

Metody uczenia maszynowego w naukach przyrodniczych

Artur Kalinowski

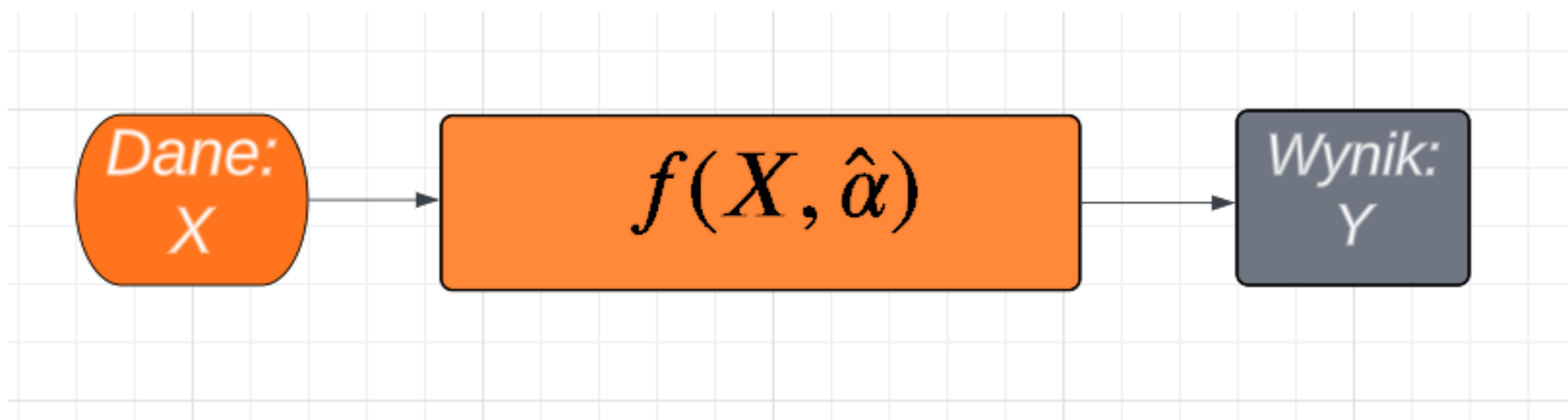
Uniwersytet Warszawski,
Wydział Fizyki

Dane:

\mathbf{R}^n – n-wymiarowa przestrzeń danych wejściowych: X

\mathbf{R}^k – k-wymiarowa przestrzeń danych wyjściowych: Y

Szukane:



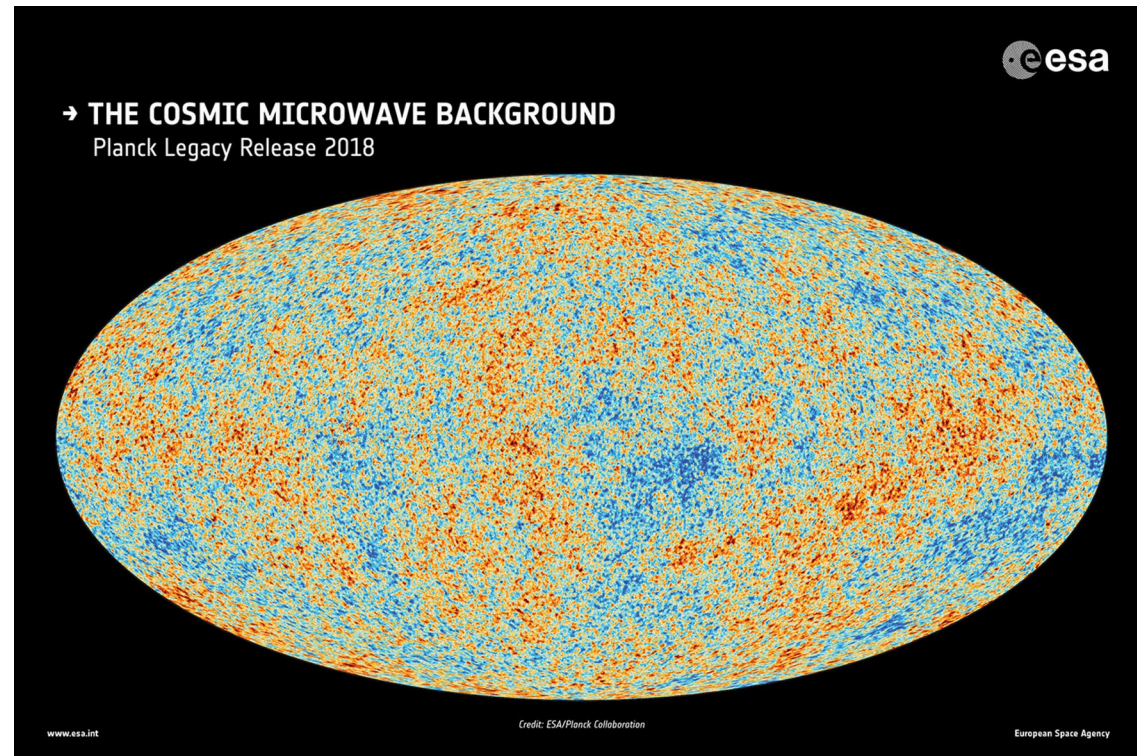
- baza w której szukamy współczynników rozwinięcia funkcji opisującej dane
- funkcje bazy są określone przez parametry α

$$f(X, \alpha)$$

Przykład:

- harmoniki sferyczne - baza funkcji zdefiniowanych na sferze:

$$T(\theta, \varphi) = \sum_{l,m} a_{l,m} Y_l^m(\theta, \varphi)$$



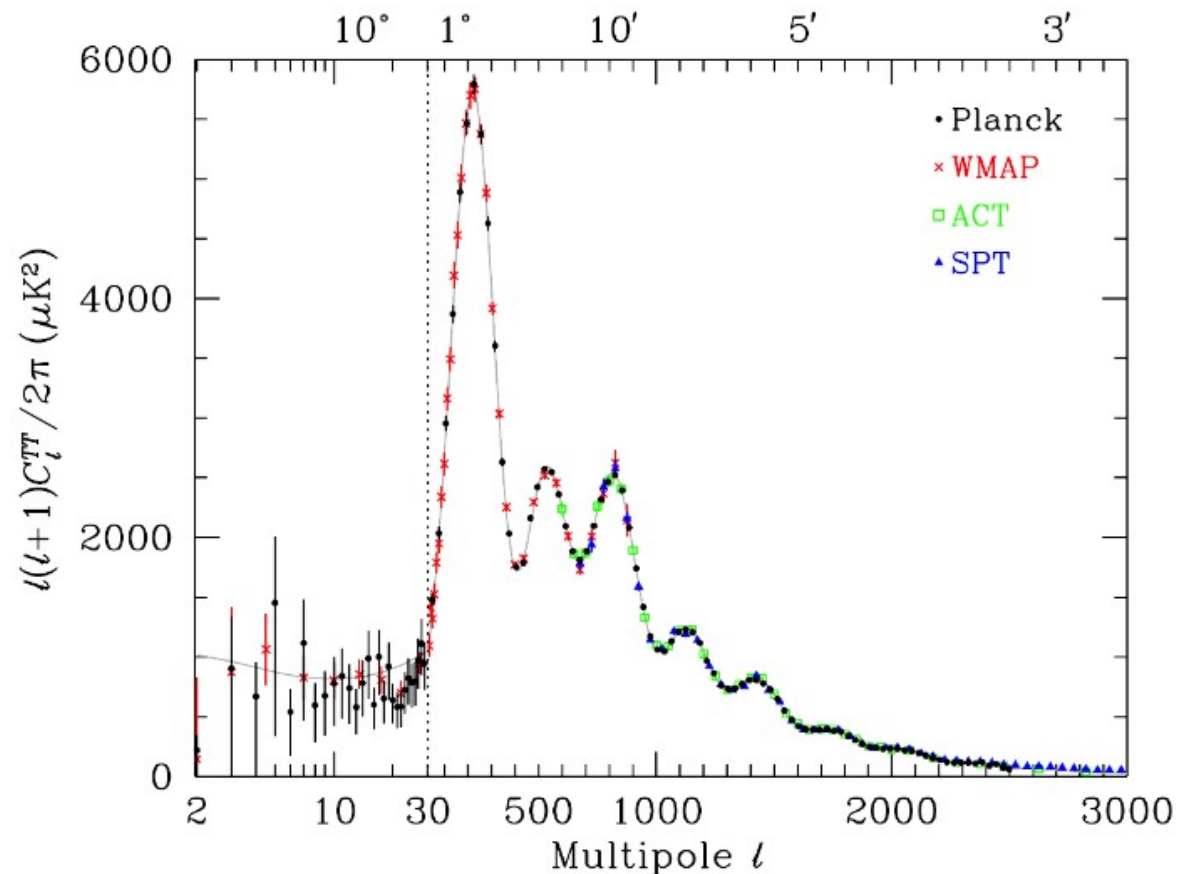
- baza w której szukamy współczynników rozwinięcia funkcji opisującej dane
- funkcje bazy są określone przez parametry α

Przykład:

- harmoniki sferyczne - baza funkcji zdefiniowanych na sferze:

$$T(\theta, \varphi) = \sum_{l,m} a_{l,m} Y_l^m(\theta, \varphi)$$

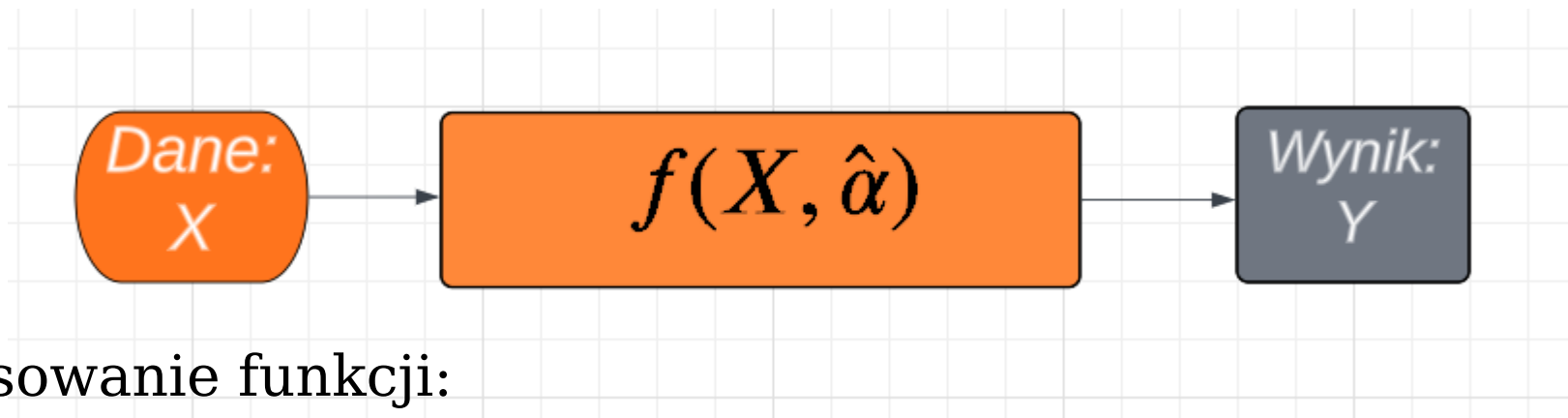
$$f(X, \alpha)$$



R.L. Workman et al. (Particle Data Group), Prog. Theor. Exp. Phys. 2022, 083C01 (2022) and 2023 update

X: n-wymiarowa przestrzeń danych wejściowych

Y: k-wymiarowa przestrzeń danych wyjściowych

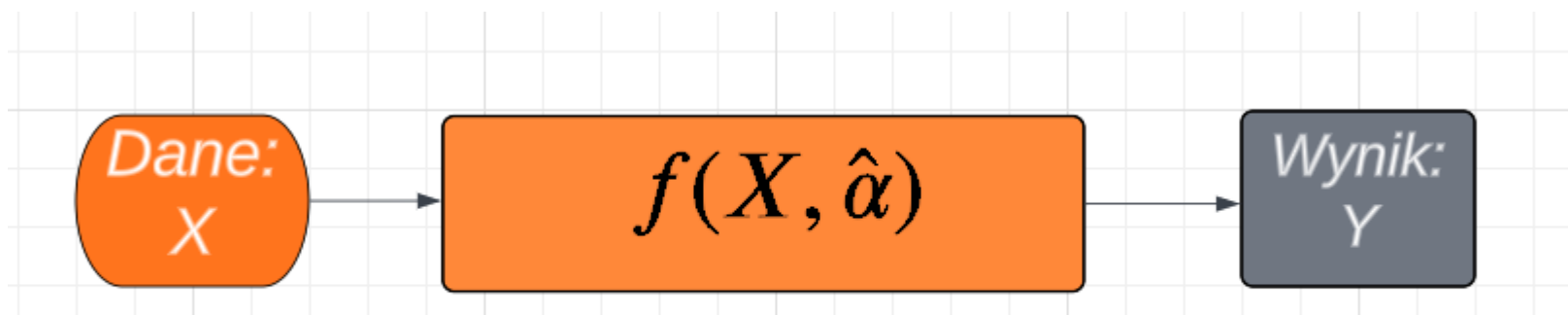


Dopasowanie funkcji:

- **$n \sim 2$, $k \sim 1$**
- **postać funkcji bazy zadana *explicite***
- **współczynniki rozwinięcia mogą być interpretowalne**

X: n-wymiarowa przestrzeń danych wejściowych

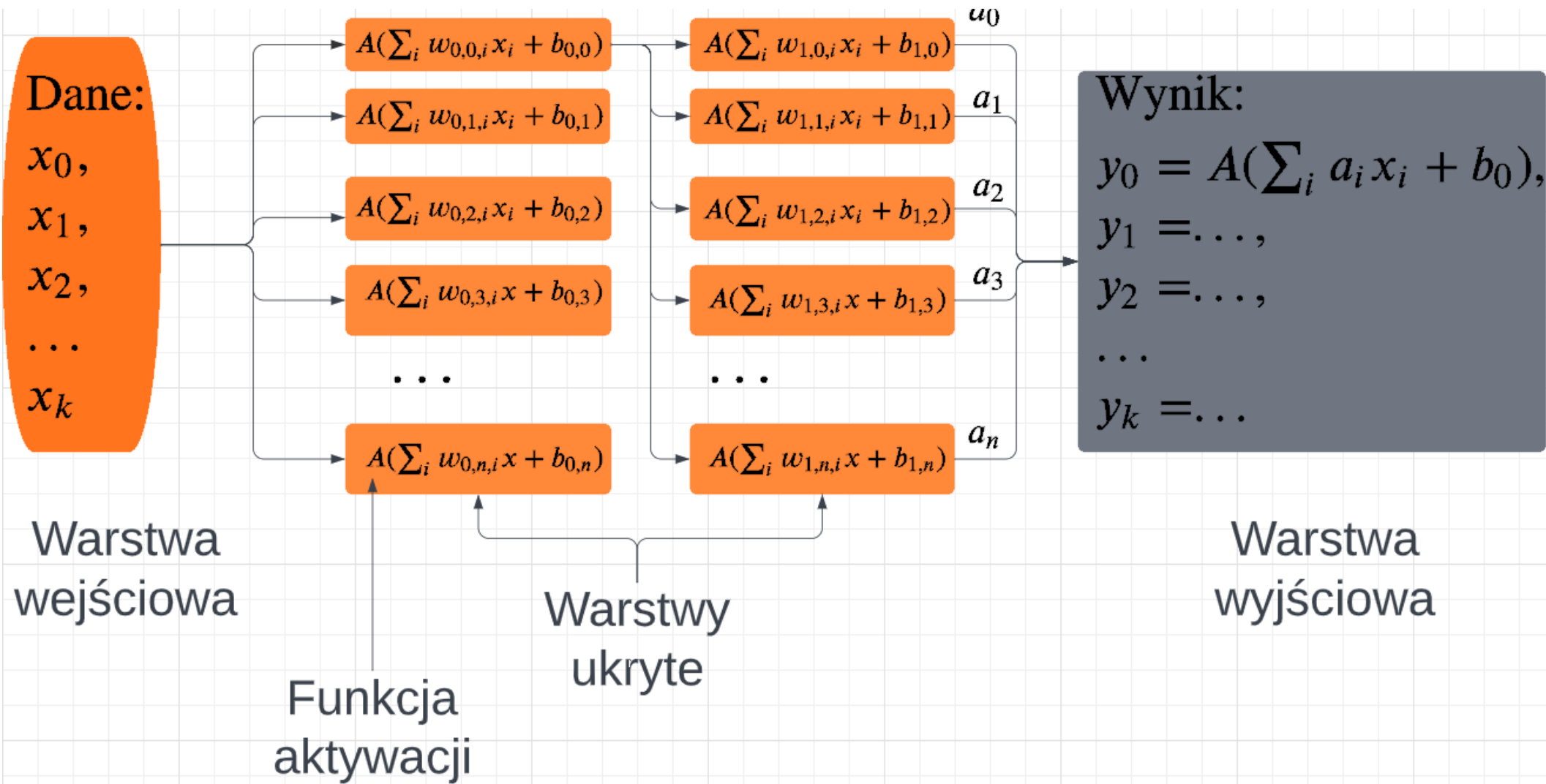
Y: k-wymiarowa przestrzeń danych wyjściowych



Uczenie maszynowe:

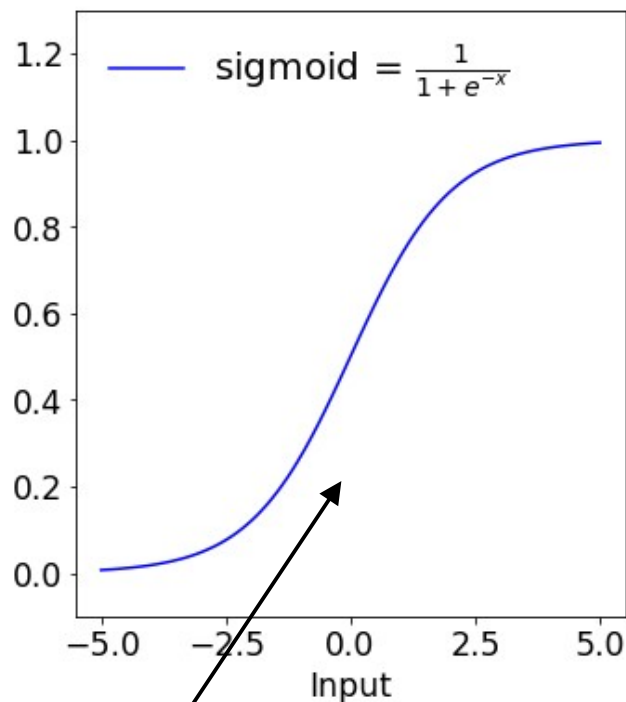
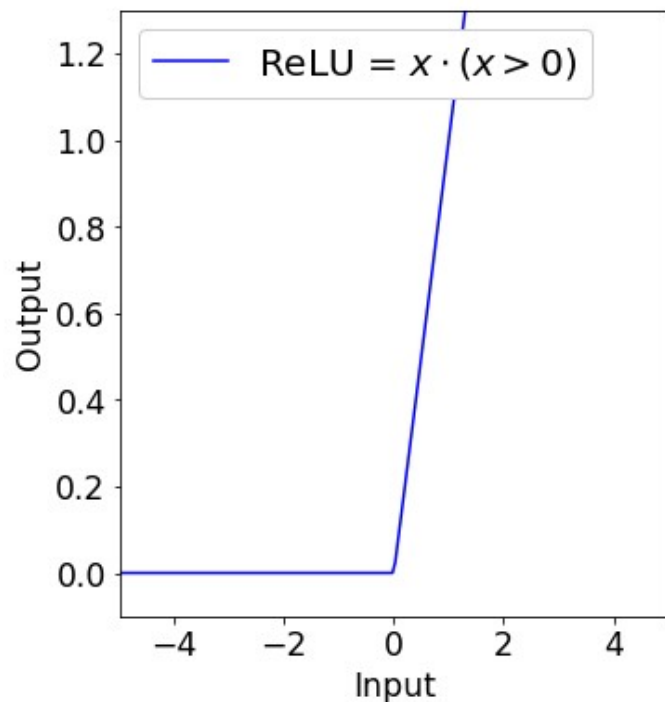
- $n \sim 10^l$, $k \sim 10^m$, $l, m \sim 6$
- baza funkcji zdefiniowana *implicitnie* przez strukturę przepływu danych - architekturę
- współczynniki rozwinięcia nie są interpretowalne

Klasyczna sieć neuronowa

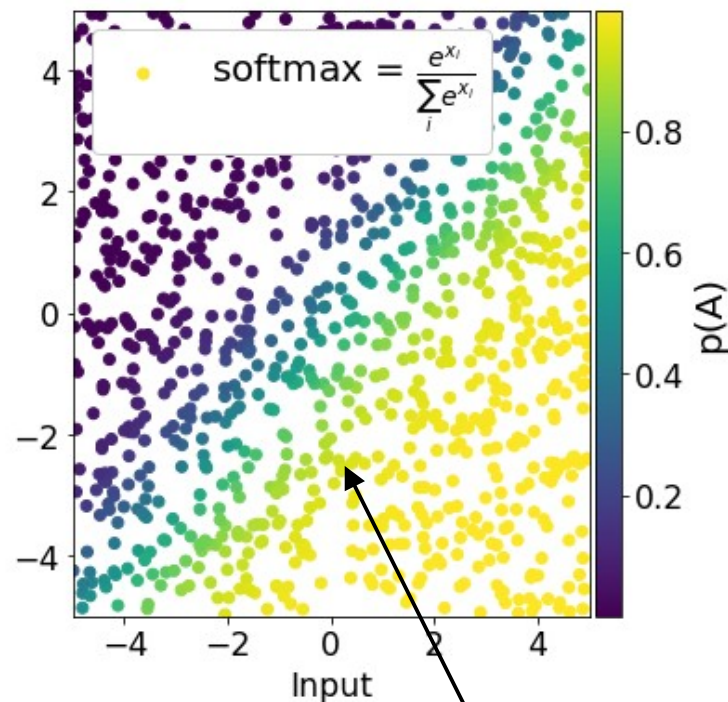


standard dla warstw ukrytych

standard dla warstwy wyjściowej:

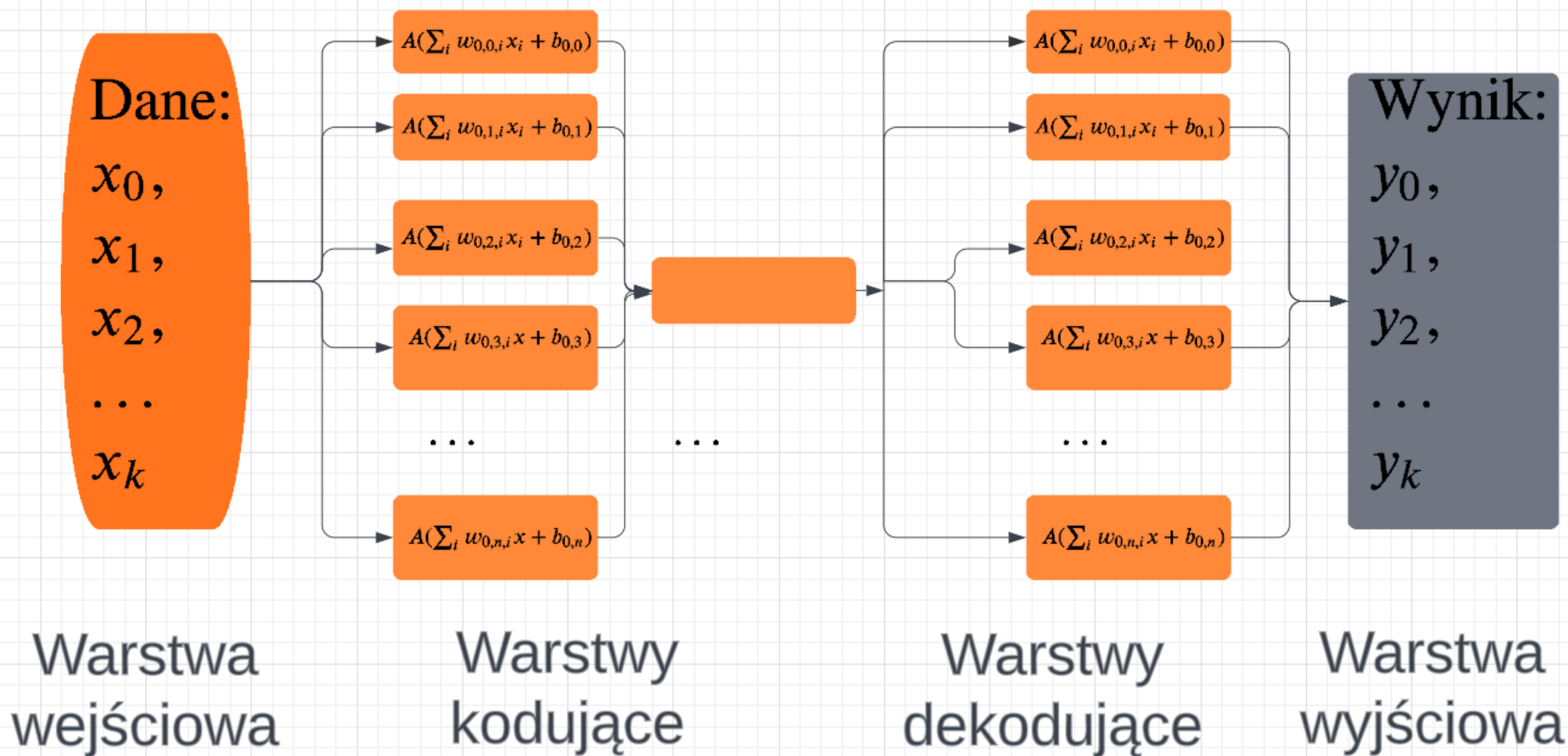


pojedyncze
prawdopodobieństwo

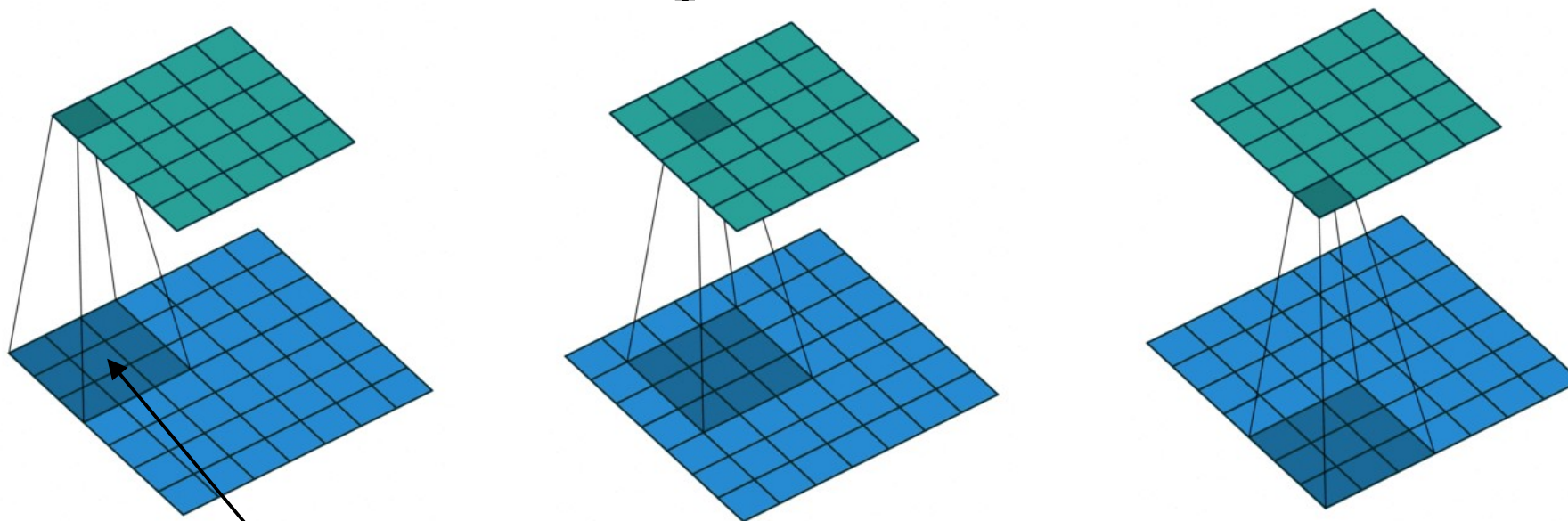


prawdopodobieństwo
dla wielu wariantów

Sieć kodująca



Sieć neuronowa ze splotem



dane z okienka
3x3 są sumowane
z tymi samymi
wagami dla
całego obrazu

3 ₀	3 ₁	2 ₂	1	0
0 ₂	0 ₂	1 ₀	3	1
3 ₀	1 ₁	2 ₂	2	3
2	0	0	2	2
2	0	0	0	1

12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0

3	3 ₀	2 ₁	1 ₂	0
0	0 ₂	1 ₂	3 ₀	1
3	1 ₀	2 ₁	2 ₂	3
2	0	0	2	2
2	0	0	0	1

12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0

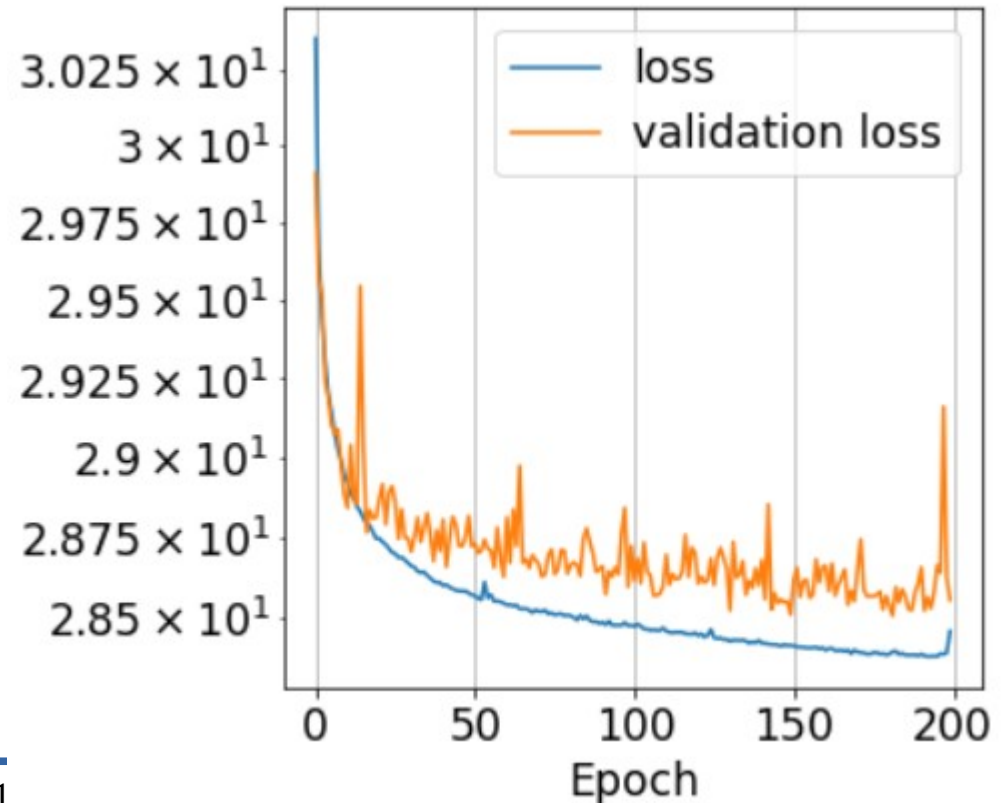
arXiv:1603.07285 [stat.ML]

- parametry znajdowane w iteracyjnym procesie minimalizacji **funkcji kosztu:**

$$L(f(X, \alpha), Y)$$

$$\alpha = \operatorname{argmin}_{\alpha} L(f(X, \alpha), Y)$$

- znajdowanie parametrów modelu jest określane mianem **treningu**



- **zagadnienie regresji:** $f(X)$ zwraca dowolną, ciągłą, wartość

$$L(f(X, \alpha), Y) = \frac{1}{N} \sum_i (f(X, \alpha) - Y)^2 + \lambda J(\alpha)$$

- **zagadnienie klasyfikacji:** $f(X)$ zwraca prawdopodobieństwo

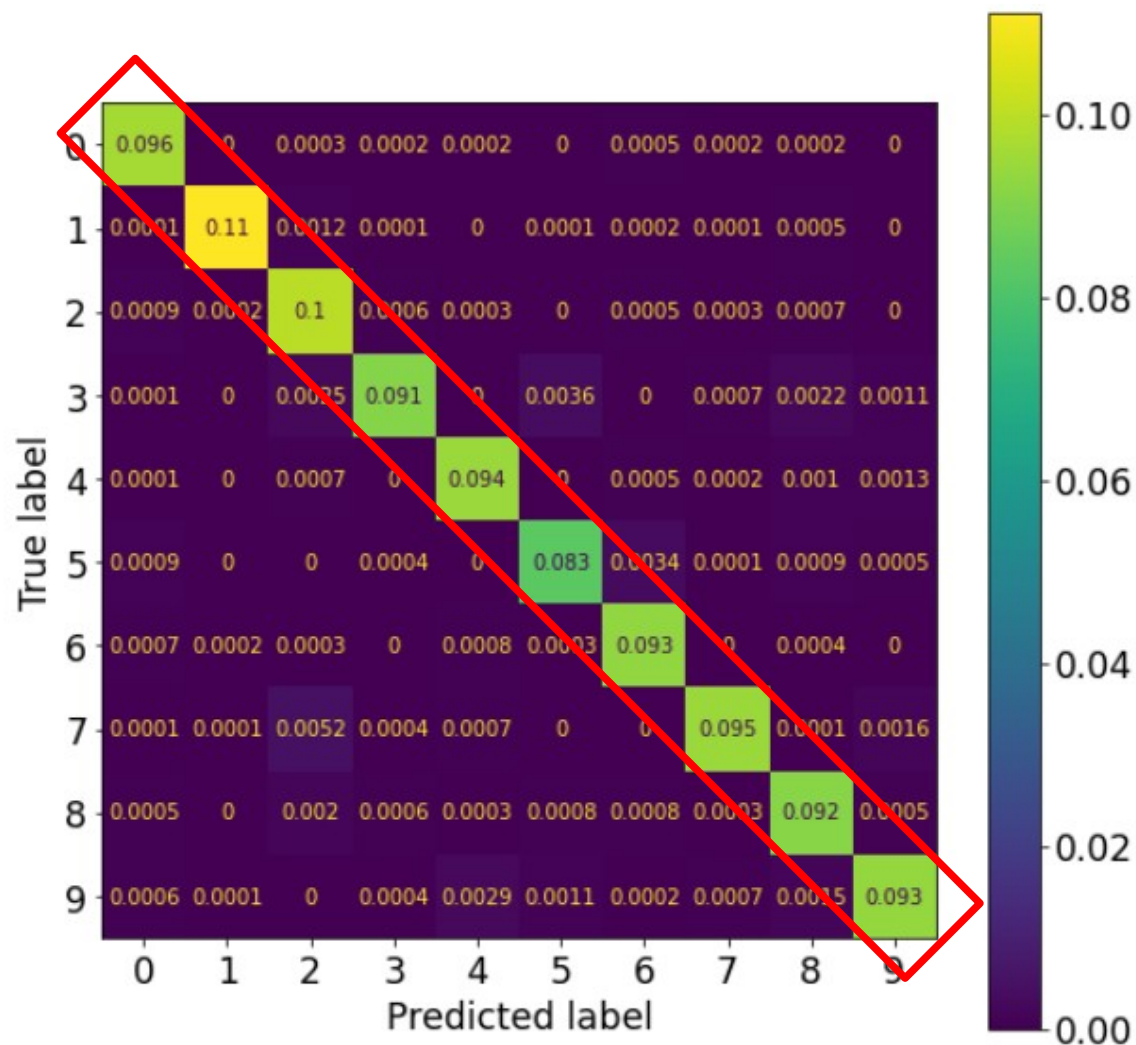
$$L(f(X, \alpha), Y) = \frac{1}{N} \sum_i \log(f_{\text{poprawna klasa}}(X, \alpha)) + \dots$$

