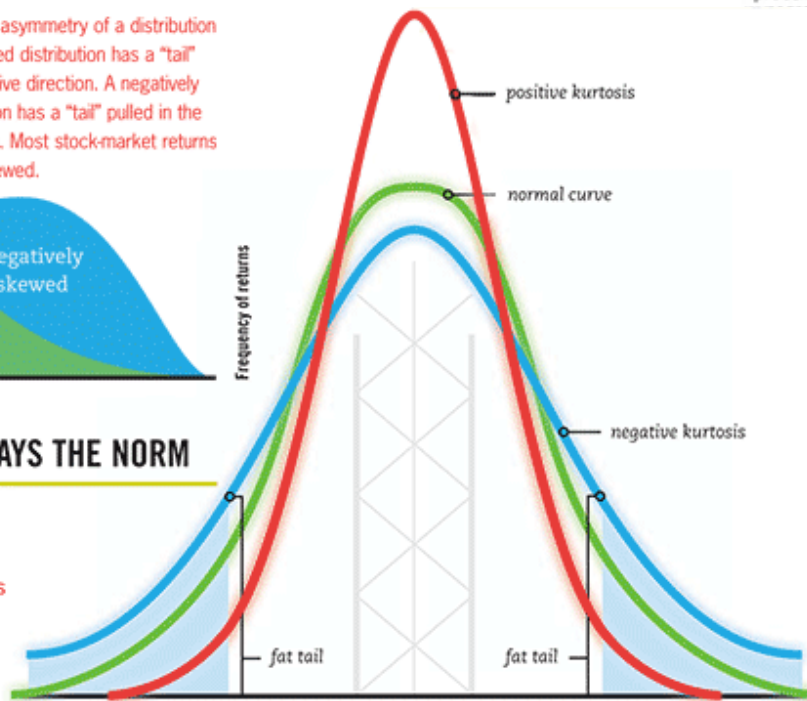


Skewness is the asymmetry of a distribution. A positively skewed distribution has a "tail" pulled in the positive direction. A negatively skewed distribution has a "tail" pulled in the negative direction. Most stock-market returns are negatively skewed.



### NORMAL NOT ALWAYS THE NORM

Kurtosis refers to how peaked the curve is: steeper means positive kurtosis and flatter means negative kurtosis. Fat tails occur when there are more outside returns on the downside or upside, or both, than the normal curve suggests.



[http://lh3.ggpht.com/-UhjcSGuME9Q/UgCqCj00\\_nI/AAAAAAAAAWXU/-OZlMA9pPnU/image\\_thumb%25255B2%25255D.png?imgmax=800](http://lh3.ggpht.com/-UhjcSGuME9Q/UgCqCj00_nI/AAAAAAAAAWXU/-OZlMA9pPnU/image_thumb%25255B2%25255D.png?imgmax=800)

dominanta

$$P(X = x_{max}) = \max$$

$$\frac{df(x)}{dx} = 0 \quad \frac{d^2 f(x)}{dx^2} < 0$$

**Momenty to uśrednienia danych podniesione do kolejnych potęg**

<http://www.advisor.ca/wp-content/uploads/2012/07/normal-not-always-the-norm.gif>

# Rozkłady 1D – kwantyle

- **Mediana** dzieli rozkład prawdopodobieństwa na dwa obszary o równym prawdopodobieństwie

$$F(x_{0,5}) = P(X < x_{0,5}) = 0,5$$

- Mediana  $x_{0,5}$  jest **kwantylem** (*ang. quantile*) rzędu 0,5

- Ogólna definicja **kwantylu rzędu  $q$** ,  $x_q$ :  $F(x_q) = P(X < x_q) = q$

- **kwartył dolny**  $x_{0,25}$

- **kwartył górny**  $x_{0,75}$

- **decyle**  $x_{0,1}, x_{0,2}, \dots, x_{0,9}$

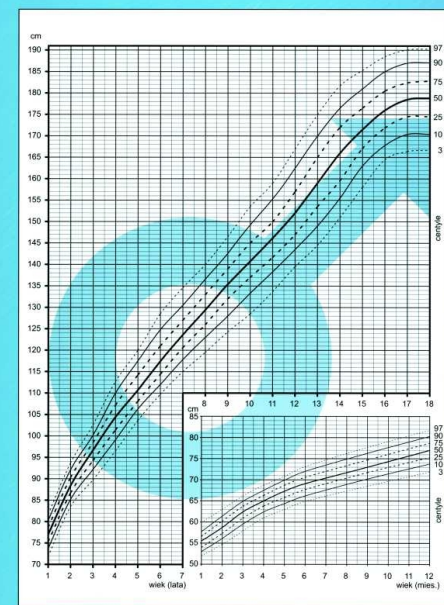
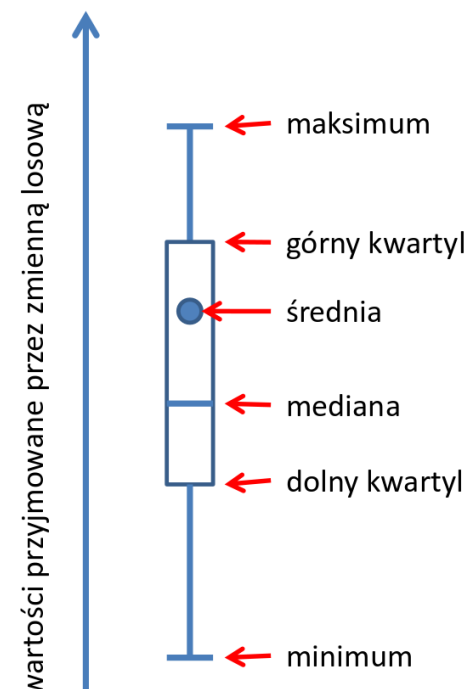
- funkcja  $x_q(q)$  jest funkcją odwrotną do dystrybuanty

$$F(x_q) = \int_{-\infty}^{x_q} f(x) dx = q, q \in \langle -1; 1 \rangle$$

- Kwantyl rzędu  $q$  jest taką liczbą  $x_q$ , że  $q \cdot 100\%$  elementów w danej próbkce (populacji) ma wartość pomiaru (badanej cechy) nie większą niż  $x_q$

- W przypadku ogonów rozkładu kwantyle mogą być lepszą wielkością niż momenty

- **Momenty i kwantyle to dwa najczęstsze opisy numeryczne danych liczbowych**



# Rozkłady 2D

## momenty

$$\lambda_{10} = E(x) = \hat{x}$$

$$\lambda_{01} = E(y) = \hat{y}$$

$$\mu_{11} = E((X - \hat{x})(Y - \hat{y})) = cov(X, Y)$$

$$\mu_{20} = E((X - \hat{x})^2) = \sigma^2(X)$$

$$\mu_{02} = E((Y - \hat{y})^2) = \sigma^2(Y)$$

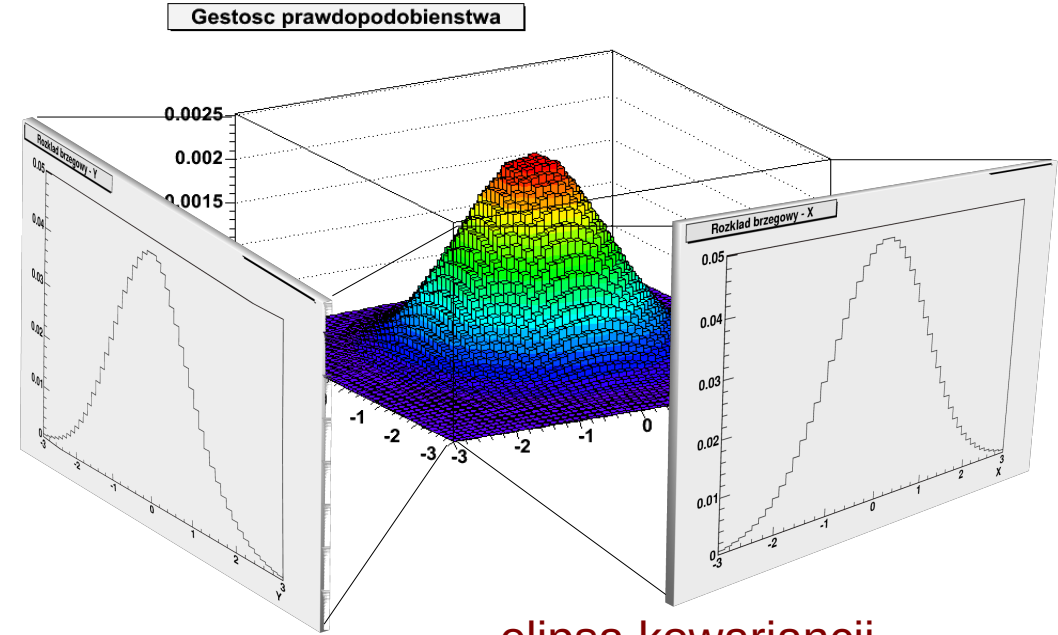
## kowariancja

$$cov(X, Y) = \mu_{11} = E(X \cdot Y) - E(X) \cdot E(Y)$$

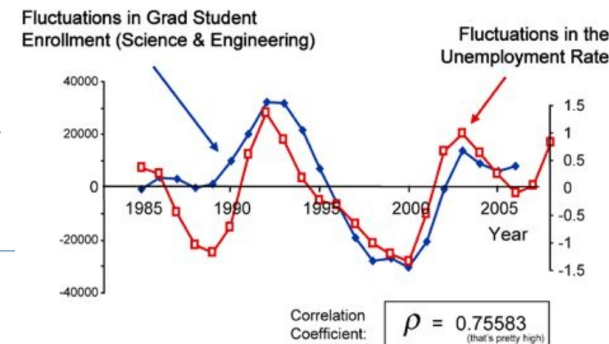
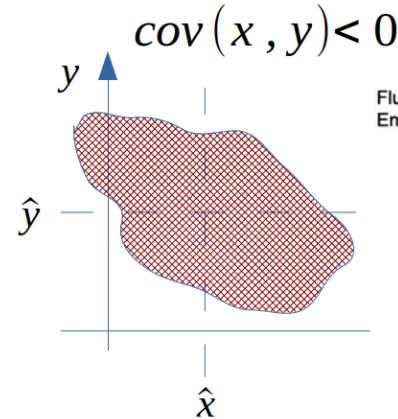
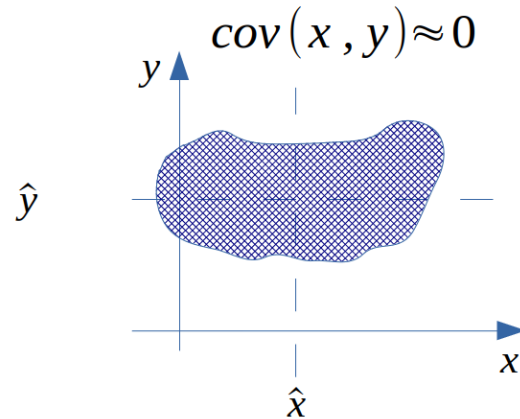
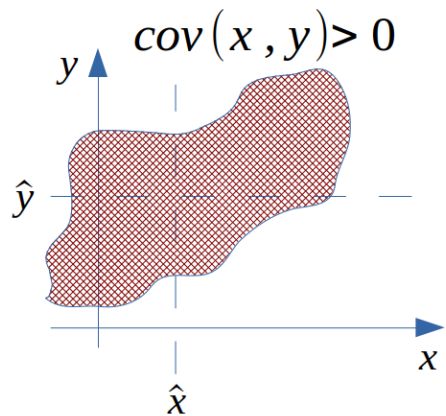
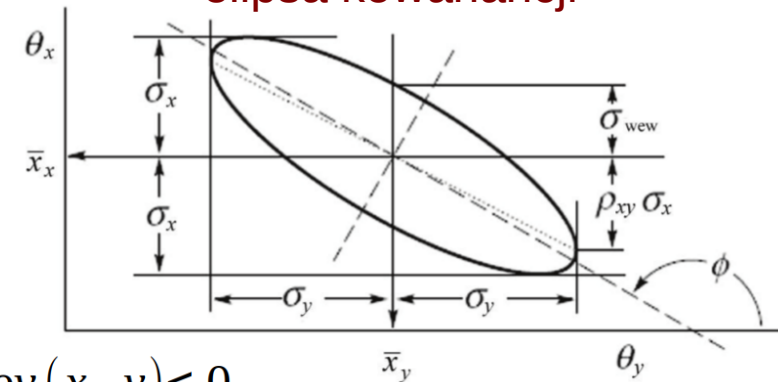
## wsp. korelacji (Pearsona)

$$\rho(X, Y) = \frac{cov(X, Y)}{\sigma(X)\sigma(Y)} \quad -1 \leq \rho(X, Y) \leq 1$$

Współczynnik korelacji a liniowa zależność zmiennych (niekoniecznie przyczyna-skutek)



## elipsa kowariancji



# Generatory liczb pseudolosowych

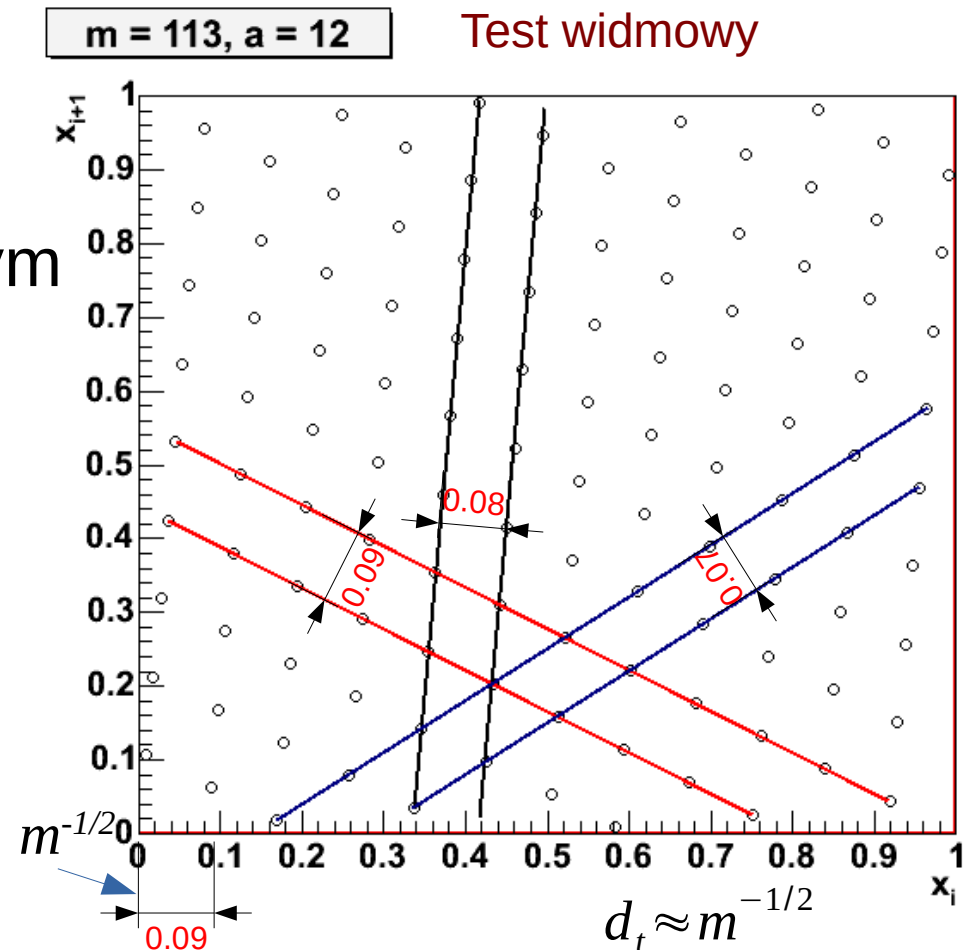
- Komputer, urządzenie deterministyczne, może generować tylko liczby pseudolosowe
  - kolejna generowana liczba jest funkcją liczb wcześniej wygenerowanych
  - *ziarno (seed)* – wartość początkowa (można ją ustalić np. z zegara systemowego)
- Generujemy liczby z jednakowym prawdopodobieństwem (rozkład jednorodny)

