

Uniwersytet Wrocławski
Wydział Fizyki i Astronomii
Fizyka komputerowa

PRACA MAGISTERSKA

TYTUŁ POLSKI

TYTUŁ ANGIELSKI

Autor:

RAFAŁ SKRZYPIEC

Promotor:

dr hab. KRZYSZTOF GRACZYK

Wrocław, 2018

Streszczenie

Tekst streszczenia

Abstract

Tekst streszczenia

Co poprawić z obecnego tekstu według priorytetu:

- analiza1
 - skomentować otrzymane funkcje postaci, zwrócić uwagę jak było liczone, skąd błędy, co się zgadza/nie zgadza
- analiza2
 - skomentować wyniki
 - wkleić tabelkę z wartościami η ?
- napisać o więzach, cross validacja w inny sposób dla tych dwóch zbiorów, dodajemy 2 punkty i jest 520
- Twierdzenie Cybenki
 - bardziej dokładnie i logicznie przeprowadzić dowód twierdzenia?
 - opisać wizualne przedstawienie twierdzenia i stworzyć rysunki z polskimi podpisami
- tekst o sigmoidzie też do ponownego przejrzania
- uporządkować kolejność podrozdziałów nt. metodologii
- spis treści mniej szczegółowy, tak aby zmieścił się na jednej stronie

Spis treści

1	Wstęp	7
1.1	Fizyka zagadnienia i cel pracy - bardzo ogólnie	7
1.2	Sieci neuronowe, wysokopoziomwy opis, zastosowania - regresja i klasyfikacja - bardzo ogólnie	7
2	Sieci Neuronowe	9
2.1	Historia sieci - perceptron, biologia	9
2.2	Funkcje aktywacji, dlaczego sigmoidy	9
2.2.1	Interpretacja probabilistyczna sigmoidy	11
2.3	Opis algorytmu uczenia prostej sieci	13
2.4	Uniwersalne twierdzenie aproksymacyjne (Twierdzenie Cybenki)	16
2.4.1	Dowód matematyczny	16
2.4.2	Przedstawienie wizualne działania	19
2.5	Problem Bias - Variance	19
3	Metodologia analizy	21
3.1	Fizyka zjawiska	21
3.2	Keras	23
3.3	Generowanie sztucznych danych	25
3.4	Walidacja krzyżowa	28
3.5	Wczesne zatrzymanie	29
3.6	Ilość neuronów	31
3.7	Algorytm uczący	33
4	Wyniki analizy	39
4.1	Analiza nr 1	39
4.2	Analiza nr 2	46

Rozdział 1

Wstęp

- 1.1 Fizyka zagadnienia i cel pracy - bardzo ogólnie
- 1.2 Sieci neuronowe, wysokopoziomwy opis, zastosowania - regresja i klasyfikacja - bardzo ogólnie

Rozdział 2

Sieci Neuronowe

2.1 Historia sieci - perceptron, biologia

2.2 Funkcje aktywacji, dlaczego sigmoidy

Funkcja aktywacji to funkcja, która działa na każdy neuron w sieci neuronowej, jako argument przyjmuje sumę iloczynów wartości neuronów z warstwy poprzedzającej i odpowiadających im wag. Każda z warstw sieci neuronowej może mieć zdefiniowaną inną funkcję aktywacji.

Perceptron, który był inspiracją powstania sieci neuronowych został skonstruowany jako uproszczony model biologicznego neuronu. W neurobiologii, neuron jest komórką, która odbiera, przetwarza i przesyła informacje wykorzystując elektryczne i chemiczne sygnały. Neurony połączone są ze sobą przez synapsy, jeden neuron może otrzymywać informacje od wielu komórek nerwowych. Jeśli suma sygnałów elektrycznych z wejściowych synaps przekroczy pewien próg, wtedy neuron transmituje dalej sygnał elektryczny. Perceptron naśladował ten mechanizm stosując przedstawioną w lewym górnym rogu na Rys. 2.1 funkcję Heaviside'a jako funkcję aktywacji. Funkcja przyjmuje wartość jeden jeśli suma wartości wejściowych jest większa od zera, w innym przypadku funkcja przyjmuje wartość zero i neuron nie propaguje sygnału. Perceptron jest najprostszym przykładem sieci neuronowej.

Wyniki badań przeprowadzonych przez [publikacja, publikacja] pokazały, że wśród pożądanych cech funkcji aktywacji znajdują się atrybuty, których funkcja Heaviside'a nie posiada, z tego powodu nie jest w praktyce często stosowana.

Koniecznym wymaganiem jest nieliniowość stosowanej funkcji, jest to cecha, która pozwala sieci neuronowej odwzorować nieliniowe zależności [Le-

Cun, Cybenko?, Hornik]. Jedynym wyjątkiem od reguły jest stosowanie w problemach regresyjnych funkcji tożsamościowej w ostatniej warstwie wyjściowej. Dobrze gdy funkcja posiada ciągłą pochodną, pozwala to na stosowanie metod optymalizacji opartych o obliczanie gradientu. Tu wyjątkiem jest stosowana poprawiona jednostka liniowa (ReLU), również przedstawiona na Rys. 2.1. Zakładając, że w zerze jej gradient równy jest zero możemy skorzystać z jej wielu zalet. Wśród nich wymienia się dokładniejsze odwzorowanie obserwowanego w neurobiologii zjawiska – tylko neurony, które otrzymały odpowiednio silny sygnał są aktywowane. Brak podatności na przeuczenie, podczas inicjalizacji sieci losowymi wagami, tylko około 50% ukrytych neuronów jest aktywowanych. Brak problemu znikającego gradientu uniemożliwiającego uczenie, w porównaniu do sigmoidy, u której wysycha się on w obu kierunkach. Jest to również funkcja często wykorzystywana w metodach głębokiego uczenia. W warstwach spłotowych sieci która służy do rozpoznawania obrazów wykorzystamy ReLU poszukując atrybutów, które nie zmieniają się podczas jej użycia.

Funkcje sigmoidalne

Częstym wyborem funkcji aktywacji są funkcje sigmoidalne. Jest to grupa monotonicznie rosnących funkcji, których zbiór wartości jest ograniczony przez asymptoty o skończonych wartościach, do których wartość funkcji dąży w $\pm\infty$ [lecun98]. Jednym z najczęściej wykorzystywanych przykładów funkcji sigmoidalnych jest sigmoida zdefiniowana równaniem

$$\sigma(x) = \frac{1}{1 + e^{-x}}. \quad (2.1)$$

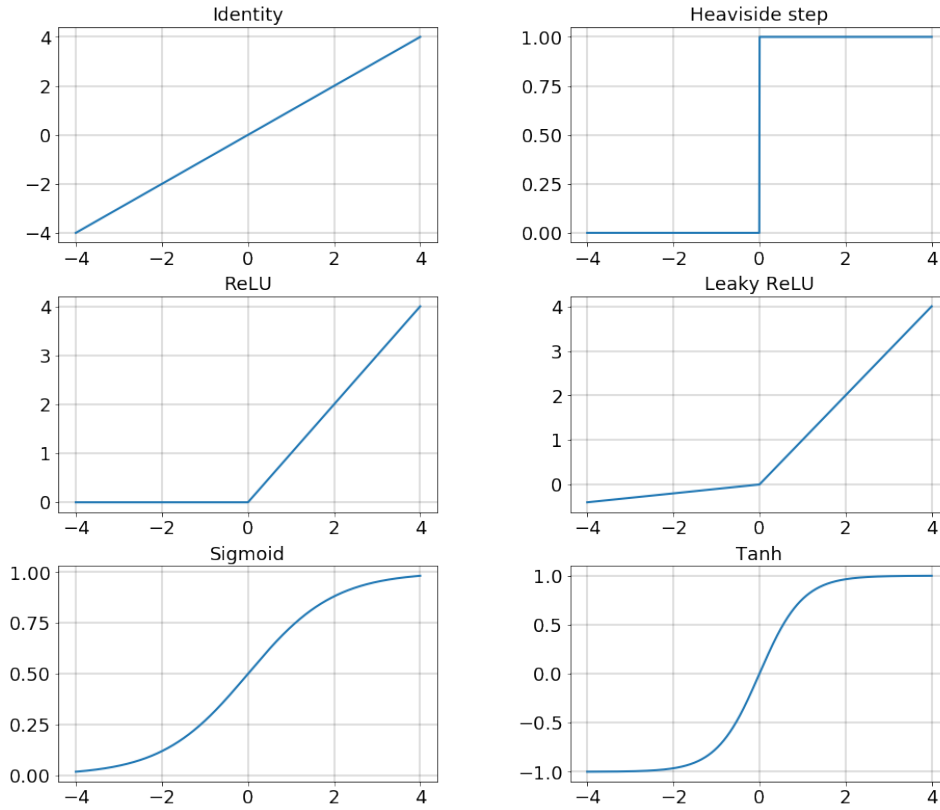
Sigmoida jest różniczkowalna w każdym punkcie co pozwala używać podczas procesu uczenia metod optymalizacji wykorzystujących gradient. Ponadto pochodna względem argumentu x wyraża się prostą relacją

$$\frac{d}{dx}\sigma(x) = \sigma(x)(1 - \sigma(x)). \quad (2.2)$$

Innym przykładem często wykorzystywanej w sztucznych sieciach neuronowych funkcji sigmoidalnej jest tangens hyperboliczny (prawy dolny róg Rys. 2.1). Wzór tej funkcji możemy wyrazić korzystając z definicji sigmoidy

$$\operatorname{tgh}(x) = 2\sigma(2x) - 1 \quad (2.3)$$

Jedną z zalet tej funkcji jest symetryczność względem początku układu współrzędnych.



Rysunek 2.1: Kilka przykładów często stosowanych funkcji aktywacji.

2.2.1 Interpretacja probabilistyczna sigmoidy

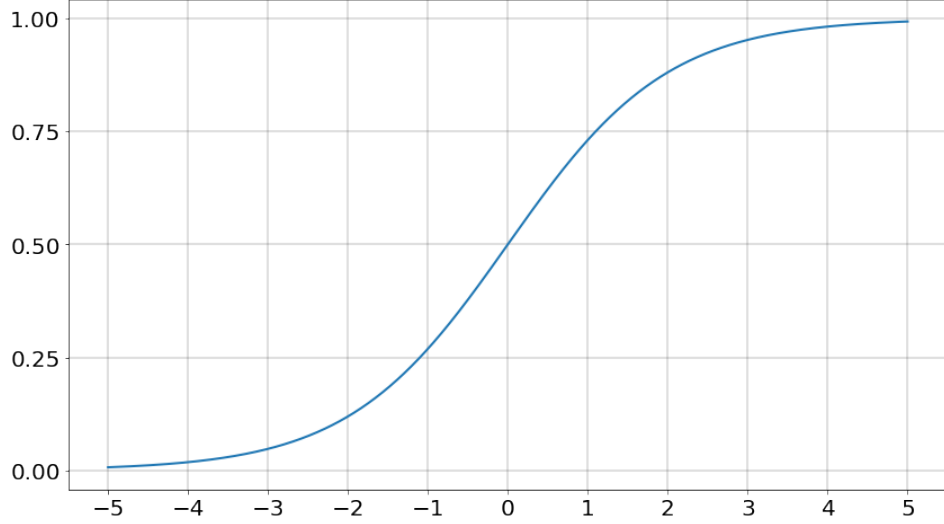
Zastosowanie sigmoidy jako funkcji aktywacji naturalnie wynika z postaci prawdopodobieństwa a posteriori w Bayesowskim podejściu do problemu klasyfikacji dwóch klas. Rozważmy sztuczną sieć neuronową z jedną warstwą ukrytą oraz funkcję dyskryminacyjną $y(\mathbf{x})$ taką, że wektor \mathbf{x} jest przypisany do klasy C_1 jeśli $y(\mathbf{x}) > 0$ i do klasy C_2 jeśli $y(\mathbf{x}) < 0$.

