

Metody inteligencji obliczeniowej w analizie danych

Sieci neuronowe II

Jan Karwowski

Wydział matematyki i Nauk Informacyjnych PW

9 marca 2021



**Fundusze
Europejskie**

Wiedza Edukacja Rozwój



**Rzeczpospolita
Polska**

**Politechnika
Warszawska**

Unia Europejska
Europejski Fundusz Społeczny

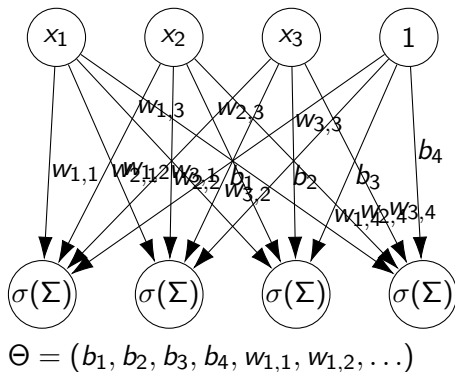


Zadanie 10 pn. „Modyfikacja programów studiów na kierunkach prowadzonych przez Wydział Matematyki i Nauk Informatycznych” realizowane jest w ramach projektu „NERW 2 PW. Nauka – Edukacja – Rozwój – Współpraca” współfinansowanego ze środków Unii Europejskiej w ramach Europejskiego Funduszu Społecznego

Ustalić jak z oddawaniem sprawozdań.

Sieć neuronowa w zadaniu regresji

- Wektor cech: $X \in [0, 1]^N$,
wektor odpowiedzi $Y \in [0, 1]^M$
- k -elementowy zbiór testowy
 $\subset X \times Y$
- Model (sieć) N i parametry
(wagi) Θ
- $MSE = \sum_i \frac{(N(x_i, \Theta) - y_i)^2}{k}$
- Skąd wziąć wektor Θ ?



$$G(x) = \sum_{j=1}^N \alpha_j \sigma \left(\sum_i^n x_i w_{i,j} + b_j \right)$$

Tw. Dla funkcji σ – sigmoidalnej – zbiór funkcji w postaci G (dla dowolnych N, α_j, w_i, b_j) jest gęsty w przestrzeni wszystkich funkcji ciągłych $[0, 1]^n \rightarrow [0, 1]$.

G. Cybenko. “Approximation by superpositions of a sigmoidal function”. In: *Mathematics of Control, Signals, and Systems* 2.4 (1989), 303–314. ISSN: 1435-568X. DOI: 10.1007/bf02551274. URL: <http://dx.doi.org/10.1007/bf02551274>

- Warstwa
 - (Wejściowa)
 - Ukryta
 - Wyjściowa
- Aktywacja neuronu
- Uczenie sieci
- Wzorzec uczący – para (\mathbf{x}, \mathbf{y}) ze zbioru uczącego

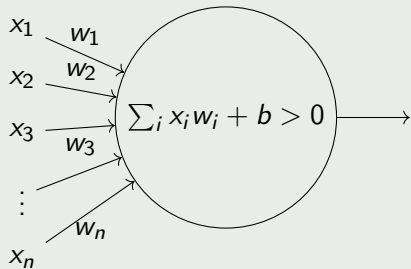
Postulat: neruony, które często są pobudzone w tym samym momencie mają silne połączenia synaptyczne.

Donald O Hebb. *The organization of behavior*. 1949

Interpretacja dla perceptronu (pojedynczego)

Zakładamy, że wszystkie wartości w zbiorze uczącym są ze zbioru $\{-1, 1\}$.

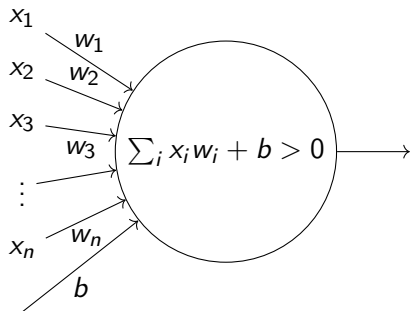
$$w_i = \frac{\sum_{k=1}^N x_i^k y^k}{N}$$



Efekt

Dla bardzo małej liczby wzorców uczących efekt jest zadowalający, w większych zbiorach sieć traci zdolność do uczenia się.

Reguła perceptronu I



Skokowa funkcja aktywacji

$$f(x) = \begin{cases} 0 & x \leq 0 \\ 1 & x > 0 \end{cases}$$

- Problem: wejście – zmienne ciągłe, wyjście – zmienne binarne.

Reguła perceptronu II

- Intuicja: zmienić wagi tak, że jeśli neuron nie jest aktywowany, a powinien być, to zmieniamy wagę tak, żeby zwiększyć jego szansę aktywacji.

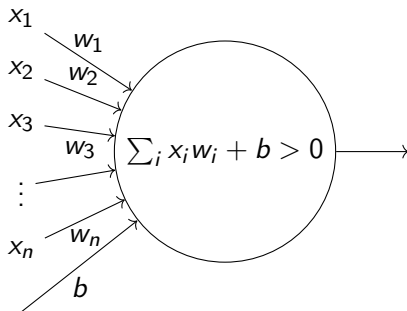
Metoda uczenia

y – klasa, \hat{y} – wartość wyliczona przez perceptron, x_i i -te wejście. Rozważamy pojedynczy wzorec uczący, prezentujemy kolejne wzorce, często wielokrotnie.

$$\Delta w_i = \begin{cases} x_i & \hat{y} < y \\ -x_i & \hat{y} > y \\ 0 & \hat{y} = y \end{cases}$$

$$\Delta b = \begin{cases} 1 & \hat{y} < y \\ -1 & \hat{y} > y \\ 0 & \hat{y} = y \end{cases}$$

Jak można uprościć ten zapis?



Reguła delta I

- Wejście – ciągłe, wyjście ciągłe, funkcja aktywacji *różniczkowalna*.
- η – parametr: krok uczenia (także współczynnik uczenia)
- Ulepszenie reguły perceptronu
- W aktualizacji uwzględniamy wpływ nachylenia funkcji aktywacji na zmianę wyniku

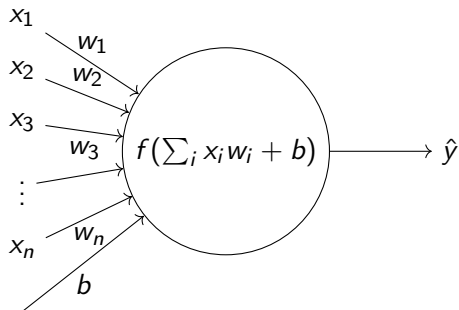
Reguła uczenia

y – wyjście wg. wzorca uczącego, \hat{y} – wartość wyliczona przez perceptron, x_i i -te wejście. Rozważamy pojedynczy wzorzec uczący, prezentujemy kolejne wzorce, często wielokrotnie.