Liniowy kongruentny generator liniowy (LCG)

$$x_{j+1} = (a \cdot x_j + c) \bmod m$$

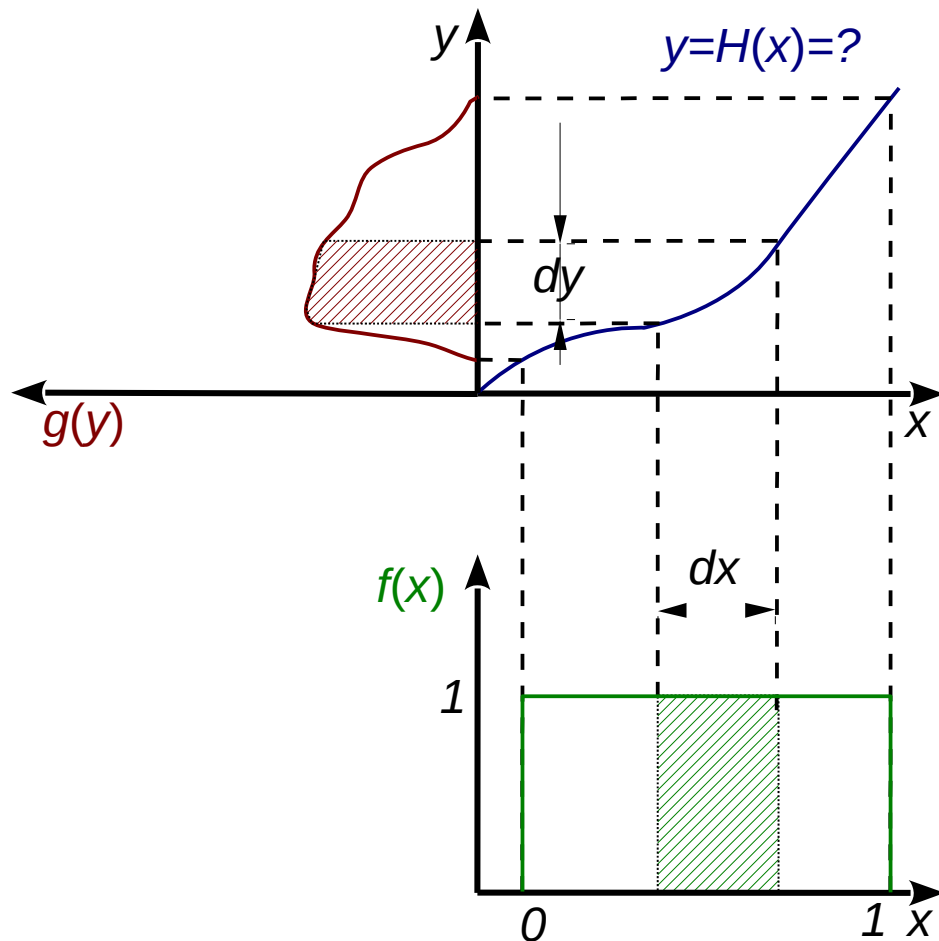
Multiplikacyjny liniowy kongruentny generator liniowy (MLCG)

$$x_{j+1} = (a \cdot x_j) \bmod m$$



# Transformacja rozkładu jednorodnego

- **Metoda odwrotnej dystrybuanty**
- Transformację rozkładu jednostajnego możemy wykorzystać do generowania liczb losowych o skomplikowanych gęstościach prawdopodobieństwa



$$f(x) dx = g(y) dy$$

$$\text{gdy } f(x) \equiv 1 \Rightarrow dx = g(y) dy = dG(y)$$

$$g(y) = G'(y)$$

↑  
dystrybuanta

$$\int dx = \int dG(y)$$

$$x = G(y)$$

↑  
musi istnieć

$$y = G^{-1}(x) \equiv H(x)$$

$$x_{min} = G(y_{min}), x_{max} = G(y_{max})$$

$$y_{min} = G^{-1}(0), y_{max} = G^{-1}(1)$$

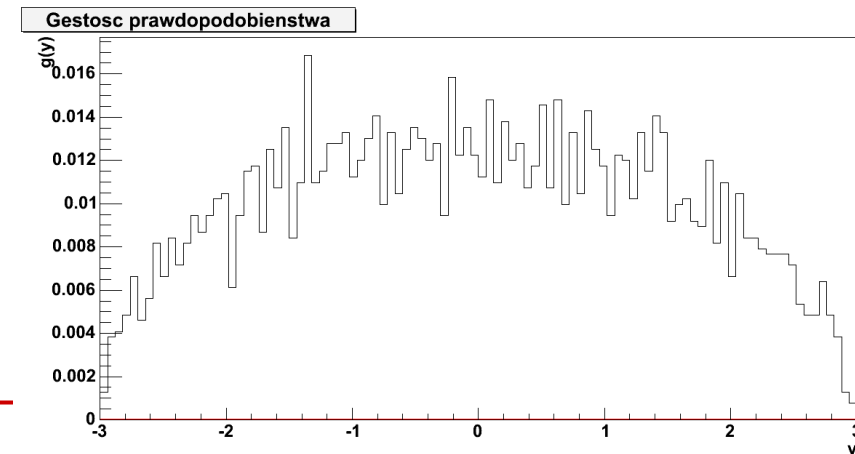
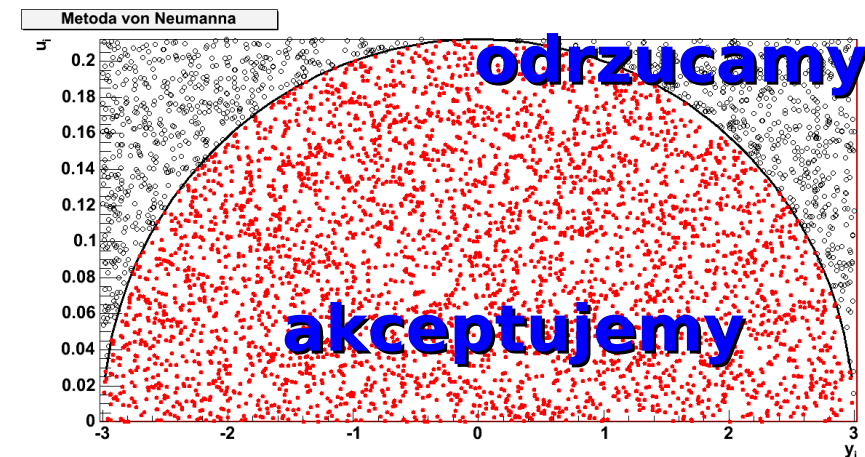
Czyli: liczymy dystrybuantę  $G(y)$   
a następnie funkcję odwrotną  $G^{-1}(y)$

Zmienna losowa  $X$  po transformacji  
 $y = G^{-1}(x)$  ma rozkład  $g(y)$



# Metoda (akceptacji) von Neumanna

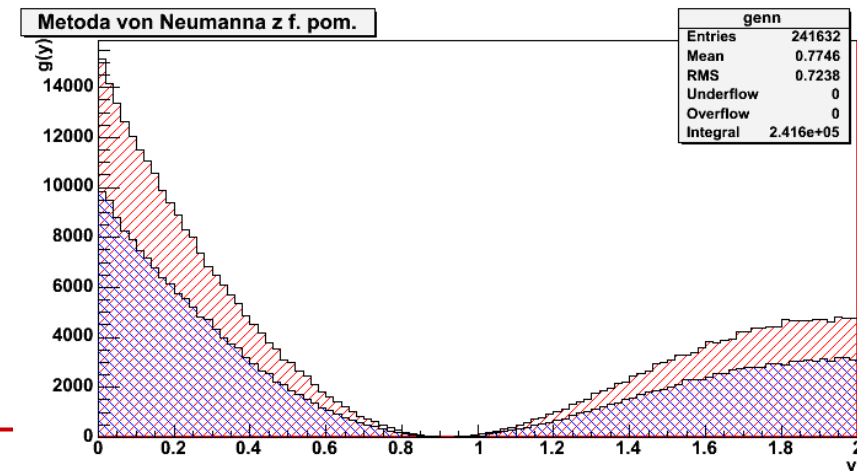
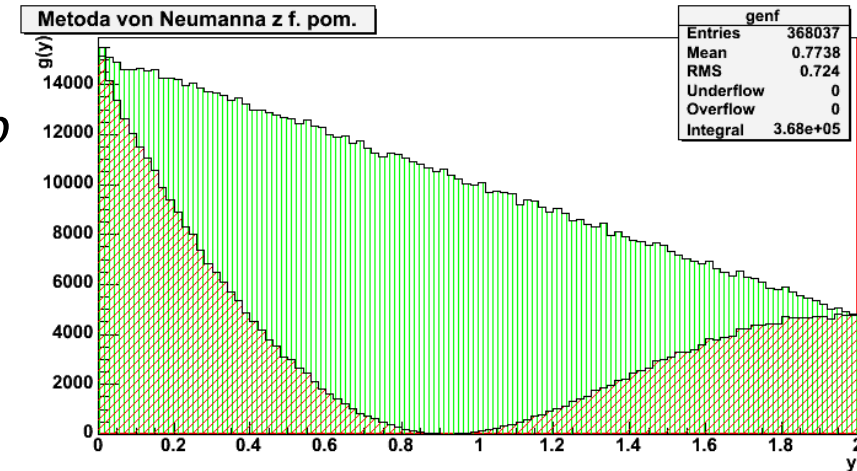
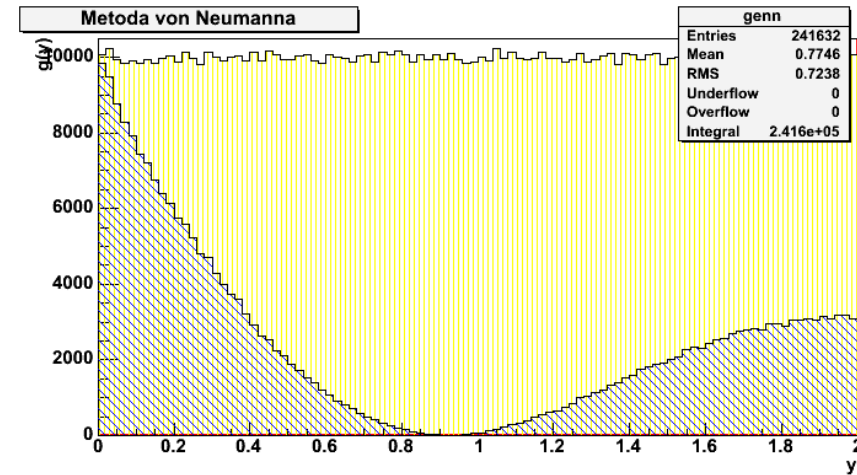
- Metoda generacji liczb metodą odwrotnej dystrybuanty ma swoje ograniczenia – dystrybuanta musi być znana (i być funkcją wzajemnie jednoznaczną), czyli musi istnieć funkcja odwrotna
- Metoda (akceptacji-odrzuceń) von Neumanna wymaga znajomości jedynie gęstości prawdopodobieństwa i pozwala na otrzymanie liczb z praktycznie dowolnego rozkładu
- Jak to działa?
  - generujemy parę liczb z rozkładu jednorodnego:  $(y_i, u_i)$
  - sprawdzamy, czy  $u_i < g(y_i)$
  - jeśli warunek jest spełniony, akceptujemy liczbę  $y_i$ ,  
jeśli nie - odrzucamy



# Metoda von Neumanna z funkcją pom.

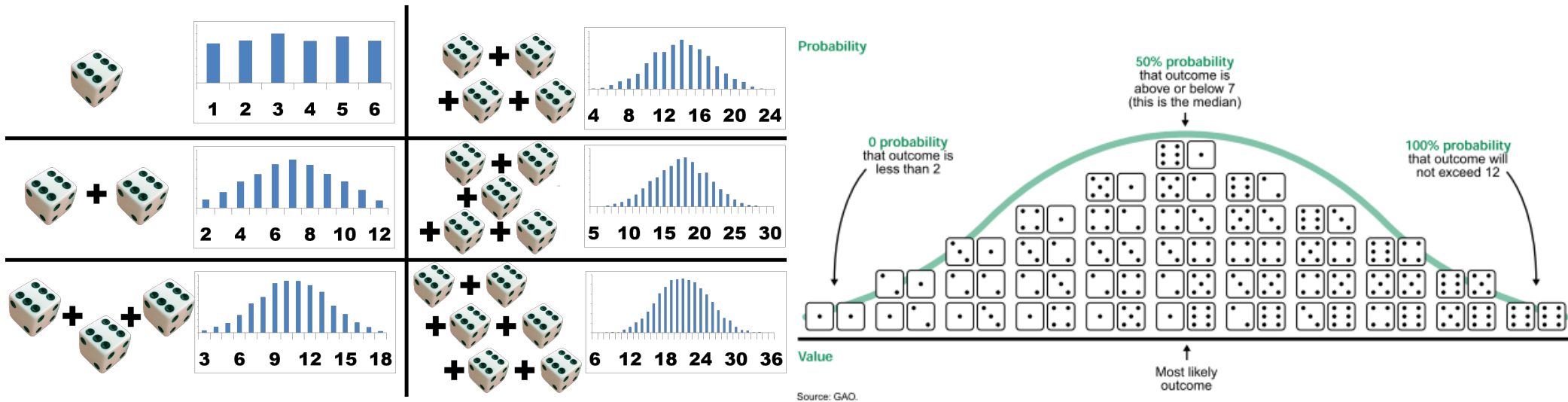
- Wydajność metody von Neumanna można poprawić, jeśli odpowiednio zawężymy obszar losowania:
  - wprowadzamy funkcję pomocniczą  $s(y)$ , z której “łatwo” wygenerować zmienne losowe (np. metodą odwrotnej dystrybuanty), i która spełnia warunek:  $g(y) \leq c \cdot s(y)$ ,  $a < y < b$
  - generujemy liczbę losową  $y_i$  z rozkładu  $s(y)$  na przedziale  $a < y_i < b$  oraz liczbę  $u_i$  z rozkładu jednorodnego na przedziale  $0 < u_i < 1$
  - odrzucaamy liczbę  $y_i$ , jeżeli:  $u_i \geq \frac{g(y_i)}{c \cdot s(y_i)}$
  - wydajność metody:

$$E = \frac{\int_a^b g(y) dy}{c \int_a^b s(y) dy}$$

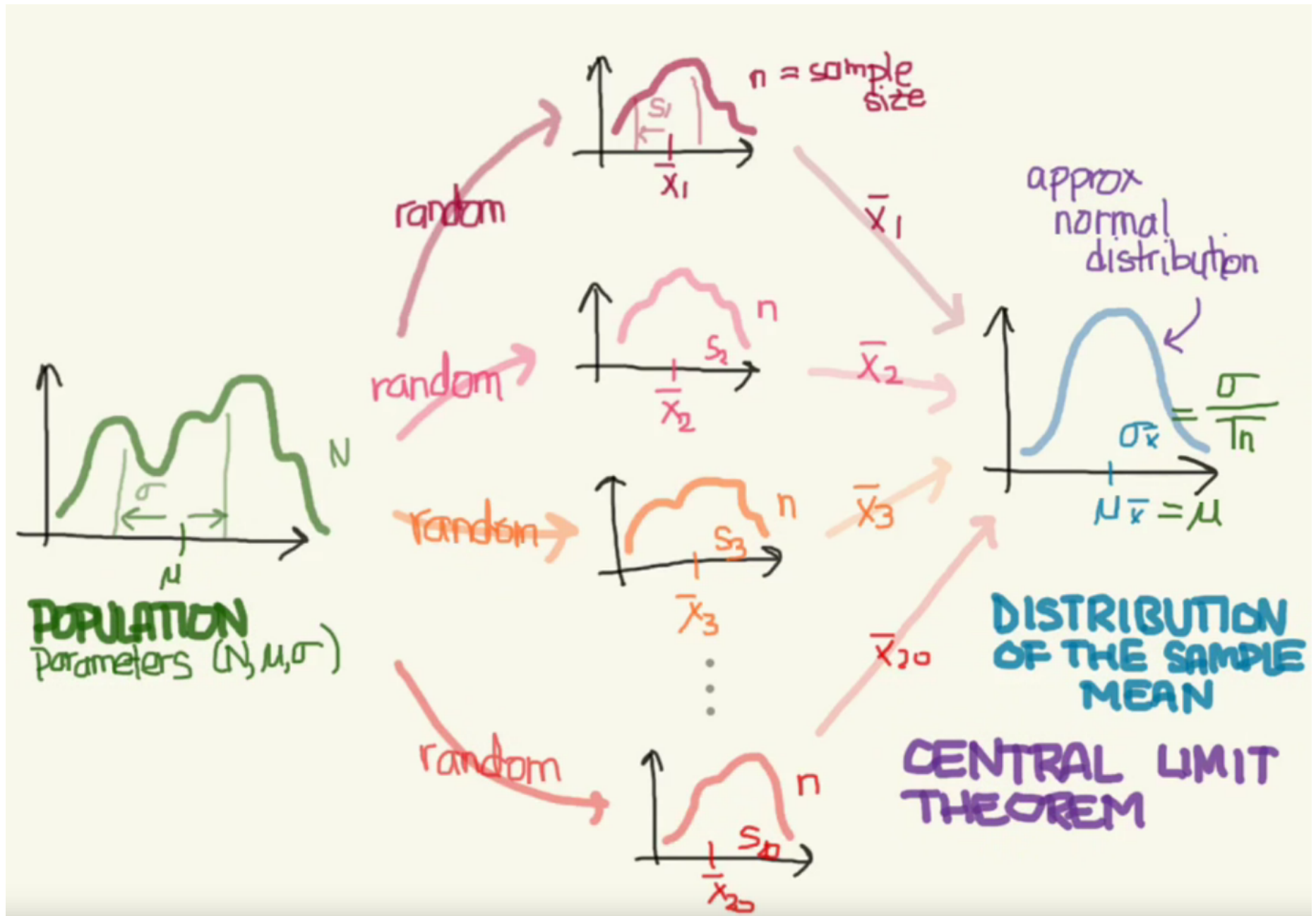


# Centralne twierdzenie graniczne

- Wyobraźmy sobie eksperyment polegający na rzucie kostką (kostkami) i obserwowaniu całkowitej liczby oczek:
  - kolejne rzuty kostką (kostkami) są niezależne
  - jeśli rzucamy kostką jednokrotnie (albo 1 kostką), to prawdopodobieństwo uzyskania danej wartości jest jednakowe
  - jeśli rzucamy kostką dwukrotnie (albo 2 kostkami), to prawdopodobieństwo uzyskania sumy oczek nie jest już jednakowe
  - jeśli rzucimy kostką  $n$ -krotnie ( $n$ -kostkami) → rozkład normalny



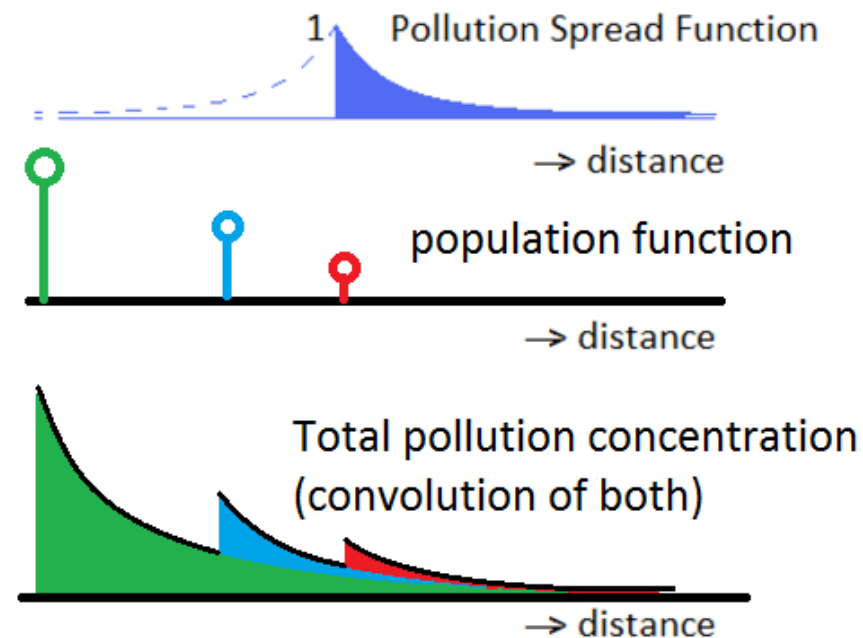
# Centralne twierdzenie graniczne



# Suma zmiennych losowych jako splot

<https://www.quora.com/The-density-function-of-the-sum-of-two-random-variables-is-the-convolution-of-their-respective-densities-What-is-the-intuition-behind-this>

- Wyobraźmy sobie taką sytuację:
  - Mieszkaś w wiosce obok rzeki
  - Mieszkańcy wioski wrzucają do rzeki odpady biologiczne
  - Koncentracja odpadów w funkcji odległości od miejsca zrzutu (*Pollution Spread Function, PSF*) jest zależna od ich rozkładu przez mikroorganizmy w rzece
  - Ilość wrzucanych odpadów zależy od populacji miejscowości na rzece



- Jaka jest pełna funkcja opisująca poziom zanieczyszczeń w rzece?
- Jest to **splot** dwóch rozkładów – funkcji populacji oraz funkcji koncentracji odpadów
- Innymi słowy, zastępujemy każdy punkt w funkcji populacji przez funkcję koncentracji przeskalowaną przez wagę funkcji populacji

# Suma zmiennych losowych jako splot

- Przypadek splotu dwóch rozkładów jednorodnych:

$$f_x(x) = \begin{cases} 1, & 0 \leq x < 1 \\ 0, & \text{w przeciwnym razie} \end{cases} \quad f_y(y) = \begin{cases} 1, & 0 \leq y < 1 \\ 0, & \text{w przeciwnym razie} \end{cases}$$

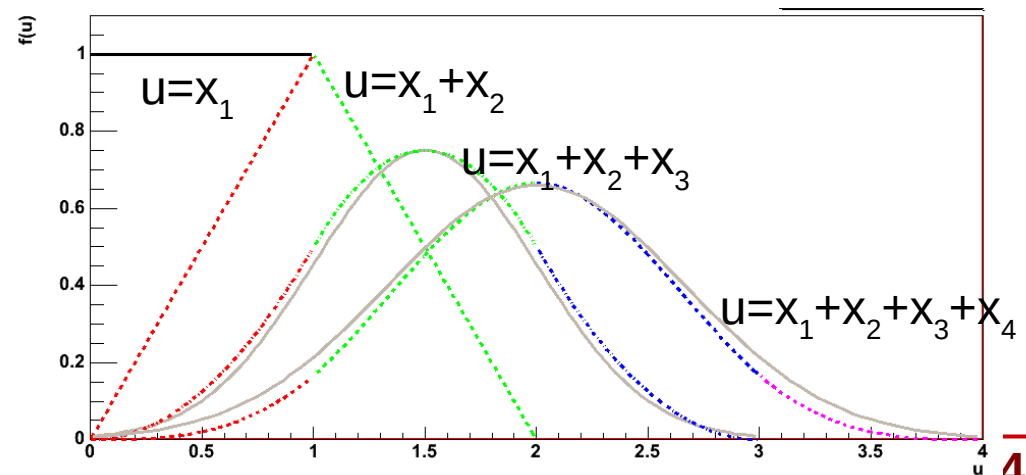
$$f(u) = \int_0^1 f_x(x) f_y(u-x) dx = \int_0^1 f_y(u-x) dx \quad \begin{matrix} v = u-x \\ dv = -dx \end{matrix} \Rightarrow f(u) = - \int_u^{u-1} f_y(v) dv = \int_{u-1}^u f_y(v) dv$$

- Zmienna  $u$  zmienia się od 0 do 2, zatem rozważmy 2 przypadki:

$$(a) \quad 0 \leq u < 1: f_1(u) = \int_0^u f_y(v) dv = \int_0^u 1 dv = u$$

$$(b) \quad 1 \leq u < 2: f_2(u) = \int_{u-1}^1 f_y(v) dv = \int_{u-1}^1 1 dv = 2 - u$$

- Zgodnie z CTG – im więcej rozkładów w splocie, tym bardziej rozkład sumy przypomina rozkład Gaussa:



# Estymatory

- Typowy problem analizy danych: znamy (np. z prawa fizycznego) ogólną postać gęstości prawdopodobieństwa w danej populacji, należy “jedynie” wyznaczyć parametry tego rozkładu. Przykład:
  - mierzymy rozpad radioaktywny w czasie:  $N(t) = N_0(1 - \exp(-\lambda t))$
  - parametr  $\lambda$  wyznaczamy na podstawie próby – mierząc skończoną ilość razy ilość rozpadów w czasie → wynik nigdy nie będzie dokładny, bo próba jest skończona, mamy problem **estymacji parametrów**
  - poszukiwana wielkość uzyskiwana jest funkcją elementów próby (**statystyką**) i jest nazywana **estymatorem**:  $S = S(X_1, X_2, \dots, X_n)$
  - estymator jest **nieobciążony**, jeżeli niezależnie od liczebności próby jego wartość oczekiwana jest równa wartości estymowanego parametru:

$$E(S(X_1, X_2, \dots, X_n)) = \lambda, \text{ dla każdego } n$$

- estymator jest **zgodny**, jeżeli jego wariancja znika:

$$\lim_{n \rightarrow \infty} \sigma(S(X_1, X_2, \dots, X_n)) = 0$$

# Estymator wartości oczekiwanej

## Populacja

- opisana funkcją gęstości:

$$f(x) = P(X=x)$$

- posiada **wartość oczekiwaną**:

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

- wartość oczekiwana rozkładu to **jedna liczba**
  - nie jest zmienną losową
  - chcemy ją zmierzyć doświadczalnie

- np. dla rozkł. Gaussa:

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

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

## Próba losowa

- zakładamy, że **średnia arytmetyczna** z elementów próby jest estymatorem wartości oczekiwanej

- **średnia arytmetyczna** jest statystyką:

$$\bar{X} = \frac{1}{n} (X_1 + X_2 + \dots + X_n)$$

- **jest zmienną losową** (zależy od elementów próby)
- posiada swoją **wartość oczekiwaną** oraz **wariancję**
- oczekujemy, że będzie ona *estymatorem nieobciążonym* i *zgodnym* wartości oczekiwanej populacji:

$$E(\bar{X}) = E(X) = \hat{x}, \text{ dla każdego } n$$

$$\lim_{n \rightarrow \infty} \sigma(\bar{X}) = 0$$

- **Na wykładzie pokazane jest jak to sprawdzić**



# Estymatory - podstawowe

- Przykładowe **estymatory nieobciążone**:

- **wartości oczekiwanej populacji** → **średnia arytmetyczna z próby**:

$$\bar{X} = \frac{1}{n} (X_1 + X_2 + \dots + X_n)$$

- **wariancji populacji** → **średnia odchyłeń kwadratowych**:

$$S^2(X) = \frac{1}{n-1} ((X_1 - \bar{X})^2 + (X_2 - \bar{X})^2 + \dots + (X_n - \bar{X})^2)$$

-1 wynika z jednego równania więzów (istnienie średniej)

- **Wariancje (niepewności) estymatorów**:

- **wariancja średniej arytmetycznej**:

$$\sigma^2(\bar{X}) = \frac{1}{n} S^2(X) = \frac{1}{n(n-1)} \sum_{i=1}^n (X_i - \bar{X})^2$$

- **wariancja średniej odchyłeń kwadratowych**:

$$\sigma^2(S^2(X)) = S^4(X) \left( \frac{2}{n-1} \right)$$

- **Uwaga!** Wariancje estymatorów są również estymatorami – możemy więc liczyć np. wariancję wariancji średniej arytmetycznej, itd.

