

Spis treści

| | | |
|----------|--|-----------|
| 1 | Wstęp | 2 |
| 1.1 | Notacja | 2 |
| 1.2 | Uczenie nadzorowane | 2 |
| 1.3 | Uczenie nienadzorowane | 4 |
| 1.4 | Praktyka uczenia maszynowego | 4 |
| 1.4.1 | Przygotowanie danych | 4 |
| 1.4.2 | Metryki do oceny regresji i klasyfikacji | 5 |
| 1.4.3 | Tuning hiperparametrów i walidacja skrośna | 9 |
| 2 | Probabilistyka | 10 |
| 2.1 | Zmienne losowe | 10 |
| 2.2 | Ważne rozkłady jednowymiarowe | 13 |
| 2.3 | Wielowymiarowy rozkład normalny | 16 |
| 2.4 | Wnioskowanie statystyczne | 17 |
| 2.5 | Liczby losowe w komputerze | 20 |
| 2.6 | Monte Carlo | 21 |
| 2.6.1 | Algorytm Importance Sampling | 22 |
| 2.6.2 | Algorytm Metropolisa–Hastingsa | 23 |
| 2.7 | Estymator jądrowy gęstości | 26 |
| 3 | Podstawy statystycznego uczenia maszynowego | 27 |
| 3.1 | Regresja liniowa | 27 |
| 3.2 | Regularyzacja | 28 |
| 3.3 | Regresja kwantylowa | 29 |
| 3.4 | Procesy gaussowskie | 31 |
| 3.5 | Naiwny klasyfikator bayesowski | 35 |
| 3.6 | Klasyfikator najbliższych sąsiadów | 36 |
| 3.7 | Regresja logistyczna | 37 |
| 3.8 | Regresja softmax | 39 |
| 4 | Uczenie głębokie i sieci neuronowe | 41 |
| 4.1 | Sieci w pełni połączone | 41 |
| 4.2 | Funkcje kosztu | 45 |
| 4.3 | Aspekty optymalizacyjne | 47 |
| 4.3.1 | Metoda pędu | 50 |
| 4.3.2 | Metoda Nesterova | 50 |
| 4.3.3 | Adam | 51 |
| 4.4 | Algorytm wstecznej propagacji błędów | 52 |
| 4.5 | Inicjalizacja | 56 |
| 4.6 | Regularyzacja | 58 |
| 4.7 | Hiperparametry | 61 |
| 4.7.1 | Dobór hiperparametrów | 61 |
| 4.7.2 | Zjawisko double descent | 62 |
| 4.8 | Sieci konwolucyjne | 64 |
| 4.9 | Sieci rezydualne | 64 |
| 4.10 | Sieci transformer | 64 |

1 Wstęp

Celem tych notatek jest zwięźle przedstawienie kompletu zagadnień związanych z szeroko pojętym uczeniem głębokim jako podejściem do Sztucznej Inteligencji (SI). Zaczynamy od minimalnego zbioru wymaganych tematów z zakresu rachunku prawdopodobieństwa i statystyki matematycznej. Następnie opisujemy podstawowe metody uczenia maszynowego z probabilistycznego punktu widzenia. W końcu przechodzimy do zasadniczej części związanej z uczeniem głębokim i sieciami neuronowymi. W każdej części staramy się przedstawiać opisywane tematy w sposób minimalistyczny, skupiając się głównie na matematycznej i ideowej, a nie implementacyjnej stronie zagadnień. Liczymy, iż takie podejście zapewni odpowiednio głębokie zrozumienie tematu, dzięki któremu dalsze studiowanie całej gamy specyficznych technicznych tematów nie sprawi żadnego problemu.

1.1 Notacja

W dalszej części tekstu będziemy stosować przedstawioną tutaj pokrótce notację. Wektory, które traktujemy jako elementy przestrzeni \mathbb{R}^d ze standardowo zdefiniowanymi operacjami dodawania i mnożenia przez skalar będziemy oznaczali wytłuszczonymi małymi lub wielkimi literami np. \mathbf{x} , \mathbf{X} . Wielkość \mathbf{x}^i będzie oznaczać dany element wektora (w tym przypadku i -ty element \mathbf{x}). Wielkość \mathbf{x}_μ będzie oznaczać pewien (w tym przypadku μ -ty) element pewnego zbioru wektorów. Macierze oraz wielowymiarowe tablice (zwane również niefortunnie tensorami) będziemy oznaczać (jedyńie) wytłuszczonymi wielkimi literami np. \mathbf{X} , Φ . Analogicznie jak w przypadku wektorów przez $\mathbf{X}^{i_1 i_2 \dots i_k}$ będziemy oznaczać (i_1, i_2, \dots, i_k) element k -wymiarowej tablicy \mathbf{X} , natomiast \mathbf{X}_μ będzie oznaczać μ -ty element pewnego zbioru tablic.

1.2 Uczenie nadzorowane

Uczenie nadzorowane jest jednym z dwóch podstawowych (pomijając tzw. uczenie ze wzmocnieniem) paradygmatów w uczeniu maszynowym, którego ogólną ideą jest zdefiniowanie pewnego modelu odwzorowującego dane wejściowe na wyjściowe predykcje. Zakładamy w nim, iż mamy dostępny zbiór obserwacji $\mathcal{X} = \{y_i(\mathbf{x}_i)\}_{i=1}^n$, gdzie $\mathbf{x} \in \mathbb{R}^m$ nazywamy wektorem cech a $y(\mathbf{x})$ jest prawidłową wartością odpowiedzi dla tych cech. Dwa najbardziej podstawowe przypadki zagadnień tego rodzaju to regresja oraz klasyfikacja. W przypadku regresji zmienna y przyjmuje wartości z pewnego podzbioru

liczb rzeczywistych. W przypadku klasyfikacji zmienna y przyjmuje wartości ze skończonego zbioru kategorii, przy czym wartości z tego zbioru nie powinny posiadać naturalnej tj. wynikającej z natury problemu, relacji porządku.

W jaki sposób tworzymy model odwzorowujący \mathbf{x} na y ? W dalszych paragrafach poznamy różne metody, ale najczęściej modelem jest pewna rodzina funkcji postaci $\phi(\mathbf{x}; \mathbf{w})$ parametryzowana skończoną liczbą parametrów, które możemy łącznie zapisać jako pewien wektor \mathbf{w} . Aby znaleźć parametry \mathbf{w} , dzięki którym dla konkretnego zagadnienia model będzie zadowalająco odwzorowywał cechy na predykcje (innymi słowy aby nauczyć model) wprowadzamy dodatkowo funkcjonal kosztu (z ang. *loss function*) $L(\mathcal{X}; \mathbf{w})$, który kwantyfikuje odpowiedzi modelu ϕ w stosunku do znanych prawidłowych odpowiedzi y dla danych ze zbioru \mathcal{X} . Najczęściej ma on postać

$$L(\mathcal{X}; \mathbf{w}) = -\frac{1}{n} \sum_{i=1}^n \log p(y_i(\mathbf{x}_i) | \mathbf{w}) ,$$

gdzie $p(y(\mathbf{x}) | \mathbf{w})$ jest warunkową gęstością prawdopodobieństwa danej obserwacji $y(\mathbf{x})$ warunkowaną przez wartość parametrów \mathbf{w} (oczywiście gęstość ta zależy od estymowanych parametrów \mathbf{w} przez funkcję $\phi(\mathbf{x}; \mathbf{w})$). Do powyższego funkcjonu możemy również dodawać tzw. człony regularyzujące (z ang. *regularizers*). Trening modelu polega wówczas na znalezieniu parametrów \mathbf{w}^* , które minimalizują funkcjonal kosztu na zbiorze treningowym \mathcal{X}

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} L(\mathcal{X}; \mathbf{w}) .$$

Zauważmy, że takie podejście ma jedną zasadniczą wadę – istotnie nie interesuje nas tak naprawdę, jak model radzi sobie na zbiorze treningowym (tzn. zagadnienie uczenia jest czymś więcej niż numeryczną minimalizacją funkcji), tylko jak będzie radził sobie na nowych, niewidzianych wcześniej danych (zależy nam przede wszystkim na generalizacji). Sytuację, w której model bardzo dobrze modeluje dane w zbiorze treningowym, ale słabo radzi sobie na nowych danych nazywamy przeuczeniem lub nadmiernym dopasowaniem (z ang. *overfitting*). Sytuację, w której model słabo radzi sobie zarówno na zbiorze treningowym, jak i na nowych danych nazywamy niedouczeniem lub niedopasowaniem (z ang. *underfitting*). Występowanie *overfittingu* i *underfittingu* jest powiązane z pojemnością (z ang. *capacity*) modelu. Złożony model o dużej pojemności potrafi dopasować się do bardzo skomplikowanych obserwacji (jest elastyczny), ale istnieje ryzyko jego przeuczenia (mówimy wówczas o *high variance*). Dla prostego modelu o małej

pojemności istnieje z kolei ryzyko, iż nie ma on wystarczająco ekspresywności (mówimy wówczas o *high bias*).

1.3 Uczenie nienadzorowane

W przypadku uczenia nienadzorowanego naszym celem nie jest znalezienie modelu odwzorowującego cechy na predykcje. Chcemy raczej zrozumieć wewnętrzną strukturę danych oraz odkryć zależności między zmiennymi lub grupami zmiennych. Modele tego rodzaju znajdują zastosowanie w analizie biznesowej, gdzie pozwalają, chociażby na analizę ważności poszczególnych wskaźników, czy wizualizację wysoko-wymiarowych danych.

1.4 Praktyka uczenia maszynowego

W dalszej części skupiamy się przede wszystkim na matematycznej stronie prezentowanych zagadnień, ale należy pamiętać, iż dowolną próbę wdrożenia modelu uczenia maszynowego należy zacząć od dokładnej inspekcji danych, dla których przygotowujemy ów model („become one with the data”, A. Karpathy), a zakończyć dogłębną analizą metryk pozwalających na ewaluację wytrenowanego modelu. Warunkiem koniecznym udanego wdrożenia modelu jest więc odpowiednie zebranie, analiza i przygotowanie danych, które trafiają następnie jako wejście do modelu ML, a następnie odpowiedni dobór i dogłębną analiza wyników ewaluacji modelu. W tym paragrafie pokrótce opisujemy elementarne praktyki, o których należy pamiętać przy wdrażaniu modeli ML.