$$\Delta w_i = \eta(y - \hat{y})f'(\sum_j x_j w_j + b)x_i$$

$$\Delta b = \eta(y - \hat{y})f'(\sum_j x_j w_j + b)$$

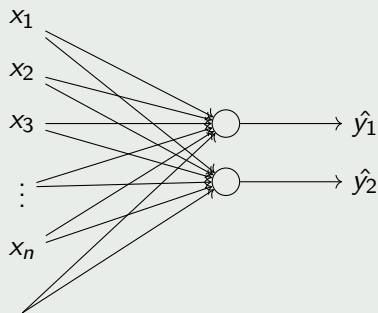
Reguła delta II



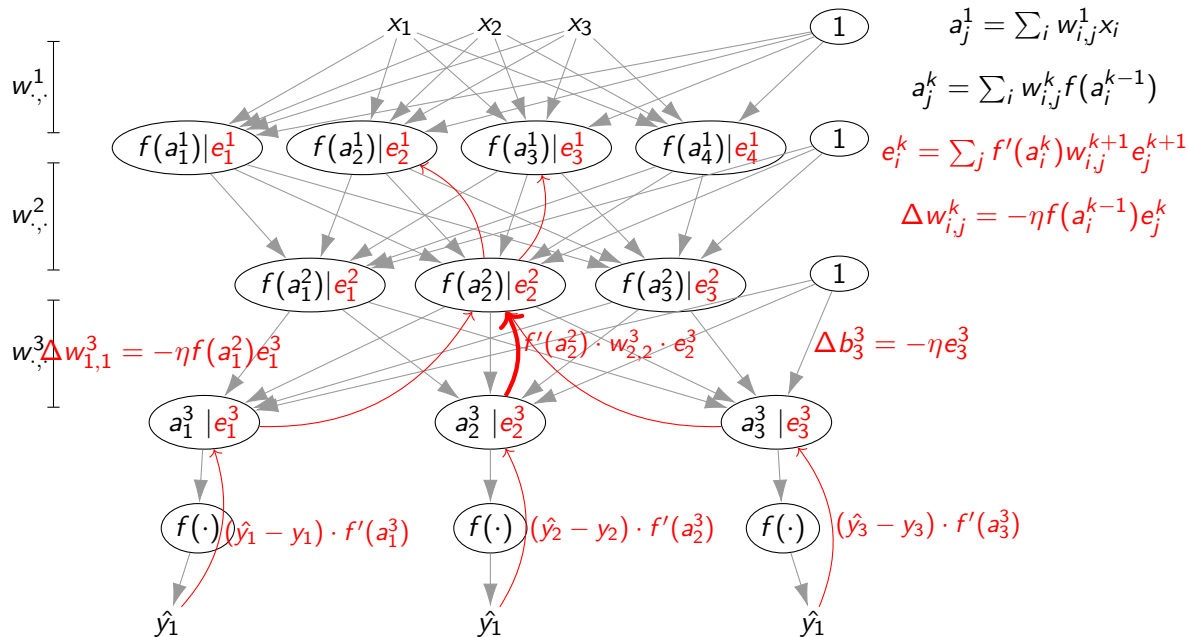
Perceptron – wiele wyjść

$$\Delta w_{i,j} = \alpha(y_j - \hat{y}_j) f'(\sum_k x_k w_{k,j} + b) x_i$$

$$\Delta b_j = \alpha(y_j - \hat{y}_j) f'(\sum_K x_K w_{K,j} + b)$$



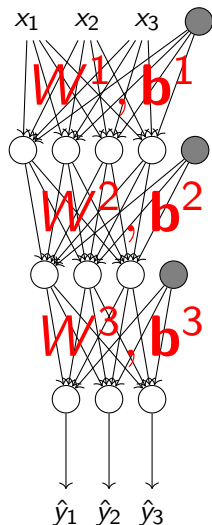
Propagacja wsteczna błędu – MLP



Algorytm

```
1  $\Theta \leftarrow \text{InicjujLosowo};$ 
2 while  $\neg \text{StopCondition}$  do
3   for  $(\mathbf{X}, \mathbf{Y}) \in \text{Zbiór uczący}$  do
4      $\hat{\mathbf{Y}} \leftarrow \text{FeedForward}(\Theta, \mathbf{X});$ 
5      $\Delta\Theta \leftarrow \text{Backpropagate}(\hat{\mathbf{Y}}, \mathbf{Y});$ 
6      $\Theta \leftarrow \Theta + \Delta\Theta;$ 
```


Uczenie gradientowe I



η – krok uczenia (niektórzy używają α)

$$\Theta = [W^1, \mathbf{b}^1, W^2, \mathbf{b}^2, \dots, \mathbf{b}^3]$$

$$MSE(\Theta) = \frac{1}{2N} \sum_i (\hat{y}_1(\Theta, \mathbf{x}^i) - y_1^i)^2 + (\hat{y}_2(\Theta, \mathbf{x}^i) - y_2^i)^2 + (\hat{y}_3(\Theta, \mathbf{x}^i) - y_3^i)^2$$

$$\Delta \Theta = -\eta \nabla MSE(\Theta)$$

$$MSE(\Theta) = \frac{1}{N} \sum_i \frac{1}{2} \left((\hat{y}_1(\Theta, \mathbf{x}^i) - y_1^i)^2 + (\hat{y}_2(\Theta, \mathbf{x}^i) - y_2^i)^2 + (\hat{y}_3(\Theta, \mathbf{x}^i) - y_3^i)^2 \right)$$

Pochodna po wadze w ostatniej warstwie $w_{1,1}^3$

- $\frac{1}{2} \left((\hat{y}_1(\Theta, \mathbf{x}) - y_1)^2 + (\hat{y}_2(\Theta, \mathbf{x}) - y_2)^2 + (\hat{y}_3(\Theta, \mathbf{x}) - y_3)^2 \right)$
- $\frac{1}{2} \frac{\partial (\hat{y}_1(\Theta, \mathbf{x}) - y_1)^2}{\partial w_{1,1}^3} = \frac{1}{2} 2 (\hat{y}_1(\Theta, \mathbf{x}) - y_1) \cdot \frac{\partial \hat{y}_1(\Theta, \mathbf{x}^i)}{\partial w_{1,1}^3} =$
 $(\hat{y}_1(\Theta, \mathbf{x}) - y_1) \cdot \frac{\partial f(a_1^3)}{\partial w_{1,1}^3} =$

$$\begin{aligned}\hat{y}_i &= f(a_i^3) \\ a_j^k &= \\ \sum_i w_{i,j}^k f(a_i^{k-1}) \\ a_j^1 &= \sum_i w_{i,j}^1 x_i\end{aligned}$$

$$MSE(\Theta) = \frac{1}{N} \sum_i \frac{1}{2} \left((\hat{y}_1(\Theta, \mathbf{x}^i) - y_1^i)^2 + (\hat{y}_2(\Theta, \mathbf{x}^i) - y_2^i)^2 + (\hat{y}_3(\Theta, \mathbf{x}^i) - y_3^i)^2 \right)$$

Pochodna po wadze w ostatniej warstwie $w_{1,1}^3$

- $\frac{1}{2} \left((\hat{y}_1(\Theta, \mathbf{x}) - y_1)^2 + (\hat{y}_2(\Theta, \mathbf{x}) - y_2)^2 + (\hat{y}_3(\Theta, \mathbf{x}) - y_3)^2 \right)$
- $\frac{1}{2} \frac{\partial (\hat{y}_1(\Theta, \mathbf{x}) - y_1)^2}{\partial w_{1,1}^3} = \frac{1}{2} 2 (\hat{y}_1(\Theta, \mathbf{x}) - y_1) \cdot \frac{\partial \hat{y}_1(\Theta, \mathbf{x})}{\partial w_{1,1}^3} =$
 $(\hat{y}_1(\Theta, \mathbf{x}) - y_1) \cdot \frac{\partial f(a_1^3)}{\partial w_{1,1}^3} = (\hat{y}_1(\Theta, \mathbf{x}) - y_1) \cdot f'(a_1^3) \cdot \frac{\partial a_1^3}{\partial w_{1,1}^3} =$

$$\begin{aligned}\hat{y}_i &= f(a_i^3) \\ a_j^k &= \sum_i w_{i,j}^k f(a_i^{k-1}) \\ a_j^1 &= \sum_i w_{i,j}^1 x_i\end{aligned}$$

$$MSE(\Theta) = \frac{1}{N} \sum_i \frac{1}{2} \left((\hat{y}_1(\Theta, \mathbf{x}^i) - y_1^i)^2 + (\hat{y}_2(\Theta, \mathbf{x}^i) - y_2^i)^2 + (\hat{y}_3(\Theta, \mathbf{x}^i) - y_3^i)^2 \right)$$