# Metoda największej wiarygodności

- **Funkcją wiarygodności** nazywamy iloczyn postaci:

$$L = \prod_{j=1}^N f(\mathbf{X}^{(j)}; \boldsymbol{\lambda})$$

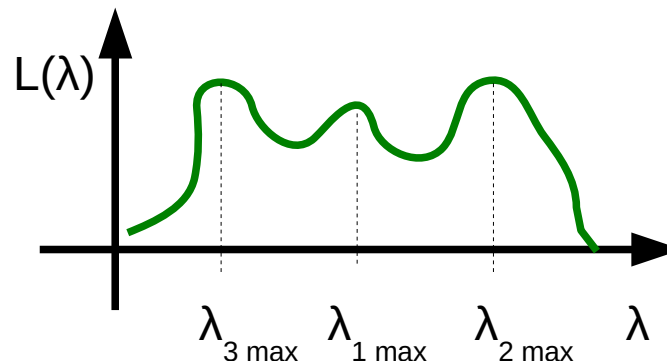
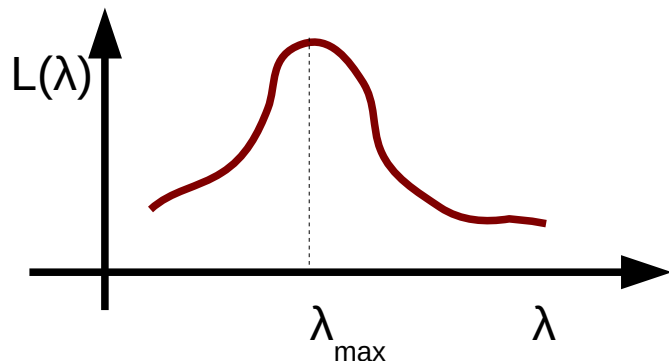
- **Jak wyznaczyć maksimum?**

– warunek konieczny: przyrównać pierwszą pochodną  $L$  do zera

- Różniczkowanie iloczynu jest jednak niewygodne, wprowadzamy więc **logarytm funkcji wiarygodności**  $l$ :

$$l = \ln L = \sum_{j=1}^N \ln f(\mathbf{x}^{(j)}; \boldsymbol{\lambda})$$

- **Funkcją wiarygodności** jest odpowiednikiem gęstości prawdopodobieństwa, tylko określona dla parametrów. Ponieważ jest funkcją próby losowej, jest również zmienną losową



równania wiarygodności

$$\frac{\partial l}{\partial \lambda_i} = 0, \quad i = 1, 2, \dots, p$$

# Przykład – estymacja rozkładu norm.

- **Przykład:** badamy rozkład wagi studentek amerykańskich college'ów
  - rozkład wagi studentek w populacji jest opisany rozkładem normalnym o wartości średniej  $\mu$  i wariancji  $\sigma^2$
  - założmy, że wybraliśmy  $N=10$  reprezentantów, których wagi (w kg)  $x^{(j)}$  układają się następująco:
    - 52,2 55,3 59,0 57,6 67,6 72,6 68,9 62,6 67,6 81,6
  - **Zadanie:** na podstawie wyniku pomiaru znajdź najbardziej wiarygodną estymację parametrów rozkładu
  - oczywiście funkcja rozkładu prawdopodobieństwa każdego wyniku (wagi studentek) dana jest wzorem:
$$f(x^{(j)}; \mu, \sigma) = \frac{1}{\sigma \sqrt{2\pi}} \exp\left(\frac{-(x^{(j)} - \mu)^2}{2\sigma^2}\right)$$
  - konstruujemy funkcję wiarygodności:

$$L(\hat{\mu}, \hat{\sigma}) = \sigma^{-n} (2\pi)^{-n/2} \exp\left[-\frac{1}{2\hat{\sigma}^2} \sum_{j=1}^N (x^{(j)} - \hat{\mu})^2\right]$$

# Przykład - estymacja rozkładu norm.

- oraz logarytmiczną funkcję wiarygodności:

$$l(\hat{\mu}, \hat{\sigma}) = \ln L(\hat{\mu}, \hat{\sigma}) = \ln \left( \hat{\sigma}^{-N} (2\pi)^{-n/2} \left[ -\frac{1}{2\hat{\sigma}^2} \sum_{j=1}^N (x^{(j)} - \hat{\mu})^2 \right] \right) = -\frac{1}{2\hat{\sigma}^2} \sum_{j=1}^N (x^{(j)} - \hat{\mu})^2 + \text{const}$$

- równania wiarygodności:

$$\frac{dl(\hat{\mu}, \hat{\sigma})}{d\hat{\mu}} = 0 \qquad \frac{dl(\hat{\mu}, \hat{\sigma})}{d\hat{\sigma}} = 0$$

$$\frac{dl(\hat{\mu}, \hat{\sigma})}{d\hat{\mu}} = \frac{-2 \sum_{j=1}^N (x^{(j)} - \hat{\mu})(-1)}{2\hat{\sigma}} = 0 \Rightarrow \sum_{j=1}^N (x^{(j)}) - N\hat{\mu} = 0 \Rightarrow \hat{\mu} = \frac{\sum_{j=1}^N x^{(j)}}{n}$$

$$\frac{dl(\hat{\mu}, \hat{\sigma})}{d\hat{\sigma}} = -\frac{n}{2\hat{\sigma}} + \frac{\sum_{j=1}^N (x^{(j)} - \hat{\mu})}{2\hat{\sigma}^2} = 0 \Rightarrow \hat{\sigma}^2 = \frac{\sum_{j=1}^N (x^{(j)} - \hat{\mu})^2}{n}$$

padło mailowo pytanie, czemu nie n-1

- estymatorem największej wiarygodności jest **obciążony** estymator wariancji

- dla formalności powinniśmy jeszcze sprawdzić drugie pochodne...

# Estymatory obciążone

- Czasem nie da się wyznaczyć estymatora nieobciążonego: jaki jest warunek osiągnięcia minimalnej wariancji?

- Można pokazać (Brandt), że zajdzie tak, gdy:  $E(S) = B(\lambda) + \lambda$

$$I' = A(\lambda)(S - E(S)) \quad I = B(\lambda)S + C(\lambda) + D \quad L = d \exp(B(\lambda)S + C(\lambda))$$

- Wtedy w przypadku estymatora nieobciążonego o minimalnej wariancji otrzymujemy:

$$\sigma^2(S) = \frac{1}{I(\lambda)} = \frac{1}{E(I'^2)}$$

$$\sigma^2(S) = \frac{1}{|A(\lambda)|}$$

- Estymatory spełniające powyższe warunki nazywamy **estymatorami o najniższej wariancji**

- Przykład – rozkład Poissona:  $f(k) = \frac{\lambda^k}{k!} e^{-\lambda}$

$$\frac{dl}{d\lambda} = I' = \sum_{j=1}^N \left\{ \frac{k^{(j)}}{\lambda} - 1 \right\} =$$

$$\frac{1}{\lambda} \sum_{j=1}^N \{k^{(j)} - \lambda\} = \frac{N}{\lambda} (\bar{K} - \lambda)$$

$$\tilde{\lambda} = \bar{K}, \quad \sigma^2(\bar{K}) = \frac{N}{\lambda}$$

$$I' = A(\lambda)(\tilde{\lambda} - \lambda) \quad \sigma^2(S) = \frac{1}{|A(\lambda)|}$$

# Rozkład $\chi^2$ – zastosowanie

- Rozkład  $\chi^2$  stosuje się jako **miarę ufności** uzyskanego wyniku (**odchylenia elementów próby od wartości średniej populacji**). Im mniejsza wartość  $\chi^2$  tym pozornie słuszniejszy wynik. Jako miary zaufania do wyniku używa się wielkości:

$$f(\chi^2) = k \cdot (\chi^2)^{\lambda-1} e^{-1/2\chi^2} \quad k = \frac{1}{\Gamma(\lambda) 2^\lambda} \quad W(\chi^2) = 1 - F(\chi^2) \equiv p = 1 - \alpha$$

nazywanej **poziomem ufności** (zwykle podawanym w % ilości odchyżeń standardowych rozkładu normalnego  $\sigma$ )

– wielkość  $\alpha$  jest nazywana **poziomem istotności**

- W rzeczywistych przypadkach mamy do czynienia z pełnym rozkładem Gaussa o dowolnym  $a$  i  $\sigma$ . Wprowadzamy wtedy odpowiednie przeskalowanie

elementy próby losowej

$$\chi^2 = X^2 = \frac{(X_1 - a)^2 + (X_2 - a)^2 + \dots + (X_n - a)^2}{\sigma^2}$$

- a w ogólnym przypadku wielowymiarowym, gdy zmienne są zależne:

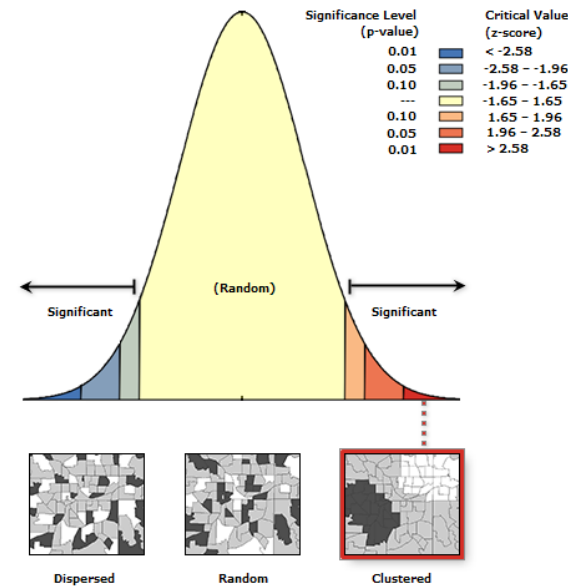
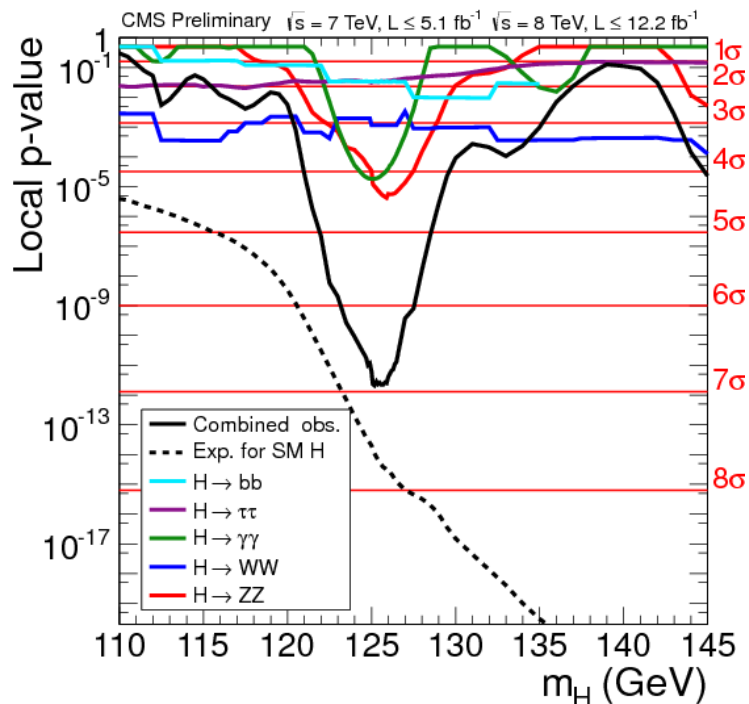
$$\chi^2 = X^2 = (\mathbf{X} - \mathbf{a})^T \mathbf{B} (\mathbf{X} - \mathbf{a})$$

# Rozkład $\chi^2$ a estymator wariancji

- Można udowodnić, że zmienna losowa:  $\frac{n-1}{\sigma^2} S^2$  ← estymator wariancji

ma rozkład  $\chi^2$  z  $f=n-1$  stopniami swobody. Wynika to stąd, że wyrażenia  $(X_i - \bar{X})^2$  nie są liniowo niezależne (co już wiemy). Każde dodatkowe równanie (więzy) pomiędzy wyrażeniami  $(X_i - \bar{X})^2$  redukuje liczbę stopni swobody o 1

- Poziomy ufności – przykład:



$$P(|Y - a| \leq \sigma) = 68,3\% \quad P(|Y - a| > \sigma) = 31,7\%$$

$$P(|Y - a| \leq 2\sigma) = 95,4\% \quad P(|Y - a| > 2\sigma) = 4,6\%$$

$$P(|Y - a| \leq 3\sigma) = 99,8\% \quad P(|Y - a| > 3\sigma) = 0,2\%$$

# Weryfikacja hipotez statystycznych

- **Przykład:** rozważamy zmienną losową  $X$  opisaną standardowym rozkładem Gaussa (średnia 0, odchylenie 1). Pobieramy 10-elementową próbę, uzyskaliśmy średnią arytmetyczną:  $\bar{X}=0,5$
- Jak na podstawie tej jednej realizacji próby (np. wyniku eksperymentu) możemy stwierdzić, czy pochodzi ona z takiej populacji? Innymi słowy, naszą **hipotezą** jest: **próba losowa pochodzi z rozkładu Gaussa o średniej 0 i odchyleniu 1**
- Procedura weryfikacji hipotezy nazywana jest **testem statystycznym**
- Jeżeli **hipoteza jest słuszna (nasze założenie)** to wartość średnia (będąca również zmienną losową)  $\bar{X}$  ma rozkład normalny ze średnią 0 i odchyleniem std.  $1/\sqrt{10}$

$$\sigma^2(\bar{X}) = \frac{1}{n} \sigma^2(X) = \frac{1}{10} \cdot 1 \Rightarrow \sqrt{\sigma^2(\bar{X})} = \frac{1}{\sqrt{10}}$$