W najprostszej, liniowej formie funkcja może być zapisana jako:

$$y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b_0. \quad (2.4)$$

Wektor \mathbf{w} , to d -wymiarowy wektor wag, natomiast parametr b_0 to bias. Rozważmy funkcję $g(\cdot)$ nazywaną dalej funkcją aktywacji, która jako argument przyjmuje jako argument sumę z równania (2.5):

$$y = g(\mathbf{w}^T \mathbf{x} + b_0) \quad (2.5)$$



Rysunek 2.2: Przykład funkcji sigmoidalnej - sigmoida, $\sigma(x) = \frac{1}{1+e^{-x}}$

Założmy, że funkcja rozkładu prawdopodobieństwa danych pod warunkiem klasy C_k zadane jest przez wielowymiarowy rozkład normalny z równymi macierzami kowariancji $\Sigma_1 = \Sigma_2 = \Sigma$

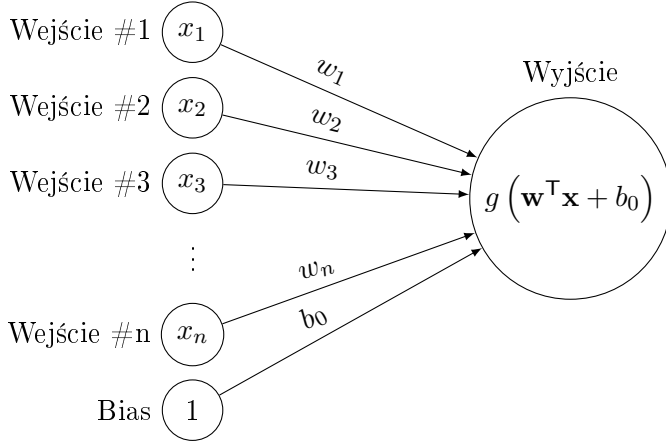
$$p(x|C_k) = \frac{1}{(2\pi)^{\frac{d}{2}} |\Sigma|^{\frac{1}{2}}} \exp \left[-\frac{1}{2} (\mathbf{x} - \mu_k)^T \Sigma^{-1} (\mathbf{x} - \mu_k) \right]. \quad (2.6)$$

Prawdopodobieństwo a posteriori klasy C_1 można zapisać używając twierdzenia Bayesa:

$$\begin{aligned} p(C_1|\mathbf{x}) &= \frac{p(\mathbf{x}|C_1)p(C_1)}{p(\mathbf{x}|C_1)p(C_1) + p(\mathbf{x}|C_2)p(C_2)} \\ &= \frac{1}{1 + \frac{p(\mathbf{x}|C_2)p(C_2)}{p(\mathbf{x}|C_1)p(C_1)}} \\ &= \frac{1}{1 + \exp(-a)}, \end{aligned} \quad (2.7)$$

gdzie

$$\begin{aligned} a &= \ln \frac{p(\mathbf{x}|C_1)p(C_1)}{p(\mathbf{x}|C_2)p(C_2)} \\ &= (\mu_1 - \mu_2)^T \Sigma^{-1} \mathbf{x} - \frac{1}{2} \mu_1^T \mu_1 + \frac{1}{2} \mu_2^T \Sigma^{-1} \mu_2 + \ln \frac{p(C_1)}{p(C_2)}, \end{aligned} \quad (2.8)$$



Rysunek 2.3: Re-prezentacja funkcji dyskryminacyjnej $y(x)$ w postaci diagramu sieci neuronowej, mającej n wejść, parametr bias i jedno wyjście.

pamiętając o tym, że macierz kowariancji jest symetryczna otrzymujemy

$$\mathbf{x} = \Sigma^{-1}(\mu_1 - \mu_2) \quad (2.9a)$$

$$b_0 = -\frac{1}{2}\mu_1^T \mu_1 + \frac{1}{2}\mu_2^T \Sigma^{-1} \mu_2 + \ln \frac{p(C_1)}{p(C_2)} \quad (2.9b)$$

Zatem widzimy, że użycie funkcji aktywacji w postaci sigmoidy pozwala nie tylko dokonać decyzji klasyfikacji ale również interpretować wynik funkcji dyskryminacyjnej jako prawdopodobieństwa a posteriori.

2.3 Opis algorytmu uczenia prostej sieci

Dane

Zbiór danych treningowych zawiera m jednowymiarowych próbek zadanych przez wektory $X \in \mathbb{R}^{1 \times m}$ i odpowiadające im wyniki $Y \in \mathbb{R}^{1 \times m}$.

Parametry

Sieć ma dwie warstwy: 1) ukryta, zawierająca L neuronów i 2) wyjściowa, składająca się z 1 neuronu. Warstwy są zdefiniowane przez:

1. parametry warstwy ukrytej, które odwzorowują 1-wymiarowe wektory wejściowe w aktywacje L neuronów: macierz wag $W^h \in \mathbb{R}^{L \times 1}$ i wektor parametru bias $b^h \in \mathbb{R}^{L \times 1}$,

2. parametry warstwy wyjściowej, które odwzorowują L -wymiarowy wektor aktywacji neuronów ukrytych w jeden neuron warstwy wyjściowej: macierz wag $W^o \in 1 \times L$ i wektor bias $b^o \in \mathbb{R}^{1 \times 1}$.

Propagacja sygnału

Wejście każdego neuronu w warstwie ukrytej jest iloczynem danych wejściowych i odpowiadającej im wagi plus parametr bias. Na przykład dla i -tego przykładu danych wejściowych, w l -tym neuronie mamy

$$a_l^{h(i)} = W_l^h x^{(i)} + b_l^h \quad (2.10)$$

Funkcją aktywacyjną neuronów jest sigmoida $\sigma(a) = \frac{1}{1+e^{-a}}$, jako argument przyjmuje ona wejście neuronów:

$$h_l^{h(i)} = \sigma(a_l^{h(i)}) \quad (2.11)$$

Neuron warstwy wyjściowej zawiera sumę iloczynów aktywacji neuronów i odpowiadających im wag plus parametr bias. Dla i -tego przykładu mamy

$$\begin{aligned} a^{o(i)} &= \sum_l W_l^o h_l^{h(i)} + b^o \\ &= \sum_l W_l^o \sigma(a_l^{h(i)}) + b^o \\ &= \sum_l W_l^o \sigma(W_l^h x^{(i)} + b_l^h) + b^o \end{aligned} \quad (2.12)$$

Jako funkcja straty zostanie wykorzystany błąd średniokwadratowy

$$\begin{aligned} J^{(i)}(\Theta) &= \frac{1}{2} \left(y^{(i)} - a^{o(i)} \right)^2 \\ J(\Theta) &= \frac{1}{m} \sum_{i=1}^m J^{(i)}(\Theta) = \frac{1}{2m} \sum_{i=1}^m \left(y^{(i)} - a^{o(i)} \right)^2. \end{aligned} \quad (2.13)$$

Propagacja wsteczna

Użycie reguły łańcuchowej umożliwia obliczenie gradientu funkcji straty względem parametrów sieci neuronowej.

Na początku policzmy gradient względem wyniku warstwy wyjściowej.

$$\frac{\partial J}{\partial a^{o(i)}} = \frac{1}{m} \left(y^{(i)} - a^{o(i)} \right), \quad (2.14)$$

następnie policzmy gradient wyjścia neuronów ukrytych:

$$\frac{\partial J}{\partial h_l^{h(i)}} = \frac{\partial J}{\partial a^{o(i)}} \frac{\partial a^{o(i)}}{\partial h_l^{h(i)}} = \frac{\partial J}{\partial a^{o(i)}} W_{ol}^{o}, \quad (2.15)$$

co umożliwia obliczenie gradientu względem wejścia neuronów ukrytych:

$$\frac{\partial J}{\partial a_l^{h(i)}} = \frac{\partial J}{\partial h_l^{h(i)}} \frac{\partial h_l^{h(i)}}{\partial a_l^{h(i)}} = \frac{\partial J}{\partial h_l^{h(i)}} h_l^{h(i)} (1 - h_l^{h(i)}) \quad (2.16)$$

gdzie została wykorzystana relacja

$$\frac{\partial \sigma(x)}{\partial x} = \sigma(x)(1 - \sigma(x)).$$

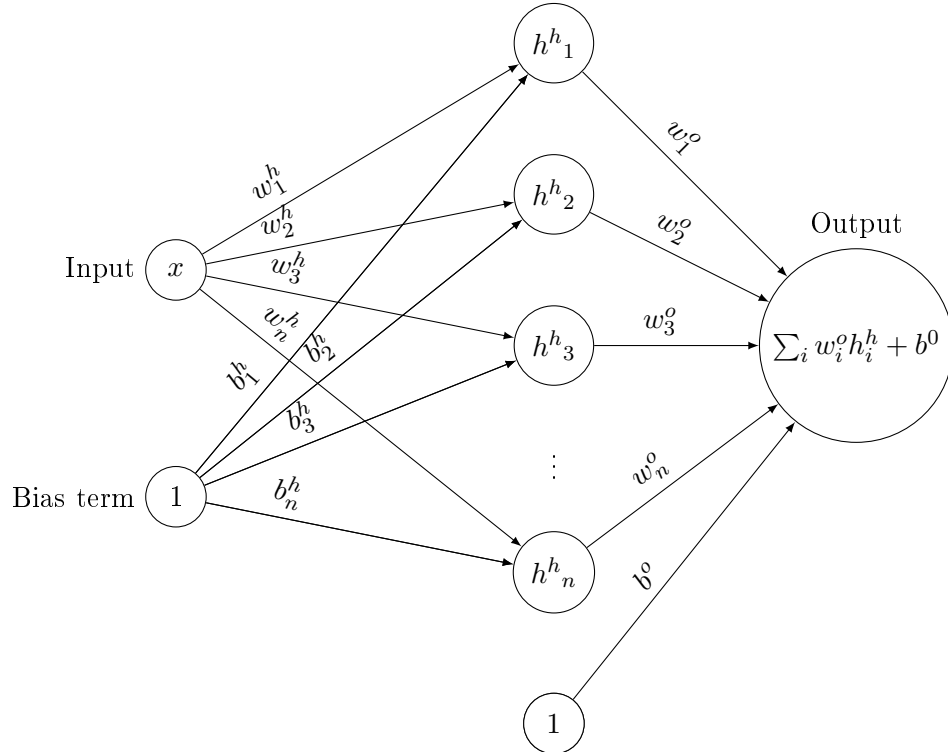
Ostatecznie możemy policzyć gradienty względem parametrów sieci, np. dla warstwy wejściowej:

$$\frac{\partial J}{\partial W_{ol}^{o}} = \sum_i \frac{\partial J}{\partial a^{o(i)}} \frac{\partial a^{o(i)}}{\partial W_{ol}^{o}} = \sum_i \frac{\partial J}{\partial a^{o(i)}} h_l^{h(i)}, \quad (2.17)$$

$$\frac{\partial J}{\partial b^o} = \sum_i \frac{\partial J}{\partial a^{o(i)}} \frac{\partial a^{o(i)}}{\partial b^o} = \sum_i \frac{\partial J}{\partial a^{o(i)}}. \quad (2.18)$$

2.4 Uniwersalne twierdzenie aproksymacyjne (Twierdzenie Cybenki)

2.4.1 Dowód matematyczny



Według uniwersalnego twierdzenia aproksymacyjnego jednokierunkowa sieć neuronowa z jedną warstwą ukrytą i skończoną ale wystarczająco dużą liczbą neuronów, może przybliżyć z dowolną dokładnością każdą funkcję.

W 1989 roku Cybenko [cytowanie] udowodnił uniwersalne twierdzenie aproksymacyjne dla jednokierunkowej sieci neuronowej z sigmoidalną funkcją aktywacji. Jeszcze w tym samym roku, po pracy Cybenki ukazała się praca Hornika, Stinchcombe'a and White'a, którzy udowodnili prawdziwość powyższego twierdzenia dla dowolnej funkcji aktywacji.

Funkcje sigmoidalne to rodzina funkcji szeroko stosowanych w jednokierunkowych sieciach neuronowych, szczególnie tych stworzonych do celów regresji. W tej części zaprezentuję dowód uniwersalnego twierdzenia aproksymacyjnego podany przez Cybenkę w 1989 roku, następnie zademonstruję dowód wizualny posługując się sigmoidą jako funkcją aktywacji.

Niech I_n oznacza n -wymiarową jednostkową kostkę, $[0, 1]^n$. $C(I_n)$ to przestrzeń ciągłych funkcji na I_n . Dodatkowo, niech $M(I_n)$ oznacza przestrzeń skończonych, regularnych miar borelowskich na n -wymiarowej kostce jednostkowej I_n .

Definicja 2.4.1. Miara μ jest regularna jeśli dla każdego mierzalnego zbioru A , $\mu(A)$ równa się supremum miar zamkniętych podzbiorów A i infimum otwartych nadzbiorów A . [Probability measures on metric spaces K.R. Parthasarathy]

Definicja 2.4.2. Funkcja $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ jest funkcją sigmoidalną jeśli

$$\sigma(x) \rightarrow \begin{cases} 1 & \text{as } x \rightarrow +\infty \\ 0 & \text{as } x \rightarrow -\infty \end{cases}$$

Definicja 2.4.3. Funkcja σ jest funkcją dyskryminacyjną jeśli dla miary $\mu \in M(I_n)$ zachodzi

$$\int_{I_n} \sigma(w^T x + b_0) d\mu(x) = 0 \quad (2.19)$$

dla każdego $w \in \mathbb{R}$ i $b_0 \in \mathbb{R}$ co implikuje, że $\mu = 0$.

Twierdzenie 2.4.1. Każda ograniczona, mierzalna funkcja sigmoidalna σ jest funkcją dyskryminacyjną. W szczególności każda ciągła funkcja sigmoidalna jest dyskryminacyjna. [cytowanie Cybenko]

Dowód uniwersalnego twierdzenia aproksymacyjnego przy wykorzystaniu funkcji sigmoidalnych wymaga wprowadzenia kilku przydatnych definicji i twierdzeń. Pierwsze z nich to twierdzenie Hahna-Banacha, które formułuje możliwość rozszerzenia każdego ograniczonego funkcjonału liniowego z podprzestrzeni unormowanej na całą podprzestrzeń, przy zachowaniu jego właściwości.

Twierdzenie 2.4.2 (Twierdzenie Hahna-Banacha). Niech X to rzeczywista przestrzeń wektorowa, p to funkcja rzeczywista zdefiniowana na X spełniająca

$$p(\alpha x + (1 - \alpha)y) \leq \alpha p(x) + (1 - \alpha)p(y) \quad \forall \alpha \in [0, 1], x, y \in X$$

Przypuśćmy, że λ to funkcjonal liniowy zdefiniowany na zbiorze $Y \subset X$, który spełnia

$$\lambda(x) \leq p(x) \quad \forall x \in Y.$$

Wtedy istnieje funkcjonal liniowy Λ zdefiniowany na X spełniający

$$\Lambda(x) \leq p(x) \quad \forall x \in X,$$

tak, że

$$\Lambda(x) = \lambda(x) \quad \forall x \in Y.$$

Reed & Simon (1980), Methods of Modern Mathematical Physics. Functional Analysis

Definicja 2.4.4. Przestrzeń $\mathcal{L}(\mathcal{H}, \mathbb{C})$ nazywana jest przestrzenią dualną przestrzeni Hilberta \mathcal{H} i oznaczamy ją przez \mathcal{H}^* . Elementy \mathcal{H}^* nazywane są ciągłymi funkcjonalami liniowymi.

