

Analiza danych

Mariusz Przybycień

Wydział Fizyki i Informatyki Stosowanej
Akademia Górniczo-Hutnicza

Wykład 1

- Statystyka dla fizyków, R.J. Nowak, PWN 2002.
- Statystyka dla fizyków. Ćwiczenia, R.J. Nowak, PWN 2002.
- Data Analysis Techniques for Physical Scientists, C. Pruneau, Cambridge 2017.
- Statistical Methods for Data Analysis in Particle Physics, L. Lista, Springer 2017.
- Statistical Data Analysis, G. Cowan, Oxford 1998.
- Data Analysis in High Energy Physics, O. Behnke, K. Kröninger, G. Schott, T. Schrörner-Sadenius, Wiley-VCH 2013.
- Statistical Analysis Techniques in Particle Physics, I. Narsky, F.C. Porter, Wiley-VCH 2014.
- Probability and Statistics for Particle Physics, C. Mana, Springer 2017.

⇒ <http://home.agh.edu.pl/mariuszp>

⇒ **Wymagania wstępne - znajomość rachunku prawdopodobieństwa i statystyki matematycznej na poziomie kursu dla II roku I stopnia FT na WFiS AGH.**

Analiza danych w fizyce cząstek elementarnych

- Rejestrujemy przypadki zderzeń w akceleratorze (np. $e^{\pm}e^{-}$, $e^{\pm}p$, pp , $p+A$, $A+A$, $e^{\pm}+A$, ...) i dla każdego mierzymy pewien zestaw charakterystyk:
 - pędy cząstek, ich rodzaje, liczbę muonów, energie dżetów, ...
- Rozkłady mierzonych charakterystyk porównujemy z przewidywaniami teoretycznymi (modele teoretyczne + symulacja odpowiedzi detektora w środowisku Geant4).
- Na podstawie tych porównań staramy się oszacować swobodne parametry modeli, ich niepewności (statystyczne i systematyczne):
 - przekroje czynne, masy nowych cząstek, temp. przejścia do stanu QGP, ...
- ... a także ocenić jak dobrze poszczególne modele teoretyczne opisują dane eksperymentalne.
- Staramy się także stwierdzić czy w danych nie występują znaczące odstępstwa od modeli, wskazujące na zaobserwowanie nowych procesów/cząstek, itp.
- Aby to wszystko zrobić w sposób ścisły i ilościowy potrzebne jest podejście wykorzystujące rachunek prawdopodobieństwa i metody statystyczne.
- W dalszej części tego wykładu zakładam znajomość rachunku p-twa i statystyki matematycznej na poziomie kursu dla II roku FT na WFiIS AGH.

Charakterystyki rozkładów prawdopodobieństwa

Rozważamy zmienną losową ciągłą x o funkcji gęstości $f(x)$ oraz funkcję $y = g(x)$.

- Wartość oczekiwana:

$$\mu_x \equiv \mathcal{E}[x] = \int_{-\infty}^{+\infty} x f(x) dx, \quad \mu_y = \mathcal{E}[y] = \int_{-\infty}^{+\infty} g(x) f(x) dx$$

- Wariancja:

$$\sigma_x^2 \equiv \mathcal{V}[x] = \mathcal{E}[(x - \mu_x)^2] = \int_{-\infty}^{+\infty} (x - \mu)^2 f(x) dx = \mathcal{E}[x^2] - \mathcal{E}^2[x]$$

- Momenty rzędu n - zwykły m_n oraz centralny μ_n :

$$m_n \equiv \mathcal{E}[x^n] = \int_{-\infty}^{+\infty} x^n f(x) dx, \quad \mu_n = \mathcal{E}[(x - \mu_x)^n] = \int_{-\infty}^{+\infty} (x - \mu)^n f(x) dx$$

- Moment centralny mieszany $\mu_{r,s}$ rzędu $r + s$ (μ_{11} nazywamy kowariancją):

$$\mu_{r,s} \equiv \mathcal{E}[(x - \mu_x)^r (y - \mu_y)^s] = \iint_{-\infty}^{+\infty} (x - \mu_x)^r (y - \mu_y)^s f(x, y) dx dy$$

- Współczynnik korelacji: $\rho(x, y) = \frac{\text{COV}(x, y)}{\sigma_x \sigma_y}$

Wektor losowy, macierz kowariancji

Niech \vec{x} będzie wektorem losowym (tzn. takim którego składowe są zmiennymi losowymi). Wówczas wektorem wartości oczekiwanych jest $\vec{\mu} = \mathcal{E}[\vec{x}]$ natomiast **macierzą kowariancji** nazywamy:

$$V[\vec{x}] = \mathcal{E} [(\vec{x} - \vec{\mu})(\vec{x} - \vec{\mu})^T]$$

gdzie $v_{ij} = \mathcal{E} [(x_i - \mu_i)(x_j - \mu_j)]$, $i, j = 1, 2, \dots, n$.