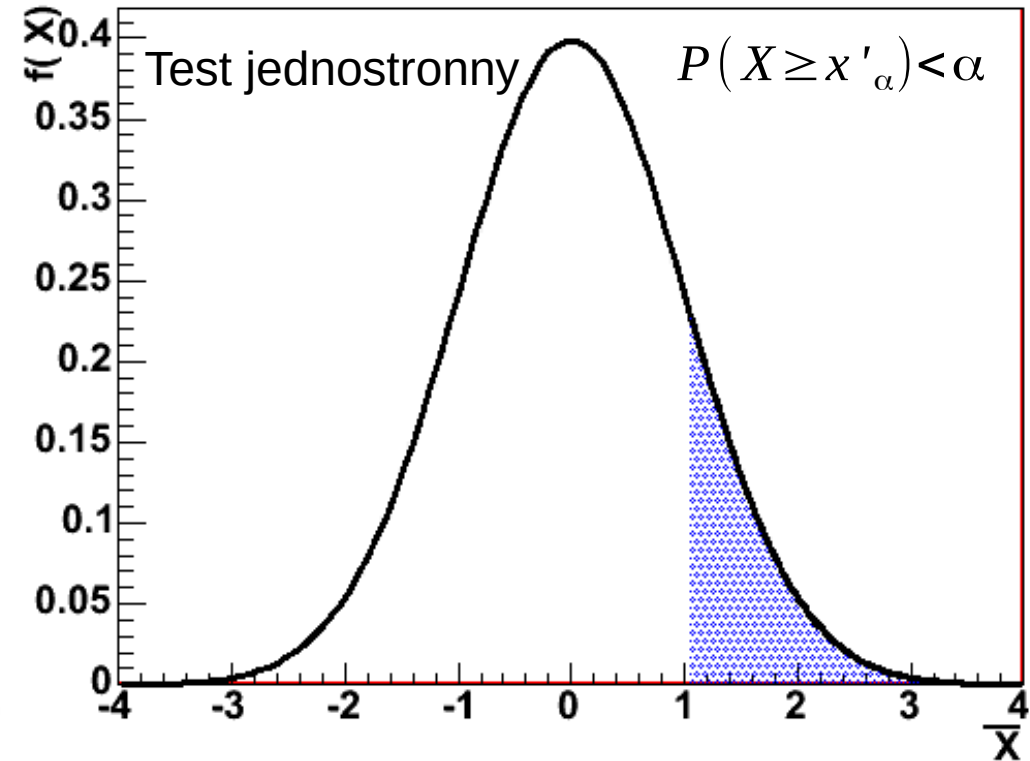
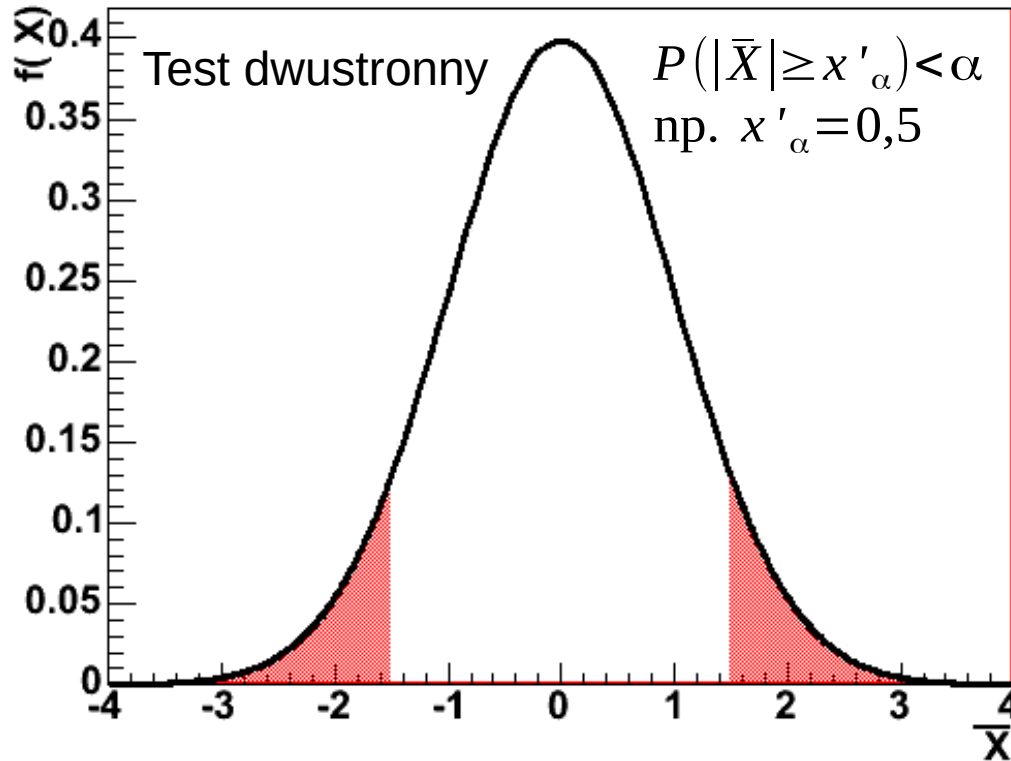


# Weryfikacja hipotez statystycznych

- Jak na podstawie **konkretnej realizacji próby** sprawdzić, czy założona hipoteza jest prawdziwa?
  - **I:** musimy ustalić pewną wartość prawdopodobieństwa  $\alpha$  (zwanego **poziomem istotności**, z reguły mała wartość, np. 0,01, albo 0,03, czy 0,05)
  - **II:** pytamy, czy prawdopodobieństwo zaobserwowania określonych wartości próby jest mniejsze niż  $\alpha$ :  $P(|\bar{X}| \geq 0,5) < \alpha$
  - **nierówność spełniona** – jest mało prawdopodobne, aby próba pochodziła z rozkładu określonego przez testowaną hipotezę → **możemy ją odrzucić**
  - **prawdopodobieństwo zaobserwowania tego, że  $|\bar{X}|$  jest duże, jest bardzo małe, ale takie nam się trafiło – więc prawdopodobnie (z prawdopodobieństwem  $1-\alpha$ ) nasza hipoteza nie jest słuszna**
  - **III:** jeśli prawdopodobieństwo jest mniejsze niż przyjęta wartość prawdopodobieństwa (poziom istotności)  $\alpha$ , odrzucamy hipotezę na zadanym poziomie istotności

# Weryfikacja hipotez statystycznych

## Rozkład wartości średniej $\bar{X}$



- Jeśli (w naszym przykładzie) wartość średnia znajduje się w zaznaczonym obszarze (nazywamy go **obszarem krytycznym**), to hipotezę odrzucamy
  - jeśli oczekujemy rozkładu normalnego o średniej 0 i małym odchyleniu (np. 10), a z próby losowej (konkretny eksperyment) mamy średnią 1000, to lądujemy w “ogonie” rozkładu średniej i na podstawie tej konkretnej próby odrzucamy hipotezę (**ale na podstawie innej próby moglibyśmy zaakceptować**)

# Weryfikacja hipotez statystycznych

- W ogólnym przypadku używamy innych wielkości niż średnia:
  - definiujemy jakąś (wygodną dla nas) statystykę testową  $T$  (np. różnicę między wynikiem eksperymentu a krzywą teoretyczną)
  - ustalamy poziom istotności  $\alpha$
  - wyznaczamy taki zbiór  $U$ , który określa obszar zmienności statystyki testowej  $T$ , taki że prawdopodobieństwo znalezienia się w nim jest ograniczone wartością  $\alpha$ :  $P(T \in U) = \alpha$
  - z pobranej próby wyznaczamy konkretną wartość statystyki testowej  $T'$ : jeżeli znajduje się ona **wewnątrz** obszaru krytycznego  $U$ , **odrzucaamy hipotezę** (mówimy: krzywa teoretyczna nie opisuje wyniku eksperymentu), czyli odrzucaamy hipotezę, jeżeli  $T' \in U$



# Test dobroci $\chi^2$ dopasowania

# Test $\chi^2$ dobroci dopasowania

- Mamy  $N$  pomiarów  $g_i$ ,  $i=1, 2, \dots, N$  oraz ich niepewności  $\sigma_i$
- Wartości  $f_i$ ,  $i=1,2,\dots,N$  określają nam prawdziwy rozkład danej wielkości mierzonej (**np. znaleziony poprzez estymację**)
- **Dla każdego pomiaru liczymy wielkość  $u_i$ :**  $u_i = \frac{g_i - f_i}{\sigma_i}$ ,  $i=1,2,\dots,N$
- Jeśli nasza teoria (wartości  $f_i$ ) jest prawdziwa, to rozkłady różnic  $u_i$  mają postać standardowego rozkładu normalnego – **nasza hipoteza**
- Jeśli tak, to rozkład  $\chi^2$  o  $N$  stopniach swobody będzie miała wielkość:  
$$T = \sum_{i=1}^N u_i^2 = \sum_{i=1}^N \left( \frac{g_i - f_i}{\sigma_i} \right)^2$$
- **(Subiektywnie)** oczekujemy małej wartości wielkości  $T$
- Gdy hipoteza jest **fałszywa**, wówczas poszczególne różnice  $u_i$  przyjmują duże wartości (wartość  $T$  jest duża)
- Jak określić granicę zmienności  $T$ ? Można zauważyć, że granica ta jest określona **kwantylem**  $\chi_{1-\alpha}^2$ , czyli:

$$P(T > \chi_{1-\alpha}^2) = \alpha$$

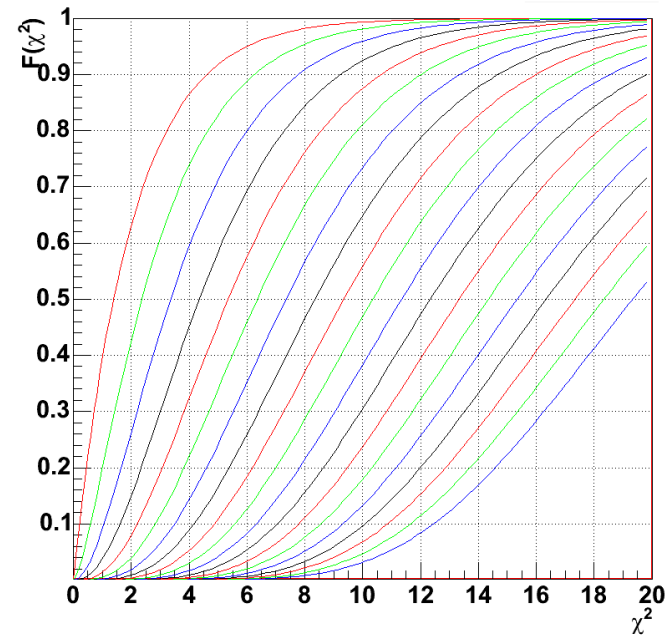
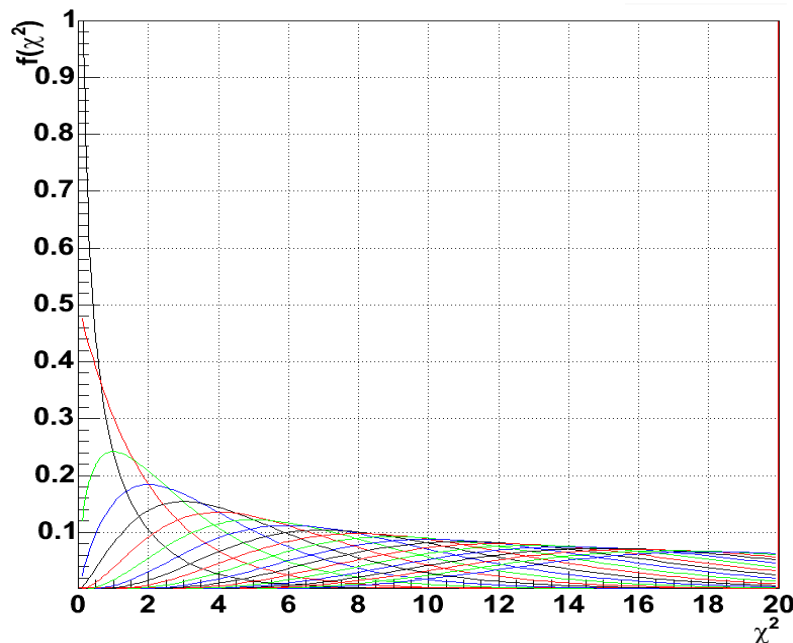
$$F(x_q) = P(X \leq x_q) = q$$
$$P(X > x_q) = 1 - q$$

# Test $\chi^2$ dobroci dopasowania

- Podsumowując, w naszym przypadku musimy dla danej realizacji próby (wyniku eksperymentu) wyznaczyć wartość testową  $T$  i porównać ją z odpowiednim kwantylem rozkładu  $\chi^2$  o odpowiedniej liczbie stopni swobody:

$$T > \chi_{1-\alpha}^2$$

- **Jeżeli ten warunek jest spełniony, to hipotezę odrzucamy** (punkty teoretyczne nie opisują danych eksperymentalnych na zadanym poziomie istotności)
- Skąd wziąć kwantyl? Z tablic lub z dystrybuanty:



# Test $\chi^2$ - przykład

## Zadanie

### Weryfikacja hipotez statystycznych (5 pkt.)

- ▶ Przeprowadzono eksperyment naświetlania wodorowej komory pęcherzykowej wiązką fotonów w celu badania oddziaływań fotonów z protonami. Fotony powodują powstawanie par elektron-pozyton, które mogą być wykorzystane do monitorowania wiązki fotonów. Częstość występowania zdjęć z 0,1,2,... parami elektron-pozyton powinna podlegać rozkładowi Poissona. Należy wczytać dane z pliku [plik](#) (w pierwszej kolumnie znajduje się liczba par elektronowych na zdjęciu k, a w drugiej liczba zdjęć zawierających k par elektronowych). Widzimy, że rozkład ten przypomina rozkład Poissona - próbujemy zatem obliczyć estymator największej wiarygodności dla parametry rozkładu Poissona (patrz [Wykład 10](#) slajd 13) (1 pkt.)
- ▶ Narysować na jednym wykresie punkty pomiarowe i dopasowanie (metodą estymatora największej wiarygodności).
- ▶ Sprawdzić jakość dopasowania za pomocą testu  $\chi^2$ . W tym celu należy zaimplementować funkcję obliczającą statystykę testową

$$\chi^2 \text{ zgodnie z wzorem } T = \sum_k \frac{(n_k - np_k)^2}{np_k}$$

gdzie:  $n_k$  - liczba obserwacji w k-tym binie,  $np_k$  - przewidywana przez teorię liczba przypadków w k-tym binie

- ▶ Określić liczbę stopni swobody i obliczyć wartość statystyki testowej. (1 pkt.)
- ▶ Zaimplementować funkcję zwracającą wynik testu  $\chi^2$  na zadanym poziomie istotności  $\alpha$

Wykorzystując zaimplementowaną funkcję zweryfikować hipotezę mówiącą, że dane pomiarowe podlegają rozkładowi Poissona. Dobrać odpowiednią wartość poziomu istotności. Uwaga! Kwanyl możemy odczytać z policzonej na ostatnich zajęciach dystrybuanty. (2 pkt.)