Pochodna po wadze w ostatniej warstwie $w_{1,1}^3$

- $\frac{1}{2} \left((\hat{y}_1(\Theta, \mathbf{x}) - y_1)^2 + (\hat{y}_2(\Theta, \mathbf{x}) - y_2)^2 + (\hat{y}_3(\Theta, \mathbf{x}) - y_3)^2 \right)$
- $\frac{1}{2} \frac{\partial (\hat{y}_1(\Theta, \mathbf{x}) - y_1)^2}{\partial w_{1,1}^3} = \frac{1}{2} 2 (\hat{y}_1(\Theta, \mathbf{x}) - y_1) \cdot \frac{\partial \hat{y}_1(\Theta, \mathbf{x})}{\partial w_{1,1}^3} =$
 $(\hat{y}_1(\Theta, \mathbf{x}) - y_1) \cdot \frac{\partial f(a_1^3)}{\partial w_{1,1}^3} = (\hat{y}_1(\Theta, \mathbf{x}) - y_1) \cdot f'(a_1^3) \cdot \frac{\partial a_1^3}{\partial w_{1,1}^3} =$
 $(\hat{y}_1(\Theta, \mathbf{x}) - y_1) \cdot f'(a_1^3) \cdot \frac{\partial \sum_j w_{j,1}^3 a_j^2}{\partial w_{1,1}^3} =$

$$\begin{aligned}\hat{y}_i &= f(a_i^3) \\ a_j^k &= \sum_i w_{i,j}^k f(a_i^{k-1}) \\ a_j^1 &= \sum_i w_{i,j}^1 x_i\end{aligned}$$

$$MSE(\Theta) = \frac{1}{N} \sum_i \frac{1}{2} \left((\hat{y}_1(\Theta, \mathbf{x}^i) - y_1^i)^2 + (\hat{y}_2(\Theta, \mathbf{x}^i) - y_2^i)^2 + (\hat{y}_3(\Theta, \mathbf{x}^i) - y_3^i)^2 \right)$$

Pochodna po wadze w ostatniej warstwie $w_{1,1}^3$

- $\frac{1}{2} \left((\hat{y}_1(\Theta, \mathbf{x}) - y_1)^2 + (\hat{y}_2(\Theta, \mathbf{x}) - y_2)^2 + (\hat{y}_3(\Theta, \mathbf{x}) - y_3)^2 \right)$
- $\frac{1}{2} \frac{\partial (\hat{y}_1(\Theta, \mathbf{x}) - y_1)^2}{\partial w_{1,1}^3} = \frac{1}{2} 2 (\hat{y}_1(\Theta, \mathbf{x}) - y_1) \cdot \frac{\partial \hat{y}_1(\Theta, \mathbf{x})}{\partial w_{1,1}^3} =$
 $(\hat{y}_1(\Theta, \mathbf{x}) - y_1) \cdot \frac{\partial f(a_1^3)}{\partial w_{1,1}^3} = (\hat{y}_1(\Theta, \mathbf{x}) - y_1) \cdot f'(a_1^3) \cdot \frac{\partial a_1^3}{\partial w_{1,1}^3} =$
 $(\hat{y}_1(\Theta, \mathbf{x}) - y_1) \cdot f'(a_1^3) \cdot \frac{\partial \sum_j w_{j,1}^3 a_j^2}{\partial w_{1,1}^3} =$
 $(\hat{y}_1(\Theta, \mathbf{x}) - y_1) \cdot f'(a_1^3) \cdot \frac{\partial (w_{1,1}^3 f(a_1^2) + w_{2,1}^3 f(a_2^2) + w_{3,1}^3 f(a_3^2))}{\partial w_{1,1}^3} =$

$$\begin{aligned}\hat{y}_i &= f(a_i^3) \\ a_j^k &= \sum_i w_{i,j}^k f(a_i^{k-1}) \\ a_j^1 &= \sum_i w_{i,j}^1 x_i\end{aligned}$$

$$MSE(\Theta) = \frac{1}{N} \sum_i \frac{1}{2} \left((\hat{y}_1(\Theta, \mathbf{x}^i) - y_1^i)^2 + (\hat{y}_2(\Theta, \mathbf{x}^i) - y_2^i)^2 + (\hat{y}_3(\Theta, \mathbf{x}^i) - y_3^i)^2 \right)$$

Pochodna po wadze w ostatniej warstwie $w_{1,1}^3$

- $\frac{1}{2} \left((\hat{y}_1(\Theta, \mathbf{x}) - y_1)^2 + (\hat{y}_2(\Theta, \mathbf{x}) - y_2)^2 + (\hat{y}_3(\Theta, \mathbf{x}) - y_3)^2 \right)$
- $\frac{1}{2} \frac{\partial (\hat{y}_1(\Theta, \mathbf{x}) - y_1)^2}{\partial w_{1,1}^3} = \frac{1}{2} 2 (\hat{y}_1(\Theta, \mathbf{x}) - y_1) \cdot \frac{\partial \hat{y}_1(\Theta, \mathbf{x})}{\partial w_{1,1}^3} =$
 $(\hat{y}_1(\Theta, \mathbf{x}) - y_1) \cdot \frac{\partial f(a_1^3)}{\partial w_{1,1}^3} = (\hat{y}_1(\Theta, \mathbf{x}) - y_1) \cdot f'(a_1^3) \cdot \frac{\partial a_1^3}{\partial w_{1,1}^3} =$
 $(\hat{y}_1(\Theta, \mathbf{x}) - y_1) \cdot f'(a_1^3) \cdot \frac{\partial \sum_j w_{j,1}^3 a_j^2}{\partial w_{1,1}^3} =$
 $(\hat{y}_1(\Theta, \mathbf{x}) - y_1) \cdot f'(a_1^3) \cdot \frac{\partial (w_{1,1}^3 f(a_1^2) + w_{2,1}^3 f(a_2^2) + w_{3,1}^3 f(a_3^2))}{\partial w_{1,1}^3} =$
 $(\hat{y}_1(\Theta, \mathbf{x}) - y_1) \cdot f'(a_1^3) \cdot f(a_1^2) =$

$$\begin{aligned} \hat{y}_i &= f(a_i^3) \\ a_j^k &= \sum_i w_{i,j}^k f(a_i^{k-1}) \\ a_j^1 &= \sum_i w_{i,j}^1 x_i \end{aligned}$$

Wyprowadzenie pochodnej

$$MSE(\Theta) = \frac{1}{N} \sum_i \frac{1}{2} \left((\hat{y}_1(\Theta, \mathbf{x}^i) - y_1^i)^2 + (\hat{y}_2(\Theta, \mathbf{x}^i) - y_2^i)^2 + (\hat{y}_3(\Theta, \mathbf{x}^i) - y_3^i)^2 \right)$$

Pochodna po wadze w ostatniej warstwie $w_{1,1}^3$

- $\frac{1}{2} \left((\hat{y}_1(\Theta, \mathbf{x}) - y_1)^2 + (\hat{y}_2(\Theta, \mathbf{x}) - y_2)^2 + (\hat{y}_3(\Theta, \mathbf{x}) - y_3)^2 \right)$
- $\frac{1}{2} \frac{\partial (\hat{y}_1(\Theta, \mathbf{x}) - y_1)^2}{\partial w_{1,1}^3} = \frac{1}{2} 2 (\hat{y}_1(\Theta, \mathbf{x}) - y_1) \cdot \frac{\partial \hat{y}_1(\Theta, \mathbf{x})}{\partial w_{1,1}^3} =$
 $(\hat{y}_1(\Theta, \mathbf{x}) - y_1) \cdot \frac{\partial f(a_1^3)}{\partial w_{1,1}^3} = (\hat{y}_1(\Theta, \mathbf{x}) - y_1) \cdot f'(a_1^3) \cdot \frac{\partial a_1^3}{\partial w_{1,1}^3} =$
 $(\hat{y}_1(\Theta, \mathbf{x}) - y_1) \cdot f'(a_1^3) \cdot \frac{\partial \sum_j w_{j,1}^3 a_j^2}{\partial w_{1,1}^3} =$
 $(\hat{y}_1(\Theta, \mathbf{x}) - y_1) \cdot f'(a_1^3) \cdot \frac{\partial (w_{1,1}^3 f(a_1^2) + w_{2,1}^3 f(a_2^2) + w_{3,1}^3 f(a_3^2))}{\partial w_{1,1}^3} =$
 $(\hat{y}_1(\Theta, \mathbf{x}) - y_1) \cdot f'(a_1^3) \cdot f(a_1^2) = \mathbf{e}_1^3 \cdot f(a_1^2)$

$$\begin{aligned}\hat{y}_i &= f(a_i^3) \\ a_j^k &= \sum_i w_{i,j}^k f(a_i^{k-1}) \\ a_j^1 &= \sum_i w_{i,j}^1 x_i\end{aligned}$$

$$MSE(\Theta) = \frac{1}{N} \sum_i \frac{1}{2} (\hat{y}_1(\Theta, \mathbf{x}^i) - y_1^i)^2 + (\hat{y}_2(\Theta, \mathbf{x}^i) - y_2^i)^2 + (\hat{y}_3(\Theta, \mathbf{x}^i) - y_3^i)^2$$