1.4.1 Przygotowanie danych

Kluczem do uzyskania dobrych wyników przy korzystaniu z algorytmów uczenia maszynowego jest odpowiednie przygotowanie danych (z ang. *pre-processing*). Typowo preprocessing składa się z:

- eksploracji danych oraz wstępnego czyszczenia, w szczególności usunięcia jawnych wartości odstających (z ang. *outliers*) oraz cech posiadających zbyt dużo wartości brakujących;
- analizy rozkładu zmiennej docelowej oraz ewentualnej transformacji logarymicznej, która poprawia stabilność numeryczną, gdy przewidywane wartości są dużymi dodatnimi liczbami rzeczywistymi, zmienia dziedzinę zmiennej objaśnianej z \mathbb{R}_+ na \mathbb{R} oraz dodatkowo jest przykładem transformacji stabilizującej wariancję;

- podziału zbioru na część treningową oraz testową;
- dokonania skalowania i imputacji brakujących wartości cech (metody `.fit()` wywołujemy jedynie dla zbioru treningowego);
- usunięcia silnie skorelowanych cech;
- zakodowania wartości kategorycznych za pomocą tzw. *one-hot encoding* pamiętając o *dummy variable trap* – jedną z k kategorii kodujemy za pomocą wektora *one-hot* długości $n - 1$, aby uniknąć zależności liniowej między cechami (opcja `drop="first"` w `OneHotEncoder` w `scikit-learn`);
- wykonania feature engineering – dodania wielomianów cech do naszych danych lub konstruowania innych cech (np. cech określających miesiąc, dzień itp.).

Podział zbioru na część treningową i testową jest najważniejszym etapem preprocessingu. Zbiór testowy wydzielamy, aby po wytrenowaniu modelu sprawdzić, jak poradzi on sobie na nowych, niewidzianych wcześniej danych. Powinniśmy go traktować jako dane, które będziemy w przyszłości dostawać po wdrożeniu modelu do realnego systemu. Takie dane również będziemy musieli przeskalować, zakodować itp., ale parametry potrzebne do wykonania tych transformacji możemy wziąć jedynie z dostępnego wcześniej zbioru treningowego. Wykorzystanie danych testowych w procesie treningu to *błąd wycieku danych* (z ang. *data leakage*). Skutkuje on niepoprawnym, nadmiernie optymistycznym oszacowaniem jakości modelu.

1.4.2 Metryki do oceny regresji i klasyfikacji

Zasadniczo, aby ocenić predykcje modelu używamy odpowiednich metryk, których wartości określają jak dobry jest model.

W przypadku regresji najczęściej używanymi metrykami są RMSE (z ang. *Root Mean Squared Error*) oraz MAE (z ang. *Mean Absolute Error*) zdefiniowane odpowiednio jako

$$\text{RMSE} := \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}, \quad \text{MAE} := \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|.$$

Metryki te mają jednakową jednostkę jak predykcje. Jeśli chcielibyśmy mieć liczbę względną określającą jakość modelu to mamy do dyspozycji metryki

MAPE (z ang. *Mean Absolute Percentage Error*) oraz SMAPE (z ang. *Symmetric Mean Absolute Percentage Error*) zdefiniowane odpowiednio jako

$$\text{MAPE} := \frac{1}{n} \sum_{i=1}^n \left| \frac{y_i - \hat{y}_i}{y_i} \right|, \quad \text{SMAPE} := \frac{1}{n} \sum_{i=1}^n \frac{|y_i - \hat{y}_i|}{|y_i| + |\hat{y}_i|}.$$

Obie te metryki mają zakres od 0 do 1, przy czym niższa wartość oznacza lepszy model. Metryki te mają jednak szereg problemów, z których najważniejsze to: problemy, gdy wartości są bliskie 0, asymetryczne traktowanie predykcji za dużych oraz za małych. Z tych powodów znacznie lepszą względną metryką jest MASE (z ang. *Mean Absolute Scaled Error*)

$$\text{MASE} := \frac{\sum_{i=1}^n |y_i - \hat{y}_i|}{\sum_{i=1}^n |y_i - \bar{y}|},$$

gdzie $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$. Metryka MASE jest zatem względnym błędem MAE jaki popełnia nasz model w stosunku do modelu naiwnego, który przewiduje zawsze wartość średnią.

W przypadku zadania klasyfikacji binarnej naszym celem dla danego wektora cech jest zwrócenie jednej z dwóch klas, które będziemy nazywać klasą pozytywną i negatywną. O ile w przypadku regresji pomiar jakości modelu był całkiem prosty, o tyle w przypadku klasyfikacji sytuacja jest nieco bardziej skomplikowana. Zauważmy bowiem, iż mamy 4 możliwości odpowiedzi klasyfikatora

- *True Positive (TP)* – poprawnie zaklasyfikowaliśmy klasę pozytywną jako pozytywną
- *True Negative (TN)* – poprawnie zaklasyfikowaliśmy klasę negatywną jako negatywną
- *False Positive (FP)* – niepoprawnie zaklasyfikowaliśmy klasę negatywną jako pozytywną
- *False Negative (FN)* – niepoprawnie zaklasyfikowaliśmy klasę pozytywną jako negatywną

Na podstawie ilości TP, TN, FP i FN w zbiorze testowym możemy wykreślić tzw. *macierz pomyłek* (z ang. *confusion matrix*) pokazującą ilość każdej z możliwości. Następnie możemy obliczyć różne stosunki tych wartości, aby uzyskać różne metryki. Najbardziej standardowymi są *accuracy*, *precision* oraz *recall* (lub inaczej *sensitivity*) zdefiniowane jako

$$\text{Accuracy} := \frac{\text{TP} + \text{TN}}{n}, \quad \text{Precision} := \frac{\text{TP}}{\text{TP} + \text{FP}}, \quad \text{Recall} := \frac{\text{TP}}{\text{TP} + \text{FN}}.$$

Wartość accuracy mówi po prostu jaki stosunek przykładów został poprawnie zaklasyfikowany (zauważmy tutaj, że $TP + TN + FP + FN = n$). Nie jest to jednak dobra miara jakości, gdy nasz zbiór jest niezbalansowany, tj. zawiera więcej przykładów określonej klasy.

| | | Predicted condition | | | | |
|------------------|---|---|--|--|--|--|
| | | Total population = P + N | Predicted Positive (PP) | Predicted Negative (PN) | Informedness, bookmaker informedness (BM) = TPR + TNR - 1 | Prevalence threshold (PT) $= \frac{\sqrt{TPR \times FPR} \cdot FPR}{TPR \cdot FPR}$ |
| Actual condition | Positive (P) [a] | True positive (TP), hit ^[b] | False negative (FN), miss, underestimation | True positive rate (TPR), recall, sensitivity (SEN), probability of detection, hit rate, power $= \frac{TP}{P} = 1 - FNR$ | False negative rate (FNR), miss rate type I error ^[c] $= \frac{FN}{P} = 1 - TPR$ | |
| | Negative (N) ^[d] | False positive (FP), false alarm, overestimation | True negative (TN), correct rejection ^[e] | False positive rate (FPR), probability of false alarm, fail-out type I error ^[f] $= \frac{FP}{N} = 1 - TNR$ | True negative rate (TNR), specificity (SPC), selectivity $= \frac{TN}{N} = 1 - FPR$ | |
| | Prevalence $= \frac{P}{P + N}$ | Positive predictive value (PPV), precision $= \frac{TP}{PP} = 1 - FDR$ | False omission rate (FOR) $= \frac{FN}{PN} = 1 - NPV$ | Positive likelihood ratio (LR+) $= \frac{TPR}{FPR}$ | Negative likelihood ratio (LR-) $= \frac{FNR}{TNR}$ | |
| | Accuracy (ACC) $= \frac{TP + TN}{P + N}$ | False discovery rate (FDR) $= \frac{FP}{PP} = 1 - PPV$ | Negative predictive value (NPV) $= \frac{TN}{PN} = 1 - FOR$ | Markedness (MK), deltaP (Δp) $= PPV + NPV - 1$ | Diagnostic odds ratio (DOR) $= \frac{LR+}{LR-}$ | |
| | Balanced accuracy (BA) $= \frac{TPR + TNR}{2}$ | F ₁ score $= \frac{2 \cdot PPV \times TPR}{PPV + TPR}$ $= \frac{2 \cdot TP}{2 \cdot TP + FP + FN}$ | Fowkes–Mallows index (FM) $= \sqrt{PPV \times TPR}$ | Matthews correlation coefficient (MCC) $= \sqrt{TPR \times TNR \times PPV \times NPV}$ $- \sqrt{FNR \times FPR \times FOR \times FDR}$ | Threat score (TS), critical success index (CSI), Jaccard index $= \frac{TP}{TP + FN + FP}$ | |

Rysunek 1: Macierz pomyłek oraz możliwe metryki oceny jakości klasyfikatora. Źródło: en.wikipedia.org/wiki/Confusion_matrix

Wartość precision określa jak pewny jest klasyfikator przy wykrywaniu klasy pozytywnej, natomiast recall mówi o tym jak dobrze klasyfikator „wyławia” przykłady pozytywne. Zauważmy jednak, iż nie możemy stosować żadnej z tych metryk w odosobnieniu. Istotnie klasyfikator, który zwraca zawsze klasę pozytywną ma maksymalny recall, a klasyfikator, który zwraca zawsze klasę negatywną ma nieokreślony precision (i jest oczywiście beznadziejnym klasyfikatorem). Musimy więc zawsze ewaluować model na obu tych metrykach i jedynie dobry wynik obu z nich mówi o jakości klasyfikatora. Oczywiście czasami chcielibyśmy określić jakość modelu za pomocą jednej liczby, a niekoniecznie sprawdzać zawsze macierz pomyłek (choć jest to bardzo użyteczne) lub podawać wartości dwóch metryk. Metryką, która