- poprawność modeli klasyfikacyjnych jest szacowana przy użyciu macierzy pomyłek, oraz **metryk** - liczb opisujących zgodność modelu z danymi, np.:

$$\epsilon = \frac{\textit{poprawne}}{\textit{normalizacja}}$$

- accuracy:**

normalizacja -
wszystkie przykłady

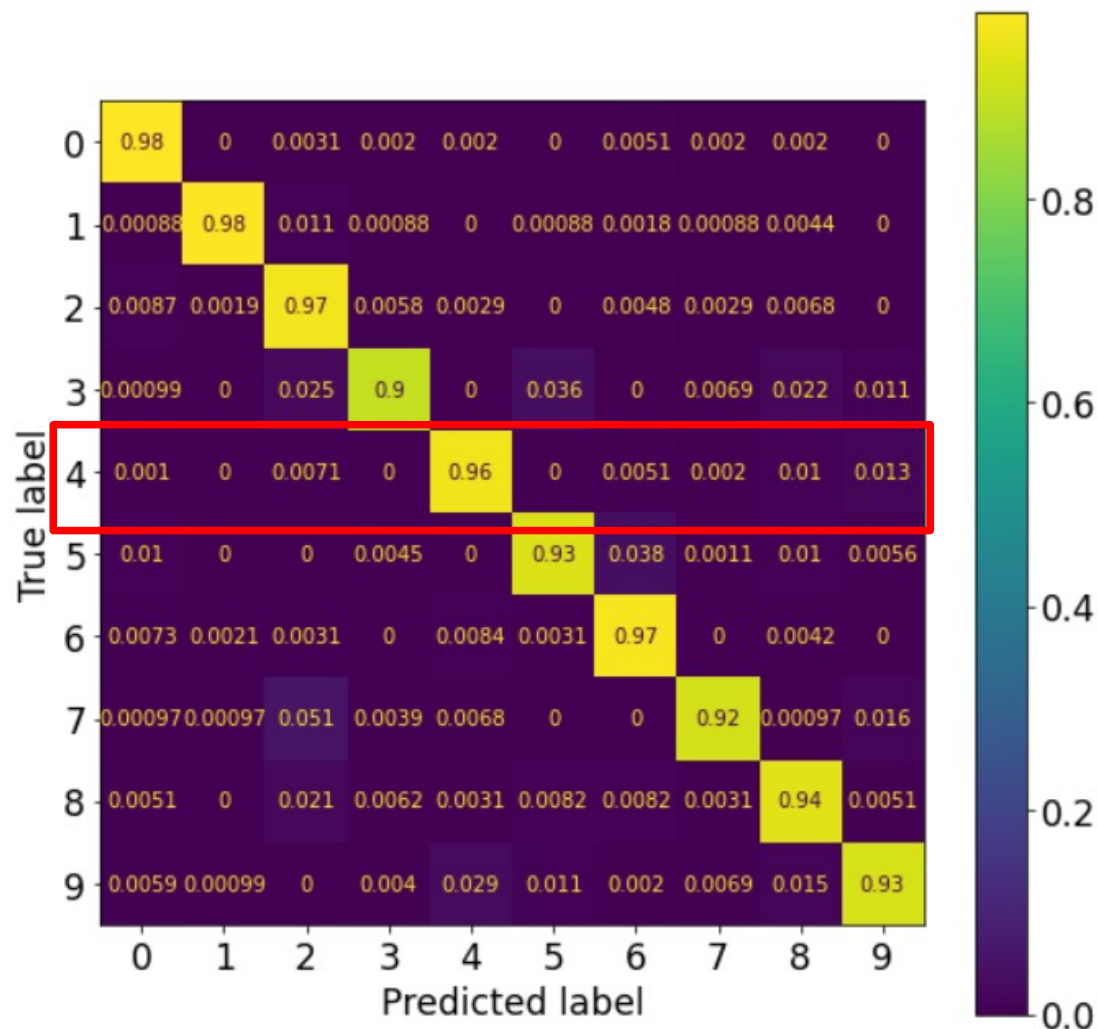


- poprawność modeli klasyfikacyjnych jest szacowana przy użyciu macierzy pomyłek, oraz **metryk** - liczb opisujących zgodność modelu z danymi, np.:

$$\epsilon = \frac{\textit{poprawne}}{\textit{normalizacja}}$$

- recall:**

normalizacja - przykłady z danej klasy

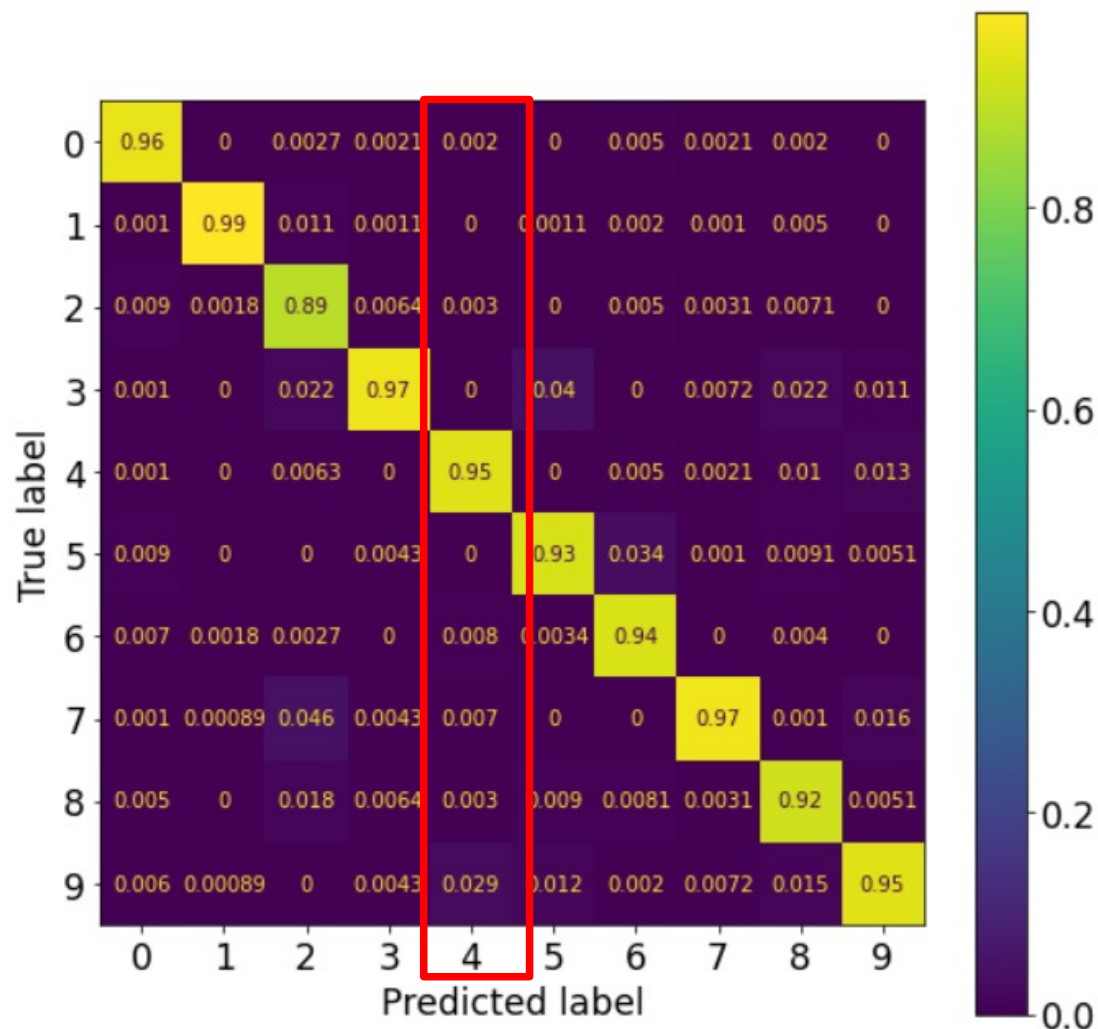


- poprawność modeli klasyfikacyjnych jest szacowana przy użyciu macierzy pomyłek, oraz **metryk** - liczb opisujących zgodność modelu z danymi, np.:

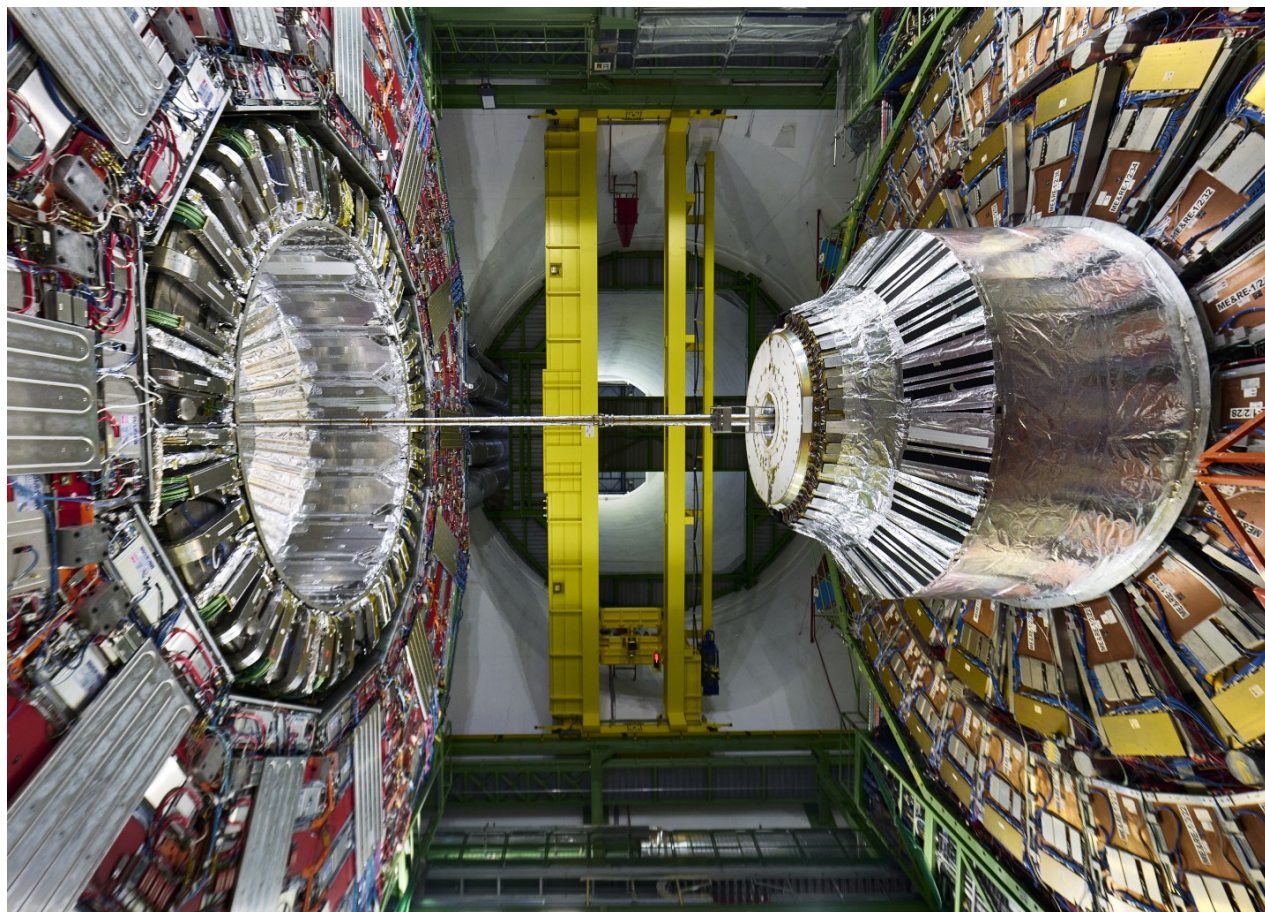
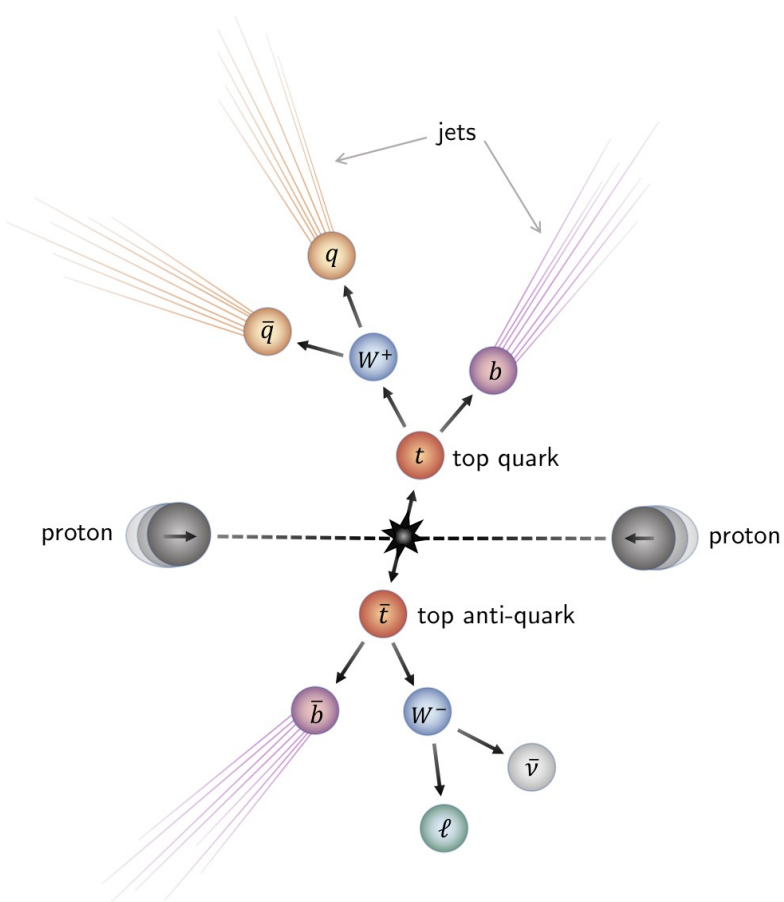
$$\epsilon = \frac{\textit{poprawne}}{\textit{normalizacja}}$$

- precision:**

normalizacja - przykłady
którym model przyznał
daną klasę



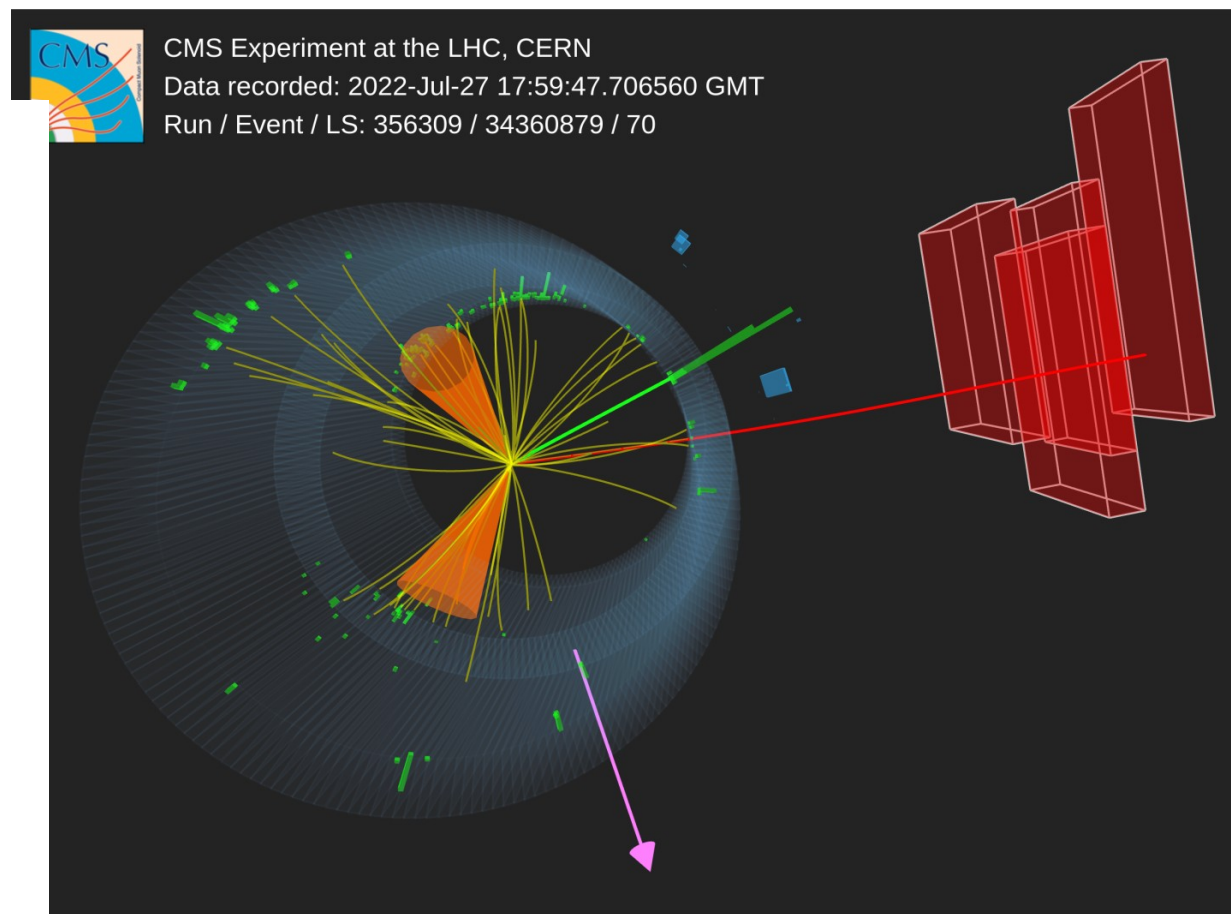
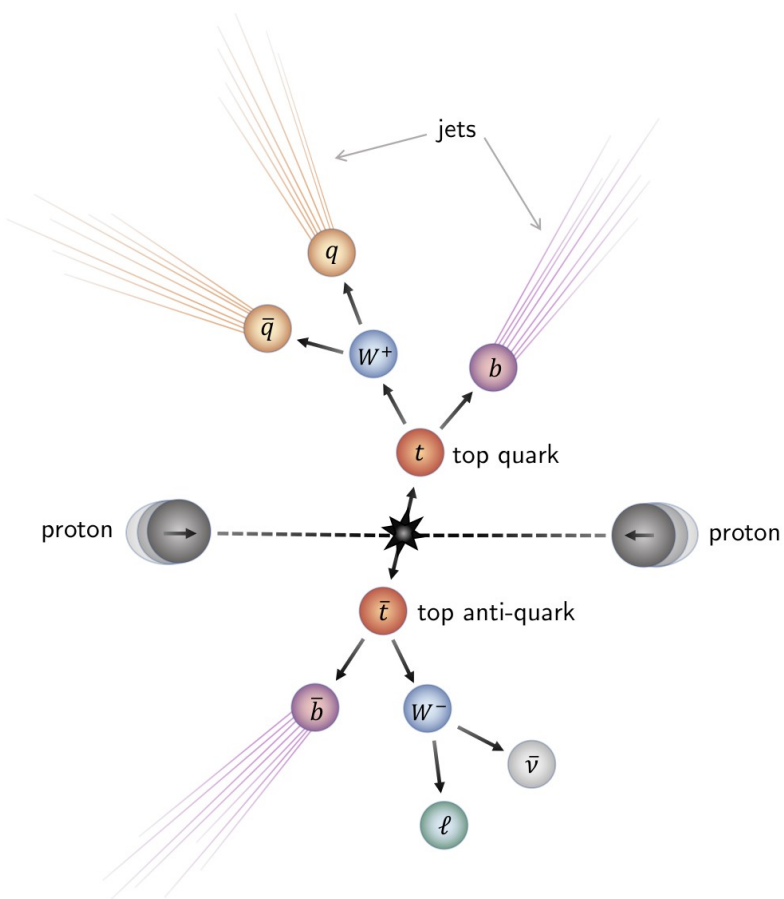
Eksperyment CMS działający przy LHC (CERN)