Pochodna po wadze w przedostatniej warstwie $w_{1,1}^2$

- $$\frac{1}{2} \frac{\partial ((\hat{y}_1(\Theta, \mathbf{x}) - y_1)^2 + (\hat{y}_2(\Theta, \mathbf{x}) - y_2)^2 + (\hat{y}_3(\Theta, \mathbf{x}) - y_3)^2)}{\partial w_{1,1}^2} = \dots$$
- $$\begin{aligned} \frac{\partial \hat{y}_i(\Theta, \mathbf{x})}{\partial w_{1,1}^2} &= \frac{\partial f\left(\sum_j w_{j,i}^3 a_j^2\right)}{\partial w_{1,1}^2} = f'(\sum_j w_{j,i}^3 a_j^2) \frac{\partial \sum_j w_{j,i}^3 a_j^2}{\partial w_{1,1}^2} = \\ &= f'(a_i^3) \frac{\partial (w_{1,i}^3 f(a_1^2) + w_{2,i}^3 f(a_2^2) + \dots)}{\partial w_{1,1}^2} = f'(a_i^3) w_{1,i}^3 \frac{\partial f(a_1^2)}{\partial w_{1,1}^2} = \\ &= f'(a_i^3) w_{1,i}^3 f'(a_1^2) \frac{\partial \sum_j w_{j,1}^2 f(a_j^1)}{\partial w_{1,1}^2} = f'(a_i^3) w_{1,i}^3 f'(a_1^2) \cdot f(a_1^1) \end{aligned}$$

$$\begin{aligned} \hat{y}_i &= f(a_i^3) \\ a_j^k &= \sum_i w_{i,j}^k f(a_i^{k-1}) \\ a_j^1 &= \sum_i w_{i,j}^1 x_i \\ e_i^3 &= (\hat{y}_i(\Theta, \mathbf{x}) - y_i^3) f'(a_i^3) \end{aligned}$$

$$MSE(\Theta) = \frac{1}{N} \sum_i \frac{1}{2} (\hat{y}_1(\Theta, \mathbf{x}^i) - y_1^i)^2 + (\hat{y}_2(\Theta, \mathbf{x}^i) - y_2^i)^2 + (\hat{y}_3(\Theta, \mathbf{x}^i) - y_3^i)^2$$

Pochodna po wadze w przedostatniej warstwie $w_{1,1}^2$

- $\frac{1}{2} \frac{\partial ((\hat{y}_1(\Theta, \mathbf{x}) - y_1)^2 + (\hat{y}_2(\Theta, \mathbf{x}) - y_2)^2 + (\hat{y}_3(\Theta, \mathbf{x}) - y_3)^2)}{\partial w_{1,1}^2} = \dots$

$$\begin{aligned}\hat{y}_i &= f(a_i^3) \\ a_j^k &= \sum_i w_{i,j}^k f(a_i^{k-1}) \\ a_j^1 &= \sum_i w_{i,j}^1 x_i \\ e_i^3 &= (\hat{y}_i(\Theta, \mathbf{x}) - y_i^3) f'(a_i^3) \\ \frac{\partial \hat{y}_i(\Theta, \mathbf{x})}{\partial w_{1,1}^2} &= \\ f'(a_i^3) w_{1,i}^3 f'(a_1^2) \cdot f(a_1^1)\end{aligned}$$

$$MSE(\Theta) = \frac{1}{N} \sum_i \frac{1}{2} (\hat{y}_1(\Theta, \mathbf{x}^i) - y_1^i)^2 + (\hat{y}_2(\Theta, \mathbf{x}^i) - y_2^i)^2 + (\hat{y}_3(\Theta, \mathbf{x}^i) - y_3^i)^2$$

Pochodna po wadze w przedostatniej warstwie $w_{1,1}^2$

$$\begin{aligned} \bullet \quad \frac{1}{2} \frac{\partial \sum_j (\hat{y}_j(\Theta, \mathbf{x}) - y_j)^2}{\partial w_{1,1}^2} &= \frac{1}{2} \cdot 2 \sum_j (\hat{y}_j(\Theta, \mathbf{x}) - y_j) \frac{\partial \hat{y}_j(\Theta, \mathbf{x})}{\partial w_{1,1}^2} = \\ &= \sum_j (\hat{y}_j(\Theta, \mathbf{x}) - y_j) f'(a_j^3) w_{1,j}^3 f'(a_1^2) f(a_1^1) = \\ &= \left(\sum_j e_j^3 w_{1,j}^3 \right) f'(a_1^2) f(a_1^1) \end{aligned}$$

$$\begin{aligned} \hat{y}_i &= f(a_i^3) \\ a_j^k &= \sum_i w_{i,j}^k f(a_i^{k-1}) \\ a_j^1 &= \sum_i w_{i,j}^1 x_i \\ e_i^3 &= (\hat{y}_i(\Theta, \mathbf{x}) - y_i^3) f'(a_i^3) \\ \frac{\partial \hat{y}_i(\Theta, \mathbf{x})}{\partial w_{1,1}^2} &= \\ f'(a_i^3) w_{1,i}^3 f'(a_1^2) \cdot f(a_1^1) \end{aligned}$$

$$MSE(\Theta) = \frac{1}{N} \sum_i \frac{1}{2} (\hat{y}_1(\Theta, \mathbf{x}^i) - y_1^i)^2 + (\hat{y}_2(\Theta, \mathbf{x}^i) - y_2^i)^2 + (\hat{y}_3(\Theta, \mathbf{x}^i) - y_3^i)^2$$