łączy precision i recall jest F_β -score zdefiniowany jako

$$F_\beta := (1 + \beta^2) \frac{\text{Precision} \cdot \text{Recall}}{\beta^2 \cdot \text{Precision} + \text{Recall}},$$

gdzie β określa ile razy bardziej ważny jest recall od precision. Typowo używa się F_1 -score

$$F_1 = 2 \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}.$$

Wiele klasyfikatorów oprócz twardych predykcji zwraca również rozkład prawdopodobieństwa nad klasami. W przypadku klasyfikacji binarnej jest to oczywiście rozkład zero-jedynkowy z parametrem p określającym prawdopodobieństwo klasy pozytywnej dla danego wektora cech. Standardowo oczywiście twardą predykcją jest ta z klas, która ma większe prawdopodobieństwo, czyli (co równoważne) predykcją jest klasa pozytywna jeśli $p > 0.5$. W niektórych problemach chcemy jednak zmienić ten próg i dokonać tzw. *threshold tuning*. Wykresem, który pozwala na dokonanie tuningu progu jest krzywa *ROC* (z ang. *Receiver Operatic Characteristic curve*), która jest krzywą parametryczną wyznaczoną przez punkty $(\text{FPR}(\text{threshold}), \text{TPR}(\text{threshold}))$ dla progów z zakresu $[0; 1]$, gdzie

$$\text{TPR} := \frac{\text{TP}}{\text{TP} + \text{FN}}, \quad \text{FPR} := \frac{\text{FP}}{\text{FP} + \text{TN}}.$$

Metryką niezależną od wybranego progu jest tzw. *AUROC* (z ang. *Area under ROC curve*) zdefiniowany jako pole powierzchni pod krzywą ROC dla danego klasyfikatora. Zauważmy, że klasyfikator losowy, który zwraca zawsze klasę pozytywną z prawdopodobieństwem równym wartości progu ma wartość AUROC równą 0.5, natomiast idealny klasyfikator, który niezależnie od wartości progu klasyfikuje wszystkie przykłady poprawnie ma AUROC równy 1.

Inną analogiczną metryką jest *AUPRC*, gdzie zamiast krzywej ROC stosujemy krzywą *PRC* (z ang. *Precision-Recall Curve*), w której zamiast TPR i FPR używamy odpowiednio Precision i Recall. Metryka AUPRC jest często wykorzystywana w przypadku klasyfikacji ekstremalnie niebalansowanej, w której mamy bardzo mało ($< 1\%$) klasy pozytywnej.

W przypadku klasyfikacji wieloklasowej używamy zasadniczo takich samych metryk jak w klasyfikacji binarnej, ale wprowadzamy mikro i makro uśrednianie (z ang. *micro/macro-averaging*). Przez TP_k będziemy rozumieć liczbę prawidłowo zaklasyfikowanych przykładów z klasy k , FP_k to liczba przykładów z innych klas, które zaklasyfikowaliśmy nieprawidłowo jako k -tą klasę, FN_k to liczba przykładów z klasy k , które zaklasyfikowaliśmy jako

inną klasę. Wówczas odpowiednie metryki mają postać

$$\text{MicroPrecision} := \frac{\sum_k \text{TP}_k}{\sum_k \text{TP}_k + \sum_k \text{FP}_k},$$

$$\text{MacroPrecision} := \frac{1}{K} \sum_{k=1}^K \frac{\text{TP}_k}{\text{TP}_k + \text{FP}_k}$$

oraz

$$\text{MicroRecall} := \frac{\sum_k \text{TP}_k}{\sum_k \text{TP}_k + \sum_k \text{FN}_k},$$

$$\text{MacroRecall} := \frac{1}{K} \sum_{k=1}^K \frac{\text{TP}_k}{\text{TP}_k + \text{FN}_k}.$$

W przypadku klasyfikacji wieloklasowej macierz pomyłek jest macierzą wymiaru $K \times K$, gdzie K jest liczbą klas.

1.4.3 Tuning hiperparametrów i walidacja skrośna

Praktycznie wszystkie modele uczenia maszynowego mają hiperparametry, często liczne, które w zauważalny sposób wpływają na wyniki, a szczególnie na underfitting i overfitting. Ich wartości trzeba dobrać zatem dość dokładnie. Proces doboru hiperparametrów nazywa się tuningiem hiperparametrów (z ang. *hyperparameter tuning*).

Istnieje na to wiele sposobów. Większość z nich polega na tym, że trenuje się za każdym razem model z nowym zestawem hiperparametrów i wybiera się ten zestaw, który pozwala uzyskać najlepsze wyniki. Metody głównie różnią się między sobą sposobem doboru kandydujących zestawów hiperparametrów. Najprostsze i najpopularniejsze to:

- pełne przeszukiwanie (z ang. *grid search*) – definiujemy możliwe wartości dla różnych hiperparametrów, a metoda sprawdza ich wszystkie możliwe kombinacje (czyli siatkę),
- losowe przeszukiwanie (z ang. *randomized search*) – definiujemy możliwe wartości jak w pełnym przeszukiwaniu, ale sprawdzamy tylko ograniczoną liczbę losowo wybranych kombinacji.

Jak ocenić, jak dobry jest jakiś zestaw hiperparametrów? Nie możemy sprawdzić tego na zbiorze treningowym – wyniki byłyby zbyt optymistyczne. Nie możemy wykorzystać zbioru testowego – mielibyśmy wyciek danych, bo

wybieralibyśmy model explicite pod nasz zbiór testowy. Trzeba zatem osobnego zbioru, na którym będziemy na bieżąco sprawdzać jakość modeli dla różnych hiperparametrów. Jest to zbiór walidacyjny (z ang. *validation set*). Zbiór taki wycina się ze zbioru treningowego.

Jednorazowy podział zbioru na części nazywa się *split validation* lub *holdout*. Używamy go, gdy mamy sporo danych, i 10-20% zbioru jako dane walidacyjne czy testowe to dość dużo, żeby mieć przyzwoite oszacowanie. Zbyt mały zbiór walidacyjny czy testowy da nam mało wiarygodne wyniki – nie da się nawet powiedzieć, czy zbyt pesymistyczne, czy optymistyczne. W praktyce niestety często mamy mało danych. Trzeba zatem jakiejś magicznej metody, która stworzy nam więcej zbiorów walidacyjnych z tej samej ilości danych. Taką metodą jest walidacja skrośna (z ang. *cross validation*, CV). Polega na tym, że dzielimy zbiór treningowy na K równych podzbiorów, tzw. foldów. Każdy podzbiór po kolei staje się zbiorem walidacyjnym, a pozostałe łączymy w zbiór treningowy. Trenujemy zatem K modeli dla tego samego zestawu hiperparametrów i każdy testujemy na zbiorze walidacyjnym. Mamy K wyników dla zbiorów walidacyjnych, które możemy uśrednić (i ewentualnie obliczyć odchylenie standardowe). Takie wyniki są znacznie bardziej wiarygodne.

2 Probabilistyka

2.1 Zmienne losowe

Zmienna losowa to formalnie odwzorowanie ze zbioru zdarzeń elementarnych Ω tj. zbioru atomowych wyników doświadczenia losowego w zbiór \mathbb{R}^n

$$\mathbf{X} : \Omega \mapsto \mathbb{R}^n.$$

Jest to zatem funkcja, która przyporządkowuje zdarzeniom losowym wartość liczbową. Każda zmienna losowa opisuje więc zmienną w klasycznym sensie, której wartości pochodzi z pewnego rozkładu. Rozkład ten jest zadany jednoznacznie przez funkcję $F : \mathbb{R}^n \mapsto [0; 1]$ taką, że

$$F(\mathbf{x}) := \Pr(\mathbf{X}^1 \leq x^1, \dots, \mathbf{X}^n \leq x^n),$$

którą nazywa się dystrybuantą (z ang. *cumulative distribution function*, *cdf*). Każdy rozkład zmiennej losowej można opisać za pomocą dystrybuanty jednak jest to niewygodne. W dwóch przypadkach: rozkładów dyskretnych i rozkładów ciągłych rozkład zmiennej losowej można opisać prościej za pomocą odpowiednio funkcji prawdopodobieństwa (z ang. *probability*

mass function, pmf) oraz gęstości prawdopodobieństwa (z ang. *probability density function, pdf*).

Definicja 2.1. Zmienna losowa \mathbf{X} ma dyskretny rozkład prawdopodobieństwa, jeśli istnieje skończony lub przeliczalny zbiór $S \subset \mathbb{R}^n$ taki, że $\Pr(\mathbf{X} \in S) = 1$. Wówczas rozkład ten jest zadany przez podanie funkcji prawdopodobieństwa $p(\mathbf{x}) = \Pr(\mathbf{X} = \mathbf{x})$ dla $\mathbf{x} \in S$.

Definicja 2.2. Zmienna losowa \mathbf{X} ma z kolei ciągły rozkład prawdopodobieństwa, jeśli istnieje funkcja $p : \mathbb{R}^n \mapsto \mathbb{R}_+$ taka, że

$$\Pr(\mathbf{X}^1 \in (a_1; b_1), \dots, \mathbf{X}^n \in (a_n; b_n)) = \int_{a_1}^{b_1} \cdots \int_{a_n}^{b_n} p(\mathbf{x}) d^n \mathbf{x}$$

dla dowolnej kostki $(a_1; b_1) \times \dots \times (a_n; b_n)$.

Definicja 2.3 (Wartości oczekiwanej). Wartością oczekiwaną funkcji zmiennej losowej $\mathbf{f}(\mathbf{X})$ nazywamy wektor $\mathbb{E}[\mathbf{f}(\mathbf{x})]$ określoną wzorem

$$\sum_{\mathbf{x} \in S} \mathbf{f}(\mathbf{x}) p(\mathbf{x}), \quad \int_{\mathbb{R}^n} \mathbf{f}(\mathbf{x}) p(\mathbf{x}) d^n \mathbf{x}$$

odpowiednio dla rozkładu dyskretnego i ciągłego.

