

UCZENIE SIECI NEURONOWYCH

Lato 2025

Kacper Borys

SPIIS TREŚCI

1. Definicje
2. Uczenie sieci

DEFINICJE

Definicja - perceptron

Niech $\vec{x} \in \mathbb{R}^n$. Perceptronem nazywamy funkcję $P(\vec{x}) = F(\vec{x} \cdot \vec{w} + b)$, gdzie F - funkcja aktywacji (np. sigmoidalna, ReLU, tanh), $\vec{w} \in \mathbb{R}^n, b \in \mathbb{R}$.

Przykładowo - mamy wektor $\vec{x} = (1, 7, 3)$, wagi $\vec{w} = (0.1, 0.5, -2)$ oraz $b = 1$. Dla tych wartości $P(\vec{x}) = \max(0, -1.4) = 0$.

Oznacza to, że dla takich danych perceptron nie aktywuje się.

Definicja - warstwa

Warstwą nazywamy funkcję $W(\vec{x}) = (P_1(\vec{x}), P_2(\vec{x}), \dots, P_k(\vec{x}))$, gdzie P_i to i -ty perceptron i $k \in \mathbb{N}^+$ to liczba perceptronów w warstwie.

Kontynuując przykład. Mamy $\vec{x} = (1, 7, 3)$, wtedy dla konkretnej warstwy $W(\vec{x}) = (0, 0.01, 65.4, 3.1, \dots, 0.422)$. Jest to nasz nowy wektor danych który możemy wykorzystać jako \vec{x}' w następnej warstwie.

Równoważnie możemy przedstawić naszą warstwę $W(\vec{x})$ jako $F(\mathcal{W}\vec{x} + \vec{b})$,

gdzie $\mathcal{W} = \begin{pmatrix} w_{1,1} & w_{1,2} & \cdots & w_{1,n} \\ w_{2,1} & w_{2,2} & \cdots & w_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ w_{k,1} & w_{k,2} & \cdots & w_{k,n} \end{pmatrix}$ to macierz wag, a $\vec{b} = (b_1, b_2, \dots, b_k)$ to wektor przesunięcia.

Definicja - sieć neuronowa

Siecią neuronową nazywamy funkcję $N(\vec{x}) = \vec{y}$, gdzie $\vec{x} \in \mathbb{R}^n$ oraz $\vec{y} \in \mathbb{R}^m$, $m, n \in \mathbb{N}^+$. Jest ona złożeniem funkcji $W^{\text{in}}(\vec{x})$, $W_n^h(\vec{h}_{n-1})$ oraz $W^{\text{out}}(\vec{h}_n)$, gdzie $\vec{h}_0 = W^{\text{in}}(\vec{x})$. Co za tym idzie $N(\vec{x}) = W^{\text{out}}(W_k^h(W_{k-1}^h(\dots(W_1^h(W^{\text{in}}(\vec{x}))))))$