W szczególności: $v_{ii} = \mathcal{V}[x_i]$ oraz $v_{ij} = \text{cov}[x_i, x_j]$.

Twierdzenie: Niech będzie dany wektor losowy \vec{x} wraz z wektorem wartości oczekiwanych $\vec{\mu}$ oraz macierzą kowariancji Λ . Wówczas wektor wartości oczekiwanych i macierz kowariancji wektora losowego $\vec{y} = B\vec{x} + \vec{b}$, gdzie B jest macierzą $m \times n$, a \vec{b} stałym wektorem, dane są przez:

$$\mathcal{E}[\vec{y}] = B\vec{\mu} + \vec{b} \quad \text{oraz} \quad V[\vec{y}] = BV[\vec{x}]B^T$$

Dowód:

$$\mathcal{E}[\vec{y}] = B\mathcal{E}[\vec{x}] + \vec{b} = B\vec{\mu} + \vec{b}$$

$$\begin{aligned} V[\vec{y}] &= \mathcal{E} [(\vec{y} - \mathcal{E}[\vec{y}])(\vec{y} - \mathcal{E}[\vec{y}])^T] = \mathcal{E} [B(\vec{x} - \vec{\mu})(\vec{x} - \vec{\mu})^T B^T] = \\ &= B\mathcal{E} [(\vec{x} - \vec{\mu})(\vec{x} - \vec{\mu})^T] B^T = BV[\vec{x}]B^T \end{aligned}$$

Rozkład jednorodny (płaski, jednostajny)

- Dyskretny: $P_k = \frac{1}{n}$, $k = 1, \dots, n$, $\mathcal{E}[k] = \frac{1}{2}(n+1)$, $\mathcal{V}[k] = \frac{n^2-1}{12}$
- Ciągły: $f(x) = \frac{1}{b-a}$, $a < x < b$, $\mathcal{E}[x] = \frac{1}{2}(b+a)$, $\mathcal{V}[x] = \frac{(b-a)^2}{12}$
- Dystrybuanta dowolnej zmiennej losowej ma rozkład płaski na przedziale $[0, 1]$:

$$F(x) = \int_{-\infty}^x f(x) dx \quad \left\{ \begin{array}{l} y = F(x) \\ \Rightarrow \end{array} \right. \quad g(y) = f(x(y)) \left| \frac{dx}{dy} \right| = f(x) \frac{1}{f(x)} = 1$$

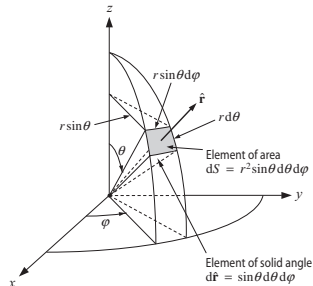
- Rozkład jednorodny na sferze o promieniu $r = 1$:

$$P(\Omega) d\Omega = \frac{1}{4\pi} \sin \theta d\theta d\varphi = P(\theta, \varphi) d\theta d\varphi$$

$$P(\theta, \varphi) = \frac{\sin \varphi}{4\pi} \Rightarrow \begin{cases} P(\varphi) = \int_0^\pi P(\theta, \varphi) d\theta = \frac{1}{2\pi} \\ P(\theta) = \int_0^{2\pi} P(\theta, \varphi) d\varphi = \frac{\sin \theta}{2} \end{cases}$$

$$u = F(\varphi) = \frac{\varphi}{2\pi} \Rightarrow \varphi = 2\pi u$$

$$v = F(\theta) = \frac{1 - \cos \theta}{2} \Rightarrow \theta = \cos^{-1}(2v - 1)$$