Definicja 2.4 (Macierzy kowariancji). Macierzą kowariancji funkcji zmiennej losowej $\mathbf{f}(\mathbf{X})$ nazywamy macierz

$$\mathbb{E} [(\mathbf{f}(\mathbf{x}) - \mathbf{m}_f)(\mathbf{f}(\mathbf{x}) - \mathbf{m}_f)^T],$$

gdzie

$$\mathbf{m}_f = \mathbb{E} [\mathbf{f}(\mathbf{x})].$$

Elementy diagonalne macierzy kowariancji nazywamy wariancjami, a elementy pozadiagonalne kowariancjami.

Definicja 2.5 (Kwantyla i mody). Kwantylem q_p rzędu $p \in (0; 1)$ zmiennej losowej jednowymiarowej o rozkładzie ciągłym z dystrybucją F nazywamy dowolne rozwiązanie równania

$$F(x) = p.$$

Modą tej zmiennej nazywamy dowolne maksimum lokalne gęstości tego rozkładu.

Twierdzenie 2.1. Niech zmienna n -wymiarowa \mathbf{X} ma rozkład ciągły o gęstości $p_{\mathbf{X}}$ i niech $\mathbf{Y}^i = \boldsymbol{\varphi}^i(\mathbf{X})$ dla $i = 1, \dots, n$. Jeśli odwzorowanie $\boldsymbol{\varphi}$ jest różniczkowalne i odwracalne, przy czym odwzorowanie odwrotne $\boldsymbol{\psi} = \boldsymbol{\varphi}^{-1}$ jest różniczkowalne, to n -wymiarowa zmienna \mathbf{Y} ma rozkład o gęstości

$$p_{\mathbf{Y}}(\mathbf{y}) = |J|p_{\mathbf{X}}(\boldsymbol{\psi}(\mathbf{y})),$$

gdzie $J := \det \left[\frac{\partial \psi^j}{\partial y^i} \right]$ jest jakobianem odwzorowania $\boldsymbol{\psi}$.

Twierdzenie 2.2. Dla dowolnych zmiennych losowych \mathbf{X} , \mathbf{Y} zachodzi

$$p_{\mathbf{X}}(\mathbf{x}) = \int_{\mathbb{R}^k} p(\mathbf{x}, \mathbf{y}) \mathrm{d}^k \mathbf{y} \quad (\text{sum rule})$$

Twierdzenie 2.3. Dla dowolnych zmiennych losowych \mathbf{X} , \mathbf{Y} zachodzi

$$p(\mathbf{x}, \mathbf{y}) = p(\mathbf{x} \mid \mathbf{y})p_{\mathbf{Y}}(\mathbf{y}) \quad (\text{product rule})$$

Twierdzenie 2.4. Zmienne losowe \mathbf{X} i \mathbf{Y} są niezależne wtedy i tylko wtedy, gdy

$$p(\mathbf{x}, \mathbf{y}) = p_{\mathbf{X}}(\mathbf{x})p_{\mathbf{Y}}(\mathbf{y}) \quad (\text{independence})$$

Twierdzenie 2.5. Dla dowolnych zmiennych losowych \mathbf{X} , \mathbf{Y} zachodzi

$$\underbrace{p(\mathbf{x} \mid \mathbf{y})}_{\text{posterior}} = \frac{\underbrace{p(\mathbf{y} \mid \mathbf{x})}_{\text{likelihood}} \underbrace{p_{\mathbf{X}}(\mathbf{x})}_{\text{prior}}}{\underbrace{\int_{\mathbb{R}^k} p(\mathbf{y} \mid \mathbf{x}) p_{\mathbf{X}}(\mathbf{x}) d^k \mathbf{x}}_{\text{evidence}}} \quad (\text{Bayes theorem})$$

Definicja 2.6 (Warunkowej wartości oczekiwanej). Warunkową wartość oczekiwaną $\mathbf{f}(\mathbf{X})$ pod warunkiem $\mathbf{Y} = \mathbf{y}$ nazywamy wielkość

$$\mathbb{E}[\mathbf{f}(\mathbf{x}) \mid \mathbf{y}] = \int_{\mathbb{R}^n} \mathbf{f}(\mathbf{x}) p(\mathbf{x} \mid \mathbf{y}) d^n \mathbf{x}$$

2.2 Ważne rozkłady jednowymiarowe

Definicja 2.7 (Rozkładu dwupunktowego). Jeśli X jest zmienną losową rzeczywistą o rozkładzie dyskretnym i $\mathcal{S} = \{x_1, x_2\}$ oraz $p(x_1) = p$, to mówimy, że X ma rozkład dwupunktowy z parametrem p . Jeśli $x_1 = 1$ i $x_2 = 0$ to taki rozkład dwupunktowy nazywamy *rozkładem zero-jedynkowym* (lub rozkładem Bernoulliego) i oznaczamy jako $X \sim \text{Ber}(p)$.

Definicja 2.8 (Schematu dwumianowego). Rozważmy doświadczenie losowe o dwu możliwych wynikach: sukces osiągamy z prawdopodobieństwem p , porażkę z prawdopodobieństwem $1 - p$. Doświadczenie tego rodzaju nazywamy *próbą Bernoulliego*. Doświadczenie takie jest modelowane zmienną losową o rozkładzie dwupunktowym z parametrem p . Schematem dwumianowym (lub schematem Bernoulliego) nazywamy doświadczenie polegające na n -krotnym powtórzeniu próby Bernoulliego, przy założeniu, iż poszczególne próby są od siebie niezależne.

Definicja 2.9 (Rozkładu dwumianowego). Niech X będzie zmienną losową taką, że X jest liczbą sukcesów w schemacie dwumianowym długości n z prawdopodobieństwem sukcesu w każdej próbie równym p . Wówczas

$$\Pr(X = k) = \binom{n}{k} p^k (1 - p)^{n-k}.$$

Rozkład prawdopodobieństwa określony powyższym wzorem nazywam się rozkładem dwumianowym o parametrach n, p . Jeśli zmienna X ma rozkład dwumianowy to stosujemy notację $X \sim \text{Bin}(n, p)$.

Definicja 2.10 (Rozkładu geometrycznego). Mówimy, że zmienna losowa X ma rozkład geometryczny z parametrem $p \in (0; 1)$, tj. $X \sim \text{Geo}(p)$, jeśli $\mathcal{S} = \mathbb{N} \setminus \{0\}$, a funkcja prawdopodobieństwa ma postać

$$p(x) = (1 - p)^{x-1} p.$$

Zmienna X opisuje czas oczekiwania na pierwszy sukces w schemacie dwumianowym o nieskończonej długości.

Definicja 2.11 (Rozkładu Poissona). Jeśli zmienna X o wartościach w \mathbb{N} opisuje liczbę wystąpień pewnego powtarzalnego zdarzenia w przedziale czasowym $[0; t]$, przy czym spełnione są następujące założenia:

- powtórzenia zdarzenia występują niezależnie od siebie;
- „intensywność” wystąpień r jest stała;
- w danej chwili (rozumianej jako odpowiednio mały przedział) może zajść co najwyżej jedno zdarzenie

to zmienna ta ma rozkład Poissona z parametrem $\lambda = rt$, tj. $X \sim \text{Pos}(\lambda)$. Jeśli $X \sim \text{Pos}(\lambda)$, to

$$\Pr(X = k) = \frac{e^{-\lambda} \lambda^k}{k!}.$$

Twierdzenie 2.6 (Poissona). Niech (X_n) będzie ciągiem zmiennych losowych takich, że $X_n \sim \text{Bin}(n, p_n)$, gdzie (p_n) jest ciągiem takim, że

$$\lim_{n \rightarrow \infty} np_n = \lambda$$

dla pewnej liczby $\lambda > 0$. Wówczas

$$\lim_{n \rightarrow \infty} \Pr(X_n = k) = \frac{e^{-\lambda} \lambda^k}{k!}.$$

Definicja 2.12 (Rozkładu jednostajnego). Mówimy, że zmienna X o rozkładzie ciągłym ma rozkład jednostajny na przedziale $[a; b]$ tzn. $X \sim \mathcal{U}(a, b)$ jeśli jej gęstość wyraża się wzorem

$$p(x) = \begin{cases} \frac{1}{b-a}, & x \in [a; b] \\ 0, & x \notin [a; b] \end{cases}.$$

Definicja 2.13 (Rozkładu wykładniczego). Niech T będzie zmienną modelującą czas oczekiwania na pierwsze zdarzenie w ciągu zdarzeń takim, że czas wystąpienia każdego z nich w przedziale $[0; t]$ jest opisany przez zmienną $X \sim \text{Pos}(\lambda t)$. wtedy

$$\Pr(T > t) = \Pr(X = 0) = e^{-\lambda t}$$

oraz

$$\Pr(T > 0) = 1.$$

Mówimy wtedy, że T ma rozkład wykładniczy z parametrem λ , tzn. $T \sim \text{Exp}(\lambda)$. Gęstość rozkładu wykładniczego ma postać

$$p(t) = \begin{cases} 0, & t \leq 0 \\ \lambda e^{-\lambda t}, & t > 0 \end{cases}.$$

Definicja 2.14 (Rozkładu normalnego). Mówimy, że zmienna losowa X o gęstości $p(x)$ ma rozkład normalny z parametrami $\mu \in$

$\mathbb{R}, \sigma^2 \in [0; +\infty)$, tzn. $X \sim \mathcal{N}(\mu, \sigma^2)$, jeśli

$$p(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right).$$

Jeśli $\mu = 0$ i $\sigma = 1$, to taki rozkład nazywamy standardowym rozkładem normalnym.