CERN/CMS Collaboration

Schemat zderzenia proton-proton

Eksperyment CMS działający przy LHC (CERN)

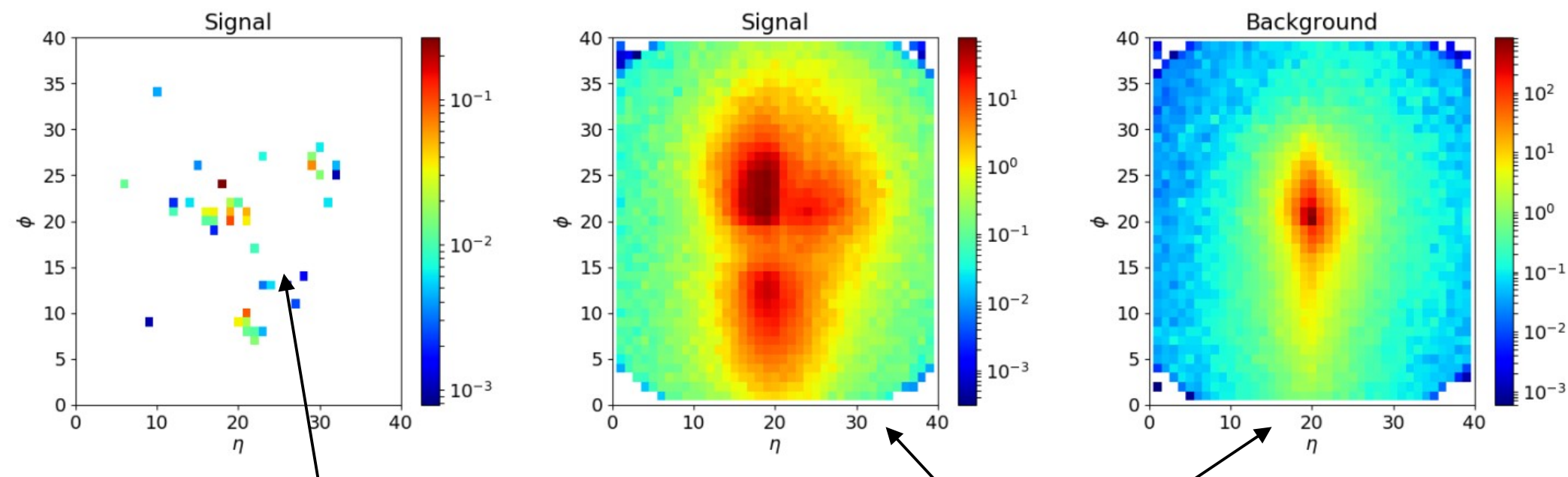


CERN/CMS Collaboration

Schemat zderzenia proton-proton

Zadanie: klasyfikacja “obrazów” zderzeń: kwark top/coś innego

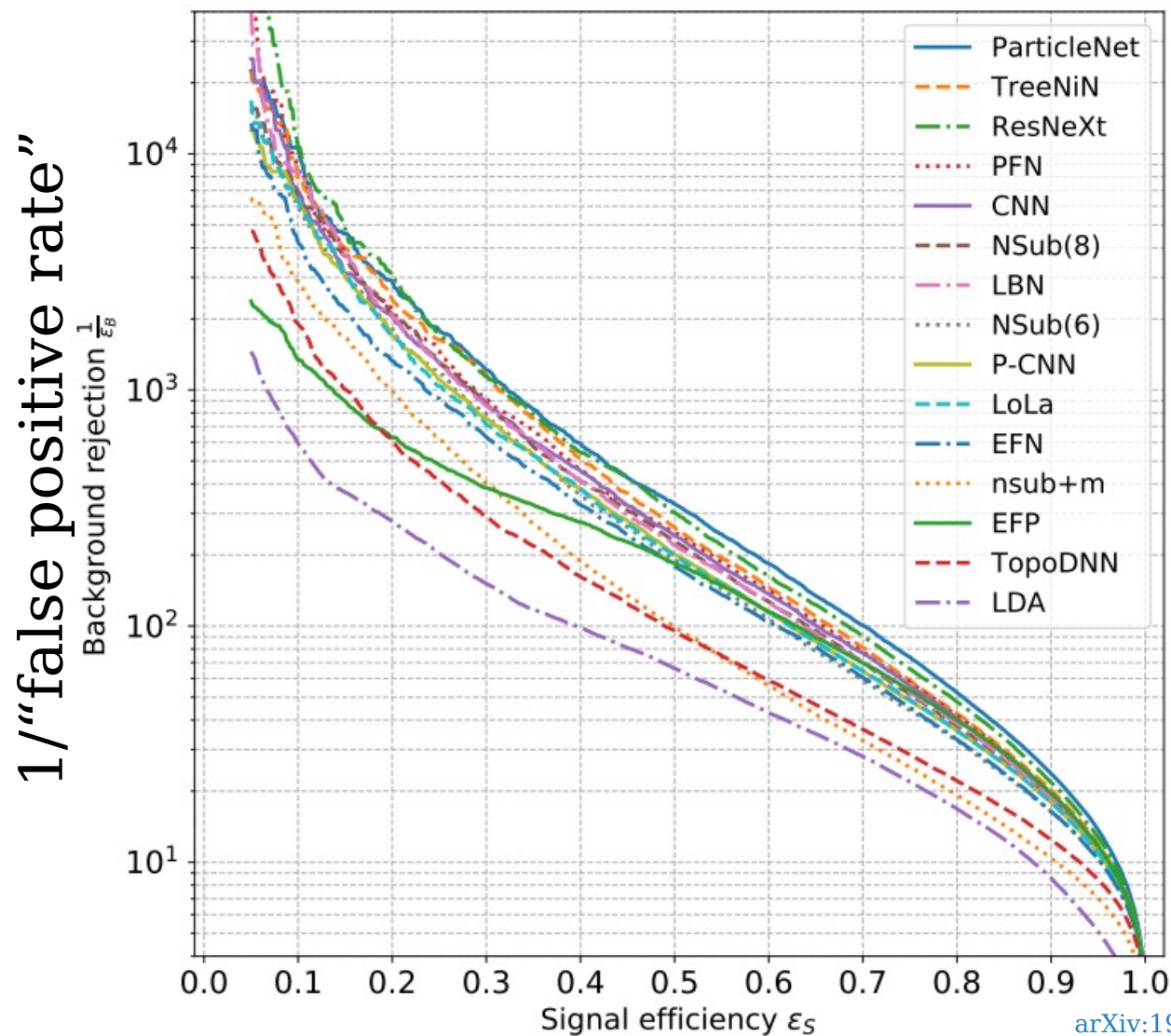
Dane wejściowe: obrazy dżetów



pojedyncze zderzenie
proton-proton

Obrazy uśrednione po
10 000 zderzeń
proton-proton

Zadanie: klasyfikacja “obrazów” zderzeń: kwark top/coś innego



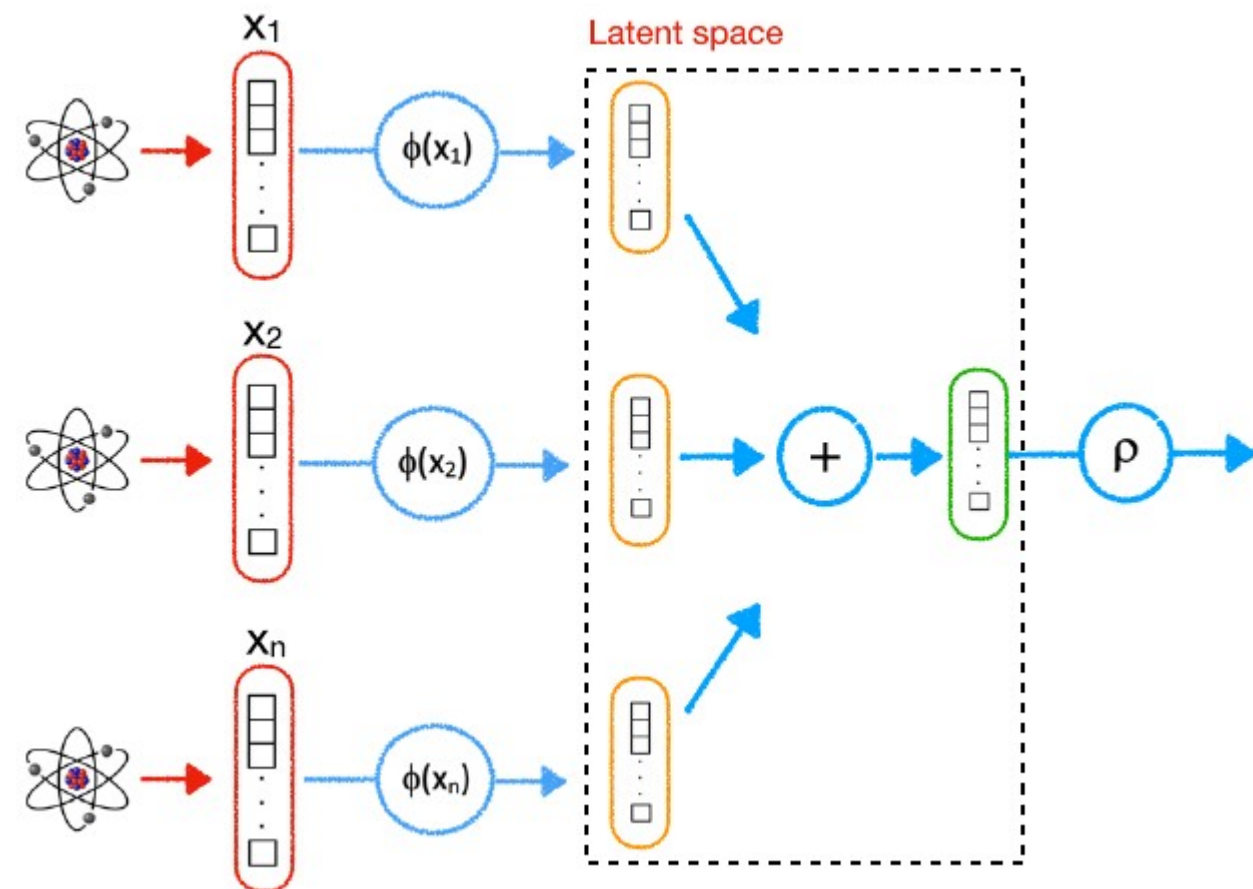
arXiv:1902.09914 [hep-ph]