Pochodna po wadze w przedostatniej warstwie $w_{1,1}^2$

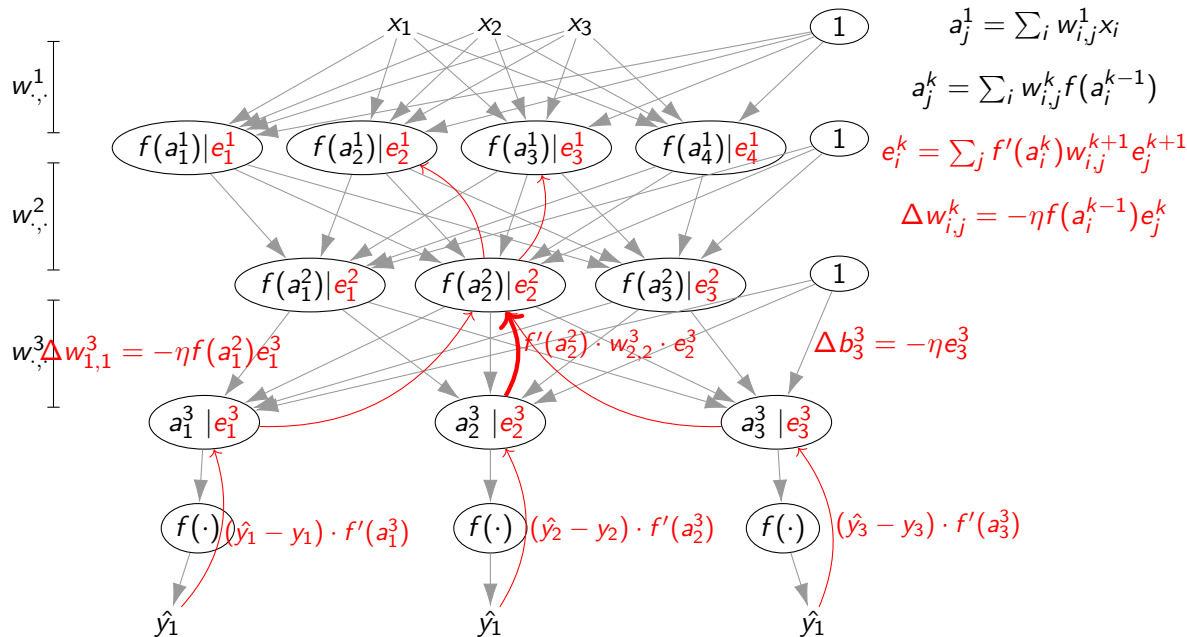
- $$\frac{1}{2} \frac{\partial \sum_j (\hat{y}_j(\Theta, \mathbf{x}) - y_j)^2}{\partial w_{1,1}^2} = \frac{1}{2} \cdot 2 \sum_j (\hat{y}_j(\Theta, \mathbf{x}) - y_j) \frac{\partial \hat{y}_j(\Theta, \mathbf{x})}{\partial w_{1,1}^2} =$$

$$\sum_j (\hat{y}_j(\Theta, \mathbf{x}) - y_j) f'(a_j^3) w_{1,j}^3 f'(a_1^2) f(a_1^1) =$$

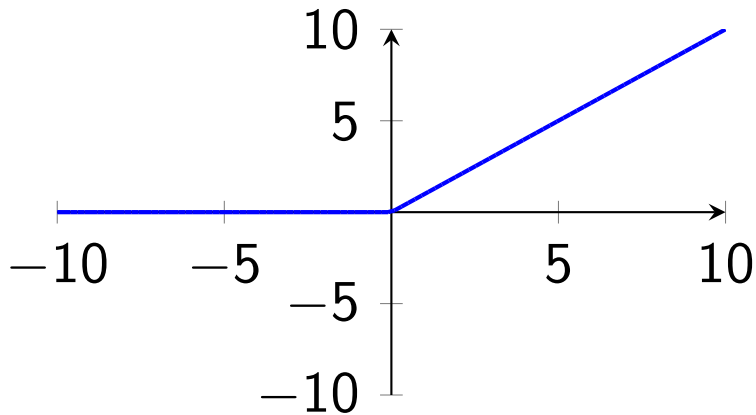
$$\left(\sum_j e_j^3 w_{1,j}^3 \right) f'(a_1^2) f(a_1^1)$$
- $$e_1^2 = \sum_j e_j^3 w_{1,j}^3 f'(a_1^2)$$

$$\begin{aligned} \hat{y}_i &= f(a_i^3) \\ a_j^k &= \sum_i w_{i,j}^k f(a_i^{k-1}) \\ a_j^1 &= \sum_i w_{i,j}^1 x_i \\ e_i^3 &= (\hat{y}_i(\Theta, \mathbf{x}) - y_i^3) f'(a_i^3) \\ \frac{\partial \hat{y}_i(\Theta, \mathbf{x})}{\partial w_{1,1}^2} &= \\ f'(a_i^3) w_{1,i}^3 f'(a_1^2) \cdot f(a_1^1) \end{aligned}$$

Propagacja wsteczna błędu – MLP



$$\text{ReLU}(x) = \max\{x, 0\}$$



$$\Delta\Theta = -\eta\nabla MSE(\Theta)$$

- Zwykle 10^{-3} lub mniej

- Mean Absolute Error $\frac{1}{N} \sum_i |\hat{y}^i - y^i|$
- Mean Squared Error $\frac{1}{N} \sum_i (\hat{y}^i - y^i)^2$
- Categorical crossentropy $\frac{1}{N} \sum_i - \sum_j y_j^i \log \hat{y}_j^i$

$$MSE(\Theta) = \frac{1}{N} \sum_i \frac{1}{2} (\hat{y}_1(\Theta, \mathbf{x}^i) - y_1^i)^2 + (\hat{y}_2(\Theta, \mathbf{x}^i) - y_2^i)^2 + (\hat{y}_3(\Theta, \mathbf{x}^i) - y_3^i)^2$$

SGD

```
1  $\Theta \leftarrow \text{InicujLosowo};$ 
2 while  $\neg \text{StopCondition}$  do
3   for  $(\mathbf{X}, \mathbf{Y}) \in \text{Zbiór uczący}$  do
4      $\hat{\mathbf{Y}} \leftarrow \text{Network}(\Theta, X);$ 
5      $\Delta\Theta \leftarrow -\eta \nabla \text{Network}(\Theta, X);$ 
6      $\Theta \leftarrow \Theta + \Delta\Theta;$ 
```

Podstawowy Gradient Descent

```
1  $\Theta \leftarrow \text{InicujLosowo};$ 
2 while  $\neg \text{StopCondition}$  do
3    $\Delta\Theta \leftarrow 0;$ 
4   for  $(\mathbf{X}, \mathbf{Y}) \in \text{Zbiór uczący}$  do
5      $\hat{\mathbf{Y}} \leftarrow \text{Network}(\Theta, X);$ 
6      $\Delta\Theta \leftarrow \Delta\Theta - \eta \nabla \text{Network}(\Theta, X);$ 
7    $\Theta \leftarrow \Theta + \Delta\Theta;$ 
```


Trajektoria i szybkość zbieżności

Mini-batch Gradient Descent

- Podstawowy wariant bywa nazywany Batch Gradient Descent

```
1  $\Theta \leftarrow \text{InicjujLosowo};$ 
2 while  $\neg \text{StopCondition}$  do
3    $\text{MiniBatches} \leftarrow \text{PodzielLosowoNaRozczceZbiory}(\text{ZbiorUczcy}, k);$ 
4   for  $\text{batch} \in \text{MiniBatches}$  do
5      $\Delta\Theta \leftarrow 0;$ 
6     for  $(\mathbf{X}, \mathbf{Y}) \in \text{Batch}$  do
7        $\hat{\mathbf{Y}} \leftarrow \text{Network}(\Theta, \mathbf{X});$ 
8        $\Delta\Theta \leftarrow \Delta\Theta - \eta \nabla \text{Network}(\Theta, \mathbf{X});$ 
9    $\Theta \leftarrow \Theta + \Delta\Theta;$ 
```

W ostatniej warstwie (MSE)

$$\mathbf{e}^m = f'(\mathbf{a}^m) \odot (\hat{\mathbf{y}} - \mathbf{y})$$

Pozostałe warstwy

$$\mathbf{e}^{k-1} = f'(\mathbf{a}^{k-1}) \odot \left(\mathbf{e}^k (W^k)^T \right)$$