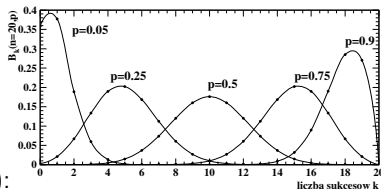


Rozkłady dwu- i wielomianowy oraz geometryczny

- Rozkład dwumianowy zmiennej losowej k :

$$\mathcal{B}_k(n, p) = \binom{n}{k} p^k (1-p)^{n-k}, \quad k = 0, 1, \dots, n$$

$$\mathcal{E}[k] = np, \quad \mathcal{V}[k] = np(1-p)$$



- Rozkład wielomianowy ($k_1 + k_2 + \dots + k_j = n$):

$$\mathcal{W}_{k_1 k_2 \dots k_j}(n, p_1, p_2, \dots, p_j) = \frac{n!}{k_1! k_2! \dots k_j!} p_1^{k_1} p_2^{k_2} \dots p_j^{k_j}$$

$$\mathcal{E}[k_i] = np_i, \quad \mathcal{V}[k_i] = np_i(1-p_i)$$

- Rozkład geometryczny:

$$\mathcal{G}_k(p) = p(1-p)^{k-1}, \quad k = 1, 2, \dots \quad \mathcal{E}[k] = \frac{1}{p}, \quad \mathcal{V}[k] = \frac{1-p}{p^2}$$

- Rozkład ujemny dwumianowy:

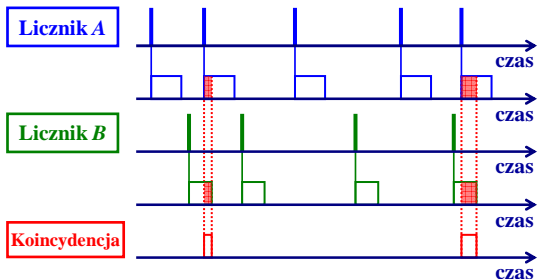
$$\mathcal{U}_k(n, p) = \binom{k-1}{n-1} p^n (1-p)^{k-n}, \quad n = 1, 2, \dots \quad k = n, n+1, \dots$$

$$\mathcal{E}[k] = \frac{n}{p}, \quad \mathcal{V}[k] = \frac{n(1-p)}{p^2}$$

Koincydencje przypadkowe

Rozważmy dwa liczniki Geigera-Müllera (A i B) umieszczone w poziomie w pewnej odległości od siebie i rejestrujące promieniowanie kosmiczne.

- f_A, f_B - typowe częstotliwości rejestracji promieni przez każdy z liczników.
- $\tau_A = 1/f_A, \tau_B = 1/f_B$ - typowe odstępy czasu pomiędzy kolejnymi cząstkami.
- Impulsy wyjściowe mają standardowy kształt prostokątny o czasach trwania T_A i T_B które są znacznie krótsze niż τ_A i τ_B .
- Impulsy wyjściowe wysyłamy na układ koincydencyjny, który daje sygnał gdy impulsy wejściowe przekrywają się choć minimalnie.



Niech czas obserwacji $\tau \gg \tau_A, \tau_B$, wtedy p-two sygnału na wyjściu detektorów wynosi odpowiednio:

$$\frac{\tau}{T_A} T_A = \tau f_A T_A \Rightarrow P_A = f_A T_A$$

$$\frac{\tau}{T_B} T_B = \tau f_B T_B \Rightarrow P_B = f_B T_B$$

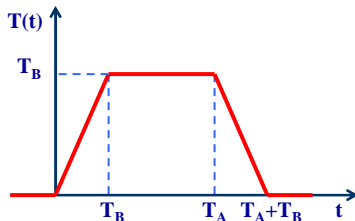
Zdarzenia są niezależne, więc:

$$P = P_A P_B = f_A T_A f_B T_B$$

Koincydencje przypadkowe

- Oszacowanie średniego czasu trwania pojedynczej koincydencji ($T_A > T_B$):

$$T(t) = \begin{cases} t & 0 \leq t \leq T_B \\ T_B & T_B \leq t \leq T_A \\ T_A + T_B - t & T_A \leq t \leq T_A + T_B \end{cases}$$



$$\bar{T} = \frac{1}{T_A + T_B} \left(\int_0^{T_B} t dt + \int_{T_B}^{T_A} T_B dt + \int_{T_A}^{T_A+T_B} (T_A + T_B - t) dt \right) = \frac{T_A T_B}{T_A + T_B}$$

- Częstość koincydencji przypadkowych dana jest przez:

$$f_{AB} = \frac{P\tau/\bar{T}}{\tau} = f_A T_A f_B T_B \left(\frac{1}{T_A} + \frac{1}{T_B} \right) = f_A f_B (T_A + T_B)$$