Definicja - gradient

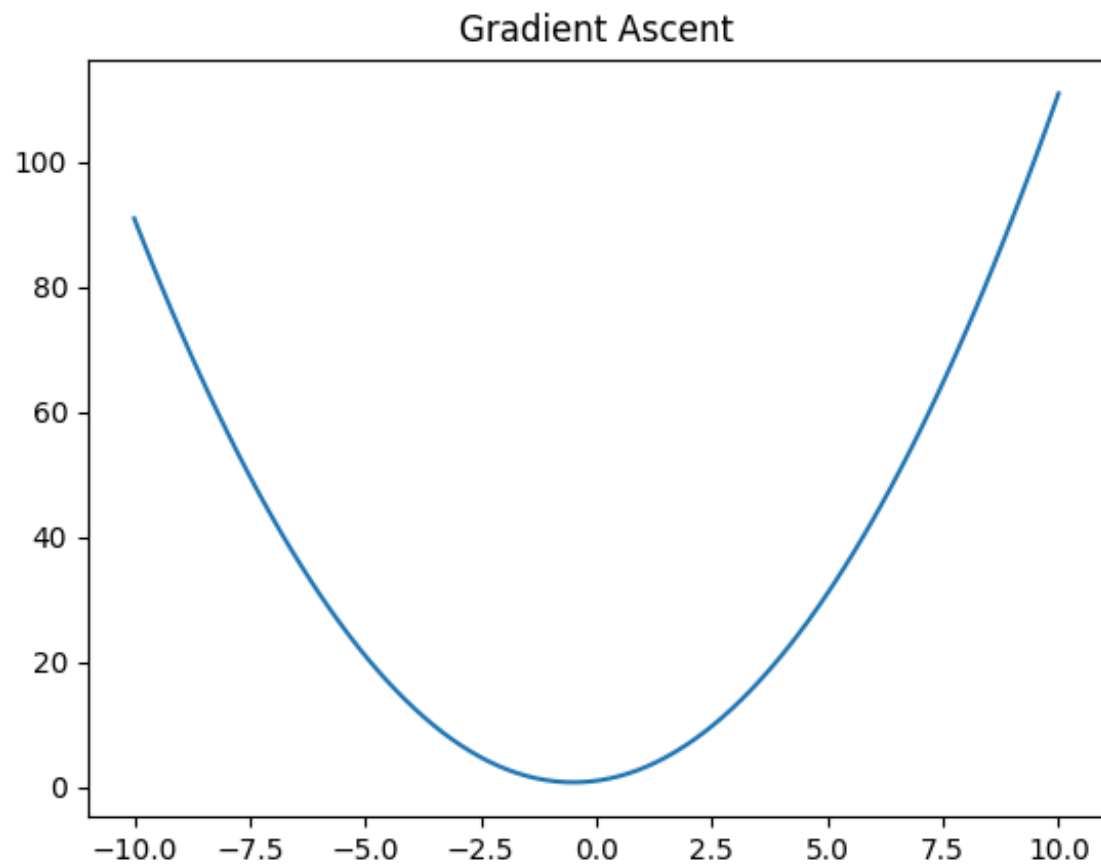
Gradientem funkcji f w punkcie x nazywamy wektor $\nabla f(x) = \left(\frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, \dots, \frac{\partial f}{\partial x_n} \right)$.

Mając jakąś funkcję która przyjmuje wektor, możemy zastosować na niej gradient - zróżniczkować ją względem każdego z parametrów. W ten sposób uzyskujemy wektor, który pokazuje w którą stronę powinniśmy iść aby zwiększyć wartość funkcji.

Przykładowo - mamy funkcję $f(x) = x^2 + 2x + 1$, wtedy $\nabla f(x) = (2x + 2)$, a dla $x = 1$ mamy $\nabla f(1) = 4$. Czyli jeśli jesteśmy w $x = 1$ to powinniśmy iść w prawo na osi liczbowej.

Uwaga

Gradient wskazuje **kierunek** w którym powinniśmy się poruszać. Długość wektora obrazuje **szybkość** zmiany funkcji.



UCZENIE SIECI

Niech $p(x) = ax^2 + bx + c$. Aby uzyskać wartość $p(1) = a + b + c = 2$ musimy dobrać odpowiednie parametry. Możemy wyliczyć je analitycznie, ale możemy również wykorzystać gradient.

Skoro gradient wskazuje kierunek największego wzrostu funkcji, to możemy iść w przeciwną stronę aby zmniejszyć wartość funkcji, innymi słowy - minimalizować ją.

To co chcemy minimalizować to $L(a, b, c) = (2 - p(a, b, c))^2$ - nazywamy to funkcją kosztu. Pokazuje nam ona jak bardzo funkcja oddalona jest od zadanej wartości.

Przykładowo - dobierzmy parametry $a = 1, b = 1, c = 0$, wtedy $p(1) = 2$, a funkcja kosztu wynosi 0. Jeśli dobierzemy $a = 1, b = 1, c = 1$, wtedy $p(1) = 3$, a funkcja kosztu wynosi 1.