Reed & Simon (1980), Methods of Modern Mathematical Physics. Functional Analysis

Twierdzenie Riesz opisuje przestrzeń \mathcal{H}^* .

Twierdzenie 2.4.3 (Twierdzenie Riesz (znaleźć polskie źródło)). *Dla każdego $T \in \mathcal{H}^*$, istnieje unikalne $y_T \in \mathcal{H}$ takie, że*

$$T(x) = \langle y_T, x \rangle \quad \forall x \in \mathcal{H}$$

Ponadto

$$\|y_T\|_{\mathcal{H}} = \|T\|_{\mathcal{H}^*}$$

Reed & Simon (1980), Methods of Modern Mathematical Physics. Functional Analysis

Twierdzenie 2.4.4. *Niech σ będzie ciągłą funkcją dyskryminacyjną, wtedy skończona suma*

$$G(x) = \sum_{i=1}^N w_i^o \sigma(w_i^h{}^\top x + b_i^h) \quad (2.20)$$

jest gęsta w $C(I_n)$. Innymi słowy, dla danej funkcji $f \in C(I_n)$ i $\epsilon > 0$, istnieje suma

$G(x)$ mająca powyższą postać, dla której

$$|G(x) - f(x)| < \epsilon \quad \forall x \in I_n$$

Dowód. Niech $S \subset C(I_n)$ będzie zbiorem funkcji w postaci $G(x)$ lub w innych słowach - zbiorem sieci neuronowych. Z pewnością S jest podprzestrzenią liniową $C(I_n)$. Jeśli S jest gęsty, domknięcie S jest całą przestrzenią $C(I_n)$.

Przyjmijmy, że domknięcie S nie jest całą przestrzenią $C(I_n)$. Wtedy domknięcie $S - S'$ jest domkniętą podprzestrzenią $C(I_n)$. Przez twierdzenie Hahna-Banacha, istnieje ograniczony funkcjonal liniowy na $C(I_n)$, nazwijmy go L , z własnością, że $L \neq 0$ ale $L(S') = L(S) = 0$.

Przez twierdzenie Riesz, ograniczony funkcjonal liniowy L ma postać

$$L(h) = \int_{I_n} h(x) d\mu(x)$$

dla $\mu \in M(I_n)$, dla każdego $h \in C(I_n)$. W szczególności, odkąd $\sigma(w^\top x + b) \in S'$ dla każdego w i b , musi zachodzić

$$\int_{I_n} \sigma(w^\top x + b) d\mu(x) = 0$$

Jednakże, założyliśmy, że σ jest funkcją dyskryminacyjną, ten warunek implikuje, że $\mu = 0$ co jest sprzeczne z naszym założeniem. Stąd, podprzestrzeń S jest gęsta w $C(I_n)$.

Pokazuje to, że suma

$$G(x) = \sum_{i=1}^N w_i^o \sigma(w_i^h{}^\top x + b_i^h)$$

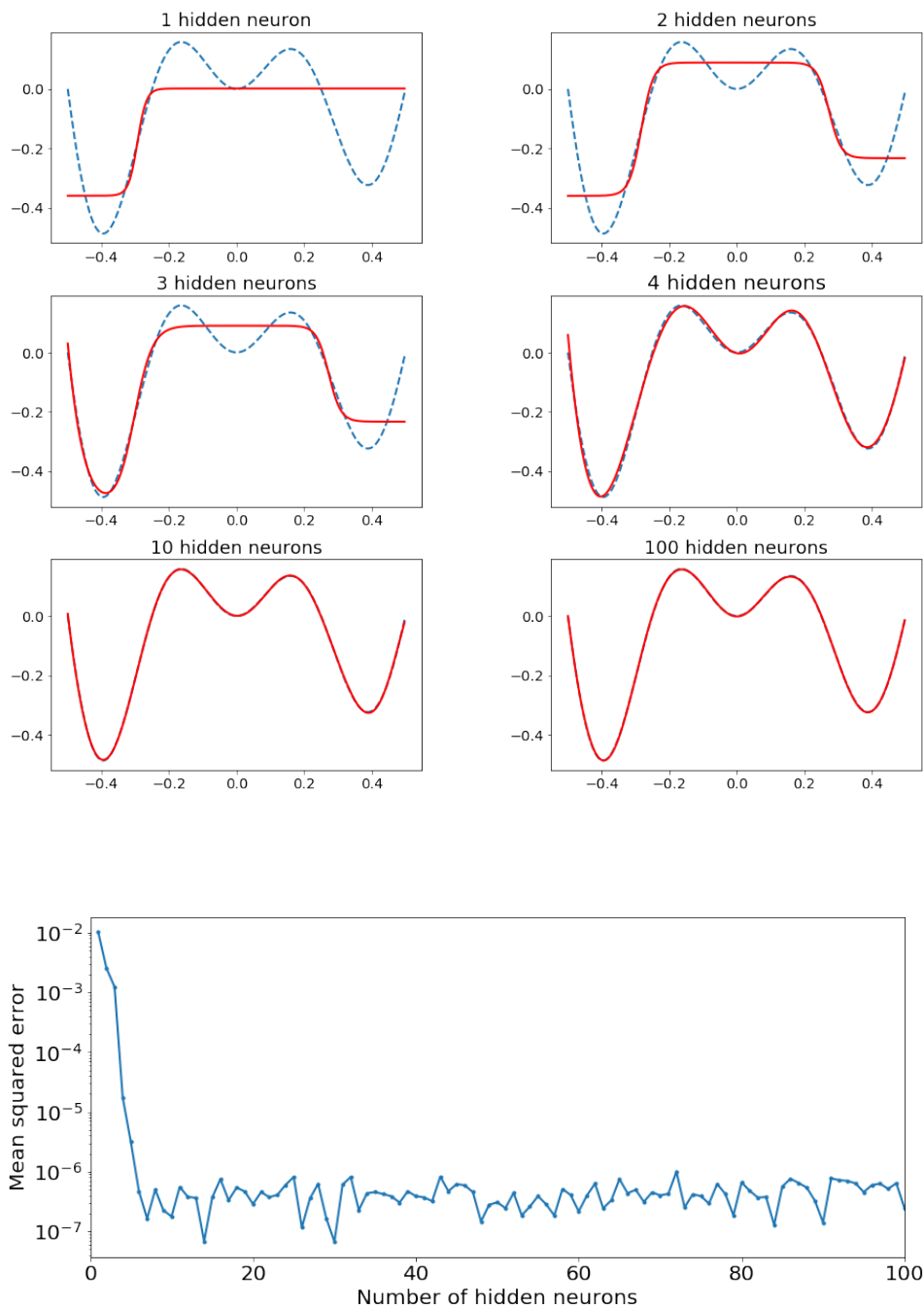
jest gęsta w $C(I_n)$ pod warunkiem, że σ jest ciągła i dyskryminacyjna.

Z twierdzenia wynika, że każda sieć neuronowa o wystarczająco dużej liczbie neuronów w jednej warstwie ukrytej i sigmoidalną funkcją aktywacyjną może z dowolną dokładnością przybliżyć przebieg każdej funkcji.

□

2.4.2 Przedstawienie wizualne działania

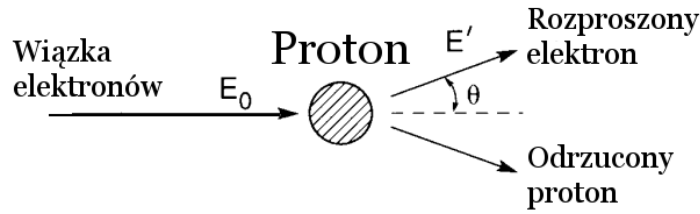
2.5 Problem Bias - Variance



Rozdział 3

Metodologia analizy

3.1 Fizyka zjawiska



Rysunek 3.1: Rozpraszanie elastyczne elektronu o energii początkowej E_0 i energii końcowej E' na jądrze atomu wodoru, θ to kąt rozpraszania.

Energia rozproszonego elektronu E' jest mniejsza niż energia początkowa E_0 o wartość energii przekazanej protonowi o masie M i wynosi

$$E' = \frac{E_0}{1 + \frac{2E_0}{M} \sin^2 \theta/2} \quad (3.1)$$

Zdefiniujmy $Q^2 \equiv -q^2$, q to przekaz czteropędu.

$$\begin{aligned} Q^2 \equiv -q^2 &= -(p^\mu - p'^\mu)^2 = 2M(E_0 - E') \\ &= 4E_0 E' \sin^2 \frac{\theta}{2} \end{aligned} \quad (3.2)$$

Struktura protonu opisywana jest przez dwie funkcje postaci, elektryczną G_{Ep} i magnetyczną G_{Mp} które są transformatami Fouriera, odpowiednio rozkładu ładunku elektrycznego i rozkładu momentu magnetycznego protonu.

Wartości tych istotnych funkcji są wyznaczane eksperymentalnie. Zmierzony w laboratorium przekrój czynny rozpraszania można w przybliżeniu jedno-fotonowym wyrazić formułą:

$$\frac{d\sigma}{d\Omega} = \left(\frac{d\sigma}{d\Omega} \right)_M \times \left[G_{Ep}^2 + \frac{\tau}{\epsilon} G_{Mp}^2 \right] \frac{1}{(1 + \tau)} \quad (3.3)$$

gdzie $\left(\frac{d\sigma}{d\Omega} \right)_M$ to przekrój czynny Motta, który wynosi

$$\left(\frac{d\sigma}{d\Omega} \right)_M = \frac{\pi \alpha^2 E' \cos^2(\theta/2)}{2E^3 \sin^4(\theta/2)}, \quad (3.4)$$

ponadto

$$\tau = \frac{Q^2}{4M^2},$$

$$\epsilon = \left[1 + 2(1 + \tau) \operatorname{tg}^2 \left(\frac{\theta}{2} \right) \right]^{-1},$$

ϵ jest czynnikiem kinematycznym i także polaryzacją wirtualnego fotonu. Metoda Rosenblutha to pierwsza poznana technika pozwalająca na otrzymanie wartości funkcji G_E i G_M dla protonu. Wymaga ona pomiarów przekroju czynnego rozpraszania elektron-proton podczas wielu eksperymentów, dla różnych parametrów θ i Q^2 . Zmianę tych parametrów można otrzymać poprzez korygowanie energii wiązki oraz kąta rozpraszania elektronu w tak dużym zakresie, jak to jest wykonalne eksperymentalnie. Metoda separacji Rosenblutha pozwala zapisać zredukowany przekrój czynny jako kombinację liniową funkcji postaci G_{Ep} oraz G_{Mp} [7]

$$\begin{aligned} \sigma_R(\epsilon, Q^2) &\equiv \epsilon(1 + \tau) \left(\frac{d\sigma}{d\Omega} \right) / \left(\frac{d\sigma}{d\Omega} \right)_M \\ &= \tau G_{Mp}^2(Q^2) + \epsilon G_{Ep}^2(Q^2). \end{aligned} \quad (3.5)$$

Kolejna metody oparte są na pomiarach transferu polaryzacji odbitego podczas rozpraszania protonu lub na pomiarach asymetrii rozpraszania, pozwalają one określić wzajemny stosunek elektrycznego oraz magnetycznego czynnika postaci

$$\mathcal{R}(Q^2) \equiv \mu_p \frac{G_{Ep}(Q^2)}{G_{Mp}(Q^2)},$$

gdzie $\mu_p = 2,793$ to moment magnetyczny protonu. W przybliżeniu jedno-fotonowym, otrzymujemy tylko dwa niezerowe składniki wektora polaryzacji,

poprzeczny P_t oraz podłużny P_l . Stosunek czynników postaci możemy otrzymać bezpośrednio ze stosunku składowych polaryzacji, otrzymujemy [7]:

$$\mathcal{R}(Q^2) \equiv -\mu_p \frac{P_t}{P_l} \frac{E + E'}{2M} \operatorname{tg}^2\left(\frac{\theta}{2}\right)$$

gdzie P_l i P_t to podłużny i poprzeczny składnik wektora polaryzacji odrzuconego protonu. E oraz E' to początkowa i końcowa energia elektronu, θ to kąt rozpraszania elektronu i M to masa protonu. Współczynnik $\mathcal{R}(Q^2)$ może także zostać wyznaczony na podstawie pomiaru asymetrii podczas sprężystego rozpraszania elektron-proton [1, 7]

$$\frac{\sigma_+ - \sigma_-}{\sigma_+ + \sigma_-} = -2\mu_p \sqrt{\tau(1+\tau)} \operatorname{tg}\left(\frac{\theta}{2}\right) \frac{\mathcal{R} \sin \theta^* \cos \phi^* + \mu_p \sqrt{\tau[1 + (1+\tau) \operatorname{tg}^2(\frac{\theta}{2})]} \cos \theta^*}{\mathcal{R}^2 + \mu_p \tau / \epsilon},$$

gdzie σ_+ i σ_- to przekroje czynne dla dodatniej i ujemnej skrętności, θ^* i ϕ^* to kąty polarny i azymutalny polaryzacji protonu względem wektora przekazu pędu \vec{q} i płaszczyzny rozpraszania. Podczas analizy wykorzystano stosunki funkcji postaci opublikowane w pracach: 25-36 z pracy [1] (uzupełnić)