# Test $\chi^2$ - przykład

## Zadanie

### Weryfikacja hipotez statystycznych (5 pkt.)

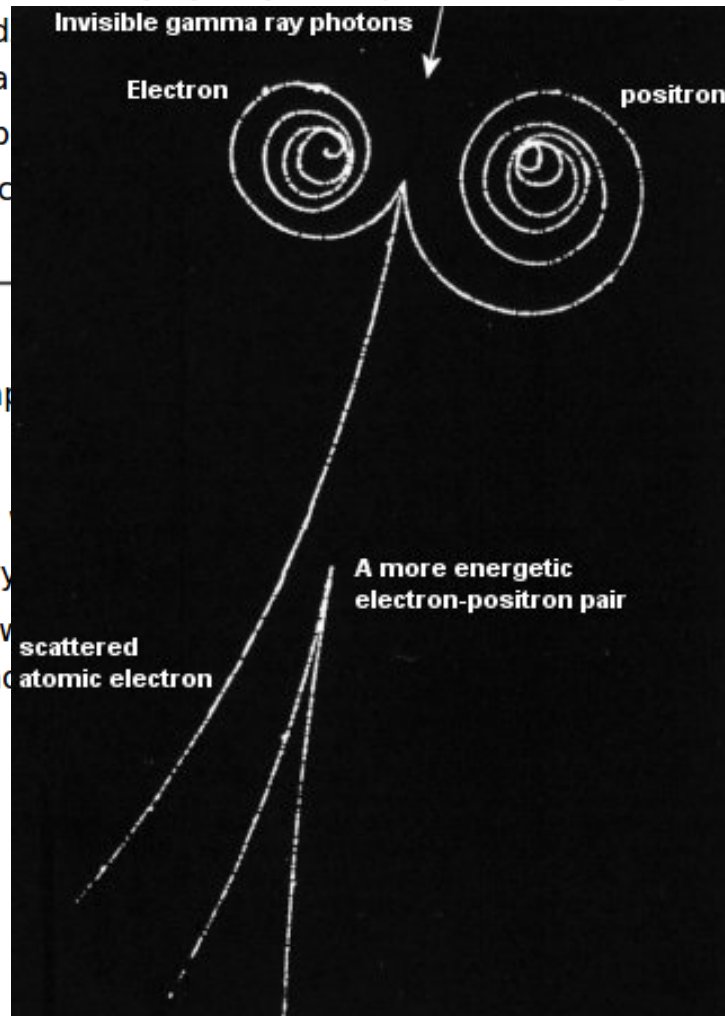
- ▶ Przeprowadzono eksperyment naświetlania wodorowej komory pęcherzykowej wiązką fotonów w celu badania oddziaływań fotonów z protonami. Fotony powodują powstawanie par elektron-pozyton, które mogą być wykorzystane do monitorowania wiązki fotonów. Częstość występowania zdjęć z 0,1,2,... parami elektron-pozyton powinna podlegać rozkładowi Poissona. Należy wczytać dane z pliku (w pierwszej kolumnie znajduje się liczba par elektronowych na zdjęciu k, a w drugiej liczba zdjęć zawierających k par elektronowych). Widujemy dane (na - próbujemy zatem obliczyć estymator największej wiarygodności dla (na slajd 13) (1 pkt.)
- ▶ Narysować na jednym wykresie punkty p (na największej wiarygodności).
- ▶ Sprawdzić jakość dopasowania za pomocą (na obliczając funkcję obliczającą statystykę testową

$\chi^2$  zgodnie z wzorem 
$$T = \sum_k \frac{(n_k - \mu_k)^2}{\mu_k}$$

gdzie:  $n_k$  - liczba obserwacji w k-tym binie,  $\mu_k$  - oczekiwana liczba obserwacji w k-tym binie

- ▶ Określić liczbę stopni swobody i obliczyć wartość testową
- ▶ Zaimplementować funkcję zwracającą wartość testową

Wykorzystując zaimplementowaną funkcję zwracającą wartość testową, obliczamy wartość testową. Dobrać odpowiednią wartość poziomu istotności (2 pkt.)



na - próbujemy zatem obliczyć estymator największej wiarygodności dla (na slajd 13) (1 pkt.)

ów w k-tym binie

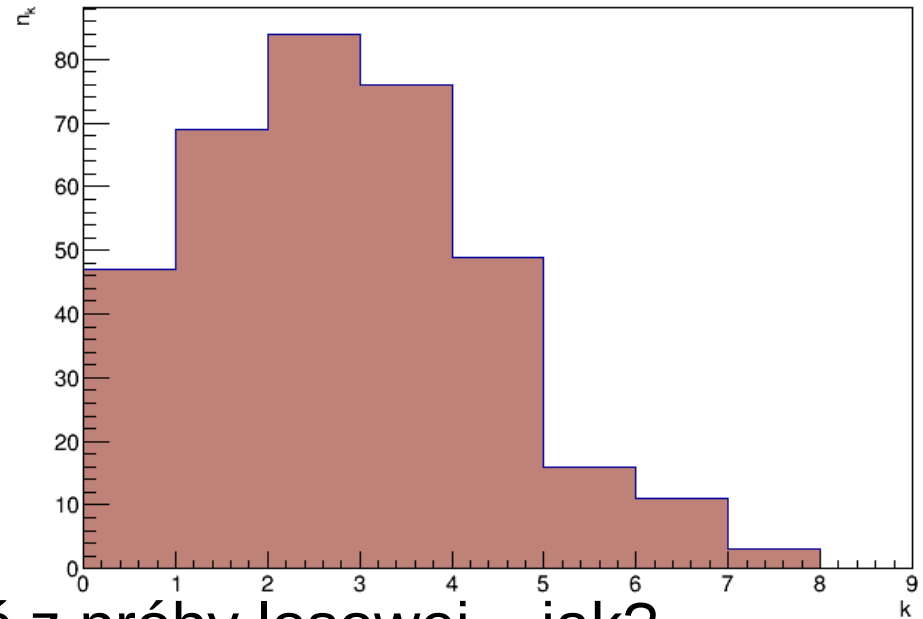
owe podlegają rozkładowi Poissona. (na ostatnich zajęciach dystrybuanty.



# Test $\chi^2$ - przykład

- Po wczytaniu danych z pliku histogram eksperymentalny wygląda następująco (nasza próba losowa):

Wynik eksperymentu



- Zakładamy hipotezę:** teoria mówi to jest rozkład Poissona (“na oko” zresztą tak wygląda)
- Rozkład Poissona ma tylko jeden parametr (wartość średnią):

$$f(k) = \frac{\lambda^k}{k!} e^{-\lambda}$$

- Musimy go zatem jakoś wyznaczyć z próby losowej – jak?  
**Na przykład metodą największej wiarygodności** – szukamy estymatora nieobciążonego największej wiarygodności o minimalnej wariancji (kilka slajdów temu):

$$\frac{dl}{d\lambda} = l' = \sum_{j=1}^N \left\{ \frac{k^{(j)}}{\lambda} - 1 \right\} =$$

$$\frac{1}{\lambda} \sum_{j=1}^N \{k^{(j)} - \lambda\} = \frac{N}{\lambda} (\bar{K} - \lambda)$$

$$\tilde{\lambda} = \bar{K}, \quad \sigma^2(\bar{K}) = \frac{N}{\lambda}$$

**Przypomnienie – definicja estymatora o min. wariancji:**

$$l' = A(\lambda)(\tilde{\lambda} - \lambda)$$

$$\sigma^2(S) = \frac{1}{|A(\lambda)|}$$

# Test $\chi^2$ - przykład

- Czyli estymatorem największej wiarygodności o minimalnej wariancji dla rozkładu Poissona jest średnia arytmetyczna z próby
- Oczywiście w naszym przypadku mamy histogram, który zawiera jakąś całkowitą liczbę wejść (całka z histogramu nie jest równa 1), wobec tego do średniej dodajemy wagi w postaci liczby wejść w danym binie i średnia staje się średnią ważoną:

$$\tilde{\lambda} = \frac{\sum_k k \cdot n_k}{\sum_k n_k}$$

- W naszym przypadku wartość ta wynosi mniej więcej:  $\tilde{\lambda} \approx 2,33$

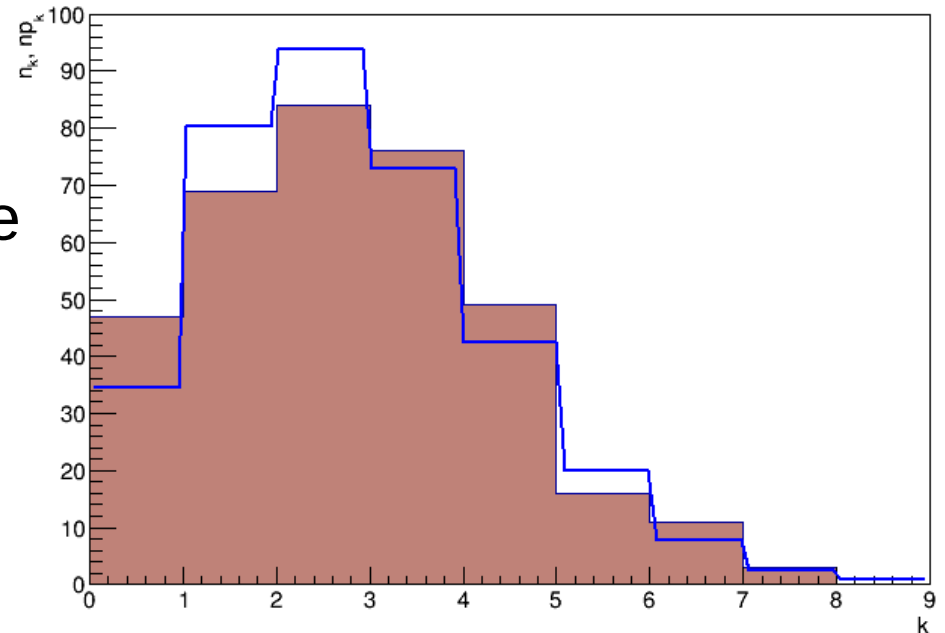
- Rysujemy więc funkcję:

$$n \cdot p_k = n \cdot f(k) = n \cdot \frac{\tilde{\lambda}^k}{k!} e^{-\tilde{\lambda}}, \text{ gdzie } n = \sum_k n_k$$

- Jak teraz sprawdzić, czy faktycznie nasza hipoteza jest słuszna?

- **Testujemy dobroć dopasowania**

Wynik eksperymentu



# Test $\chi^2$ - przykład

- Musimy zatem wyznaczyć wartość statystyki testowej  $T$ :

$$T = \sum_{k=0}^7 u_i^2 = \sum_{k=0}^7 \frac{(n_k - np_k)^2}{np_k} \approx 10,53$$

- Co dalej? Zakładamy poziom istotności, na przykład:  $\alpha = 0,01$
- Musimy jeszcze określić liczbę stopni swobody – ile ich jest?
  - liczba binów (8) minus 1 minus liczba parametrów (1)

$$r - 1 - p = 8 - 1 - 1 = 6$$

- Teraz szukamy odpowiedniego kwantyla rozkładu  $\chi^2$  o 6 stopniach swobody:  $\chi_{1-\alpha}^2 = \chi_{0,99}^2 \approx 16,81$
- Porównujemy statystykę z kwantylem:  $T = 10,51 < \chi_{0,99}^2 = 16,81$
- **Warunek  $T > \chi_{1-\alpha}^2$  nie jest spełniony, zatem na poziomie istotności  $\alpha = 0,01$  nie ma podstaw do odrzucenia hipotezy**



# Hipotezy zerowa i alternatywna

## Błędy I i II rodzaju

Na podstawie:

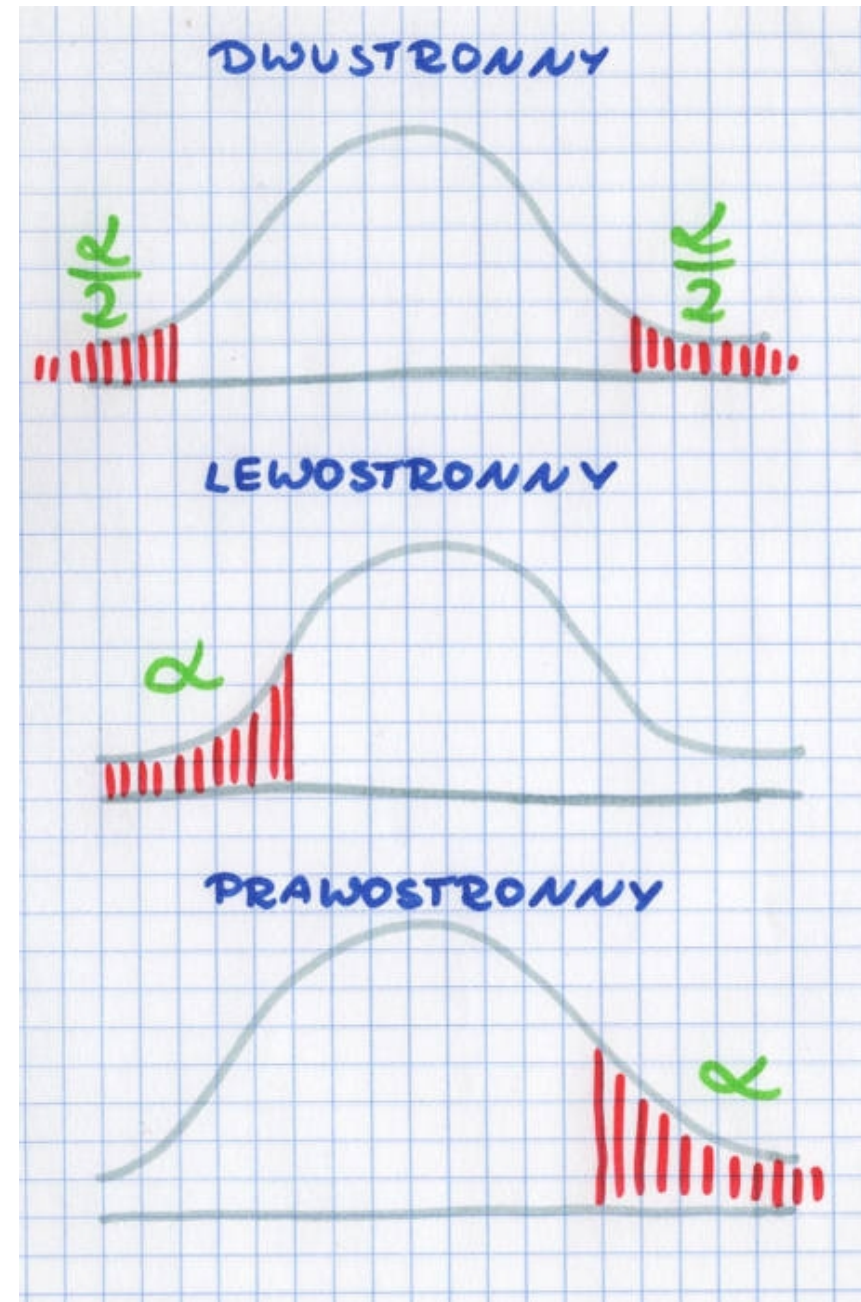
<https://www.statystyczny.pl/hipotezy-statystyczne/>

# Hipoteza zerowa i alternatywna

- Hipoteza może być różna, przykłady hipotez:
  - średni wzrost Polaków to 175 cm
  - 2% dzieci w wieku szkolnym nie lubi czekolady
  - poziom szczęścia dwie minuty po zjedzeniu dużej porcji lodów jest wyższy niż przed zjedzeniem tejże dużej porcji lodów
- **Hipoteza zerowa** ( $H_0$ ) to w uproszczeniu taka, gdy nie widzimy różnicy – np. czujemy się tak samo po zjedzeniu lodów jak przed
  - w przypadku pomiarów, np. wartość  $\chi^2$  jest mała → teoria opisuje dane
- **Hipoteza alternatywna** ( $H_A$ ) to przeciwieństwo hipotezy zerowej, którą możemy zdefiniować na kilka sposobów, np.:
  - np. jest różnica w szczęściu w zjedzeniu lodów (test dwustronny)
  - poziom szczęścia po zjedzeniu lodów jest mniejszy (test jednostronny)
  - poziom szczęścia po zjedzeniu lodów jest większy (test jednostronny)
  - w przypadku pomiarów  $\chi^2$  jest duże → teoria nie opisuje danych (test jednostronny – rozkład  $\chi^2$  jest niesymetryczny)
- **Statystyka testowa** - to funkcja próby, na podstawie której wnioskuje się o odrzuceniu lub nie hipotezy statystycznej – wielkość mająca swój rozkład prawd.

