

Uczenie głębokie — zarys

Bartosz Hanc

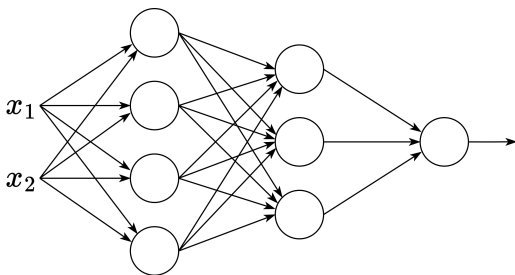
10 lutego 2025

1 Wprowadzenie

Historia sztucznych sieci neuronowych sięga lat 50. XX wieku, kiedy to Frank Rosenblatt zaproponował model perceptronu – prostej 3-warstwowej sieci, której zadaniem miało być rozpoznawanie znaków alfanumerycznych. Podstawową jednostką przetwarzającą w takiej sieci był neuron McCullocha–Pittsa, który realizował następujące przekształcenie $F : \mathbb{R}^n \mapsto \mathbb{R}$

$$F(\mathbf{x}) = \varphi(\mathbf{x}\mathbf{w}^\top + b),$$

gdzie $\mathbf{x} = (x_1, \dots, x_n)$ określa sygnały wchodzące do neuronu, φ jest pewną nieliniową funkcją zwaną funkcją aktywacji, a \mathbf{w} i b to pewne parametry (zwane również wagami) określające, które z sygnałów mają być wzmocnione, a które osłabione. Perceptron jest zbudowany z wielu takich neuronów połączonych ze sobą. Neurony te są zorganizowane w warstwy, w obrębie których nie ma połączeń, natomiast między sąsiednimi warstwami neurony są połączone „każdy z każdym”. Na Rysunku 1 przedstawiono schemat takiego prostego perceptronu wielowarstwowego. Strzałki pokazują przepływ sygnałów w sieci, natomiast każdy wierzchołek reprezentuje odwzorowanie nieliniowe F z właściwymi dla siebie wagami. Efektywnie perceptron wielowarstwowy realizuje więc pewne przekształcenie nieliniowe $\mathbb{R}^n \mapsto \mathbb{R}^m$.



Rysunek 1: Schemat prostego perceptronu wielowarstwowego realizującego odwzorowanie $\mathbb{R}^2 \mapsto \mathbb{R}$

Podany powyżej opis konstrukcji sztucznych sieci neuronowych, choć często przywoływany we wprowadzeniach, nie jest zbyt użyteczny w nowoczesnym

podejściu do uczenia głębokiego. Bardziej przydatnym sformułowaniem, czym właściwie są sieci neuronowe jest następująca definicja: sieć neuronowa to dowolny skierowany graf acykliczny, którego wierzchołki reprezentują dowolną sparametryzowaną funkcję o argumentach i wartościach tensorowych, a skierowane krawędzie reprezentują przepływ tych tensorów. Sieć neuronowa realizuje więc pewne przekształcenie wejściowych tensorów (przykładowo może to być 4-wymiarowy tensor reprezentujący zbiór kolorowych obrazków). Trening sieci neuronowej będzie polegał na takiej zmianie jej parametrów, aby sieć realizowała określone zadanie.

Pojęcie tensora w uczeniu maszynowym nie ma wiele wspólnego z tensorami używanymi w algebrze, czy w fizyce. W naszym kontekście tensor to po prostu wielowymiarowa tablica liczb. Standardowo w tym tekście tensory będziemy oznaczać literami wytłuszczonymi, np. \mathbf{A} , natomiast odnosząc się do ustalonego elementu tensora \mathbf{A} będziemy zapisywać $A_{\alpha_1, \dots, \alpha_k}$ lub A_α , gdzie α będzie oznaczać wielowskażnik $(\alpha_1, \dots, \alpha_k)$. Liczbę k nazywamy wymiarem tensora \mathbf{A} . Zbiór wszystkich tensorów wymiaru k o liczbie elementów wzdłuż osi i równej n_i będziemy oznaczać $\mathbb{R}^{(n_1, \dots, n_k)}$.