Gdy nałożymy gradient na funkcję kosztu i zmienimy znak otrzymamy wektor który wskazuje na największy spadek funkcji.

Obliczmy gradient funkcji kosztu.

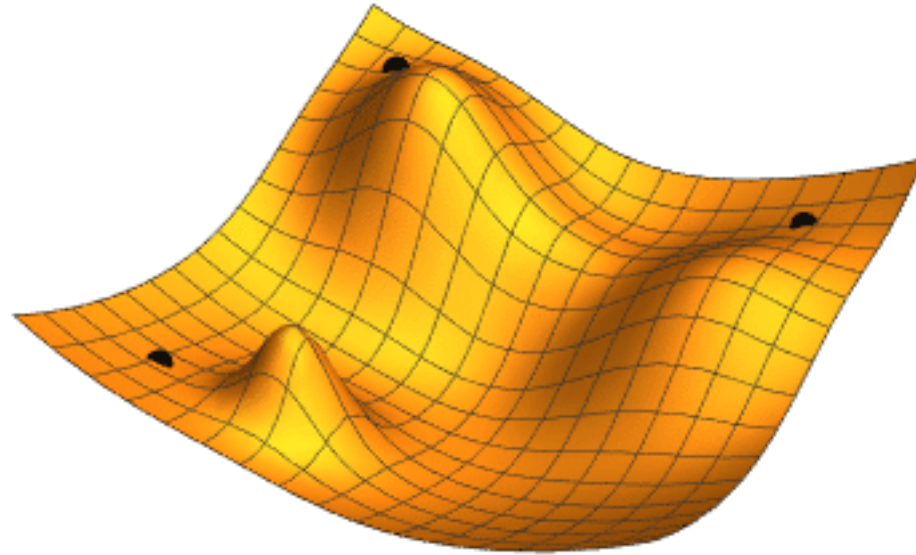
$$\nabla L(a, b, c) = \left(\frac{\partial L}{\partial a}, \frac{\partial L}{\partial b}, \frac{\partial L}{\partial c} \right)$$

Korzystając z reguły łańcuchowej mamy:

$$\nabla L(a, b, c) = 2(p(a, b, c) - 2)(\nabla p(a, b, c))$$

Skoro mamy gradient, to możemy go wykorzystać do zmiany parametrów. Zbudujmy algorytm który będzie aktualizował nasze parametry:

- 1 zainicjuj parametry a, b, c
- 2 zainicjuj krok α
- 3 zainicjuj ε
- 4 **dopóki** $L(a, b, c) > \varepsilon$
 - 5 | oblicz gradient $\nabla L(a, b, c)$
 - 6 | $(a, b, c) = (a, b, c) - \alpha * \nabla L(a, b, c)$
- 7 **koniec**



Definicja - propagacja wsteczna

Propagacją wsteczną nazywamy proces aktualizowania parametrów sieci neuronowej przy pomocy $\nabla L(\vec{y} - N(\vec{p}))$ gdzie \vec{p} jest wektorem złożonym z wag oraz biasów.

Znając gradient W^{out} , dzięki regule łańcuchowej, możemy obliczyć gradient W_n^{h} , a znając gradient W_n^{h} możemy obliczyć gradient W_{n-1}^{h} i tak dalej aż do W^{in} . Tym sposobem możemy uczyć naszą sieć neuronową.

Te równania dokładnie opisują algorytm propagacji wstecznej:

$$\delta^{\text{out}} = \nabla_{\sigma(\vec{h})} L \odot \sigma'(\vec{h})$$

$$\delta_{n-1}^h = (W_n)^T \delta_n \odot \sigma'(\vec{h}_n)$$

$$\delta_n = \frac{\partial L}{\partial \vec{b}_n}$$

$$\sigma(\vec{h}_{n-1}) \delta_n = \frac{\partial L}{\partial \vec{w}_n}$$

Gdzie \odot jest mnożeniem Hadamarda (*element-wise*), σ jest funkcją aktywacji, a L jest funkcją kosztu.