- Podobnie, dla potrójnych koincydencji otrzymujemy:

$$\frac{1}{\bar{T}_{ABC}} = \frac{1}{\bar{T}_{AB}} + \frac{1}{T_C} = \frac{1}{T_A} + \frac{1}{T_B} + \frac{1}{T_C}$$

$$f_{ABC} = f_A T_A f_B T_B f_C T_C \left(\frac{1}{T_A} + \frac{1}{T_B} + \frac{1}{T_C} \right)$$

Rozkład wykładniczy i Weibulla

- Rozkład wykładniczy opisuje czas oczekiwania na zdarzenie:

$$\text{Exp}(t; \lambda) = \lambda e^{-\lambda t} \quad \text{lub} \quad \text{Exp}(t; \tau) = \frac{1}{\tau} e^{-t/\tau} \quad \text{gdzie} \quad \tau = \frac{1}{\lambda}$$

$$\mathcal{E}[t] = \tau = \frac{1}{\lambda} \quad \mathcal{V}[t] = \tau^2 = \frac{1}{\lambda^2}$$

- Częstość koincydencji przypadkowych:

λ_A, λ_B - intensywności rejestracji promieniowania przez liczniki A i B ,

T_A, T_B - czasy trwania impulsów generowanych przez te liczniki.

Niech koincydencja będzie inicjowana przez licznik A , wtedy:

$$f(t_A, t_B; \lambda_A, \lambda_B) = \lambda_B (\lambda_A + \lambda_B) \exp(-\lambda_A t_A - \lambda_B t_B) \quad 0 \leq t_A \leq t_B < \infty$$

P-two koincydencji:

$$P(t_A - t_B \leq T_A) = \lambda_B (\lambda_A + \lambda_B) \int_0^{\infty} dt_A e^{-\lambda_A t_A} \int_{t_A}^{t_A + T_A} dt_B e^{-\lambda_B t_B} = 1 - e^{-\lambda_B T_A}$$

Łączna liczba koincydencji przypadkowych inicjowanych przez detektory A i B :

$$f_{AB} = \lambda_A (1 - e^{-\lambda_B T_A}) + \lambda_B (1 - e^{-\lambda_A T_B}) \approx \lambda_A \lambda_B (T_A + T_B)$$

- Rozkład Weibulla: $f(x; \alpha, \sigma) = \frac{\alpha}{\sigma} \left(\frac{x}{\sigma}\right)^{\alpha-1} \exp\left(-\left(\frac{x}{\sigma}\right)^{\alpha}\right)$

Rozkład Poissona

- Rozkład Poissona opisuje p -two wystąpienia k zdarzeń w zadanym przedziale czasu Δt jeśli czas oczekiwania na każde kolejne zdarzenie podlega rozkładowi wykładniczemu $Exp(\lambda)$.

- Rozkład Poissona:

$$\mathcal{P}_k(\mu) = \frac{\mu^k}{k!} e^{-\mu}, \quad \mu > 0, \quad k = 0, 1, \dots$$

$$\mathcal{E}[k] = \mu, \quad \mathcal{V}[k] = \mu$$

Dla $n \rightarrow \infty, p \rightarrow 0, np \rightarrow \mu$:

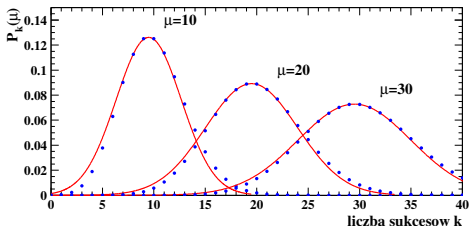
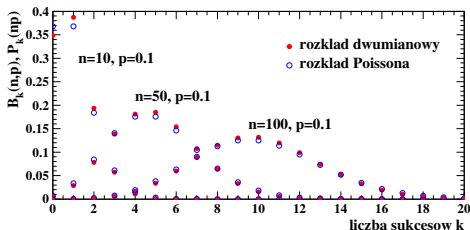
$$\mathcal{B}_k(n, p) \longrightarrow \mathcal{P}_k(\mu)$$

- Parametr $\mu = \lambda \Delta t$ określa średnią liczbę zdarzeń w przedziale czasu Δt .
- Rozkład Poissona może także odnosić się do zdarzeń w przestrzeni (λ wtedy oznacza gęstość $[1/m^n]$).