3.2 Keras

Modele sieci neuronowych opisywane w tej pracy zostały zaprogramowane przy użyciu biblioteki Keras. Keras jest interfejsem API wysokiego poziomu służącym do tworzenia i szkolenia modeli głębokiego uczenia. Początkowo Keras został opracowany dla naukowców, którzy mogli dzięki niemu dokonywać szybkich eksperymentów i symulacji. Dzięki temu, że jest rozpowszechniany pod licencją MIT, co oznacza, że może być za darmo wykorzystywany w projektach komercyjnych, zdobył dużą popularność. Dziś ma on kilka set tysięcy użytkowników, od nauczycieli akademickich po inżynierów oprogramowania pracujących zarówno w start-upach jak i dużych firmach, i hobbyistów. Jego zalety są wykorzystywane między innymi w wiodących ośrodkach naukowych takich jak Europejska Organizacja Badań Jądrowych CERN i setkach firm, z których największe to Google, Netflix, Uber, Yelp, Opera Software [4]. Kaggle to platforma internetowa, która organizuje konkursy na najlepsze modele służące do przewidywania i opisywania zbiorów danych przesyłanych przez firmy i użytkowników. Jednym z najpopularniejszych narzędzi wykorzystywanych przez analityków jest Keras, wiele z konkursów zostało wygranych przez modele zbudowane przy użyciu wspomnianego interfejsu API. Do największych zalet Keras należą:

- posiada przyjazny użytkownikowi interfejs, który ułatwia szybkie prototypowanie modeli sieci neuronowych
- prosty i spójny interfejs zoptymalizowany pod kątem typowych przypadków użycia
- zapewnia przejrzyste informacje zwrotne dotyczące błędów użytkownika
- obsługuje dowolne architektury sieciowe: modele z wieloma wejściami lub wieloma wyjściami
- posiada wbudowane wsparcie dla splotowych sieci neuronowych oraz rekurencyjnych sieci neuronowych
- pozwala na bezproblemowe działanie tego samego kodu na CPU oraz GPU

Keras jest biblioteką, o której można powiedzieć, że zapewnia cegły służące do zbudowania modelu głębokiego uczenia natomiast w minimalnym stopniu pozwala użytkownikom na ingerencję w ich strukturę. W zamian wykorzystuje wyspecjalizowaną i dobrze zoptymalizowaną bibliotekę wyspecjalizowaną w operacjach na tensorach. Szczególnie szybko wykonują się obliczenia numeryczne typowe dla algorytmów uczenia maszynowego takich jak mnożenie macierzy i obliczanie gradientu. Można wybierać wśród trzech istniejących implementacji, każda z nich ma otwarte źródło. Pierwsza z nich wykorzystuje Tensorflow opracowany i rozwijany przez Google'a, druga korzysta z Theano opracowanego i rozwijanego przez LISA Lab w Uniwersytecie Montrealskim, ostatnia i najmniej popularna wykorzystuje CNTK opracowane i rozwijane przez Microsoft. W przyszłości prawdopodobnie pojawi się więcej możliwości wyboru, między innymi niedawno powstały, zdobywający coraz większą popularność projekt Torch finansowany przez Facebooka. Obecnie najczęściej wykorzystywany jest TensorFlow, został on także wykorzystany w tej pracy. Poniżej zaprezentuję jak proste jest zbudowanie i wytrenowanie bardzo podstawowego przykładu sieci neuronowej przy użyciu biblioteki Keras. Cały proces wymaga wykonania kilku kroków:

1. Zdefiniuj swoje dane treningowe: dane wejściowe i dane wyjściowe
2. Zdefiniuj warstwy swojej sieci neuronowej, które przekształcają dane wyjściowe w wyjście
3. Skonfiguruj proces uczenia poprzez wybranie funkcji straty, algorytmu szukającego minimum funkcji straty
4. Przeprowadź odpowiednią do wytrenowania sieci ilość iteracji

Zdefiniowana poniżej sieć składa się z dwóch warstw ukrytych o odpowiednio 10 i 5 neuronach ukrytych. Funkcją aktywacji w pierwszej warstwie jest sigmoida, dane wejściowe zawierają dwie cechy, które posłużą do zbudowania modelu, druga warstwa wykorzystuje tangens hiperboliczny jako funkcję aktywacji. Model podczas nauki minimalizuje błąd średniokwadratowy, wykorzystuje do tego algorytm rmsprop, trenowanie modelu skończy się po 100 pełnych iteracjach zbioru danych.

```
#Zaimportuj wymagane pliki
from keras import models
from keras import layers

#Zainicjalizuj model
model = models.Sequential()

#Dodaj pierwszą warstwę
model.add(layers.Dense(units = 10, activation =
    'sigmoid', input_shape = 2))

#Dodaj drugą warstwę
model.add(layers.Dense(units = 5, activation = 'tanh'))

#Dodaj warstwę wyjściową
model.add(layers.Dense(units = 1))

#Skompiluj model
model.compile(optimizer = 'rmsprop', loss='mse')

#Trenuj model
model.fit(inputs = X, outputs = Y, epochs = 100)
```

3.3 Generowanie sztucznych danych

Każdy punkt pomiarowy oprócz zmiennej objaśnianej zawiera przypisaną do niej niepewność pomiaru. Bazując na idei zaproponowanej w [3], wykorzystując niepewność pomiaru możemy wygenerować następne zestawy danych, które będą zawierały wartości przekrojów czynnych z zakresów, w których były możliwe do zmierzenia dla zadanych wartości Q^2 i ε . Następnie każda z replik posłuży do treningu osobnej sieci neuronowej, co pozwoli otrzymać

rozkład funkcji $\sigma^{(net)}(Q^2, \varepsilon)$.

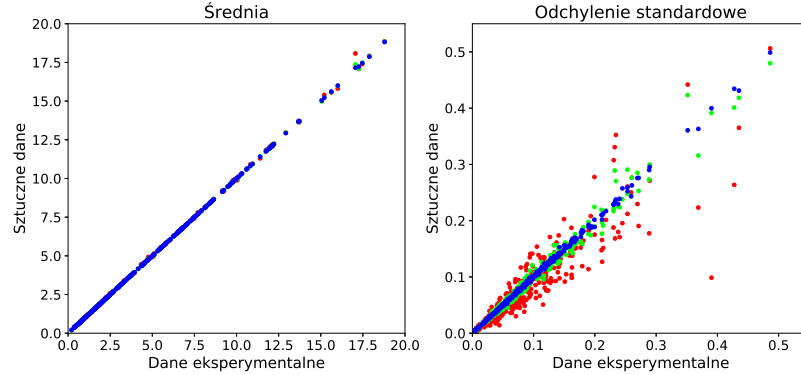
$$\sigma_{a,i}^{(art)(k)} = \sigma_i^{(exp)} + \mathcal{N}\left(0, \Delta\sigma_i^{(exp)}\right)^{(k)} \quad (3.6)$$

Ważne jest aby wybrać optymalną wartość liczby replik N_{rep} tak aby rozkład wygenerowanych danych zawierał charakterystyki zgodne z danymi eksperymentalnymi. Aby dokonać wyboru odpowiedniej wartości N_{rep} porównano wartości średnie oraz odchylenie standardowe próbek po wygenerowaniu sztucznych danych. Rysunek 3.2 przedstawia dwa wykresy punktowe powyższych wartości dla 10, 100 oraz 1000 replik. Szczególny wpływ ilości wygenerowanych danych widoczny jest w części przedstawiającej porównanie odchyleń standardowych. Większa liczba klonów powoduje, że charakterystyki rozkładów wygenerowanych oraz eksperymentalnych danych są bardziej zgodne, co na wykresach prezentuje się jako ułożenie punktów wzdłuż prostej $y = x$. Wartość średnia oraz odchylenie standardowe sztucznych danych zostały zdefiniowane w równaniach 3.7a oraz 3.7b.

$$\langle \sigma_{a,i}^{(art)} \rangle_{rep} = \frac{1}{N_{rep}} \sum_{k=1}^{N_{rep}} \sigma_{a,i}^{(art)(k)} , \quad (3.7a)$$

$$\langle \Delta\sigma_{a,i}^{(art)} \rangle_{rep} = \sqrt{\langle \sigma_{a,i}^{(art)2} \rangle_{rep} - \langle \sigma_{a,i}^{(art)} \rangle_{rep}^2} . \quad (3.7b)$$

Aby wskazać jak bardzo wygenerowane dane różnią się od danych ekspe-



Rysunek 3.2: $\langle \sigma_{a,i}^{(art)} \rangle$ vs. $\sigma_i^{(exp)}$ po lewej oraz $\langle \Delta\sigma_{a,i}^{(art)} \rangle$ vs. $\Delta\sigma_i^{(exp)}$ po prawej dla $N_{rep} = 10$ (czerwony), 100 (zielony), 1000 (niebieski).

rymentalnych zdefiniowano średnią wariancję oraz średni błąd względny dla

Tabela 3.1: Porównanie pomiędzy danymi eksperymentalnymi i danymi sztucznie wygenerowanymi

N_{rep}	10	100	1000
$\left\langle V \left[\left\langle \sigma_{a,i}^{(art)} \right\rangle_{rep} \right] \right\rangle_{dat}$	3.6×10^{-3}	4.3×10^{-4}	4.0×10^{-5}
$\left\langle PE \left[\left\langle \sigma_{a,i}^{(art)} \right\rangle_{rep} \right] \right\rangle_{dat}$	0.60%	0.20%	0.06%
$\left\langle V \left[\left\langle \Delta \sigma_{a,i}^{(art)} \right\rangle_{rep} \right] \right\rangle_{dat}$	1.5×10^{-3}	9.4×10^{-5}	1.4×10^{-5}
$\left\langle PE \left[\left\langle \Delta \sigma_{a,i}^{(art)} \right\rangle_{rep} \right] \right\rangle_{dat}$	19.4%	5.8%	1.7%

wszystkich punktów pomiarowych ($N_{dat} = 426$):

$$\left\langle V \left[\left\langle \sigma_{a,i}^{(art)} \right\rangle_{rep} \right] \right\rangle_{dat} = \frac{1}{N_{dat}} \sum_{i=1}^{N_{dat}} \left(\left\langle \sigma_{a,i}^{(art)} \right\rangle_{rep} - \sigma_i^{(exp)} \right)^2, \quad (3.8a)$$

$$\left\langle PE \left[\left\langle \sigma_{a,i}^{(art)} \right\rangle_{rep} \right] \right\rangle_{dat} = \frac{1}{N_{dat}} \sum_{i=1}^{N_{dat}} \left| \frac{\left\langle \sigma_{a,i}^{(art)} \right\rangle_{rep} - \sigma_i^{(exp)}}{\sigma_i^{(exp)}} \right|. \quad (3.8b)$$

Analogicznie możemy zdefiniować $\left\langle V \left[\left\langle \Delta \sigma_{a,i}^{(art)} \right\rangle_{rep} \right] \right\rangle_{dat}$ oraz $\left\langle PE \left[\left\langle \Delta \sigma_{a,i}^{(art)} \right\rangle_{rep} \right] \right\rangle_{dat}$.

Tabela 3.1 przedstawia różnice między zbiorami danych dla 10, 100 oraz 1000 replik danych. Wariancja wartości średniej zachowuje się zgodnie z przewidywaniami wynikającymi z teorii Monte Carlo i jest proporcjonalna do $1/N_{rep}$. Podobnie jest w przypadku wariancji odchylenia standardowego, które powinno maleć wraz ze wzrostem N_{rep} proporcjonalnie do $1/\sqrt{N_{rep}}$ [3]. Aby osiągnąć ponad 99% zgodność w wartości średniej oraz około 99% zgodność w niepewności pomiarowej należy wygenerować około 1000 replik danych. Ponadto, każdy z 14 niezależnych zbiorów danych ma określoną procentową niepewność systematyczną $\Delta\eta$, która powinna zostać uwzględniona podczas następnego etapu generowania replik danych. Dla każdego z 14 zbiorów losowana jest jedna wartość $\mathcal{N}(0, \Delta\eta)$ i ostatecznie generowane punkty przyjmują postać

$$\begin{aligned}
\sigma_i^{(art)(k)} &= \sigma_{a,i}^{(art)(k)} \times \left(1 + \mathcal{N}(0, \Delta\eta)^{(k)}\right) \\
&= \left(\sigma_i^{(exp)} + \mathcal{N}\left(0, \Delta\sigma_i^{(exp)}\right)^{(k)}\right) \times \left(1 + \mathcal{N}(0, \Delta\eta)^{(k)}\right). \quad (3.9)
\end{aligned}$$

Przedstawiona w dalszej części pracy analiza opiera się na aż trzech różnych zbiorach danych. Oprócz powyższego zbioru dysponujemy zależnością przekroju czynnego od Q^2 i ϵ poprawioną o poprawkę dwufotonową oraz stosunkiem funkcji postaci w zależności od Q^2 . Do tych zbiorów danych zostały zastosowane analogiczne kroki służące wygenerowaniu sztucznych danych.