# Statystyka testowa

- **Statystyka testowa** - to funkcja próby, na podstawie której wnioskuje się o odrzuceniu lub nie hipotezy statystycznej – wielkość mająca swój rozkład prawdopodobieństwa
- Z naszej próby losowej (eksperymentu) dostajemy jedną wartość – ona znajduje się gdzieś w tym rozkładzie
- Obszar krytyczny (obszar odrzuceń) jest zawsze na końcu rozkładu
  - jeśli hipoteza mówi, że coś jest różne – dwustronny
  - mniejsze lub większe – jednostronny
- **Statystyka testowa ma swój (różny) rozkład zarówno dla  $H_0$  jak i  $H_A$ !!!**



# Błąd I rodzaju

- **Błąd I rodzaju** to taki, gdzie odrzucamy hipotezę zerową a była ona prawdziwa
  - $H_0$ : poziom szczęścia po zjedzeniu dużej porcji lodów jest taki sam jak przed
  - $H_A$ : poziom szczęścia się zmienia
  - my na podstawie doświadczenia (próby losowej) zjedliśmy lody i np. przez te okropne wyrzuty sumienia, że znowu za dużo kalorii :) odrzucamy hipotezę zerową na rzecz alternatywnej
  - jeżeli w wyniku wielu prób losowych wynika, że jednak lądujemy w ogonie rozkładu owego szczęścia, to popełnimy właśnie błąd pierwszego rodzaju – bo odrzuciliśmy hipotezę zerową, która była prawdziwa
- Prawdopodobieństwo popełnienia błędu I rodzaju określa **poziom istotności**  $\alpha$ , stąd oczekujemy, by był jak najmniejszy

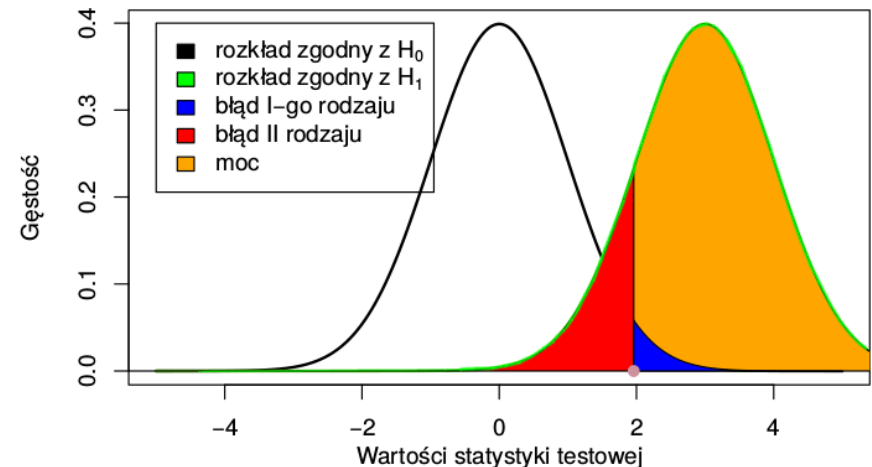
# Błąd II rodzaju

- **Błąd II rodzaju** to taki, gdy nie odrzucimy hipotezy fałszywej
  - ma on miejsce w sytuacji, kiedy jednak ten poziom szczęścia przed zjedzeniem lodów i po zjedzeniu się różni. Jeśli te wyrzuty sumienia z powodu przyswojenia dużej dawki kalorii są na tyle poważne, że pomimo zjedzenia czegoś dobrego i smacznego, wcale nie czujemy się lepiej. I jeśli poziom szczęścia zdecydowanie spada nam po zjedzeniu lodów, a my stwierdzimy, że nie ma żadnej różnicy, to wtedy popełniamy błąd II rodzaju. **Nie odrzucamy hipotezy zerowej, mimo że jest ona fałszywa.**
- **Prawdopodobieństwo** popełnienia błędu II rodzaju określamy jako  $\beta$



# Moc testu

- Moc testu (prawdopodobieństwo, że prawidłowo odrzucimy hipotezę zerową) to  $1-\beta$ . Inaczej mówiąc jest to prawdopodobieństwo niepopelnienia błędu II rodzaju.
- Moc testu zależy od kilku czynników:
  - Wielkości próby użytej w badaniu (im większa próba, tym większa moc testu).
  - Rzeczywistej wielkości efektu na tle losowej zmienności w populacji.
  - Przyjętego poziomu istotności  $\alpha$  (między błędem I i II rodzaju jest taka zależność, że jeżeli zwiększamy prawdopodobieństwo popelnienia danego błędu, jednocześnie zmniejszamy je dla drugiego).

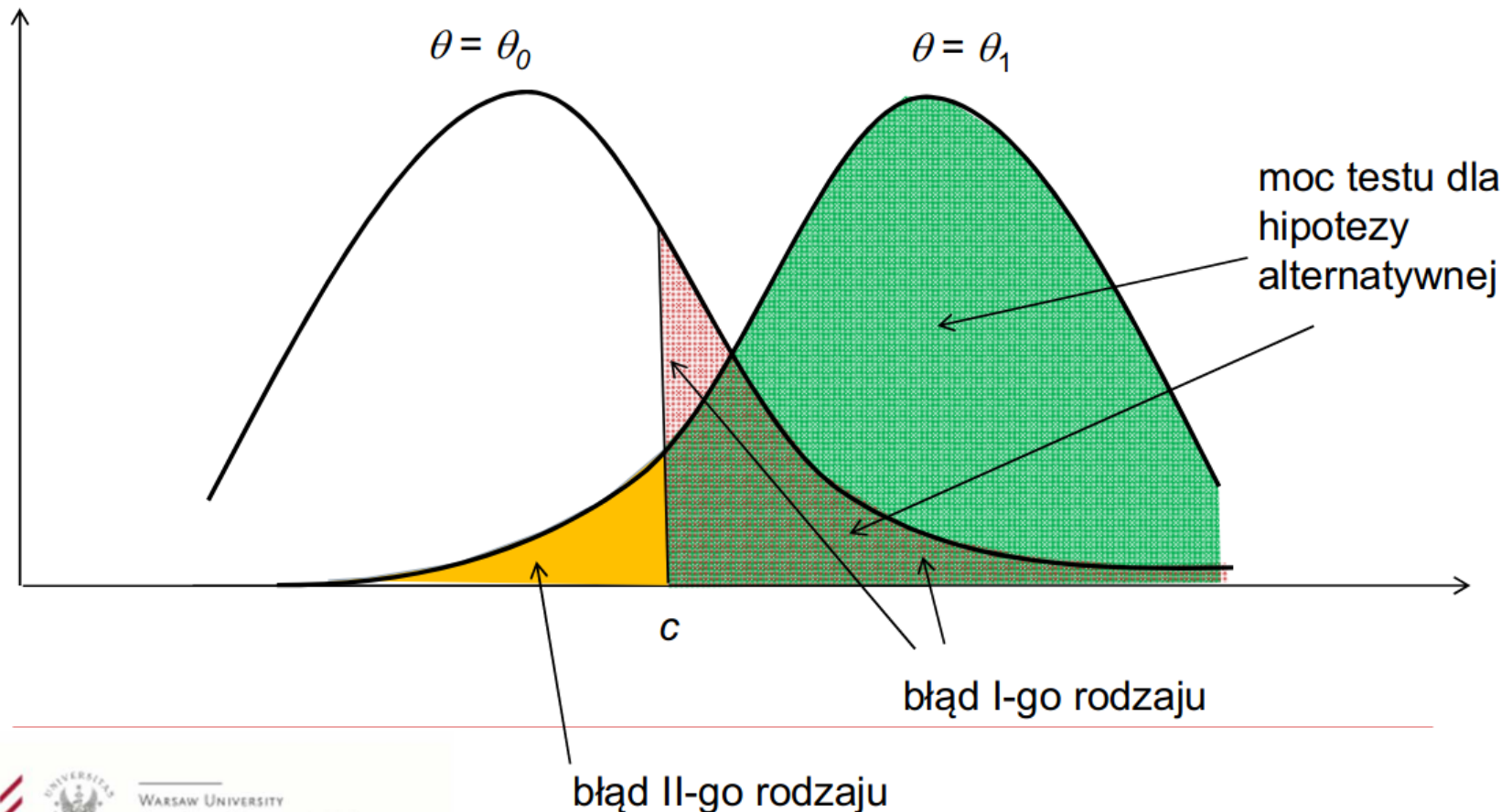


Rys.2 Relacje między błędem I, II rodzaju oraz mocą

# Moc testu

## Moc testu: interpretacja graficzna (1)

rozkłady statystyki testowej przy założeniu prawdziwości  
hipotezy zerowej i alternatywnej



# Przykład: winny czy inny

- **Przykład** działanie sądów może pomóc zrozumieć nam, kiedy przyjmujemy hipotezę alternatywną i dlaczego nieprzyjęcie hipotezy alternatywnej nie oznacza, że przyjmujemy hipotezę zerową.
- Wyobraźmy sobie Pana X oskarżonego o kradzież diamentu Pani Y. Pan X staje przed sądem.
  - **Hipotezą zerową** w tym przypadku jest niewinność Pana X.
    - Zakładamy, że Pan X wcale tych diamentów nie ukradł, że zrobiła to sprzątaczką, kucharką albo Pani Y schowała je w innej szufladzie i zupełnie o tym zapomniała.
  - **Hipoteza alternatywna**, to oczywiście wina Pana X. Skoro nie jest niewinny, to znaczy, że jednak ukradł diamenty i powinien zostać przez sąd skazany.
  - Założeniem jednak podstawowym jest brak winy – sąd musi znaleźć przekonujące argumenty, żeby móc Pana X oskarżyć.
    - **Jeśli je znajdzie** – to znaczy, że **odrzuca hipotezę zerową** (tę o niewinności) na rzecz hipotezy alternatywnej (że Pan X jest winny przestępstwa).
    - **Jeśli ich nie znajdzie**, to (nawet jeśli sąd wciąż będzie podejrzewać, że coś się w zachowaniu Pana X nie zgadza i że ta niewinność wcale nie jest zbyt pewna) **nie ma podstaw do odrzucenia hipotezy zerowej. To wcale nie znaczy, że Pan X nie ukradł diamentów. To tylko oznacza, że sąd nie znalazł wystarczającego dowodu.**
- A gdzie tu błędy?
  - Jeśli Pan X ukradł diamenty i został skazany to sąd się nie pomylił. Jeśli diamentów nie ukradł i nie został uznany winnym, to również nie ma żadnego błędu.
  - Jeśli nie ukradł, a poszedł do więzienia, mamy do czynienia z błędem pierwszego rodzaju. Natomiast jeśli Pan X ukradł diamenty i został uniewinniony, to mamy do czynienia z błędem drugiego rodzaju.

# Błąd II rodzaju

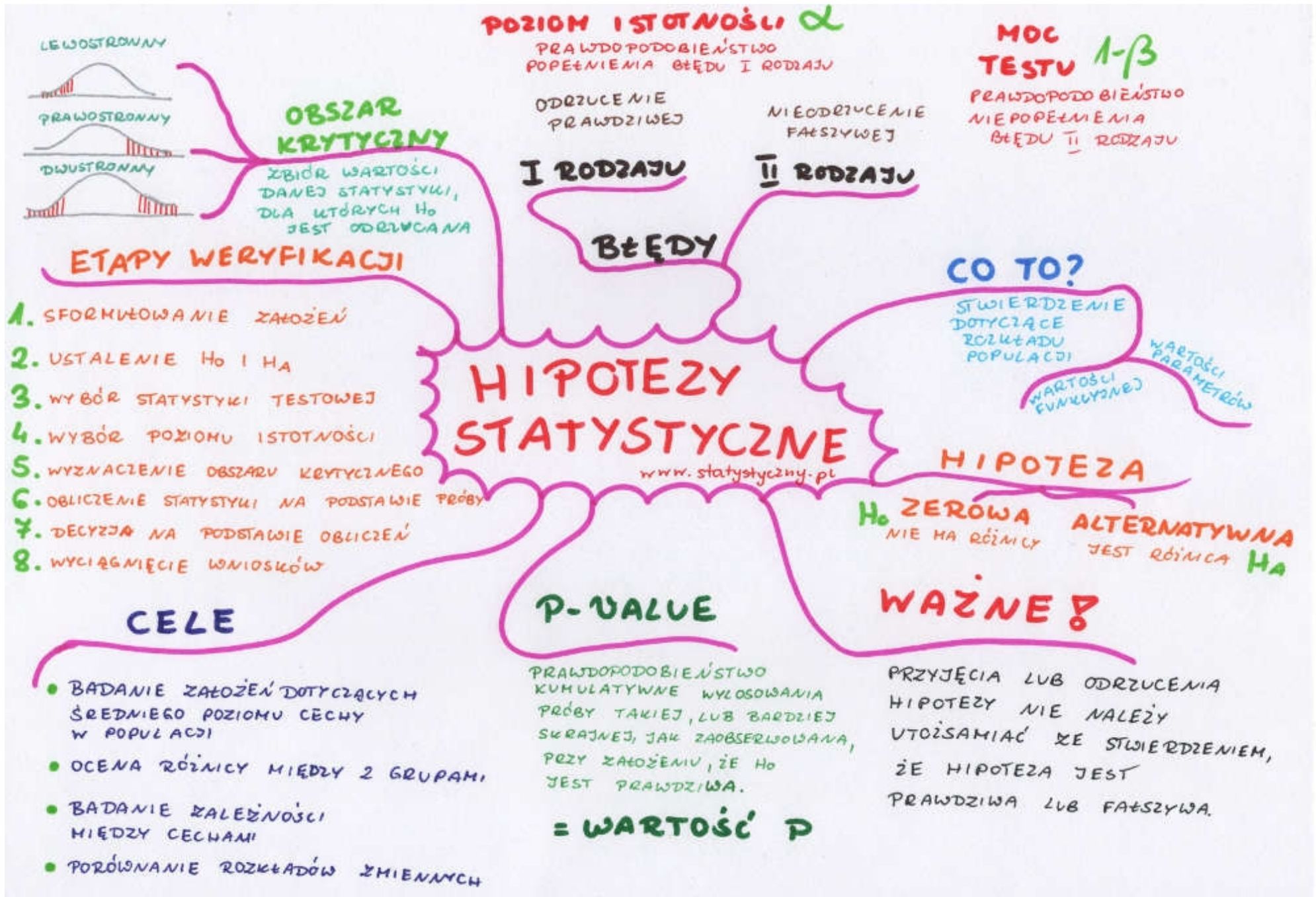
---

- Ograniczenie błędów drugiego rodzaju jest bardzo istotne w niektórych testach. Np. w przypadku medycyny lepiej powiedzieć zdrowemu pacjentowi, żeby zrobił dodatkowe badanie (kiedy w przypadku zerowej hipotezy – pacjent jest zdrowy – wyszło nam błędnie podejrzenie choroby) niż chorego pacjenta odesłać do domu bez leczenia z informacją, że jest zdrowy (błąd drugiego rodzaju).