Pochodna wagi

$$\frac{\partial MSE}{\partial w_{i,j}^k} = e_j^k f'(a_i^{k-1})$$

Początkowe wagi

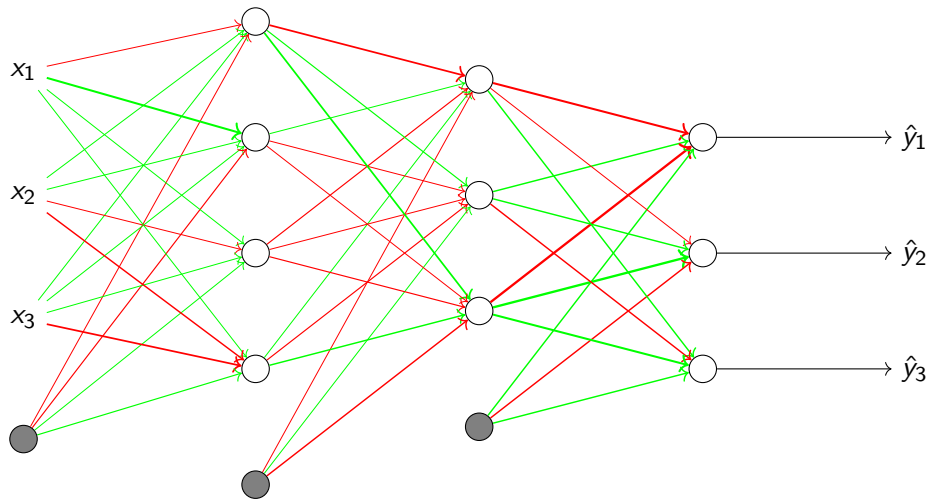
- Duży wpływ na zbieżność
- Małe wartości dookoła zera
- Najprostsze warianty
 - Rozkład gaussowski
 - Rozkład jednostajny
- Związek z funkcją aktywacji
- **Xavier**

$$w_{i,j}^k \in \left[-\frac{\sqrt{6}}{\sqrt{IN + OUT}}, \frac{\sqrt{6}}{\sqrt{IN + OUT}} \right]$$

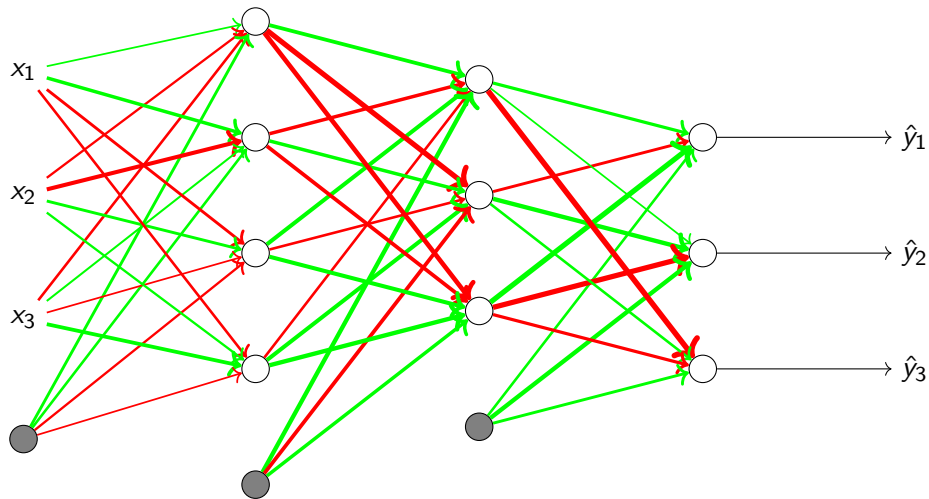
Xavier Glorot and Yoshua Bengio. "Understanding the difficulty of training deep feedforward neural networks". In: *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics, AISTATS 2010, Chia Laguna Resort, Sardinia, Italy, May 13-15, 2010*. Ed. by Yee Whye Teh and D. Mike Titterton. Vol. 9. JMLR Proceedings. JMLR.org, 2010, pp. 249–256

- **He** – inny wariant normalizacji wewnątrz warstwy Kaiming He et al. "Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet Classification". In: *2015 IEEE International Conference on Computer Vision (ICCV) (2015)*. DOI: 10.1109/iccv.2015.123

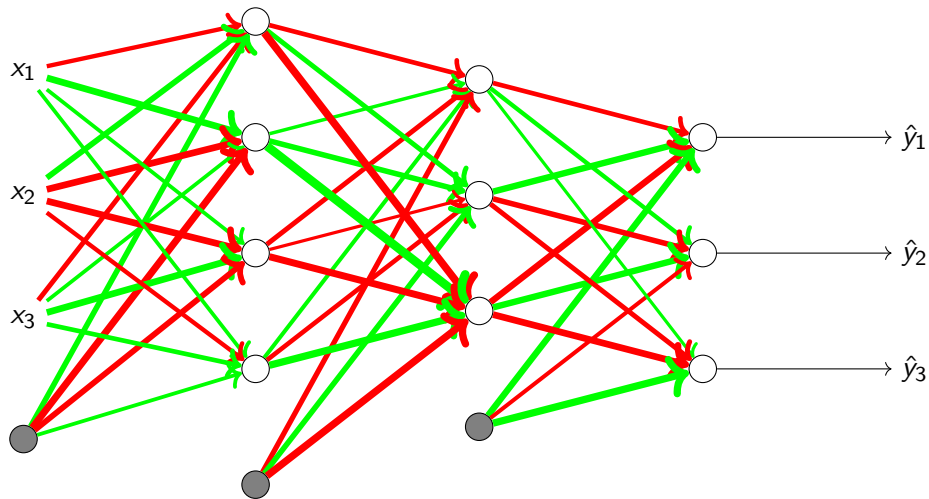
Eksploracja wag



Eksploracja wag



Eksploracja wag



Uczenie z momentem

- η – krok uczenia
- $\lambda \in (0, 1)$ – współczynnik wygaszania momentu

```
1  $\Theta \leftarrow \text{InicjujLosowo};$   
2  $\text{Momentum} \leftarrow [0, 0, \dots 0];$   
3 while  $\neg \text{StopCondition}$  do  
4    $\Delta\Theta \leftarrow [0, 0, \dots 0];$   
5   for  $(\mathbf{X}, \mathbf{Y}) \in \text{Zbiór uczący}$  do  
6      $\hat{Y} \leftarrow \text{Network}(\Theta, X);$   
7      $\Delta\Theta \leftarrow \Delta\Theta - \nabla \text{Network}(\Theta, X);$   
8    $\text{Momentum} \leftarrow \Delta\Theta + \text{Momentum} \cdot \lambda;$   
9    $\Theta \leftarrow \Theta + \eta \text{Momentum};$ 
```


RMSProp

- η – krok uczenia
- β – współczynnik wygaszania

```
1  $\Theta \leftarrow \text{InicjujLosowo};$ 
2  $\mathbb{E}[g^2] \leftarrow [0, 0, \dots, 0];$ 
3 while  $\neg \text{StopCondition}$  do
4    $g \leftarrow 0;$ 
5   for  $(\mathbf{X}, \mathbf{Y}) \in \text{Zbiór uczący}$  do
6      $\hat{Y} \leftarrow \text{Network}(\Theta, X);$ 
7      $g \leftarrow g + \nabla \text{Network}(\Theta, X);$ 
8    $\mathbb{E}[g^2] \leftarrow \beta \mathbb{E}[g^2] + (1 - \beta)g^2 // \text{ Kwadrat po współrzędnych}$ 
9    $\Theta \leftarrow \Theta - \eta \left[ \frac{g_i}{\sqrt{\mathbb{E}[g^2]_i}} \right]_{\forall i};$ 
```

Także inne warianty, np. Adagrad, Adadelta.