2.3 Wielowymiarowy rozkład normalny

Definicja 2.15 (Standardowego wielowymiarowego rozkładu normalnego). Zmienna losowa \mathbf{X} ma standardowy n -wymiarowy rozkład normalny jeśli jej składowe są niezależne i dla każdego $i = 1, \dots, n$ $\mathbf{X}^i \sim \mathcal{N}(0, 1)$. Jest to rozkład ciągły o gęstości

$$p(\mathbf{x}) = \frac{1}{\sqrt{(2\pi)^n}} \exp\left(-\frac{1}{2}\mathbf{x}^T\mathbf{x}\right).$$

Definicja 2.16 (Wielowymiarowego rozkładu normalnego). Zmienna losowa \mathbf{X} ma n -wymiarowy rozkład normalny (z ang. *Multivariate Normal Distribution, MVN*), tzn. $\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ jeśli istnieje k -wymiarowa zmienna losowa \mathbf{Z} o standardowym rozkładzie normalnym dla pewnego $k \leq n$ oraz istnieje $\boldsymbol{\mu} \in \mathbb{R}^n$ i macierz $\mathbf{A} \in \mathbb{R}^{n \times k}$ takie, że $\boldsymbol{\Sigma} = \mathbf{A}\mathbf{A}^T$ oraz

$$\mathbf{X} = \mathbf{A}\mathbf{Z} + \boldsymbol{\mu}.$$

Jeśli macierz $\boldsymbol{\Sigma}$ jest dodatnio określona, to rozkład $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ jest ciągły, a jego gęstość jest dana przez

$$p(\mathbf{x}) = \frac{1}{\sqrt{(2\pi)^n \det \boldsymbol{\Sigma}}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right)$$

Macierz $\boldsymbol{\Sigma}^{-1}$ nazywa się *macierzą precyzji*.

Poziomice gęstości niezdegenerowanego wielowymiarowego rozkładu normalnego są elipsoidami, których półosie są skierowane wzdłuż wektorów własnych macierzy $\boldsymbol{\Sigma}$ i mają długości proporcjonalne do pierwiastka z war-

tości własnych.

Twierdzenie 2.7 (Własności niezdegenerowanego rozkładu normalnego). Niech $\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ dla dodatnio określonej macierzy $\boldsymbol{\Sigma}$, wówczas

1. Wszystkie rozkłady brzegowe i warunkowe \mathbf{X} są rozkładami normalnymi. Istotnie jeśli

$$\begin{bmatrix} x \\ y \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \mu_x \\ \mu_y \end{bmatrix}, \begin{bmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_{yy} \end{bmatrix} \right)$$

to można pokazać, iż

$$x \sim \mathcal{N}(\mu_x, \Sigma_{xx}), \quad y \sim \mathcal{N}(\mu_y, \Sigma_{yy}).$$

oraz $x | y \sim \mathcal{N}(\mu_{x|y}, \Sigma_{x|y})$, gdzie

$$\mu_{x|y} = \mu_x + \Sigma_{xy} \Sigma_{yy}^{-1} (y - \mu_y), \quad \Sigma_{x|y} = \Sigma_{xx} - \Sigma_{xy} \Sigma_{yy}^{-1} \Sigma_{yx}.$$

2. Zmienne składowe X^1, \dots, X^n są niezależne wtedy i tylko wtedy, gdy $\boldsymbol{\Sigma}$ jest macierzą diagonalną.

2.4 Wnioskowanie statystyczne

Niech zmienna losowa \mathbf{X} określa model rozkładu pewnej cechy (cech) w ustalonym zbiorze instancji (tzw. *populacji generalnej*). Innymi słowy, przyjmujemy, że wartości cech zachowują się jakby zostały wybrane losowo zgodnie z rozkładem zmiennej \mathbf{X} . Do podstawowych zagadnień wnioskowania statystycznego należą:

- oszacowanie wielkości charakteryzujących rozkład \mathbf{X} (np. wartości średniej albo wariancji);
- weryfikacja hipotez dotyczących rozkładu \mathbf{X} (tym nie będziemy się zajmować).

Definicja 2.17 (Modelu statystycznego). Modelem statystycznym nazywamy parę $(\mathcal{P}, \mathcal{X})$, gdzie \mathcal{P} jest rodziną rozkładów prawdopo-

dobieństwa na zbiorze \mathcal{X} . Zazwyczaj przyjmuje się

$$\mathcal{P} = \{p(\cdot \mid \boldsymbol{\theta}) : \boldsymbol{\theta} \in \Theta\}$$

dla pewnego zbioru parametrów Θ . Model statystyczny nazywamy *parametrycznym* jeśli $\Theta \subset \mathbb{R}^k$.

Definicja 2.18 (Prostej próby losowej). Prostą próbą losową o liczności n nazywamy ciąg niezależnych zmiennych losowych $\mathbf{X}_1, \dots, \mathbf{X}_n$ o tym samym rozkładzie $p(\cdot \mid \boldsymbol{\theta}) \in \mathcal{P}$ (z ang. *independent and identically distributed, i.i.d.*).

Definicja 2.19 (Estymatora). Estymatorem nazywa się statystykę $\hat{\theta}(\mathbf{X}_1, \dots, \mathbf{X}_n)$ służącą do oszacowania wartości parametru θ . Liczbę $\hat{\theta}(\mathbf{x}_1, \dots, \mathbf{x}_n)$ dla konkretnej realizacji prostej próby losowej nazywa się wartością estymatora albo estymatą.

Definicja 2.20 (Funkcji wiarygodności). Funkcją wiarygodności (z ang. *likelihood function*) dla modelu $\mathcal{P} = \{p(\cdot \mid \boldsymbol{\theta}) : \boldsymbol{\theta} \in \Theta\}$ nazywamy funkcję

$$\mathcal{L} : \mathbb{R}^n \times \Theta \ni (\mathbf{x}, \boldsymbol{\theta}) \mapsto \mathcal{L}(\mathbf{x}; \boldsymbol{\theta}) \in [0; +\infty)$$

wyznaczającą rozkład łączny obserwowanych danych jako funkcję parametru $\boldsymbol{\theta}$.

Niech $\mathbf{X}_1, \dots, \mathbf{X}_n$ będzie prostą próbą losową. Jeśli $p(\cdot \mid \boldsymbol{\theta})$ opisuje rozkład warunkowy, z którego pochodzą obserwacje, to

$$\mathcal{L}(\mathbf{x}_1, \dots, \mathbf{x}_n; \boldsymbol{\theta}) = \prod_{i=1}^n p(\mathbf{x}_i \mid \boldsymbol{\theta}).$$

Dla wygody obliczeń często rozważa się tzw. zanegowaną logarytmiczną funkcję wiarygodności (z ang. *Negated Log-Likelihood function, NLL*), tzn.

$$L(\mathbf{x}; \boldsymbol{\theta}) = -\log \mathcal{L}(\mathbf{x}; \boldsymbol{\theta}).$$

Wówczas dla realizacji prostej próby losowej mamy

$$L(\mathbf{x}_1, \dots, \mathbf{x}_n; \boldsymbol{\theta}) = - \sum_{i=1}^n \log p(\mathbf{x}_i \mid \boldsymbol{\theta}).$$

Definicja 2.21 (Estymatora największej wiarygodności). Estymatorem największej wiarygodności (z ang. *Maximum Likelihood Estimator*, *MLE*) nazywamy funkcję $\hat{\boldsymbol{\theta}}$, która przy ustalonych wartościach obserwacji (realizacji prostej próby losowej) $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ maksymalizuje wartość funkcji wiarygodności lub, co równoważne, minimalizuje wartość zanegowanej logarytmicznej funkcji wiarygodności tj.

$$\hat{\boldsymbol{\theta}}(\mathbf{x}_1, \dots, \mathbf{x}_n) = \arg \min_{\boldsymbol{\theta} \in \Theta} \left[- \sum_{i=1}^n \log p(\mathbf{x}_i \mid \boldsymbol{\theta}) \right].$$

Jeśli funkcja wiarygodności jest różniczkowalna względem $\boldsymbol{\theta}$ dla dowolnych \mathbf{x}^i , to MLE można czasem wyznaczyć analitycznie korzystając z warunku koniecznego optymalności, tzn. rozwiązując układ równań

$$\frac{\partial L(\mathbf{x}_1, \dots, \mathbf{x}_n; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = 0.$$

Jeśli MLE nie da się wyliczyć analitycznie, wyznacza się je przy użyciu algorytmów optymalizacji numerycznej. Estymatory MLE są asymptotycznie nieobciążone.

Definicja 2.22 (Estymatora MAP). Estymatorem MAP (z ang. *Maximum a Posteriori Estimator*, *MAP*) nazywamy funkcję $\hat{\boldsymbol{\theta}}$, która przy ustalonych wartościach obserwacji (realizacji prostej próby losowej) $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ maksymalizuje wartość iloczynu funkcji wiarygodności i priora nad wartościami parametrów lub, co równoważne, minimalizuje zregularyzowaną logarytmiczną funkcję wiarygodności tj.

$$\hat{\boldsymbol{\theta}}(\mathbf{x}_1, \dots, \mathbf{x}_n) = \arg \min_{\boldsymbol{\theta} \in \Theta} \left[- \sum_{i=1}^n \log p(\mathbf{x}_i \mid \boldsymbol{\theta}) - n \log p(\boldsymbol{\theta}) \right].$$

2.5 Liczby losowe w komputerze

Opiszemy jeszcze pokrótce metody generowania liczb pseudolosowych z dowolnych rozkładów prawdopodobieństwa w sposób algorytmiczny. Podstawowym narzędziem, którego będziemy potrzebować do generowania próbek z bardziej skomplikowanych rozkładów będzie prosty generator liczb z rozkładu jednostajnego $\mathcal{U}(0, 1)$. Moglibyśmy oczywiście wykorzystać jakieś fizyczne urządzenie lub proces, który generuje liczby prawdziwie losowe (np. detektor Geigera-Mullera, szum lamp elektronowych, ruletka), ale błędem byłaby rezygnacja z odtwarzalności. Poszukujemy zatem deterministycznej metody, która generuje sekwencje liczb, które są w przybliżeniu losowe. Podstawową metodą do algorytmicznego generowania liczb pseudolosowych jest tzw. *liniowy generator kongruentny* (z ang. *Linear Congruential Generator*, *LCG*), który jest opisany zależnością rekurencyjną

$$I_{j+1} = (aI_j + c) \mod m,$$

gdzie a, c, m to pewne ustalone dodatnie liczby całkowite, a I_0 to tzw. ziarno (z ang. *seed*). LCG generuje liczby całkowite, więc w dużym uproszczeniu liczby zmiennoprzecinkowe z rozkładu $\mathcal{U}(0, 1)$ otrzymujemy jako I_j/m (trzeba tutaj jednak uwzględnić problemy wynikające z arytmetyki zmiennoprzecinkowej).