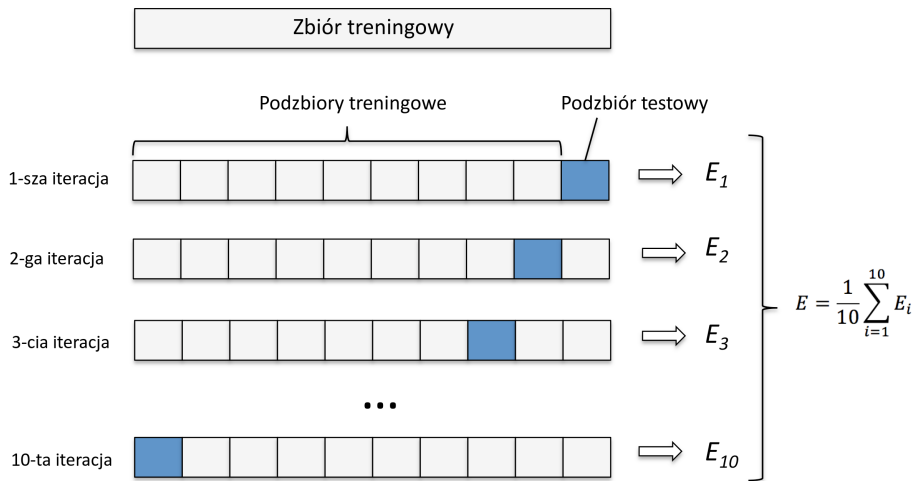
3.4 Walidacja krzyżowa

Algorytm wykorzystywany podczas nauki modelu ma za zadanie znalezienie takich parametrów, które sprawiają, że model odwzorowuje dane wykorzystane do nauki w sposób jak najlepszy z możliwych. Jeśli do walidacji modelu wykorzystamy inną, niezależną próbkę danych pochodzącą z tego samego zbioru co podzbiór uczący, zazwyczaj okaże się, że model nie działa aż tak dobrze jak przy użyciu zbioru uczącego. Rozmiar tej różnicy zwiększa się, szczególnie wtedy gdy wielkość zbioru treningowego jest niewielka, lub gdy liczba parametrów modelu jest bardzo duża. Walidacja krzyżowa to metoda statystyczna, która ma za zadanie zminimalizować tę różnicę przez co pomaga ocenić i zwiększyć trafność przewidywań modelu predykcyjnego.

W najprostszym przykładzie walidacji krzyżowej zbiór danych dzieli się na dwa podzbiory: uczący i walidacyjny. Podczas gdy zbiór uczący służy do nauki modelu, zbiór walidacyjny wykorzystuje się aby mierzyć błąd modelu na nieznanym zbiorze danych.

W algorytmie k -krotnej walidacji krzyżowej zbiór danych jest losowo dzielony na k równych wielkością podzbiorów. Jeden z k podzbiorów jest przeznaczany na zbiór walidacyjny, pozostałe $k - 1$ podzbiorów służą jako dane treningowe. Powyżej opisana procedura jest powtarzana k razy, a każdy k podzbiorów dokładnie raz zostaje wykorzystany jako zbiór testowy. Następnie k wyników modelu jest uśrednianych dając w rezultacie jeden wynik. Rysunek 3.3 przedstawia sposób działania 10-krotnej walidacji krzyżowej.

Cytując [6]: "(...) istnieje pewien kompromis między obciążeniem a wariancją, związany z wyborem parametru k w k -krotnej walidacji krzyżowej. Zazwyczaj stosuje się wartości z przedziału od 5 do 10, ponieważ pokazano empirycznie, że w takim wypadku otrzymujemy przewidywania, które nie cierpią nadmiernie ani z powodu dużego obciążenia ani dużej wariancji."



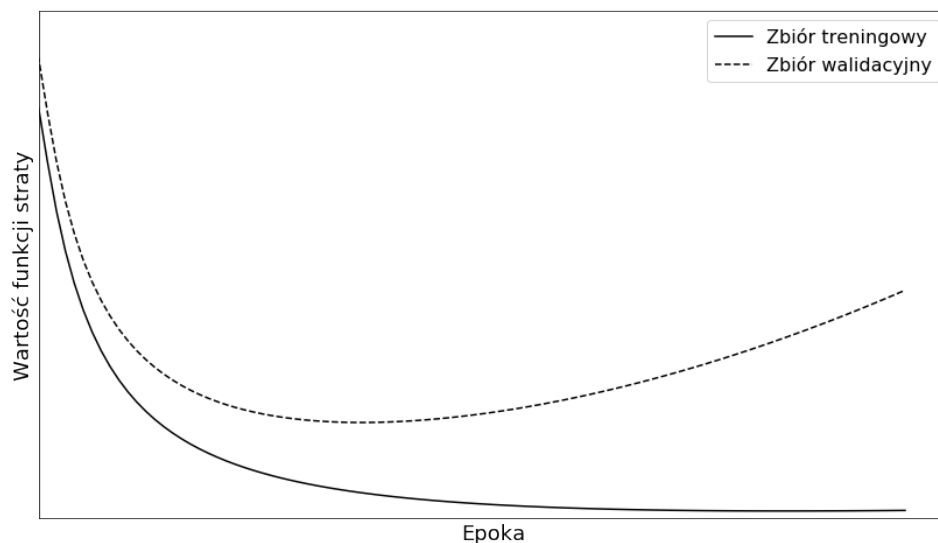
Rysunek 3.3: Przykład 10-krotnej walidacji krzyżowej, kolor niebieski oznacza podzbiór testowy, pozostała część zbioru to podzbiór treningowy. Wynik końcowy jest średnią z wyników wszystkich iteracji. Rysunek przetłumaczony na język polski z [8].

Podczas treningu modelu wybierano więc takie k z zakresu $[5, 10]$, dla którego liczba próbek w zbiorze danych jest całkowicie podzielna przez k co zapewnia równy rozmiar wszystkich zbiorów treningowych i walidacyjnych.

3.5 Wczesne zatrzymanie

Algorytmy uczenia maszynowego dopasowują parametry modelu na podstawie danych treningowych o skończonym rozmiarze. Podczas procesu szkolenia model jest oceniany na podstawie tego, jak dobrze przewiduje obserwacje zawarte w tym zbiorze. Jednak celem uczenia maszynowego jest stworzenie modelu, który ma zdolność do przewidywania uprzednio niewidzianych obserwacji. Nadmierne dopasowanie to zjawisko pojawiające się wtedy gdy model za bardzo dopasowuje się do danych w zbiorze uczącym co powoduje zmniejszenie wartości błędu na tym zbiorze lecz równocześnie jest przyczyną wzrostu błędu na zbiorze testowym. Nadmierne dopasowanie modelu to problem, który może się pojawiać gdy model zawiera więcej parametrów niż wymagałaby tego natura modelowanego zjawiska. Sieć neuronowa to struktura skłonna do przeuczenia. Podczas gdy obserwowany błąd obliczany w

oparciu o dane treningowe spada, w pewnym momencie wartość błędu dla zbioru walidacyjnego zaczyna wzrastać. Rysunek 3.4 przedstawia często zamieszczane w literaturze, wyidealizowane krzywe zmiany wartości funkcji straty w czasie, dla zbiorów treningowego i walidacyjnego. Najlepszy model predykcyjny miałby parametry, które odpowiadają momentowi globalnego minimum dla zbioru walidacyjnego.



Rysunek 3.4: Wyidealizowane przykłady krzywych przedstawiających zmianę wartości funkcji straty na zbiorach treningowym i walidacyjnym, podczas nauki modelu

W dziedzinie uczenia maszynowego, metoda wczesnego zatrzymania to forma regularyzacji, która pozwala uniknąć problemu przeuczenia, zatrzymując naukę modelu gdy wartość funkcji straty na zbiorze walidacyjnym zaczyna wzrastać. Rzeczywisty przebieg wartości funkcji straty ma wiele lokalnych minimów, dlatego na podstawie obserwacji krzywych uczenia dokonano wyboru kryteriów zatrzymania nauki modelu. Niech $\Theta_{wa}(t)$ to wartość funkcji straty na zbiorze walidacyjnym po t epokach, $\Theta_{min}(t)$ to dotychczasowe minimum funkcji straty na zbiorze walidacyjnym po t epokach, definiowane jako:

$$\Theta_{min}(t) \equiv \min_{t' < t} \Theta_{wa}(t')$$

Niech $\Theta_{sr}(t)$ będzie średnią wartością funkcji straty dla zbioru walidacyjnego z ostatnich 10 epok.

$$\Theta_{sr}(t) \equiv \frac{1}{10} \sum_{i=0}^{10} \Theta_{wa}(t-i)$$

Oraz zdefiniujemy pomocniczy parametr $GL(t)$

$$GL(t) \equiv \frac{\Theta_{sr}(t)}{\Theta_{min}} - 1$$

Podczas nauki przedstawionych modeli statystycznych oprócz wykorzystania metody wczesnego zatrzymania została ustalona minimalna wymagana liczba epok. Z powodu startu algorytmu uczącego z losowymi parametrami, szczególnie w pierwszych iteracjach nauki funkcja błędu może być poddana dużym fluktuacjom. Po przekroczeniu minimalnej liczby epok do wczesnego zatrzymania wystarczyło spełnienie jednego z dwóch obowiązujących warunków:

- $\Theta_{min}(t) = \Theta_{min}(t+200)$ dla wszystkich $t \in [t, t+200]$, brak zmniejszenia minimalnej wartości funkcji straty dla zbioru walidacyjnego przez 200 epok
- $GL(t) > 2$, względny wzrost średniej wartości funkcji straty przez ostatnie 10 epok względem osiągniętego minimum jest większy niż 200%

Po skończeniu nauki, wybierany jest model, który ma najmniejszą wartość funkcji straty na zbiorze testowym.

3.6 Ilość neuronów

Architektura sieci neuronowej, tzn. ilość warstw ukrytych oraz ilość neuronów w warstwach ukrytych jest zdeterminowana przez wymiar danych wejściowych, rodzaj rozwiązywanego problemu (klasyfikacja czy regresja) oraz relację między zmiennymi objaśniającymi i zmienną objaśnianą.

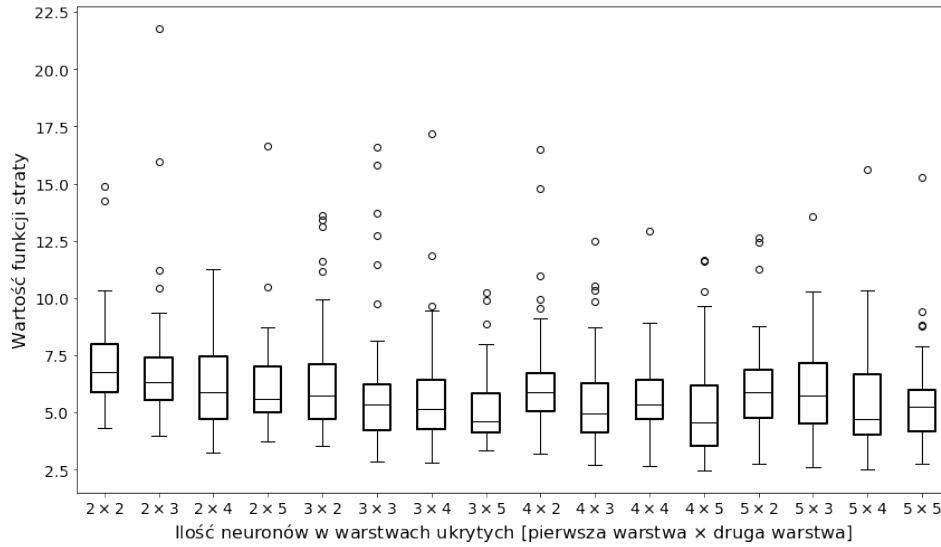
Uogólniony model liniowy przydatny w szerokim zakresie zastosowań, nie potrzebuje żadnej warstwy ukrytej. Bywa szczególnie przydatny gdy zbiór zawiera mało danych lub są one obciążone dużą niedokładnością. Nawet w przypadku gdy relacja między zmiennymi jest lekko nieliniowa, użycie prostego modelu liniowego może skutkować lepszym uogólnieniem problemu niż skomplikowany model będący wrażliwy na każdy szum znajdujący się w danych. Zgodnie z uniwersalnym twierdzeniem aproksymacyjnym jedna warstwa ukryta z wystarczająco dużą liczbą neuronów wystarcza aby z dowolną

Tabela 3.2: Liczba parametrów sieci neuronowej z dwoma warstwami ukrytymi w zależności od liczby neuronów w warstwach

I warstwa \ II warstwa	2	3	4	5
	2	3	4	5
2	14	18	22	26
3	19	24	29	34
4	24	30	36	42
5	29	36	43	50

dokładnością dowolną ciągłą funkcję [cybenko]. Jeśli zmienna objaśniająca jest jednowymiarowa, wydaje się, że nie odniesiemy żadnej korzyści z konstruowania sieci neuronowej o więcej niż jednej warstwie ukrytej. Sprawy komplikują się jednak gdy zmienna wejściowa jest dwu lub więcej wymiarowa. Dwuwarstwowa sieć neuronowa zachowuje właściwości jednowarstwowej sieci neuronowej oraz osiąga zdolność nauki każdego problemu klasyfikacyjnego [1995 Bishop 123], ponadto wielowarstwowa sieć neuronowa z dwoma warstwami może skutkować dokładniejszymi wynikami wykorzystując mniejszą ilość parametrów niż jednowarstwowa sieć [Chester (1990)]. Na tej podstawie, do rozwiązania problemu regresji gdzie wejściem jest para liczb (ε, Q^2) postanowiłem wybrać sieć neuronowa z dwoma warstwami ukrytymi.

Aby znaleźć odpowiednią liczbę neuronów w dwóch warstwach ukrytych, stworzyłem siatkę $[2, 3, 4, 5] \times [2, 3, 4, 5]$ neuronów i sprawdziłem, która konfiguracja daje najmniejszy błąd zbioru walidacyjnego. Dane zostały podzielone na zbiór treningowy i testowy w stosunku 2:1. Dla każdej konfiguracji wytrenowano 50 sieci i sprawdzono jak wygląda statystyka błędów. Tabela 3.2 zawiera porównanie liczby parametrów sieci neuronowej w zależności od liczby neuronów w warstwach ukrytych. Do eksperymentów wybrano konfiguracje charakteryzujące się rozsądną w porównaniu do rozmiaru danych wejściowych liczbą parametrów. Rysunek 3.5 przedstawia rozkłady minimalnej wartości funkcji straty uzyskanej na danych walidacyjnych uzyskanej z 50 treningów sieci dla każdej konfiguracji ilości neuronów. Wykres pudełkowy to forma graficznej prezentacji rozkładu, która pozwala w łatwy sposób ukazać położenie, rozproszenie oraz kształt empirycznego rozkładu badanej cechy statystycznej. Konfiguracja 3×5 charakteryzuje się najniższą medianą wartości funkcji straty oraz małą liczbą wartości odstających. Ta obserwacja pozwoliła zdecydować, że liczby neuronów będą wynosiły 3 i 5 w odpowiednio pierwszej i drugiej warstwie ukrytej, co za tym idzie sieć będzie miała 36 parametrów.