Mamy więc graf skierowany posiadający N wierzchołków. Dla każdego wierzchołka $i \in [N]$ wprowadzimy następujące wielkości:

- \mathcal{N}_i – zbiór następników, tj. wierzchołków $j \in [N]$, dla których istnieje krawędź skierowana (i, j) ,
- \mathcal{P}_i – zbiór poprzedników, tj. wierzchołków $j \in [N]$, dla których istnieje krawędź skierowana (j, i) ,
- $\mathbf{F}^{(i)}(\cdot; \boldsymbol{\theta}^{(i)})$ – funkcję tensorową realizowaną przez dany wierzchołek, sparametryzowaną przez $\boldsymbol{\theta}^{(i)}$,
- $\mathbf{v}^{(i)}$ – wartość funkcji $\mathbf{F}^{(i)}$ przechowywaną aktualnie w wierzchołku (w szczególności może to być tzw. null – \emptyset).

Obliczenie wartości $\mathbf{v}^{(i)}$ możemy wyrazić rekurencyjnie jako

$$\mathbf{v}^{(i)} = \mathbf{F}^{(i)} \left[\left(\mathbf{v}^{(j)} \right)_{j \in \mathcal{P}_i}; \boldsymbol{\theta}^{(i)} \right].$$

Jeśli teraz wartości \mathbf{v} będziemy obliczać w kolejności zwróconej przez sortowanie topologicznie wierzchołków grafu to przy propagacji sygnału w przód odwiedzimy każdy wierzchołek jedynie raz.

Jako przykład rozważmy wspomniany na wstępie perceptron wielowarstwowy składający się z N warstw w pełni połączonych (ang. *fully connected layers*). Sieć taka jest reprezentowana przez skierowaną ścieżkę, której każdy wierzchołek (warstwa w pełni połączona) realizuje przekształcenie $\mathbf{F} : \mathbb{R}^{(\dots, n)} \mapsto \mathbb{R}^{(\dots, m)}$

$$\mathbf{F}(\mathbf{X}; \mathbf{W}, \mathbf{b}) = \varphi(\mathbf{X}\mathbf{W} + \mathbf{b}),$$

gdzie $\mathbf{W} \in \mathbb{R}^{(n, m)}$, $\mathbf{b} \in \mathbb{R}^m$ są parametrami tego przekształcenia, a φ to pewna nieliniowa funkcja aktywacji. Warto zaznaczyć, iż \mathbf{X} jest dowolnym $k \geq 1$ wymiarowym tensorem, którego ostatnia oś ma n elementów (nie musi to być koniecznie macierz, czy wektor). Łatwo zrozumieć powyższy wzór wypisując jawne wyrażenie na element wartości funkcji \mathbf{F}

$$F_{\alpha\mu}(\mathbf{X}; \mathbf{W}, \mathbf{b}) = \varphi \left(\sum_{\beta=1}^n X_{\alpha\beta} W_{\beta\mu} + b_{\mu} \right),$$

gdzie α to wielowskaźnik, natomiast β i μ to skalarne indeksy. Obliczenie wartości $\mathbf{v}^{(i)}$ można dla takiej sieci wyrazić zatem rekurencyjnie jako

$$\mathbf{v}^{(i)} = \varphi \left(\mathbf{v}^{(i-1)} \mathbf{W}^{(i)} + \mathbf{b}^{(i)} \right)$$

gdzie $\mathbf{v}^{(0)}$ to wejście do sieci neuronowej, co wynika z faktu, iż dla ścieżki $\mathcal{N}_i = \{i+1\}$, $\mathcal{P}_i = \{i-1\}$ (zakładając naturalną numerację wierzchołków).