Adam

Adaptive moment estimation

Moment+Normalizacja gradientu

- η – Krok uczenia
- $\beta_1, \beta_2 \in [0, 1]$ – Współczynniki wygaszania

Diederik P. Kingma and Jimmy Ba. "Adam: A Method for Stochastic Optimization". In: *3rd International Conference on Learning Representations, ICLR 2015*. Ed. by Yoshua Bengio and Yann LeCun. 2015

```
1  $\Theta \leftarrow \text{InicjujLosowo}, m \leftarrow [0, 0, \dots, 0], v \leftarrow [0, 0, \dots, 0];$ 
2  $t \leftarrow 0;$ 
3 while  $\neg \text{StopCondition}$  do
4    $t \leftarrow t + 1;$ 
5    $g \leftarrow [0, 0, \dots, 0];$ 
6   for  $(\mathbf{X}, \mathbf{Y}) \in \text{Zbiór uczący}$  do
7      $\hat{Y} \leftarrow \text{Network}(\Theta, X);$ 
8      $g \leftarrow g + \nabla \text{Network}(\Theta, X);$ 
9    $m \leftarrow \beta_1 m + (1 - \beta_1)g;$ 
10   $v \leftarrow \beta_2 v + (1 - \beta_2)g^2 // \text{ Kwadrat po współrzędnych}$ 
11   $\hat{m} \leftarrow m / (1 - (\beta_1)^t);$ 
12   $\hat{v} \leftarrow v / (1 - (\beta_2)^t);$ 
13   $\Theta \leftarrow \Theta - \eta \frac{\hat{m}}{\sqrt{\hat{v}}} // \text{ Dzielenie, pierwiastek po}$ 
     $\text{współrzędnych}$ 
```

Adam

Adaptive moment estimation

Moment+Normalizacja gradientu

- η – Krok uczenia
- $\beta_1, \beta_2 \in [0, 1]$ – Współczynniki wygaszania

Diederik P. Kingma and Jimmy Ba. "Adam: A Method for Stochastic Optimization". In: *3rd International Conference on Learning Representations, ICLR 2015*. Ed. by Yoshua Bengio and Yann LeCun. 2015

```
1  $\Theta \leftarrow \text{InicjujLosowo}, m \leftarrow [0, 0, \dots, 0], v \leftarrow [0, 0, \dots, 0];$ 
2  $t \leftarrow 0;$ 
3 while  $\neg \text{StopCondition}$  do
4    $t \leftarrow t + 1;$ 
5    $g \leftarrow [0, 0, \dots, 0];$ 
6   for  $(\mathbf{X}, \mathbf{Y}) \in \text{Zbiór uczący}$  do
7      $\hat{Y} \leftarrow \text{Network}(\Theta, X);$ 
8      $g \leftarrow g + \nabla \text{Network}(\Theta, X);$ 
9    $m \leftarrow \beta_1 m + (1 - \beta_1)g;$ 
10   $v \leftarrow \beta_2 v + (1 - \beta_2)g^2 // \text{ Kwadrat po współrzędnych}$ 
11   $\hat{m} \leftarrow m / (1 - (\beta_1)^t);$ 
12   $\hat{v} \leftarrow v / (1 - (\beta_2)^t);$ 
13   $\Theta \leftarrow \Theta - \eta \frac{\hat{m}}{\sqrt{\hat{v} + \epsilon}} // \text{ Dzielenie, pierwiastek po}$ 
     $\text{współrzędnych}$ 
```

Moment+Normalizacja gradientu

- η – Krok uczenia
- $\beta_1, \beta_2 \in [0, 1]$ – Współczynniki wygaszania

Diederik P. Kingma and Jimmy Ba. "Adam: A Method for Stochastic Optimization". In: *3rd International Conference on Learning Representations, ICLR 2015*. Ed. by Yoshua Bengio and Yann LeCun. 2015

\hat{m} , \hat{v} – estymatory nieobciążone

Indeks dolny – wartość g w zadanej iteracji.

$$m = \sum_{n=0}^t \beta g_{t-n} (1 - \beta)^n$$

Przybliżenie nieskończonej sumy (mamy tylko t iteracji, więc nie możemy jej policzyć):

$$m = \sum_{n=0}^{\infty} \beta g_{t-n} (1 - \beta)^n$$

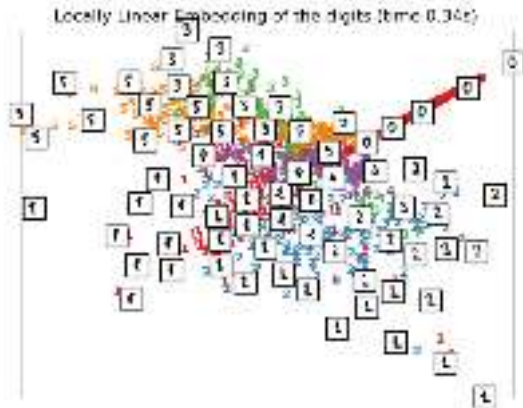
m wyliczone w algorytmie jest estymatorem obciążonym, bo „ogon” sumy jest wyzerowany. Estymator nieobciążony:

$$\hat{m} = \left(\sum_{n=0}^t \beta g_{t-n} (1 - \beta)^n \right) / (1 - \beta^t)$$

MLP, a wypukłość funkcji błędu

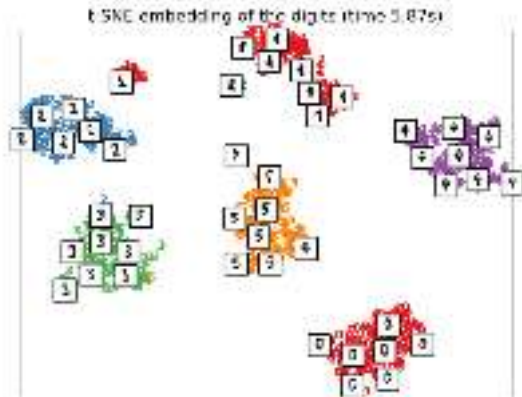
- Zbiór punktów w przestrzeni \mathbb{R}^n
- Podział punktów na k rozłącznych podzbiorów, tak żeby punkty w każdym zbiorze były możliwie „podobne”.

Locally linear embedding (LLE)



Źródło: <https://scikit-learn.org/stable/modules/manifold>

t-distributed Stochastic Neighbor Embedding (t-SNE)



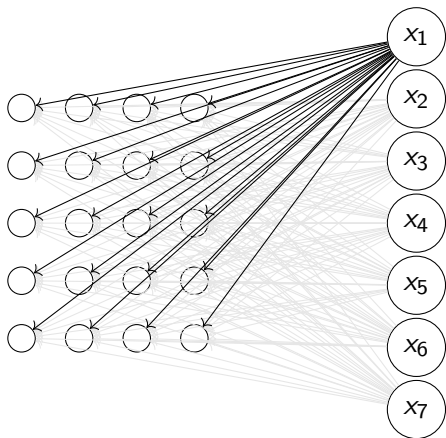
Źródło: <https://scikit-learn.org/stable/modules/manifold>

K-Means

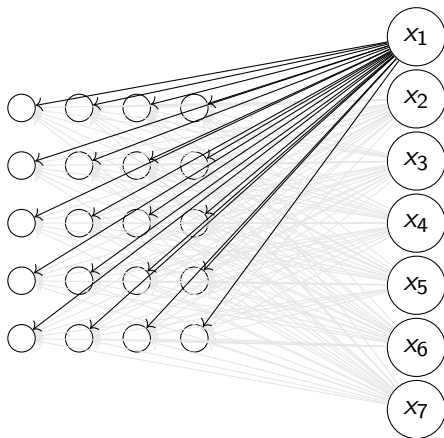
Teuvo Kohonen. “Automatic formation of topological maps of patterns in a self-organizing system”. In: *Proceedings of the 2nd scandinavian Conference on Image Analysis*. 1981, pp. 214–220

Teuvo Kohonen. “Self-organized formation of topologically correct feature maps”. In: *Biological Cybernetics* 43.1 (1982), 59–69. ISSN: 1432-0770. DOI: 10.1007/bf00337288. URL: <http://dx.doi.org/10.1007/bf00337288>

Teuvo Kohonen. “Self-Organization and Associative Memory”. In: *Springer Series in Information Sciences* (1988). ISSN: 0720-678X. DOI: 10.1007/978-3-662-00784-6. URL: <http://dx.doi.org/10.1007/978-3-662-00784-6>



- Niech podobne wektory wejścia pobudzają sąsiednie neruony
- Niech neurony dla różniących się wejść będą odseparowane
- Wagi definiują pozycję neuronu w \mathbb{R}^n