Rysunek 3.5: Wykresy pudełkowe przedstawiające rozkład wartości funkcji straty w zależności od ilości neuronów w pierwszej i drugiej warstwie ukrytej

3.7 Algorytm uczący

Bardzo istotnym elementem tworzonego modelu jest wybór algorytmu poszukującego minimum funkcji straty oznaczonej na potrzeby tego paragrafu jako $J(\theta)$. Na podstawie jego wyników aktualizowane będą parametry tworzonej sieci neuronowej. Bardzo pomocną koncepcją pozwalającą zrozumieć istotę trudności problemu jest powierzchnia błędu. "Każda z N wag i wartości progowych sieci (tzn. wszystkie wolne parametry modelu) traktowana jest jako jeden z wymiarów przestrzeni. W ten sposób każdy stan sieci, wyznaczony przez aktualne wartości jej N parametrów może być traktowany jako punkt na N -wymiarowej hiperpłaszczyźnie. $N+1$ wymiarem (zaznaczanym jako wysokość ponad wspomnianą wyżej hiperpowierzchnią) jest błąd, jaki popełnia sieć. Dla każdego możliwego zestawu wag i progów może więc zostać narysowany punkt w przestrzeni $N+1$ -wymiarowej, w taki sposób, że stan sieci wynikający z aktualnego zestawu jej parametrów lokuje ten punkt na wspomnianej wyżej N -wymiarowej hiperpłaszczyźnie zaś wartość błędu, jaki popełnia sieć dla tych właśnie wartości parametrów stanowi wysokość umieszczenia punktu ponad tą płaszczyznę. Gdybyśmy opisaną procedurę powtórzyli dla wszystkich możliwych wartości kombinacji wag i progów sieci, wówczas otrzymalibyśmy "chmurę" punktów rozciągających się

ponad wszystkimi punktami N -wymiarowej hiperpłaszczyzny parametrów sieci, tworzącą właśnie rozważaną powierzchnię błędu. Celem uczenia sieci jest znalezienie na tej wielowymiarowej powierzchni punktu o najmniejszej wysokości, czyli ustalenie takiego zestawu wag i progów, który odpowiada najmniejszej wartości błędu. Przy stosowaniu modeli liniowych z funkcją błędu opartą na sumie kwadratów powierzchnia błędu ma kształt paraboloidy (funkcji kwadratowej), ma więc kształt kielicha o gładkich powierzchniach bocznych i o jednym wyraźnym minimum. Z tego powodu wyznaczenie w tym przypadku wartości minimalnej nie stwarza większych problemów." [11]

Jeżeli dysponujemy niewielkim zbiorem danych treningowych, do znalezienia optimum funkcji doskonale sprawdzają się metody quasi-Newtonowskie. Ich zaletą jest bardzo szybka zbieżność, niestety obliczenie hesjanu funkcji wielu zmiennych charakteryzuje się dużą złożonością pamięciową $O(n^2)$ i jeszcze większą złożonością obliczeniową $O(n^3)$. Z tego powodu możliwość ich zastosowania ogranicza się do niewielu przypadków. Najbardziej znane algorytmy quasi-Newtonowskie to m.in.: *LM-BFGS*, *Levenberg-Marquardt*. Dysponując dużym zbiorem danych należy wybrać inny algorytm. Po za losowym poszukiwaniem parametrów, najłatwiejszym z nich i bardzo intuicyjnym jest metoda gradientu prostego (*gradient descent*). Parametry θ aktualizowane są w następujący sposób:

$$\theta^{k+1} = \theta^k - \alpha \nabla J(\theta^k) \quad (3.10)$$

gdzie α to wybrany odpowiednio parametr szybkości uczenia (*learning rate*) odpowiedzialny za stopień zmiany parametrów w kolejnych iteracjach. Jeśli θ^0 znajduje się odpowiednio blisko minimum funkcji, i parametr α jest wystarczająco niewielki, algorytm osiąga liniową zbieżność [2]. W ogólności metoda gradientu prostego gwarantuje zbieżność do globalnego minimum w przypadku funkcji błędu o wypukłej powierzchni i do lokalnego minimum dla funkcji błędu o powierzchni nie wypukłej. Algorytm jednak jest bardzo wolny, co jest jego największą słabością. Ze względu na częstotliwość aktualizacji wag, metodę gradientu prostego możemy podzielić na *batch gradient descent* oraz *stochastic gradient descent*. W pierwszym przypadku wagi są dostosowywane po przetworzeniu pełnego zbioru danych, w metodzie stochastycznej zbiór uczący dzielony jest na podzbiory a wagi aktualizowane są po przetworzeniu każdego z podzbiorów. Druga metoda jest szczególnie użyteczna dla dużych zbiorów danych. Spodziewamy się, że dla dobrze przygotowanych danych kierunek podążania wartości wag będzie podobny jeśli policzymy gradient zarówno dla 10% jak i dla 100% zbioru treningowego.

Wyobraźmy sobie, że poszukiwanie minimum powierzchni błędu to prze-

mierzenie przestrzeni pełnej dolin, pagórków, wąwozów. W kolejnych iteracjach przeskakujemy między tymi obszarami, w pewnym momencie może się zdarzyć, że gradient zaniknie lub będzie bardzo słaby a nasze poszukiwania zatrzymają się nie osiągając wystarczającego minimum. Idea pędu inspirowana zjawiskami fizycznymi to nadanie gradientowi krótkotrwałej pamięci. Posługując się kolejną analogią, popchnięta w dół piłka nabierając prędkości zwiększa swój pęd. To samo dzieje się z parametrami sieci, wartość pędu wzrasta dla wymiarów, których gradienty wskazują te same kierunki i zmniejsza modyfikacje wartości dla wymiarów, w których gradienty zmieniają kierunki. W rezultacie otrzymujemy szybszą zbieżność i mniejsze oscylacje.

$$v^{k+1} = \beta v^k + \nabla J(\theta^k) \quad (3.11)$$

$$\theta^{k+1} = \theta^k - \alpha v^{k+1} \quad (3.12)$$

Zmiana jest niewielka, gdy $\beta = 0$, otrzymujemy zwykłą metodę gradientu prostego, zazwyczaj jednak ustala się wartość parametru β , zwanego pędem na około 0.9 [10].

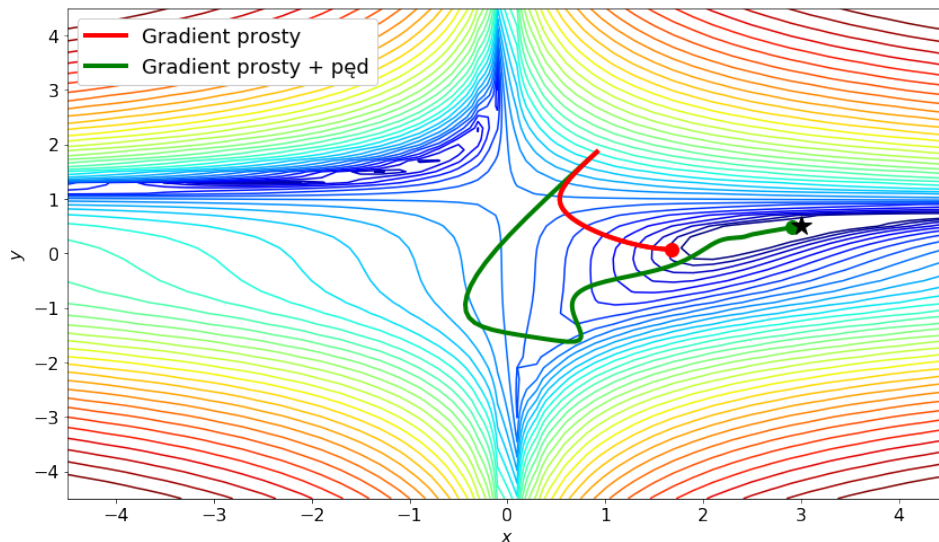
Porównanie efektywności przedstawionych wyżej algorytmów znajduje się na Rysunek 3.6, w zaprezentowanym przykładzie metoda gradientu prostego potrzebuje około 10 razy więcej iteracji od modyfikacji z pędem aby dotrzeć do minimum zaprezentowanej funkcji. Jest to przykład świadczący o tym jak duży wpływ na szybkość działania algorytmu wywiera ta niewielka modyfikacja.

Wykorzystany podczas treningu modelu algorytm korzysta jednak z jeszcze z jednej modyfikacji. Nie chcielibyśmy aby piłka spuszczone w dół ślepo podążała za zboczem widząc, że za niedługo mocno się ono podniesie. Przyspieszenie Nesterova (*NAG*) jest sposobem na uwzględnienie podczas obliczania gradientu przybliżonej przyszłej pozycji parametrów sieci. Algorytm opisują równania 3.13, 3.14 [9].

$$v^{k+1} = \beta v^k + \nabla J(\theta^k - \beta v^k) \quad (3.13)$$

$$\theta^{k+1} = \theta^k - \alpha v^{k+1} \quad (3.14)$$

Niezwyczajnie istotnym parametrem algorytmu jest α , jego niezmiennosc wraz z postępem iteracji powoduje bardzo niską efektywność algorytmu. Ze względu na metodę zmiany tego parametru, który może być indywidualnie ustalany dla każdej wagi powstało wiele szeroko wykorzystywanych algorytmów. Do najpopularniejszych należą między innymi *Adam*, *Nadam*, *Adagrad*, *Adadelta*, *AMSGrad*, *RMSprop*.



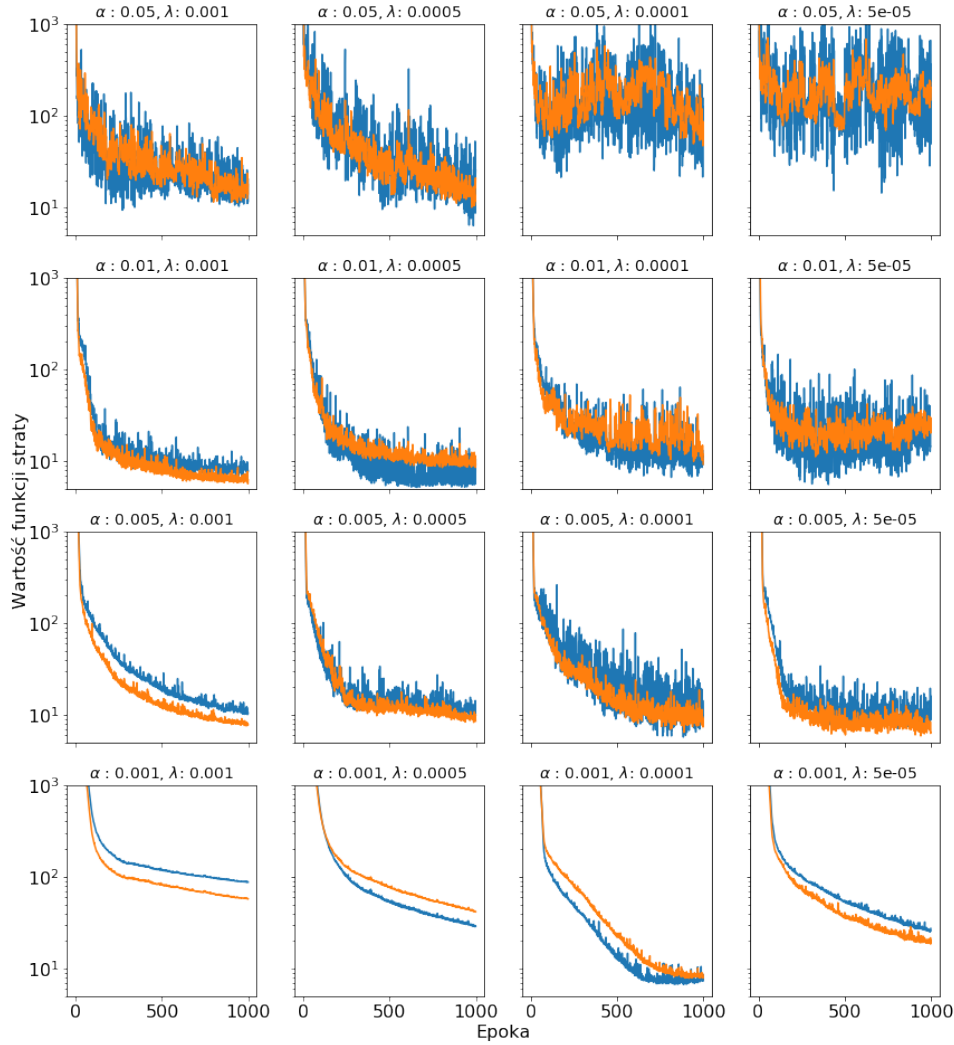
Rysunek 3.6: Funkcja $f(x, y) = (1.5 - x + xy)^2 + (2.25 - x + xy^2)^2 + (2.625 - x + xy^3)^2$, osiąga minimum równe 0, w punkcie $(3, 0.5)$ oznaczonym czarną gwiazdą. Grafika przedstawia porównanie działania metody gradientu prostego oraz jego modyfikacji poprzez dodanie pędu. Przyjmując, że punkt początkowy to $(2, 1)$, $\alpha = 0.001$ i $\beta = 0.9$, możemy prześledzić trajektorie algorytmów przez pierwsze 500 iteracji działania.