2 Funkcje straty

Idea treningu sieci neuronowych polega na znalezieniu takich wartości parametrów sieci Θ , aby sieć realizowała określone zadanie. Zadanie to jest specyfikowane w postaci zbioru przykładów \mathcal{X} , który zawiera wzorce odpowiedzi modelu. W najbardziej typowym przypadku uczenia nadzorowanego zbiór \mathcal{X} składa się z par zawierających wektor cech i prawidłową etykietę dla nich. Zbiór ten może mieć jednak również inną postać, np. w przypadku trenowania modeli generatywnych dla obrazków, składa się on jedynie z pewnych zdjęć i nie zawiera żadnych etykiet. Nie wchodząc w żadne szczegóły dotyczące struktury samego zbioru \mathcal{X} , w każdym z nich trening sieci

wygląda podobnie. Przede wszystkim wprowadzamy funkcję $L : 2^{\mathcal{X}} \times \Theta \mapsto \mathbb{R}$ zwaną funkcją kosztu lub straty (ang. *cost, loss function*), która w sposób ilościowy określa jak sieć neuronowa z danymi parametrami $\boldsymbol{\theta}$ radzi sobie na dowolnym podzbiórze danych. Następnie wykorzystujemy algorytmy optymalizacji, aby znaleźć parametry które minimalizują stratę (na całym zbiorze lub na jego części)

$$\boldsymbol{\theta}^* = \arg \min_{\boldsymbol{\theta} \in \Theta} L(\mathcal{X}, \boldsymbol{\theta}).$$

Zauważmy, że w takim podejściu istnieje pewna pułapka. Otóż sieć neuronową (czy też dowolny inny model uczenia maszynowego) uczymy na zbiorze danych treningowych. Nie interesuje nas jednak tak naprawdę jego wydajność na tym zbiorze, tylko na zbiorze nowych, niewidzianych w czasie treningu danych (czyli interesuje nas generalizacja). Szukanie minimum globalnego funkcji kosztu na zbiorze treningowym może wówczas doprowadzić do nadmiernego dopasowania (ang. *overfitting*), czyli sytuacji, w której model bardzo dobrze radzi sobie na danych treningowych natomiast fatalnie radzi sobie na danych testowych, których nie były dostępne podczas treningu.

Pozostaje pytanie, jak właściwie konstruujemy funkcje kosztu. Jednym z lepiej umotywowanych podejść (choć nie jedynym i chyba nawet nie najbardziej rozpowszechnionym) jest konstrukcja w oparciu o kryterium maksymalnej wiarygodności. Zasadnicza idea jest taka, iż modelujemy nasze dane jako pochodzące z pewnej sprametryzowanej rodziny rozkładów prawdopodobieństwa, a parametry tego rozkładu są zależne od parametrów sieci. Następnie maksymalizujemy funkcję wiarygodności dla zbioru danych treningowych, przy czym ową maksymalizację wyrażamy najczęściej jako minimalizację zanegowanej logarytmicznej funkcji wiarygodności (ang. *Negated Log-Likelihood*). Poniżej rozpatrzmy szczegółowo konstrukcję dwóch najczęściej używanych funkcji straty odpowiednio dla zadań regresji i klasyfikacji.

2.1 Błąd średniokwadratowy

Rozpatrujemy zagadnienie regresji, którego celem jest przewidzenie dla danego wektora cech \mathbf{x} ciągłej wartości y . Zbiór treningowy \mathcal{X} składa się z n par $(\mathbf{x}^{(i)}, y^{(i)})$. Zakładamy, że wartość y pochodzi z rozkładu normalnego z parametrami $f(\mathbf{x}; \boldsymbol{\theta})$, σ^2 , gdzie wariancja jest znana, a $f(\mathbf{x}; \boldsymbol{\theta})$ to wyjście sieci neuronowej

$$y(\mathbf{x}) \sim \mathcal{N}(f(\mathbf{x}; \boldsymbol{\theta}), \sigma^2).$$