“recall”

Zadanie: klasyfikacja związków: nadprzewodzący? **TAK/NIE**

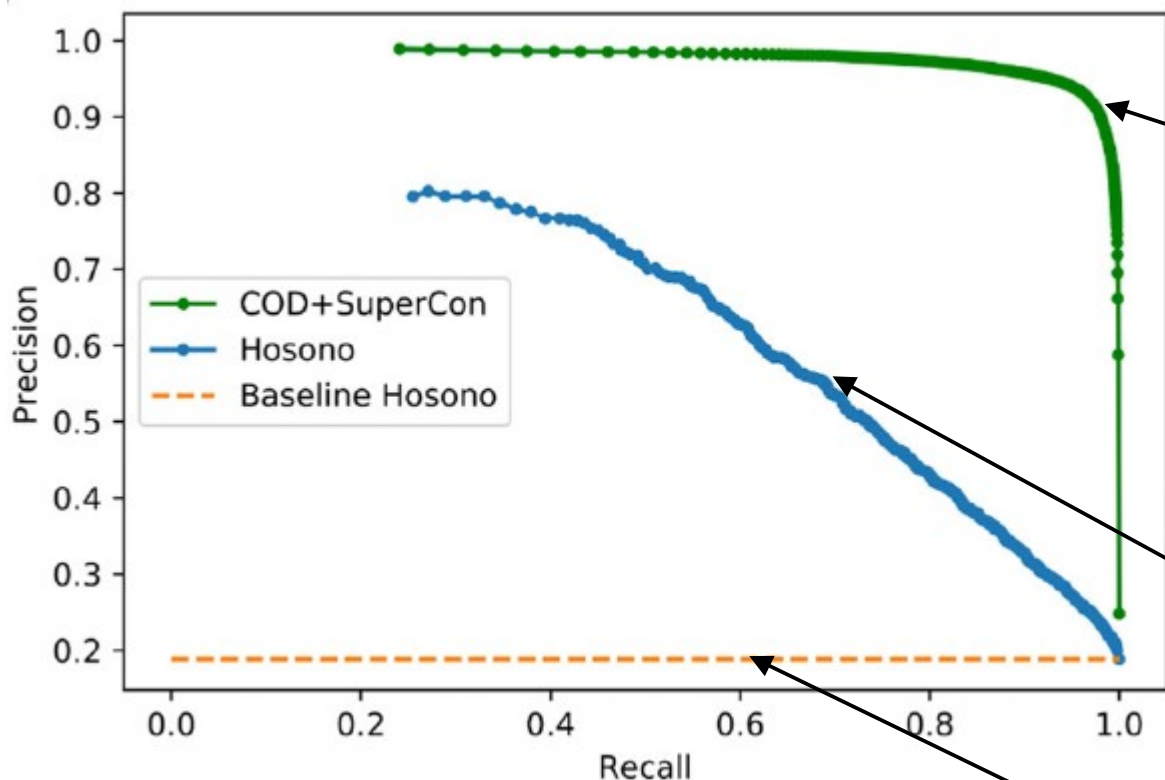
Dane wejściowe:

- zestaw pierwiastków tworzących analizowany związek,
- 22 cechy na pierwiastek + współczynnik stoichiometryczny



$\phi(x), \rho(x)$ - sieci neuronowe

K., Guizouarn, T. et al. From individual elements to materials: in search of new superconductors via machine put Mater 9, 71 (2023).
1038/s41524-023-01023-6



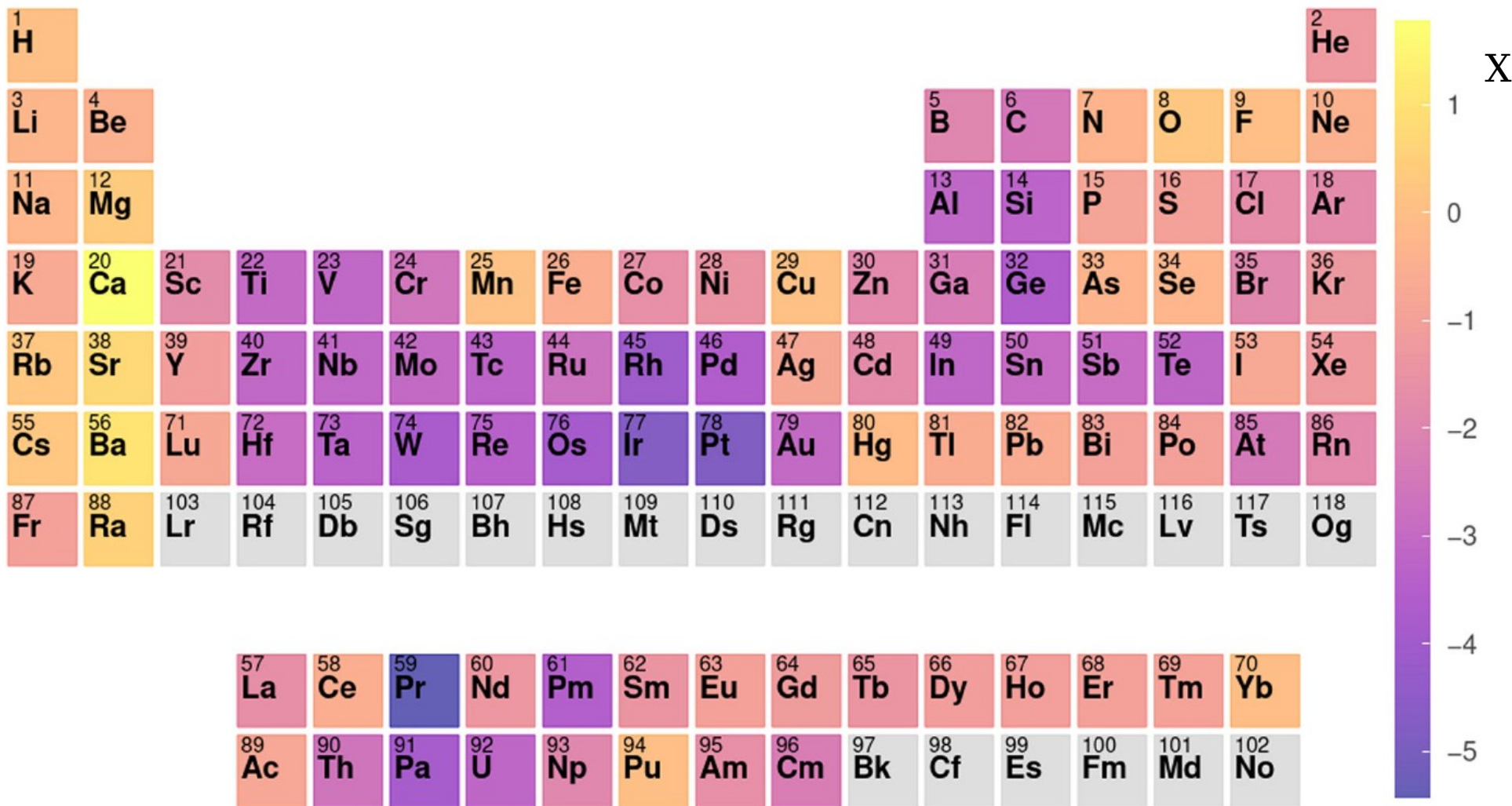
zbiór o strukturze identycznej ze zbiorem użytym do treningu

niezależny zbiór, nie używany do trenowania

zbiór Hosono zawiera 20% nadprzewodników

Pereti, C., Bernot, K., Guizouarn, T. et al. From individual elements to macroscopic materials: in search of new superconductors via machine learning. npj Comput Mater 9, 71 (2023).
<https://doi.org/10.1038/s41524-023-01023-6>

Mapa wpływu na zwiększenie temperatury krytycznej:
duże X → zwiększona T_c



Pereti, C., Bernot, K., Guizouarn, T. et al. From individual elements to macroscopic materials: in search of new superconductors via machine learning. npj Comput Mater 9, 71 (2023).
<https://doi.org/10.1038/s41524-023-01023-6>

Współczesne uczenie maszynowe ma swoje korzenie w metodach statystycznej analizy danych rozwijanych w drugiej połowie XX wieku.

Błyskawiczny rozwój metod obliczeniowych pozwala na zastosowanie znanego podejścia – „dopasowanie funkcji” w całkiem nowych kontekstach