W swoim algorytmie postanowiłem dokonywać zmiany parametru α wraz ze wzrostem iteracji. Ponadto szybkość uczenia zależna jest od wybranego parametru λ decydującego o tym z jaką szybkością maleje.

$$\alpha(i) = \alpha_0 \times \frac{1}{1 + \lambda \times i} \quad (3.15)$$

Rysunek 3.7 przedstawia porównanie przykładowych krzywych zmian wartości funkcji straty w czasie dla różnych wartości α i λ . Na ich podstawie widać jak duży wpływ wnosi parametr α w proces nauki modelu. Zbyt duża szybkość uczenia powoduje bardzo duże oscylacje krzywej funkcji straty, za mała wartość α bardzo mocno spowalnia proces nauki. Pewien kompromis przynosi wybranie odpowiednio dużej początkowej wartości szybkości uczenia, co przynosi szybkie przejście algorytmu w obszar minimum i następnie zmniejszenie go do wartości potrafiącej efektywnie dalej poszukiwać optimum. Zadowalający przebieg mają krzywe o parametrach $\alpha = 0.005$, $\lambda = 0.001$, które przedstawiają porządkany, eksponencjalny kształt o niewiel-

kiej oscylacji. Na podstawie powyższej analizy to właśnie te hiperparametry zostały wykorzystane w modelu, dodatkowo parametr pędu β został ustalony na wartość 0.9



Rysunek 3.7: Porównanie przykładowych krzywych zmiany wartości funkcji straty w czasie dla zbiorów treningowego (kolor pomarańczowy) i walidacyjnego (kolor niebieski) ze względu na parametry α (*learning rate*) oraz λ (*decay*)

Rozdział 4

Wyniki analizy

4.1 Analiza nr 1

Celem pierwszej analizy jest modelowanie elektrycznego i magnetycznego czynnika postaci przy wykorzystaniu wyłącznie dane przekrojów czynnych rozpraszania elektron-proton w eksperymentach (lista referencji).

Dane wejściowe i funkcja straty

Na zbiór analizowanych danych składa się 24 niezależnych zbiorów danych z eksperymentów, w których dokonywano rozpraszania elektron-proton, razem daje to 426 punktów pomiarowych. Zestaw danych składa się z 4 kolumn, które zawierają kolejno zmienną objaśnianą σ - przekrój czynny, niepewność pomiaru zmiennej objaśnianej $\Delta\sigma$ oraz dwie zmienne objaśniające Q^2 - kwadrat przekazanego czteropędu i czynniki kinematyczny ε . Ponadto każdy z niezależnych zbiorów ma określoną niepewność systematyczną $\Delta\eta$. Dodatkowo do każdego ze zbiorów dodano sztuczny punkt pomiarowy, który korzysta z założenia, że $\sigma(Q^2 = 0, \varepsilon = 1) = 1$, niepewność pomiarowa punktu wynosi $\Delta\sigma = 0.01$, zwiększa to liczbę wszystkich punktów pomiarowych do 450. Funkcja straty to z definicji funkcja przyporządkowująca nieujemną wielkość kary poprzez porównanie zmiennej objaśnianej do wyliczonego estymatora. W przedstawionym modelu, wykorzystana została zmodyfikowana postać funkcji chi-kwadrat, która bierze pod uwagę zarówno niepewność pomiarową oraz systematyczną

$$\chi^2 = \frac{1}{n} \chi_\sigma^2 \quad (4.1)$$

$$\chi_\sigma^2 = \sum_{k=1}^{N_\sigma} \left[\sum_{i=1}^{n_k} \left(\frac{\eta_k \sigma_{ki}^{th} - \sigma_{ki}^{ex}}{\Delta \sigma_{ki}} \right)^2 + \left(\frac{\eta_k - 1}{\Delta \eta_k} \right)^2 \right], \quad (4.2)$$

gdzie N_σ to liczba zbiorów danych z niezależnych eksperymentów, n_k to liczba punktów w k -tym zbiorze danych, $n = \sum_{k=1}^{N_\sigma} n_k$ to liczba wszystkich punktów pomiarowych, η_k to parametr normalizacyjny dla k -tego zbioru danych, $\Delta \eta_k$ to błąd systematyczny. σ_{ki}^{ex} to wartość eksperymentalna przekroju czynnego i -tego pomiaru z k -tego zbioru danych, zmierzona dla określonego Q_{ki}^2 , ε_{ki} , $\Delta \sigma_{ki}^{ex}$ oznacza odpowiadającą niepewność pomiaru, σ_{ki}^{th} to przewidywanie modelu statystycznego. η_k , $k = 1, 2, \dots, N_\sigma$ to parametry normalizacyjne. Ich wartości są aktualizowane podczas każdej iteracji nauki modelu [5], powinny one spełniać warunek

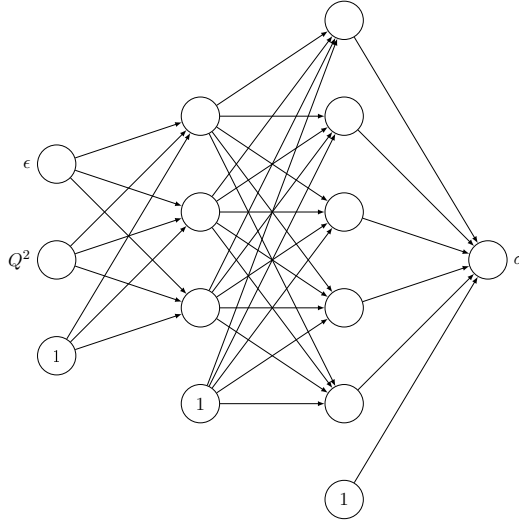
$$\frac{\partial \chi_\sigma^2}{\partial \eta_k} = 0, \quad k = 1, \dots, N_\sigma, \quad (4.3)$$

co można zapisać jako

$$\eta_k = \frac{\sum_{i=1}^{n_k} \frac{\sigma_{ki}^{th} \sigma_{ki}^{ex}}{(\Delta \sigma_{ki})^2} + \frac{1}{(\Delta \eta_k)^2}}{\sum_{i=1}^{n_k} \frac{(\sigma_{ki}^{th})^2}{(\Delta \sigma_{ki})^2} + \frac{1}{(\Delta \eta_k)^2}}. \quad (4.4)$$

Parametry i nauka sieci

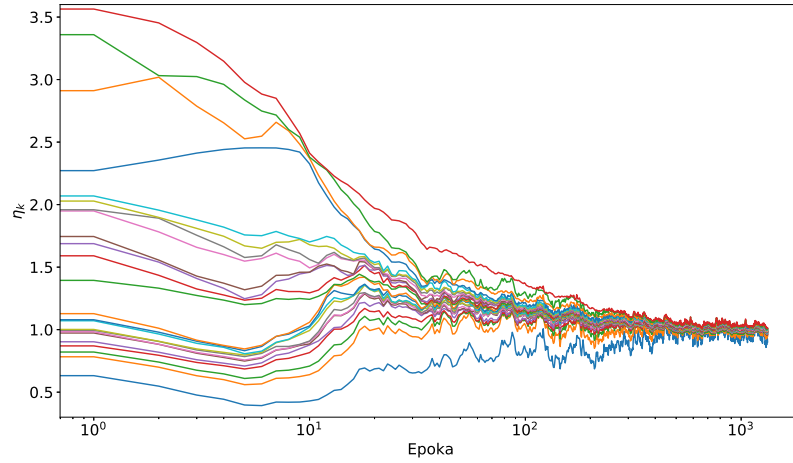
Sieć neuronowa o dwóch warstwach ukrytych i hipparametrach przedstawionych w Tabeli 4.1 daje w rezultacie wartość σ w zależności od zmiennych Q^2 oraz ϵ . Podczas treningu elementy ze zbioru parametrów η_k dążą do wartości bliskich 1. Rysunek 4.2 przedstawia ewolucję parametrów η_k wraz z nauką sieci dla każdego z 24 niezależnych zbiorów danych. Łącznie zostało wytrenowanych 2500 ($N_{rep} \times k$) modeli. Przykładowy przebieg wartości funkcji straty dla zbioru treningowego i walidacyjnego został przedstawiony na Rysunku 4.3, możemy zauważyć, że około 1100 epoki funkcja straty osiąga minimum, zgodnie z opisanymi wcześniej zasadami wczesnego zatrzymania model kończy naukę po następnych 200 epokach.



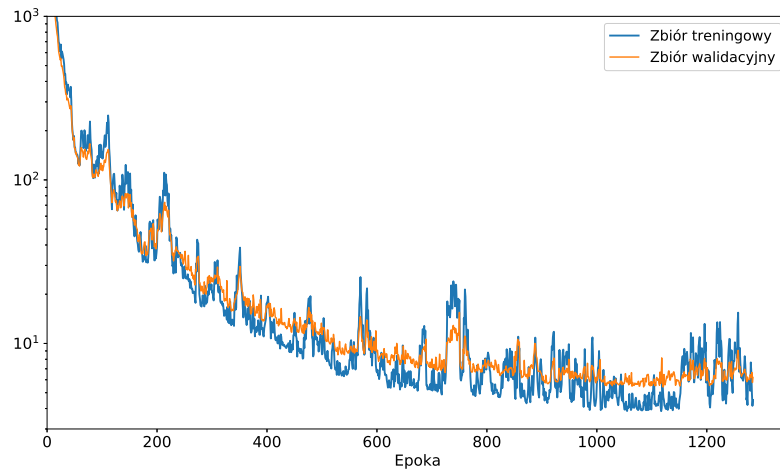
Rysunek 4.1: Schemat sieci neuronowej zastosowanej w pierwszej analizie, która składa się z: i) warstwy wejściowej z dwoma neuronami, ii) dwóch warstw ukrytych z odpowiednio trzema i pięcioma neuronami, iii) warstwy wyjściowej z jednym neuronem. Linie zakończone strzałką oznaczają wagę odpowiadającą każdej z par neuronów

Tabela 4.1: Hiperparametry modelu

Kategoria	Parametr	Wartość
Generowanie danych	N_{rep}	500
k -krotna walidacja krzyżowa	k	5
Algorytm uczący	α (<i>learning rate</i>)	0,003
	λ (<i>decay</i>)	0,0005
	β (<i>pęd</i>)	0,9
Sieć neuronowa	Liczba warstw	2
	Ilość neuronów	(3,5)



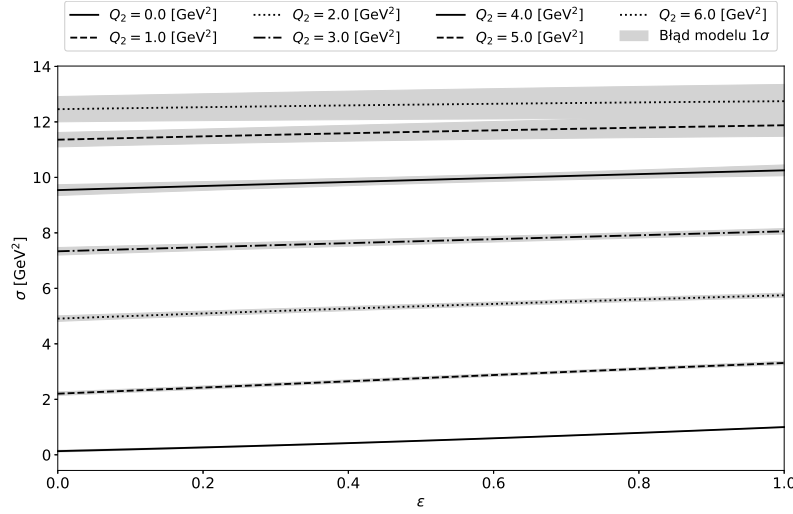
Rysunek 4.2: Ewolucja parametrów η_k podczas jednego z treningów modelu. Wraz ze wzrostem epok, wartości parametrów ustalane na podstawie równania 4.4 zbiegają do wartości bliskich 1.



Rysunek 4.3: Zmiana wartości funkcji straty podczas nauki modelu. Wartość funkcji obliczana na podstawie zbioru treningowego oznaczona jest kolorem niebieskim, dla zbioru walidacyjnego - pomarańczowym.