Szkic algorytmu

- 1 Weź wektor ze zbioru danych
- 2 Wybierz neuron, który jest najbliżej tego wzorca
- 3 Przesuń najbliższy neuron w kierunku wzorca
- 4 Przesuń, z pewnym osłabieniem, **sąsiadów** neuronu w kierunku wzorca
 - Sąsiedzi w przestrzeni siatki 2D

```
1  $\forall w_{i,j}^k = RAND$  // Losowa inicjacja
2 for  $t \leftarrow 1 \dots \lambda$  do
3   for  $\mathbf{x} \in RandomShuffle(data)$  do
4      $(i^*, j^*) \leftarrow \arg \min_{i,j} d(\mathbf{w}_{i,j}, \mathbf{x});$ 
5     for  $(i, j) \in N \times M$  do
6        $\mathbf{w}_{i,j} \leftarrow$   

          $\mathbf{w}_{i,j} + \theta((i^*, j^*), (i, j), t) \alpha(t) (\mathbf{x} - \mathbf{w}_{i,j});$ 
```

- $w_{i,j}^k$ – waga między wejściem k , a neuronem na pozycji (i, j)
- λ – liczba iteracji uczenia
- d – metryka w przestrzeni danych
- N, M – wymiary siatki
- $\alpha(t)$ – wygaszanie w czasie
- $\theta(n_1, n_2, t)$ – waga sąsiedztwa wygaszana w czasie

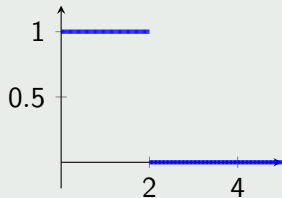
Sąsiedztwo

Względem odległości euklidesowej między punktami

$$\theta(n_1, n_2, t) = f(d(n_1, n_2), t)$$

Koło

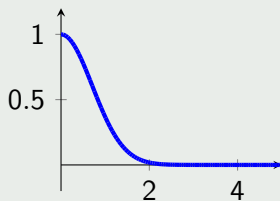
$$f(x, r) = \begin{cases} 1 & x \leq r \\ 0 & x > r \end{cases}$$



Uwaga: skalowanie szerokości.

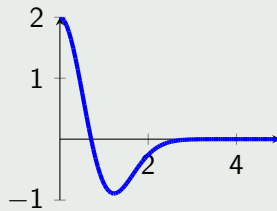
Funkcja Gaussa

$$f(x, t) = e^{-x^2}$$



„Meksykański kapelusz”

$$f(x, r) = (2 - 4x^2)e^{-x^2}$$



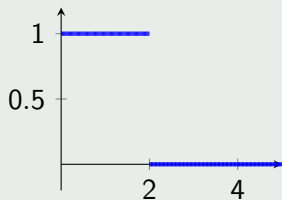
Sąsiedztwo

Względem odległości euklidesowej między punktami

$$\theta(n_1, n_2, t) = f(d(n_1, n_2), t)$$

Koło

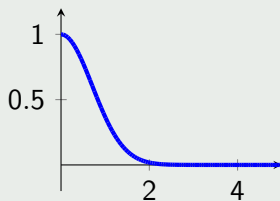
$$f(x, r) = \begin{cases} 1 & x \leq re^{-t/\lambda} \\ 0 & x > re^{-t/\lambda} \end{cases}$$



Uwaga: skalowanie szerokości.

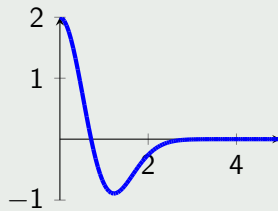
Funkcja Gaussa

$$f(x, t) = e^{-(xt)^2}$$



„Meksykański kapelusz”

$$f(x, r) = (2 - 4x^2)e^{-x^2}$$

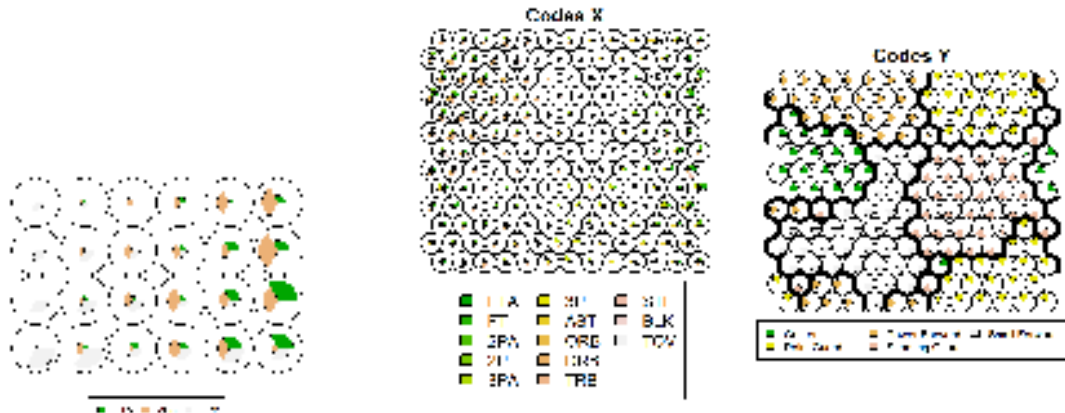


- $\alpha(t) = \eta e^{-t\lambda}$, gdzie η to parametr $\in \mathbb{R}$
 - $\alpha(t) = \eta^t$, $\eta \in (0, 1)$
- $\alpha(t) = \eta$ – brak redukcji w czasie

- Niekoniecznie siatka 2D
- **Siatka sześciokątna 2D**

Przykład działania

https://clarkdatalabs.github.io/soms/SOM_NBA





G. Cybenko. “Approximation by superpositions of a sigmoidal function”. In: *Mathematics of Control, Signals, and Systems* 2.4 (1989), 303–314. ISSN: 1435-568X. DOI: 10.1007/bf02551274. URL: <http://dx.doi.org/10.1007/bf02551274>.



Xavier Glorot and Yoshua Bengio. “Understanding the difficulty of training deep feedforward neural networks”. In: *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics, AISTATS 2010, Chia Laguna Resort, Sardinia, Italy, May 13-15, 2010*. Ed. by Yee Whye Teh and D. Mike Titterton. Vol. 9. JMLR Proceedings. JMLR.org, 2010, pp. 249–256.



Kaiming He et al. “Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet Classification”. In: *2015 IEEE International Conference on Computer Vision (ICCV)* (2015). DOI: 10.1109/iccv.2015.123.



Donald O Hebb. *The organization of behavior*. 1949.



Diederik P. Kingma and Jimmy Ba. “Adam: A Method for Stochastic Optimization”. In: *3rd International Conference on Learning Representations, ICLR 2015*. Ed. by Yoshua Bengio and Yann LeCun. 2015.



Teuvo Kohonen. "Automatic formation of topological maps of patterns in a self-organizing system". In: *Proceedings of the 2nd scandinavian Conference on Image Analysis*. 1981, pp. 214–220.



Teuvo Kohonen. "Self-Organization and Associative Memory". In: *Springer Series in Information Sciences* (1988). ISSN: 0720-678X. DOI: 10.1007/978-3-662-00784-6. URL: <http://dx.doi.org/10.1007/978-3-662-00784-6>.



Teuvo Kohonen. "Self-organized formation of topologically correct feature maps". In: *Biological Cybernetics* 43.1 (1982), 59–69. ISSN: 1432-0770. DOI: 10.1007/bf00337288. URL: <http://dx.doi.org/10.1007/bf00337288>.



Fundusze Europejskie
Wiedza Edukacja Rozwój



**Rzeczpospolita
Polska**

**Politechnika
Warszawska**

Unia Europejska
Europejski Fundusz Społeczny



Zadanie 10 pn. „Modyfikacja programów studiów na kierunkach prowadzonych przez Wydział Matematyki i Nauk Informacyjnych” realizowane jest w ramach projektu „NERW 2 PW. Nauka – Edukacja – Rozwój – Współpraca” współfinansowanego ze środków Unii Europejskiej w ramach Europejskiego Funduszu Społecznego