Mając już generator liczb z rozkładu jednostajnego $\mathcal{U}(0, 1)$ i znając jawny wzór na dystrybuantę $F(x)$ innego rozkładu jednowymiarowego \mathcal{D} możemy generować liczby z tego rozkładu korzystając z tzw. *metody odwrotnej dystrybuanty*. Istotnie, jeśli $U \sim \mathcal{U}(0, 1)$, to $F^{-1}(U) \sim \mathcal{D}$. Istotnie

$$\Pr(F^{-1}(U) \leq x) = \Pr(U \leq F(x)) = F(x).$$

Metoda ta ma jedną zasadniczą wadę – musimy znać jawny wzór na dystrybuantę $F(x)$. W przypadku np. tak ważnych rozkładów, jak rozkład normalny dystrybuanta nie jest funkcją elementarną i metoda ta nie jest najlepsza. W przypadku rozkładu normalnego znacznie lepszą metodą jest tzw. *metoda Boxa-Mullera*. Weźmy rozkład łączny dwóch niezależnych zmiennych losowych X, Y pochodzących ze standardowego rozkładu normalnego

$$p(x, y) = \frac{1}{2\pi} \exp\left(-\frac{1}{2}(x^2 + y^2)\right).$$

Skorzystamy ze wzoru na transformację zmiennych losowych. Istotnie niech

$$X = \sqrt{Z} \cos \Phi, \quad Y = \sqrt{Z} \sin \Phi,$$

dla $0 < Z$ oraz $0 \leq \Phi < 2\pi$, wówczas

$$q(z, \phi) = p(x(z, \phi), y(z, \phi)) \begin{vmatrix} \frac{1}{2\sqrt{z}} \cos \phi & \frac{1}{2\sqrt{z}} \sin \phi \\ -\sqrt{z} \sin \phi & \sqrt{z} \cos \phi \end{vmatrix} = \frac{1}{4\pi} e^{-\frac{z}{2}}.$$

Zauważmy, że otrzymaliśmy sferycznie symetryczny rozkład wykładniczy. Możemy zatem wylosować z rozkładu jednostajnego kąt ϕ oraz z rozkładu wykładniczego wartość z korzystając z metody odwrotnej dystrybuanty. Wówczas wartości $x = \sqrt{z} \cos \phi$, $y = \sqrt{z} \sin \phi$ będą pochodzić ze standardowego rozkładu normalnego. Aby wygenerować próbki z ogólnego wielowymiarowego rozkładu normalnego, korzystamy z definicji, tj. najpierw generujemy n próbek ze standardowego rozkładu normalnego, a następnie korzystamy z przekształcenia afinicznego $\mathbf{x} = \mathbf{A}\mathbf{z} + \boldsymbol{\mu}$, gdzie $\boldsymbol{\Sigma} = \mathbf{A}\mathbf{A}^T$.

2.6 Monte Carlo

Umiemy już generować próbki z wielowymiarowego rozkładu normalnego. Chcemy teraz poznać metodę, która umożliwi generowanie próbek ze skomplikowanych, wielowymiarowych rozkładów prawdopodobieństwa, których gęstość znamy jedynie z dokładnością do stałej normalizującej, tj. znamy jedynie $\tilde{p}(\mathbf{x}) = Zp(\mathbf{x})$. Ograniczenie to wynika z chęci próbkowania z posteriora $p(\mathbf{x} \mid \mathbf{y})$ w sytuacji, gdy znamy jedynie rozkład łączny $p(\mathbf{x}, \mathbf{y}) = p(\mathbf{y} \mid \mathbf{x})p_{\mathbf{X}}(\mathbf{x})$. Okazuje się, iż znajomość rozkładu jedynie z dokładnością do stałej normalizującej jest wystarczająca do generowania próbek z tego rozkładu. Generowanie próbek z kolei wystarcza natomiast, na mocy silnego prawa wielkich liczb, do szacowania wartości średnich dowolnych funkcji zmiennej \mathbf{x} . Przypomnijmy, iż na mocy silnego prawa wielkich liczb ciąg średnich częściowych $(\bar{\mathbf{X}}_n)$ ciągu zmiennych losowych (\mathbf{X}_n) i.i.d. z rozkładu $\mathbf{X} \sim \mathcal{D}$ jest zbieżny z prawdopodobieństwem 1 do wartości oczekiwanej $\mathbb{E}[\mathbf{X}]$ tj.

$$\Pr \left(\lim_{n \rightarrow \infty} \bar{\mathbf{X}}_n = \mathbb{E}[\mathbf{X}] \right) = 1.$$

Wartość oczekiwaną $\mathbb{E}[\mathbf{X}]$ możemy zatem przybliżyć średnią $\bar{\mathbf{X}}_n$ z dużej ilości próbek.

Pozostaje pytanie w jaki sposób generować próbki ze skomplikowanych rozkładów prawdopodobieństwa, których gęstości znamy jedynie z dokładnością do stałej normalizującej. Poniżej przedstawimy dwa algorytmy próbkowania: algorytm IS oraz Metropolisa–Hastingsa będący szczególną realizacją całej rodziny algorytmów próbkowania zwanych Markov Chain Monte Carlo (MCMC).

2.6.1 Algorytm Importance Sampling

Załóżmy, iż chcemy obliczyć wartość oczekiwaną pewnej funkcji zmiennej losowej \mathbf{x} względem skomplikowanego rozkładu prawdopodobieństwa $p(\mathbf{x})$, który znamy jedynie z dokładnością do stałej normalizującej

$$p(\mathbf{x}) = \frac{1}{Z_p} \tilde{p}(\mathbf{x})$$

tj. szukamy

$$\mathbb{E}_p[f(\mathbf{x})] = \int f(\mathbf{x}) p(\mathbf{x}) d^n \mathbf{x}.$$

Jeśli umiemy generować próbki \mathbf{x} z innego (prostsze) rozkładu $q(\mathbf{x})$ (np. wielowymiarowego rozkładu normalnego), który nazywamy rozkładem proponującym kandydatów (z ang. *proposal distribution*) to możemy zapisać

$$\begin{aligned} \mathbb{E}_p[f(\mathbf{x})] &= \int_{\mathbb{R}^n} f(\mathbf{x}) p(\mathbf{x}) d^n \mathbf{x} = \int_{\mathbb{R}^n} f(\mathbf{x}) \frac{p(\mathbf{x})}{q(\mathbf{x})} q(\mathbf{x}) d^n \mathbf{x} \\ &= \mathbb{E}_q \left[f(\mathbf{x}) \frac{p(\mathbf{x})}{q(\mathbf{x})} \right] = \frac{Z_q}{Z_p} \mathbb{E}_q \left[f(\mathbf{x}) \frac{\tilde{p}(\mathbf{x})}{\tilde{q}(\mathbf{x})} \right]. \end{aligned}$$

Zakładamy tutaj, iż nośnik rozkładu p zawiera się w nośniku q tj. $\text{supp } p \subseteq \text{supp } q$. Stosunek stałych Z_p/Z_q również możemy oszacować z próbek z q , gdyż mamy

$$Z_p = \int_{\mathbb{R}^n} \tilde{p}(\mathbf{x}) d^n \mathbf{x} = Z_q \int_{\mathbb{R}^n} \frac{\tilde{p}(\mathbf{x})}{\tilde{q}(\mathbf{x})} q(\mathbf{x}) d^n \mathbf{x} = Z_q \mathbb{E}_q \left[\frac{\tilde{p}(\mathbf{x})}{\tilde{q}(\mathbf{x})} \right],$$

skąd ostatecznie

$$\mathbb{E}_p[f(\mathbf{x})] = \frac{\mathbb{E}_q \left[f(\mathbf{x}) \frac{\tilde{p}(\mathbf{x})}{\tilde{q}(\mathbf{x})} \right]}{\mathbb{E}_q \left[\frac{\tilde{p}(\mathbf{x})}{\tilde{q}(\mathbf{x})} \right]}.$$

Jeśli z rozkładu q wygenerowaliśmy próbki $X = \{\mathbf{x}_1, \dots, \mathbf{x}_m\}$ to na mocy silnego prawa wielkich liczb mamy

$$\mathbb{E}_p[f(\mathbf{x})] \approx \frac{\sum_{i=1}^m f(\mathbf{x}_i) \frac{\tilde{p}(\mathbf{x}_i)}{\tilde{q}(\mathbf{x}_i)}}{\sum_{j=1}^m \frac{\tilde{p}(\mathbf{x}_j)}{\tilde{q}(\mathbf{x}_j)}} = \sum_{i=1}^m \lambda_i f(\mathbf{x}_i),$$

gdzie

$$\lambda_i = \frac{\tilde{p}(\mathbf{x}_i)/\tilde{q}(\mathbf{x}_i)}{\sum_{j=1}^m \tilde{p}(\mathbf{x}_j)/\tilde{q}(\mathbf{x}_j)}.$$

Algorytm Importance Sampling jest prostym algorytmem Monte Carlo, który ma jeden zasadniczy problem. W jaki sposób mamy wybrać rozkład proponujący kandydatów q ? Pewną odpowiedź na to pytanie sugeruje analiza wariancji statystyki

$$\bar{f}_m(\mathbf{x}_1, \dots, \mathbf{x}_m) = \frac{1}{m} \sum_{i=1}^m \frac{f(\mathbf{x}_i)p(\mathbf{x}_i)}{q(\mathbf{x}_i)}$$

dla $\mathbf{x}_i \sim q$ mamy

$$\mathbb{V}[\bar{f}_m] = \frac{1}{m} \mathbb{V}_q \left[f(\mathbf{x}) \frac{p(\mathbf{x})}{q(\mathbf{x})} \right] = \frac{1}{m} \int_{\mathbb{R}^n} \frac{(f(\mathbf{x})p(\mathbf{x}) - \mu_f q(\mathbf{x}))^2}{q(\mathbf{x})} d^n \mathbf{x} .$$

Chcemy oczywiście, aby wariancja była jak najmniejsza, gdyż wówczas mała liczba próbek da dobre przybliżenie wartości oczekiwanej. Rozkład proponujący kandydatów powinien być zatem proporcjonalny do $f(\mathbf{x})p(\mathbf{x})$, co może być trudne do praktycznego zrealizowania.