# Testy statystyczne

		DECYZJA	
		PRZYJAĆ <del><math>H_0</math></del> (NIE ODRZUCIĆ)	ODRZUCIĆ $H_0$
W RZECZYWISTOŚCI	$H_0$ PRAWDZIWA	DECYZJA PRAWIDŁOWA	BŁĄD I RODZAJU
	$H_0$ FAŁSZYWA	BŁĄD II RODZAJU	DECYZJA PRAWIDŁOWA

# Testy statystyczne



# Dalsza lektura

---

- Polecam poniższe strony:
  - <https://www.statystyczny.pl/hipotezy-statystyczne/>
  - <http://pogotowiestatystyczne.pl/istotnosc-statystyczna/>
  - [http://coin.wne.uw.edu.pl/azylicz/sm/sm09\\_2016.pdf](http://coin.wne.uw.edu.pl/azylicz/sm/sm09_2016.pdf)

### Przykładowe zadania na kolokwium z KADD (wykład):

1) Dana jest dystrybuanta zmiennej X:

x	(-inf, -2>	(-2, 1>	(1, 3>	(3, inf)
F(x)	0	0.2	0.8	1.0

Wyznacz funkcję prawdopodobieństwa.

2) Zmienna losowa X ma funkcję prawdopodobieństwa postaci:

x	-3	-1	3	5
p <sub>i</sub>	0.1	0.2	0.5	0.2

Wyznaczyć funkcje prawdopodobieństwa zmienne U:

a)  $U = 2X + 3$

b)  $U = X^3$

c)  $U = X^2 - 5$

Ponadto wyznaczyć wartości oczekiwane i wariancję.

3) Zmienna losowa ma rozkład prawdopodobieństwa:

$$f(x) = c \cdot \sin(x) \quad x \in [0, \pi]$$

a) Dobrać stałą c, aby f(x) była gęstością prawdopodobieństwa

b) Wyznaczyć funkcję dystrybuanty

c) Wyznaczyć wartość oczekiwaną, wariancję, odchylenie standardowe

d) Wyznaczyć medianą

4) Zaproponować metodą do generacji liczb z rozkładu potęgowego

5) Zaproponować metodę do generacji liczb z rozkładu sinusoidalnego

6) Dwuwymiarowa zmienna losowa (X,Y) ma rozkład gęstości:  $f(x,y) = 1/(20\pi) \cdot \exp(-0.5 \cdot (x^2/4 + y^2/25))$ . Zbadać, czy zmienne X, Y są niezależne, Podać sposób na obliczenie (obliczyć):  $P(-1 < X < 2, 0 < Y < 3)$

7) Dane są dwie próby:

A) 21, 19, 14, 27, 25, 23, 22, 18, 21 (N = 9)

B) 16, 24, 22, 21, 25, 21, 18 (N=7)

Czy przy poziomie istotności 5% wariancja próby (B) jest mniejsza niż próby (A)?

8) Zweryfikuj hipotezę, że 10 pomiarów zostało wylosowane z populacji o wartości średniej 25.6? Przyjmij, że poziom istotności wynosi 10%, załóż, że populacji ma rozkład normalny. Pomiary: 22.03, 27.05, 27.5, 22.7, 25.3, 26.2, 27.9, 21.2, 23.3, 25.5.



# Zaliczanie

- W tym roku nie możemy przeprowadzić testu tego typu ani kolokwiów na laboratorium
- **Test wykładowy** będzie na platformie Moodle (30% oceny):
  - 10-20 pytań **zamkniętych**
  - skończony czas na wykonanie testu
  - pytania losowane z bazy pytań (każdy student będzie miał różne pytania w losowej kolejności)
  - Pytania jeszcze nie istnieją – musimy je stworzyć ...
    - ... dopiero wczoraj OKNO PW założyło mi konto na platformie Moodle
- **Projekt** (70% oceny):
  - analiza danych z testem statystycznym (dopasowanie modelu i przeprowadzenie weryfikacji)
    - **przykład**: sprawdzamy dopasowanie funkcji logistycznej do danych zakażeń koronawirusem
    - **inwencja własna**
  - krótki konspekt (max 1 strona) przed realizacją projektu → akceptacja przez prowadzącego
  - krótkie sprawozdanie i kod źródłowy na koniec
  - w ciągu kilki dni przygotowujemy przykładowe sprawozdanie I przykłady projektów
- Warunki zaliczenia:
  - Maksymalnie 2 niezaliczone programy z laboratoriów (odpowiednik nieobecności)
  - zaliczony test Moodle
  - zaliczony projekt



**KONIEC**

# Test t-Studenta

- Mamy zmienną losową  $X$  o rozkładzie normalnym. Pobieramy próbę losową o liczebności  $N$  i wartości średniej  $\bar{X}$
- Wariancja wartości średniej:  $\sigma^2(\bar{X}) = \sigma^2(X) / N$
- Dla dostatecznie dużych prób wartość średnia z próby (na mocy centralnego twierdzenia granicznego) ma rozkład normalny  $(\hat{x}, \sigma(\bar{X}))$
- Zmienna  $y = \frac{\bar{X} - \hat{x}}{\sigma(\bar{X})}$  ma standardowy rozkład normalny
- Na ogół nie znamy jednak odchylenia standardowego  $\sigma^2(X)$
- Posługujemy się estymatorem wariancji:  
$$s_X^2 = \frac{1}{N-1} \sum_{j=1}^N (X_j - \hat{x})^2 \qquad s_{\bar{X}}^2 = \frac{1}{N(N-1)} \sum_{j=1}^N (x_j - \bar{x})^2$$
- **Pytanie:** jak bardzo będziemy odbiegać od rozkładu Gaussa, jeżeli we wzorze na  $y$  zastąpimy odchylenie estymatorem?
- Dla uproszczenia, przyjmiemy, że  $\hat{x} = 0$  (każdy rozkład Gaussa możemy przesunąć o wartość średnią)

# Test t-Studenta

- Rozpatrzmy zmienną losową  $T$  zdefiniowaną następująco:

$$T = \bar{X} / s_{\bar{X}} = \bar{X} \cdot \sqrt{N} / S_X$$

- Wielkość  $(N-1)s_X^2 = fs_X^2$  ma rozkład  $\chi^2$  o liczbie stopni swobody  $f = N-1$
- Wzór na zmienną  $T$  zmieni się nam zatem następująco:

$$T = \bar{X} / s_{\bar{X}} = \bar{X} \cdot \sqrt{N} \cdot \sqrt{f} / \chi$$

- Dystrybuanta zmiennej  $T$  będzie określona wzorem:

$$F(t) = P(T < t) = P\left(\frac{\bar{X} \sqrt{N} \sqrt{f}}{\chi} < t\right) = \frac{\Gamma\left(\frac{1}{2}(f+1)\right)}{\Gamma\left(\frac{1}{2}f\right) \sqrt{\pi} \sqrt{f}} \int_{-\infty}^t \left(1 + \frac{t^2}{f}\right)^{\frac{1}{2}(f+1)} dt$$

- A odpowiadająca jej funkcja gęstości, nosząca nazwę **rozkładu t-Studenta**:

$$f(t) = \frac{\Gamma\left(\frac{1}{2}(f+1)\right)}{\Gamma\left(\frac{1}{2}f\right) \sqrt{\pi} \sqrt{f}} \left(1 + \frac{t^2}{f}\right)^{\frac{1}{2}(f+1)}$$

# Rozkład t-Studenta

- $f=1$
- 2
- 3
- 4
- 5
- 6
- 7
- 8
- 9
- 10
- Jak widać, rozkład t-Studenta jest symetryczny względem 0
- Rozkład dąży do rozkładu Gaussa gdy  $f \rightarrow \infty$
- Z symetrii względem 0 mamy związek (analogicznie jak dla Gaussa):  $P(|t| \leq t) = 2F(|t|) - 1$

- Możemy wyznaczyć graniczne wartości  $\pm t'_\alpha$  odpowiadające poziomowi istotności  $\alpha$  poprzez całkę:

$$\int_0^{t'_\alpha} f(t) dt = \frac{1}{2}(1-\alpha), \text{ gdzie } t'_\alpha = t_{1-\frac{1}{2}\alpha}$$

- Kwantyle  $t'_\alpha = t_{1-\frac{1}{2}\alpha}$  są stacjonarne dla różnych poziomów istotności  $\alpha$  oraz liczby stopni swobody  $f$
- Jest to dwustronny test t-Studenta (jednostronny analogicznie)

# Zastosowanie testu t-Studenta

- **Hipoteza:** zakładamy, że nasza populacja przewiduje wartość oczekiwaną z populacji mającej rozkład normalny równą  $\lambda_0$
- Pobieramy próbę o liczebności  $N$  i wyznaczamy wartość średnią  $\bar{X}$  oraz wariancję  $S_X$
- Jeżeli przy założonym poziomie istotności  $\alpha$  zachodzi nierówność:

$$|t| = \frac{|\bar{X} - \lambda_0| \sqrt{N}}{S_X} > t'_\alpha = t_{1 - \frac{1}{2}\alpha}$$

- Wtedy **odrzucaamy** naszą hipotezę
- W przypadku testu jednostronnego  $t = \frac{\bar{X} - \lambda_0}{S_X} > t_{2\alpha} = t_{1-\alpha}$

# Zastosowanie testu t-Studenta

- Powyższe rozważania możemy **uogólnić** na porównanie wartości średnich dwóch prób losowych z populacji  $X$  oraz  $Y$  o liczebnościach  $N_1$  i  $N_2$

- **Hipoteza:** równość wartości średnich z obu populacji:  $\hat{x} = \hat{y}$

- Zakładamy (z centralnego twierdzenia granicznego), że wartości średnie z prób mają rozkład normalny z wariancjami:

$$\sigma^2(\bar{X}) = \sigma^2(X)/N_1, \quad \sigma^2(\bar{Y}) = \sigma^2(Y)/N_2$$

- Wariancje są estymowane przez estymatory:

$$s_{\bar{X}}^2 = \frac{1}{N_1(N_1-1)} \sum_{j=1}^{N_1} (X_j - \bar{X})^2 \quad s_{\bar{Y}}^2 = \frac{1}{N_2(N_2-1)} \sum_{j=1}^{N_2} (Y_j - \bar{Y})^2$$

- Różnica wartości średnich z próby również ma rozkład zbliżony do normalnego:  $\Delta = \bar{X} - \bar{Y} \Rightarrow \sigma^2(\Delta) = \sigma^2(\bar{X}) + \sigma^2(\bar{Y})$

- Jeśli hipoteza jest prawdziwa, wówczas oczywiste jest, że  $\hat{\Delta} = 0$  oraz iloraz  $\hat{\Delta}/s_{\Delta}$  powinien podlegać rozkładowi Gaussa

- Tak postawiona hipoteza cicho zakłada, że  $X$  i  $Y$  to te same populacje

# Test różnic t-Studenta

- Skoro tak, to oczywiście  $\sigma^2(X)=\sigma^2(Y)$ , zatem można je estymować za pomocą jednego estymatora jako średnią ważoną z dwóch prób:

$$s^2 = \frac{(N_1 - 1)s_X^2 + (N_2 - 1)s_Y^2}{(N_1 - 1) + (N_2 - 1)}$$

- Wtedy możemy zdefiniować estymatory:

$$s_{\bar{X}}^2 = s^2 / N_1, \quad s_{\bar{Y}}^2 = s^2 / N_2, \quad s_{\Delta}^2 = s_{\bar{X}}^2 + s_{\bar{Y}}^2 = \frac{N_1 + N_2}{N_1 \cdot N_2} s^2$$

- Można udowodnić, że zmienna  $\Delta/s(\Delta)$  podlega rozkładowi t-Studenta z liczbą stopni swobody  $f = N_1 + N_2 - 2$
- Równość wartości średnich można więc weryfikować posługując się **testem różnic Studenta**
- $\Delta/s(\Delta)$  obliczana jest na podstawie wyników dwóch prób. Jej wartość bezwzględną porównujemy z kwantylem rozkładu Studenta o liczbie stopni swobody  $f$  dla ustalonego poziomu istotności  $\alpha$ . Sprawdzamy nierówność (**spełniona – odrzucamy hipotezę**):

$$|t| = \frac{|\Delta|}{s_{\Delta}} = \frac{|\bar{X} - \bar{Y}|}{s_{\Delta}} > t'_{\alpha} = t_{1 - \frac{1}{2}\alpha}$$



# Rozkład t-Studenta

- $f=1$
- 2
- 3
- 4
- 5
- 6
- 7
- 8
- 9
- 10
- Jak widać, rozkład t-Studenta jest symetryczny względem 0
- Rozkład dąży do rozkładu Gaussa gdy  $f \rightarrow \infty$
- Z symetrii względem 0 mamy związek (analogicznie jak dla Gaussa):  $P(|t| \leq t) = 2F(|t|) - 1$

- Możemy wyznaczyć graniczne wartości  $\pm t'_\alpha$  odpowiadające poziomowi istotności  $\alpha$  poprzez całkę:

$$\int_0^{t'_\alpha} f(t) dt = \frac{1}{2}(1 - \alpha), \text{ gdzie } t'_\alpha = t_{1 - \frac{1}{2}\alpha}$$

- Kwantyle  $t'_\alpha = t_{1 - \frac{1}{2}\alpha}$  są stabicowane dla różnych poziomów istotności  $\alpha$  oraz liczby stopni swobody  $f$
- Jest to dwustronny test t-Studenta (jednostronny analogicznie)

# Test różnic t-Studenta - przykład

Numer pomiaru	Przyrząd 1	Przyrząd 2				
1	100	97	-0,21	0,04	-3,8	14,44
2	101	101	0,79	0,62	0,2	0,04
3	102	102	1,79	3,2	1,2	1,44
4	100	99	-0,21	0,04	-1,8	3,24
5	98	101	-2,21	4,89	0,2	0,04
6	97	108	-3,21	10,31	7,2	51,84
7	100	101	-0,21	0,04	0,2	0,04
8	101	102	0,79	0,62	1,2	1,44
9	99	96	-1,21	1,47	-4,8	23,04
10	100	101	-0,21	0,04	0,2	0,04
11	98		-2,21	4,89		
12	101		0,79	0,62		
13	100		-0,21	0,04		
14	102		1,79	3,2		
15	103		2,79	7,78		
16	101		0,79	0,62		
17	99		-1,21	1,47		
18	100		-0,21	0,04		
19	102		1,79	3,2		
Ilość pomiarów	19	10				
Średnia	100,21	100,8	-0,59			
Stopnie swobody	18	9				
S <sup>2</sup>	43,16	95,6				
S <sup>2</sup> /f	2,4	10,62				
S <sup>2</sup>	49,1					
S <sup>2</sup> Delta	8,18					

- Mamy kwantyle:

$$t'_{0,2}(27) = t_{0,9}(27) = 1,71$$

$$t'_{0,1}(27) = t_{0,95}(27) = 2,05$$

$$t'_{0,02}(27) = t_{0,99}(27) = 2,77$$

$$t'_{0,01}(27) = t_{0,995}(27) = 3,05$$

$$t'_{0,004}(27) = t_{0,998}(27) = 3,43$$

$$t'_{0,002}(27) = t_{0,999}(27) = 3,69$$

- Hipotezy nie można odrzucić