- Dla dużych wartości μ mamy:

$$\mathcal{P}_k(\mu) \rightarrow \frac{1}{\sqrt{2\pi\sigma}} e^{-\frac{(k - (\mu - 0.5))^2}{2\sigma^2}}$$

gdzie $\sigma = \sqrt{\mu}$.



Rozkład normalny (Gaussa)

- Rozkład normalny (Gaussa):

$$\mathcal{N}(x; \mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right)$$

$$\mathcal{E}[x] = \mu, \quad \mathcal{V}[x] = \sigma^2$$

- Dystrybuanta standaryzowanego rozkładu Gaussa:

$$\Phi(x) = \frac{1}{2} \left[1 + \operatorname{erf}\left(x/\sqrt{2}\right) \right]$$

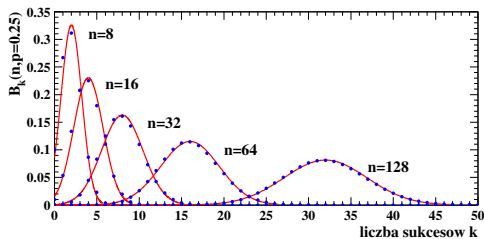
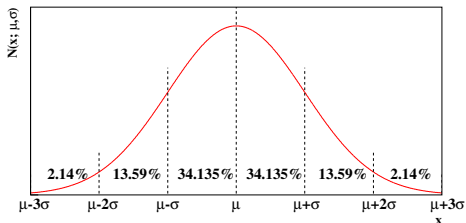
gdzie $\operatorname{erf}(x)$ to funkcja błędu:

$$\operatorname{erf}(x) = \frac{2}{\sqrt{\pi}} \int_0^x e^{-t^2} dt$$

- Dwuwymiarowy rozkład normalny:

$$\mathcal{N}(x, y; \mu_x, \mu_y, \sigma_x, \sigma_y) =$$

$$= \frac{1}{2\pi \sigma_x \sigma_y \sqrt{1-\rho^2}} \exp\left\{-\frac{1}{2(1-\rho^2)} \left[\left(\frac{x-\mu_x}{\sigma_x}\right)^2 + \left(\frac{y-\mu_y}{\sigma_y}\right)^2 - 2\rho \left(\frac{x-\mu_x}{\sigma_x}\right) \left(\frac{y-\mu_y}{\sigma_y}\right) \right]\right\}$$



Rozkład logarytmiczno-normalny

- Wielowymiarowy (m) rozkład normalny:

$$\mathcal{N}(\vec{x}; \vec{\mu}, V) = \frac{1}{(2\pi)^{m/2} |V|^{1/2}} \exp \left[-\frac{1}{2} (\vec{x} - \vec{\mu})^T V^{-1} (\vec{x} - \vec{\mu}) \right]$$

- Rozkład logarytmiczno-normalny:

$$p_{\ln}(x; \mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} \frac{1}{x} \exp \left(-\frac{(\ln x - \mu)^2}{2\sigma^2} \right), \quad x > 0$$

$$\mathcal{E}[x] = \exp \left(\mu + \frac{1}{2}\sigma^2 \right), \quad \mathcal{V}[x] = [\exp(\sigma^2) - 1] \exp(2\mu + \sigma^2)$$

$$F(x) = \Phi \left(-\frac{\ln x - \mu}{\sigma} \right)$$

- Zastosowanie rozkładu log-normalnego do iloczynów wielu zm. losowych:

Rozważmy n zmiennych losowych z_1, z_2, \dots, z_n oraz ich iloczyn $x = \prod_{i=1}^n z_i$.

Definiujemy zm. losowe w_i poprzez $z_i = e^{w_i}$ oraz $y = \sum_{i=1}^n w_i$

$$\text{Ponieważ } x = e^y \quad \Rightarrow \quad p(x) = p(y(x)) \left| \frac{dy}{dx} \right| = p(y(x)) \frac{1}{x}$$

Dla dużego n , na podstawie CTG, zmienna y ma rozkład normalny, a więc zmienna x ma rozkład log-normalny.

Centralne Twierdzenie Graniczne

- Niech będzie dany ciąg niezależnych zmiennych losowych x_1, x_2, \dots, x_n pochodzących z dowolnego rozkładu o skończonej wartości oczekiwanej μ oraz wariancji σ^2 . Rozkład zmiennej losowej:

$$z_n = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} \quad \text{gdzie} \quad \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

dąży do standaryzowanego rozkładu normalnego.