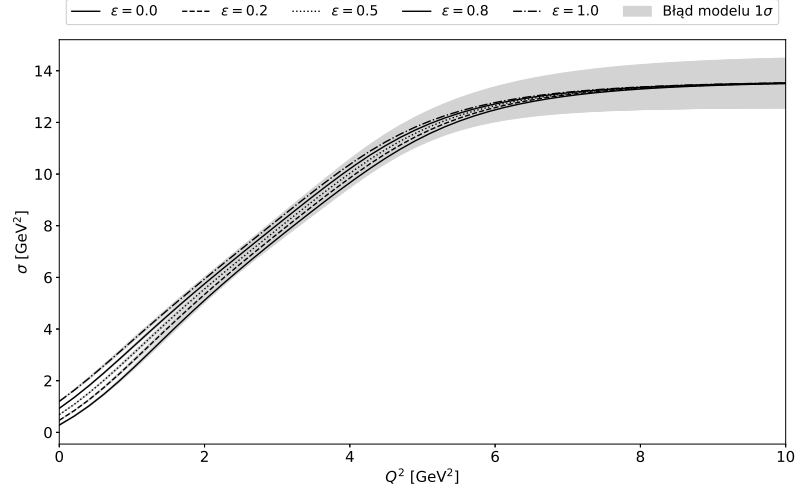
Wyniki

Wyniki wytrenowanych modeli tworzą statystykę, której najważniejszymi parametrami są średnia - często przedstawiana jako wynik analizy, oraz odchylenie standardowe, które przedstawia zakres bardzo prawdopodobnych wyników, na wykresie przedstawiany jako zacieniowany obszar. Rysunek 4.5 przedstawia zależność przekroju czynnego $\sigma(\epsilon)$, dla kilku ustalonych wartości Q^2 . Otrzymane funkcje mają przebieg liniowy o bardzo podobnym współczynniku kierunkowym a błąd modelu 1σ znacznie wzrasta wraz ze wzrostem Q^2 . Zależność przekroju czynnego $\sigma(Q^2)$ przy ustalonym para-



Rysunek 4.4: Zależność przekroju czynnego σ od czynnika kinematycznego ϵ przy ustalonym przekazie czteropędu Q^2 . Linia ciągła wyznacza średnią wartość po wszystkich wytrenowanych sieciach. Kolor szary wyznacza obszar niepewności 1σ .

metrze ϵ znajduje się na rysunku 4.5. Możemy zauważyć, że im niższa wartość ϵ tym mniejszy przekrój czynny dla $Q^2 = 0$, następnie krzywe mają bardzo podobny przebieg, niezależnie od ustalonego parametru ϵ zbiegają do tej samej maksymalnej wartości σ wraz ze wzrostem Q^2 . Ponadto, wraz ze wzrostem Q^2 rośnie niepewność otrzymanego wyniku, dla $Q^2 = 10$ GeV² obszar błędu modelu 1σ wynosi aż ± 1 GeV². Po zbadaniu podstawowych zależności estymowanej funkcji od kwadratu przekazu czteropędu oraz czynnika kinematycznego wyznaczono elektryczny i magnetyczny funkcji postaci. Wiemy, że



Rysunek 4.5: Zależność przekroju czynnego σ od przekazu czteropędu Q^2 przy ustalonym czynniku kinematycznym ϵ . Linia ciągła wyznacza średnią wartość po wszystkich wytrenowanych sieciach. Kolor szary wyznacza obszar niepewności 1σ .

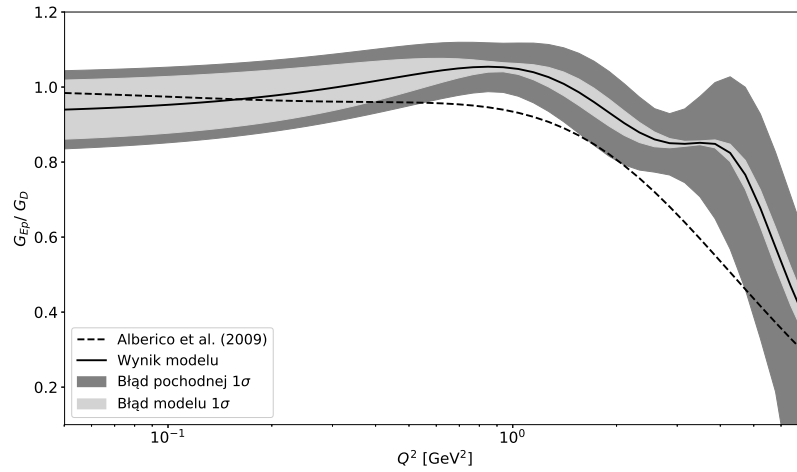
$$\sigma_R(\epsilon, Q^2) = \tau G_{Mp}^2(Q^2) + \epsilon G_{Ep}^2(Q^2), \quad (4.5)$$

obliczając pochodną po parametrze ϵ otrzymamy kwadrat elektrycznej funkcji postaci protonu, więc

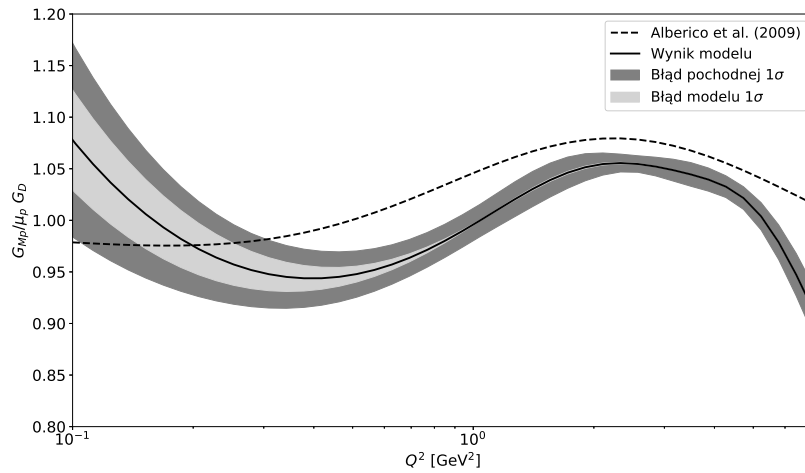
$$G_{Ep}(Q^2) = \sqrt{\frac{\partial \sigma_R(\epsilon, Q^2)}{\partial \epsilon}}. \quad (4.6)$$

Następnie magnetyczna funkcja postaci protonu wyraża się wzorem

$$G_{Mp}(Q^2) = \sqrt{\frac{\sigma_R(\epsilon, Q^2) - \epsilon \frac{\partial \sigma_R(\epsilon, Q^2)}{\partial \epsilon}}{\tau}} \quad (4.7)$$



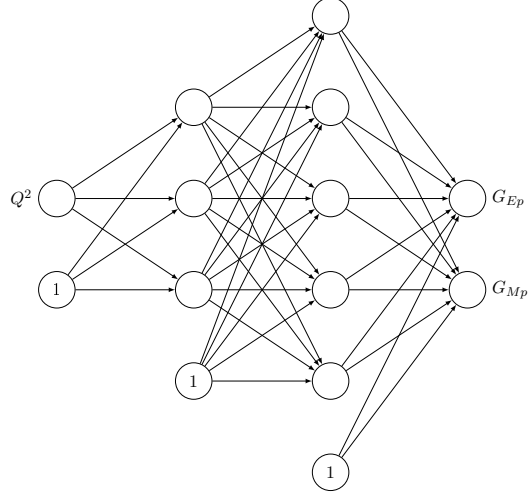
(a)



(b)

Rysunek 4.6: Elektryczna (a) i magnetyczna (b) funkcja postaci. Linia ciągła przedstawia średnią ze wszystkich modeli, zacieniowane regiony wyznaczają obszary 1σ powstałe z dwóch różnych przyczyn. Jasnoszary obszar opisany jako błąd modelu to odchylenie standardowe opisujące rozkład wyników wszystkich wytrenowanych modeli. Ciemnoszary obszar wyznacza odchylenie standardowe opisujące rozkład wyników powstałych na skutek obliczeń pochodnej dla różnych wartości ϵ z zakresu $[0, 1]$. Przerywana linia to wyniki przedstawione w [1].

Rysunek 4.7: Schemat sieci neuronowej zastosowanej w drugiej analizie, która składa się z: i) warstwy wejściowej z dwoma neuronami, ii) dwóch warstw ukrytych z odpowiednio trzema i pięcioma neuronami, iii) warstwy wyjściowej z dwoma neuronami. Linie zakończone strzałką oznaczają wagę odpowiadającą każdej z par neuronów



4.2 Analiza nr 2

Dane wejściowe i funkcja straty

Dane wejściowe w następnej analizie to 450 punktów pomiarowych analizowanych w pierwszej analizie powiększone o zbiór 68 pomiarów stosunku funkcji postaci G_{Ep}/G_{Mp} wraz z niepewnością pomiarową w zależności od kwadratu przekazu czteropędu Q^2 . Ponieważ znamy jeden z węzłów funkcji postaci, do zbioru dodany został 69 i 70 punkt: $\mathcal{R}(Q^2 = 0) = 1$, $\Delta\mathcal{R} = 0,01$. Razem otrzymujemy 520 pomiarów co sugeruje wybranie parametru k -krotnej walidacji jako $k = 5$. Wykorzystana podczas nauki drugiego modelu funkcja straty χ^2 (4.8) jest modyfikacją funkcji wykorzystanej w pierwszej analizie. Do użytej wcześniej funkcji błędu dodany został składnik (4.9) uwzględniający błąd estymacji stosunku G_{Ep}/G_{Mp} .

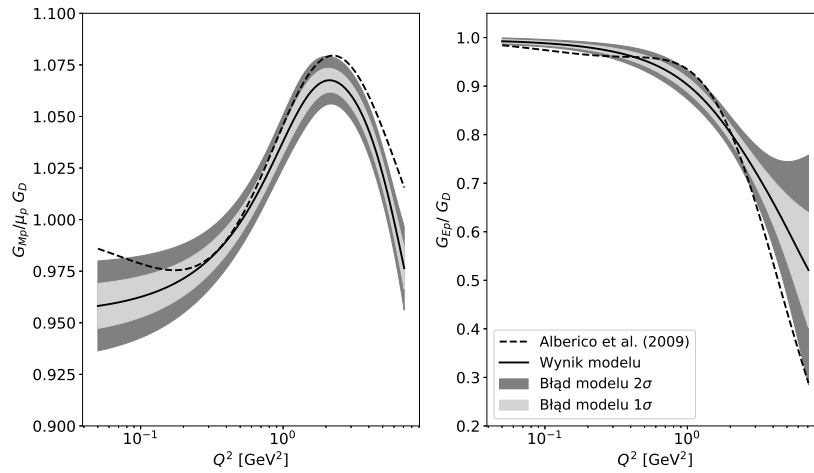
$$\chi^2 = \frac{1}{n} [\chi_\sigma^2 + \chi_{PT}^2] \quad (4.8)$$

$$\chi_{PT}^2 = \sum_{i=1}^{n_k^{PT}} \left(\frac{\mathcal{R}_i^{th} - \mathcal{R}_i^{ex}}{\Delta\mathcal{R}_i} \right)^2 \quad (4.9)$$

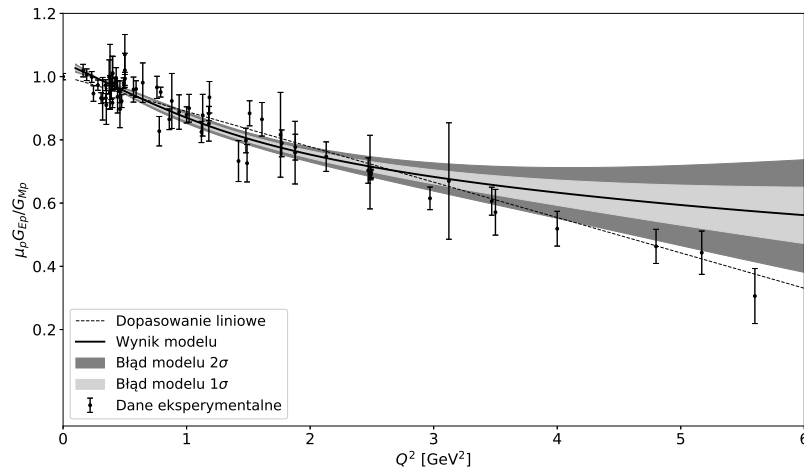
Parametry i nauka sieci

Tabela 4.2: Hiperparametry modelu

Kategoria	Parametr	Wartość
Generowanie danych	N_{rep}	500
k -krotna walidacja krzyżowa	k	5
Algorytm uczący	α (<i>learning rate</i>)	0,025
	λ (<i>decay</i>)	0,001
	β (<i>pęd</i>)	0,9
Sieć neuronowa	Liczba warstw	2
	Ilość neuronów	(3,5)



Rysunek 4.8: Magnetyczna i elektryczna funkcja postaci. Linia ciągła oznacza średnią ze wszystkich modeli, zacienione pola wyznaczają obszary 1σ (jasny) oraz 2σ (ciemny). Linia przerywana przedstawia wyniki z publikacji [1].



Rysunek 4.9: Pomiary stosunku elektrycznej i magnetycznej funkcji postaci wraz z dopasowaniem liniowym oraz dopasowaniem modelu statystycznego. Linia ciągła oznacza średni wynik ze wszystkich wytrenowanych modeli, zacienione pola wyznaczają obszary 1σ oraz 2σ .

Bibliografia

- [1] W. M. Alberico, S. M. Bilenky, C. Giunti, and K. M. Graczyk. Electromagnetic form factors of the nucleon: New fit and analysis of uncertainties. *Phys. Rev. C*, 79(6):065204, June 2009.
- [2] J. E. Dennis, Jr. and R. B. Schnabel. *Numerical Methods for Unconstrained Optimization and Nonlinear Equations (Classics in Applied Mathematics, 16)*. Soc for Industrial & Applied Math, 1996.
- [3] S. Forte, L. s. Garrido, J. I. Latorre, and A. Piccione. Neural network parametrization of deep-inelastic structure functions. *Journal of High Energy Physics*, 5:062, May 2002.
- [4] C. Francois. *Deep Learning with Python*. Manning Publications, 2017.
- [5] K. M. Graczyk. Two-photon exchange effect studied with neural networks. *Phys. Rev. C*, 84:034314, Sep 2011.
- [6] G. James, D. Witten, T. Hastie, and R. Tibshirani. *An Introduction to Statistical Learning: With Applications in R*. Springer Publishing Company, Incorporated, 2014.
- [7] C. F. Perdrisat, V. Punjabi, and M. Vanderhaeghen. Nucleon electromagnetic form factors. *Progress in Particle and Nuclear Physics*, 59:694–764, Oct. 2007.
- [8] S. Raschka. *Python Machine Learning*. Packt Publishing, 2015.
- [9] S. Ruder. An overview of gradient descent optimization algorithms. *ArXiv e-prints*, Sept. 2016.
- [10] D. E. Rumelhart, G. E. Hinton, and R. J. Williams. Learning representations by back-propagating errors. *Nature*, 323:533–536, Oct. 1986.

- [11] StatSoft. Uczenie perceptronu wielowarstwowego. *Internetowy Podręcznik Statystyki*, 2011.