2.6.2 Algorytm Metropolisa–Hastingsa

Cała klasa algorytmów próbkowania MCMC opiera się na idei wyrażenia generowania próbek jako ewolucji pewnego łańcucha Markowa.

Definicja 2.23 (Łańcucha Markowa). Łańcuchem Markowa nazwiemy ciąg zmiennych losowych (\mathbf{X}_t) o wartościach w \mathbb{R}^n taki, że spełnione jest *kryterium Markowa*

$$\begin{aligned} \forall A \subset \mathbb{R}^n : \Pr(\mathbf{X}_t \in A \mid \mathbf{X}_{t-1} = \mathbf{x}_{t-1}, \dots, \mathbf{X}_0 = \mathbf{x}_0) \\ = \Pr(\mathbf{X}_t \in A \mid \mathbf{X}_{t-1} = \mathbf{x}_{t-1}) . \end{aligned}$$

Elementy ciągu nazywamy stanami łańcucha.

Dany łańcuch jest zadany jednoznacznie przez podanie gęstości prawdopodobieństwa przejścia łańcucha ze stanu $\mathbf{x} \rightarrow \mathbf{y}$, którą będziemy oznaczać przez $\pi(\mathbf{y} \mid \mathbf{x})$ (zakładamy, iż prawdopodobieństwo przejścia jest niezależne od chwili t – łańcuch taki nazywamy jednorodnym). Funkcja π spełnia oczywiście warunek unormowania

$$\int_{\mathbb{R}^n} \pi(\mathbf{y} \mid \mathbf{x}) d^n \mathbf{y} = 1 ,$$

istotnie prawdopodobieństwo przejścia gdziekolwiek ze stanu \mathbf{x} jest równe 1. Będziemy zakładać dodatkowo, iż $\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n : \pi(\mathbf{y} | \mathbf{x}) > 0$. Rozkład $p(\mathbf{x})$ łańcucha Markowa (tj. rozkład prawdopodobieństwa z którego losujemy stan łańcucha w danej chwili t) z daną funkcją przejścia π nazwiemy rozkładem stacjonarnym tego łańcucha jeśli

$$p(\mathbf{y}) = \int_{\mathbb{R}^n} \pi(\mathbf{y} | \mathbf{x}) p(\mathbf{x}) d^n \mathbf{x} .$$

Rozkład stacjonarny danego łańcucha oznaczmy przez $p^*(\mathbf{x})$. Zauważmy, iż jeśli stan początkowy łańcucha \mathbf{X}_0 pochodzi z rozkładu stacjonarnego p^* to każdy kolejny stan \mathbf{X}_t również pochodzi z rozkładu stacjonarnego. Jeśli z kolei stan początkowy pochodzi z jakiegoś innego rozkładu p_0 to rozkład łańcucha w chwili t jest dany przez relację rekurencyjną

$$p_t(\mathbf{y}) = \int_{\mathbb{R}^n} \pi(\mathbf{y} | \mathbf{x}) p_{t-1}(\mathbf{x}) d^n \mathbf{x} , \quad \text{dla } t > 1 .$$

Rozkładem granicznym łańcucha Markowa nazwiemy granicę w sensie zbieżności punktowej

$$\lim_{t \rightarrow \infty} p_t(\mathbf{x}) .$$

Przy podanych wyżej założeniach istnieje twierdzenie, które mówi iż taki łańcuch Markowa posiada jednoznaczny rozkład stacjonarny tożsamy z rozkładem granicznym. Ponadto warunkiem wystarczającym, aby dany rozkład $p(\mathbf{x})$ był rozkładem stacjonarnym łańcucha Markowa jest

$$\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n : \pi(\mathbf{y} | \mathbf{x}) p(\mathbf{x}) = \pi(\mathbf{x} | \mathbf{y}) p(\mathbf{y}) ,$$

co wynika z scałkowania powyższego równania

$$\int_{\mathbb{R}^n} \pi(\mathbf{y} | \mathbf{x}) p(\mathbf{x}) d^n \mathbf{x} = \int_{\mathbb{R}^n} \pi(\mathbf{x} | \mathbf{y}) p(\mathbf{y}) d^n \mathbf{x} = p(\mathbf{y}) \int_{\mathbb{R}^n} \pi(\mathbf{x} | \mathbf{y}) d^n \mathbf{x} = p(\mathbf{y}) .$$

Kryterium to nazywamy *kryterium lokalnego balansu* (z ang. *detailed balance condition*).

Podstawowa idea wykorzystania łańcuchów Markowa do generowania próbek ze skomplikowanego rozkładu p jest więc następująca: tworzymy łańcuch Markowa, dla którego p jest rozkładem stacjonarnym, wówczas rozpoczynając w dowolnym dopuszczalnym stanie początkowym \mathbf{X}_0 po wykonaniu dużej liczby kroków (etap ten nazywamy okresem przejściowym z

ang. *burn-in period*) stan \mathbf{X}_t (dla $t \gg 1$) tego łańcucha będzie w przybliżeniu pochodził z rozkładu granicznego p (nie jest jednak prosto stwierdzić po jak długim okresie przejściowym przybliżenie to jest wystarczająco dobre). Aby otrzymać z takiej procedury próbki prawdziwie i.i.d. każda z próbek musiałaby pochodzić z ponownego uruchomienia takiego łańcucha. Oczywiście jest to nieefektywne, więc w praktyce generujemy próbki z jednego łańcucha po prostu odrzucając pewne z nich tak aby uniknąć znaczących korelacji. Pozostaje pytanie jak skonstruować funkcję przejścia $\pi(\mathbf{y} \mid \mathbf{x})$ dla danego rozkładu granicznego $p(\mathbf{x})$. Podstawową konstrukcję podaje algorytm Metropolisa–Hastingsa.

```

1 # Metropolis-Hastings
2 # -----
3 # Initialize state  $\mathbf{x}$  using any admissible value.
4 for i in range(max_iter):
5     # Step 1. Sample candidate  $\mathbf{y}$  from the conditional
6     # proposal distribution  $q(\mathbf{y} \mid \mathbf{x})$ .
7
8     # Step 2. Accept new candidate with probability
9      $r(\mathbf{y} \mid \mathbf{x}) = \min \left\{ 1, \frac{p(\mathbf{y})q(\mathbf{x} \mid \mathbf{y})}{p(\mathbf{x})q(\mathbf{y} \mid \mathbf{x})} \right\}$ 
10    # and change state to  $\mathbf{y}$ , otherwise stay in state  $\mathbf{x}$ .

```

Funkcja przejścia ma zatem postać

$$\pi_{\text{MH}}(\mathbf{y} \mid \mathbf{x}) = q(\mathbf{y} \mid \mathbf{x})r(\mathbf{y} \mid \mathbf{x}).$$

Pozostaje tylko wykazać, iż spełnione jest kryterium lokalnego balansu. Istotnie mamy

$$\begin{aligned} \pi_{\text{MH}}(\mathbf{y} \mid \mathbf{x})p(\mathbf{x}) &= \min \{q(\mathbf{y} \mid \mathbf{x})p(\mathbf{x}), q(\mathbf{x} \mid \mathbf{y})p(\mathbf{y})\} \\ \pi_{\text{MH}}(\mathbf{x} \mid \mathbf{y})p(\mathbf{y}) &= \min \{q(\mathbf{x} \mid \mathbf{y})p(\mathbf{y}), q(\mathbf{y} \mid \mathbf{x})p(\mathbf{x})\} \end{aligned} \quad ,$$

skąd $\pi_{\text{MH}}(\mathbf{y} \mid \mathbf{x})p(\mathbf{x}) = \pi_{\text{MH}}(\mathbf{x} \mid \mathbf{y})p(\mathbf{y})$. Zauważmy, iż nie musimy znać $p(\mathbf{x})$ z dokładnością do stałej normalizującej, gdyż

$$\frac{p(\mathbf{y})}{p(\mathbf{x})} = \frac{\tilde{p}(\mathbf{y})/Z_p}{\tilde{p}(\mathbf{x})/Z_p} = \frac{\tilde{p}(\mathbf{y})}{\tilde{p}(\mathbf{x})}.$$

Poza algorytmem Metropolisa–Hastingsa jest wiele innych algorytmów z rodziny MCMC. Większość z nich implementuje konkretny sposób generowania (zostawiając resztę struktury) tak, aby zmniejszyć korelację po okresie przejściowym i przyspieszyć zbieżność. Standardowo wykorzystywanymi algorytmami z tej klasy są algorytmy HMC (*Hamiltonian Monte Carlo*) oraz NUTS (*No U-Turn Sampler*).