Funkcja wiarygodności ma postać

$$C \cdot \prod_{i=1}^n \exp \left\{ -\frac{1}{2\sigma^2} \left[y^{(i)} - f(\mathbf{x}^{(i)}; \boldsymbol{\theta}) \right]^2 \right\},$$

skąd funkcja kosztu

$$L(\mathcal{X}; \boldsymbol{\theta}) = \frac{1}{2} \sum_{i=1}^n \left[y^{(i)} - f(\mathbf{x}^{(i)}; \boldsymbol{\theta}) \right]^2$$

pomijając człony stałe. Otrzymujemy zatem znany błąd średniokwadratowy. Wygodniej jest przyjąć, że sieć neuronowa przetwarza od razu cały *batch* przykładów uczących, tj. na wejściu otrzymujemy macierz

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}^{(1)} \\ \vdots \\ \mathbf{x}^{(n)} \end{bmatrix}$$

i zwracamy wektor odpowiedzi, tzn. sieć neuronowa realizuje przekształcenie $\mathbf{f} : \mathbb{R}^{(n,d)} \mapsto \mathbb{R}^n$. W takim przypadku błąd średniokwadratowy możemy zapisać wygodniej jako

$$L(\mathcal{X}; \boldsymbol{\theta}) = \frac{1}{2} \|\mathbf{y} - \mathbf{f}(\mathbf{X}; \boldsymbol{\theta})\|^2.$$

Ponieważ będzie to nam potrzebne później, podajemy tutaj również wzór na pochodną powyższej funkcji kosztu po odpowiedzi sieci

$$\frac{\partial L}{\partial f_\alpha} = f_\alpha - y_\alpha.$$

2.2 Entropia krzyżowa

Podobnie jak w przypadku regresji postępujemy w przypadku klasyfikacji. Rozważmy najpierw klasyfikację binarną, w której zbiór treningowy składa się z n par zawierających wektor cech $\mathbf{x}^{(i)}$ oraz etykietę $t^{(i)} \in \{0, 1\}$. Tym razem dane będziemy modelować przez rozkład dwumianmowy z pewnym parametrem $\pi(\mathbf{x}; \boldsymbol{\theta}) \in [0, 1]$ określającym prawdopodobieństwo klasy pozytywnej

$$t(\mathbf{x}) \sim \text{Bin}(\pi(\mathbf{x}; \boldsymbol{\theta})).$$

Dodatkowo założymy, że funkcja π ma postać

$$\pi(\mathbf{x}; \boldsymbol{\theta}) = \sigma(f(\mathbf{x}; \boldsymbol{\theta})),$$

gdzie $f(\mathbf{x}; \boldsymbol{\theta})$ to wyjście sieci neuronowej, a σ to funkcja logistyczna

$$\sigma(z) = \frac{1}{1 + e^{-z}}.$$

Funkcja wiarygodności ma w tym przypadku postać

$$\prod_{i=1}^n \pi(\mathbf{x}^{(i)}; \boldsymbol{\theta})^{t^{(i)}} \cdot [1 - \pi(\mathbf{x}^{(i)}; \boldsymbol{\theta})]^{1-t^{(i)}},$$

skąd

$$L(\mathcal{X}; \boldsymbol{\theta}) = - \sum_{i=1}^n \left\{ t^{(i)} \log \pi^{(i)} + (1 - t^{(i)}) \log (1 - \pi^{(i)}) \right\},$$

gdzie $\pi^{(i)} = \pi(\mathbf{x}^{(i)}; \boldsymbol{\theta})$. Taką funkcję kosztu nazywamy (binarną) entropią krzyżową.

3 Metody optymalizacji

4 Wsteczna propagacja błędu

5 Funkcje aktywacji

6 Inicjalizacja

7 Regularyzacja