2.7 Estymator jądrowy gęstości

W poprzednim paragrafie opisaliśmy w jaki sposób mając (nieznormalizowaną) funkcję gęstości prawdopodobieństwa $\tilde{p}(\mathbf{x})$ generować algorytmicznie próbki z opisanego przez nią rozkładu. Teraz zajmijmy się problemem odwrotnym tj. mając realizację prostej próby losowej z pewnego rozkładu chcemy znaleźć funkcję $\hat{p}(\mathbf{x})$, która estymuje gęstość rozkładu prawdopodobieństwa, z którego pochodzą próbki. Opiszemy tutaj jedną z najprostszych metod zwaną estymatorem jądrowym (z ang. *Kernel Density Estimator*, *KDE*). Estymatorem jądrowym gęstości funkcji p nazywamy funkcję

$$\hat{p}(\mathbf{x}) := \frac{1}{n} \sum_{i=1}^n K_{\mathbf{H}}(\mathbf{x} - \mathbf{x}_i),$$

gdzie \mathbf{H} jest symetryczną i dodatnio określoną macierzą zwaną *bandwidth matrix*, funkcja $K_{\mathbf{H}}$ ma postać

$$K_{\mathbf{H}}(\mathbf{x}) = (\det \mathbf{H})^{-1/2} K(\mathbf{H}^{-1/2} \mathbf{x}),$$

gdzie funkcja K zwana jądrem (z ang. *kernel*) jest gęstością prawdopodobieństwa pewnego sferycznie symetrycznego rozkładu wielowymiarowego. Wybór funkcji K nie jest kluczowy ze statystycznego punktu widzenia, więc możemy bez problemu założyć, iż jest to gęstość wielowymiarowego rozkładu normalnego, tj.

$$K_{\mathbf{H}}(\mathbf{x}) = \frac{1}{\sqrt{(2\pi)^m \det \mathbf{H}}} \exp\left(-\frac{1}{2} \mathbf{x}^T \mathbf{H}^{-1} \mathbf{x}\right),$$

gdzie zakładamy $\mathbf{x} \in \mathbb{R}^m$. Większym problemem w przypadku KDE jest wybór odpowiedniego parametru \mathbf{H} . Jednym z prostszych wyborów w przypadku estymatora

$$\boxed{\hat{p}(\mathbf{x}) = \frac{1}{n \sqrt{(2\pi)^m \det \mathbf{H}}} \sum_{i=1}^n \exp\left(-\frac{1}{2} (\mathbf{x} - \mathbf{x}_i)^T \mathbf{H}^{-1} (\mathbf{x} - \mathbf{x}_i)\right)}$$

jest tzw. *reguła Silvermana*, która podaje następujący przepis na macierz \mathbf{H}

$$\mathbf{H}^{ij} = 0, \quad \sqrt{\mathbf{H}^{ii}} = \left(\frac{4}{4+m}\right)^{\frac{1}{m+4}} n^{\frac{-1}{m+4}} \sigma_i,$$

gdzie σ_i jest estymatorem wariancji i -tej współrzędnej zmiennej \mathbf{X} . Inną możliwością jest tzw. *reguła Scotta*

$$\mathbf{H}^{ij} = 0, \quad \sqrt{\mathbf{H}^{ii}} = n^{\frac{-1}{m+4}} \sigma_i.$$

KDE w praktycznych zastosowaniach często przyspiesza się za pomocą odpowiednich struktur danych do wyszukiwania najbliższych sąsiadów w przestrzeni \mathbb{R}^m , tj. zamiast sumować przyczynki od wszystkich punktów \mathbf{x}_i dla danego \mathbf{x} , znajdujemy jego k najbliższych sąsiadów \mathbf{x} ze zbioru $\{\mathbf{x}_i\}_{i=1}^n$ stosując np. ANN i obliczamy przyczynki do $\hat{p}(\mathbf{x})$ tylko od nich.

3 Podstawy statystycznego uczenia maszynowego

Przechodzimy teraz do zagadnień uczenia maszynowego, w których wykorzystamy przedstawioną wcześniej teorię rachunku prawdopodobieństwa (w szczególności teorię zmiennych losowych) oraz wnioskowania statystycznego.

3.1 Regresja liniowa

Rozpatrujemy teraz zagadnienie regresji tj. predykcji ciągłej wartości $y \in \mathbb{R}$ w zależności od wektora cech \mathbf{x} . Zakładamy ponadto prosty model liniowy, tj. nasze predykcje będą miały postać

$$\phi(\mathbf{x}; \mathbf{w}) = \mathbf{w}^T \mathbf{x}$$

dla pewnych estymowanych parametrów \mathbf{w} . Zauważmy, iż taka postać modelu ϕ uwzględnia człon stały poprzez transformację $\mathbf{x} \leftarrow [\mathbf{x}^T \ 1]^T$. W przypadku regresji liniowej parametry estymujemy za pomocą metody MLE dla następującego modelu statystycznego

$$y(\mathbf{x}) \mid \mathbf{w} \sim \mathcal{N}(\phi(\mathbf{x}; \mathbf{w}), \sigma^2).$$

Niech $\mathcal{X} = \{y_i(\mathbf{x}_i)\}_{i=1}^n$ będzie naszym zbiorem obserwacji i.i.d. Wiarygodność ma zatem postać

$$\mathcal{L}(\mathcal{X}; \mathbf{w}) = \frac{1}{(2\pi\sigma^2)^{n/2}} \prod_{i=1}^n \exp\left(-\frac{1}{2\sigma^2} [y_i - \phi(\mathbf{x}_i; \mathbf{w})]^2\right).$$

Zanegowana logarytmiczna funkcja wiarygodności (funkcjonał kosztu) ma zatem postać

$$L(\mathcal{X}; \mathbf{w}) = \frac{1}{2} \sum_{i=1}^n [y_i - \phi(\mathbf{x}_i; \mathbf{w})]^2,$$

gdzie pominęliśmy multiplikatywne i addytywne człony stałe, gdyż nie wpływają one na zagadnienie optymalizacji. Minimalizując funkcję L względem \mathbf{w} otrzymamy estymatę MLE tych parametrów. Otrzymana funkcja L ma postać formy kwadratowej, a otrzymany problem optymalizacyjny nazywamy metodą najmniejszych kwadratów (z ang. *Ordinary Least Squares, OLS*). Wprowadźmy macierz $\mathbf{X}^{ij} := \mathbf{x}_i^j$ oraz wektory $\mathbf{y}^i = y_i$, $\phi^i = \phi(\mathbf{x}_i; \mathbf{w}) = \sum_j \mathbf{X}^{ij} \mathbf{w}^j$. Wówczas możemy zapisać funkcję kosztu jako

$$L(\mathcal{X}; \mathbf{w}) = \frac{1}{2} \sum_{i=1}^n (\mathbf{y}^i - \phi^i(\mathbf{X}; \mathbf{w}))^2,$$

skąd

$$\frac{\partial L}{\partial \mathbf{w}^a} = \sum_b \frac{\partial L}{\partial \phi^b} \frac{\partial \phi^b}{\partial \mathbf{w}^a}.$$

Jednocześnie

$$\frac{\partial L}{\partial \phi^b} = \mathbf{y}^b - \phi^b, \quad \frac{\partial \phi^b}{\partial \mathbf{w}^a} = \mathbf{X}^{ba},$$

skąd

$$\frac{\partial L}{\partial \mathbf{w}^a} = \sum_b \mathbf{X}^{ba} (\mathbf{y}^b - \phi^b),$$

co możemy zapisać w kompaktowej formie macierzowej

$$\frac{\partial L}{\partial \mathbf{w}} = \mathbf{X}^T (\mathbf{y} - \phi) = \mathbf{X}^T (\mathbf{y} - \mathbf{X} \mathbf{w}).$$

Przyrównując powyższe równanie do 0 otrzymujemy następujący wzór na estymatę MLE parametrów \mathbf{w}

$$\mathbf{w}_{\text{MLE}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} = \mathbf{X}^+ \mathbf{y},$$

gdzie \mathbf{X}^+ oznacza *pseudoodwrotność Moore'a-Penrose'a*, którą można efektywnie obliczyć korzystając z rozkładu SVD macierzy \mathbf{X} .

3.2 Regularyzacja

Regularyzacją nazywamy proces polegający na wprowadzeniu ad hoc do zagadnienia optymalizacji dodatkowych członów tak, aby rozwiązanie było regularne (prostsze, nieosobliwe, jednoznaczne). W przypadku funkcji kosztu L najczęściej dodajemy człon penalizujący rozwiązania o dużej normie estymowanego parametru tj. człon postaci $\lambda \|\mathbf{w}\|$ dla pewnej normy $\|\cdot\|$ i

hiperparametru λ zwanego *siłą regularyzacji*. W kontekście bayesowskim regularyzację można również rozumieć jako użycie estymacji MAP zamiast MLE dla odpowiednio dobranego priora nad estymowanymi parametrami.

Jeśli dla modelu regresji liniowej jako człon regularyzujący przyjmiemy $\lambda \|\mathbf{w}\|_2^2$, tj. kwadrat normy L2 wektora estymowanych parametrów, to otrzymana funkcja kosztu ma postać

$$L(\mathcal{X}; \mathbf{w}) = \frac{1}{2} \sum_{i=1}^n [y_i - \phi(\mathbf{x}_i; \mathbf{w})]^2 + \lambda \|\mathbf{w}\|_2^2.$$

Analogiczną postać można otrzymać rozważając estymatę MAP dla identycznego modelu statystycznego z priorem na wartości parametrów \mathbf{w} będącym rozkładem normalnym $\mathcal{N}(0, \tau^2 \mathbf{1})$ dla odpowiednio dobranego τ . Dla powyższej funkcji kosztu możemy bez problemu wyznaczyć wartość \mathbf{w} , która ją minimalizuje. Istotnie korzystając z wyników dla modelu regresji liniowej mamy

$$\mathbf{w}_{\text{MAP}} = (\lambda \mathbf{1} + \mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y},$$

Zagadnienie minimalizacji funkcji kosztu będącej formą kwadratową z dodanym członem regularyzującym w postaci normy L2 wektora nazywamy *regresją grzbietową* (z ang. *ridge regression*).

Innym przykładem regularyzacji jest tzw. regularyzacja L1, która polega na dodaniu do funkcji kosztu członu postaci $\lambda \sum_{j=1}^d |\mathbf{w}^j|$ tj. normy L1 wektora wag. Zagadnienie optymalizacji formy kwadratowej z członem regularyzującym L1 nazywamy *regresją LASSO*. W takim przypadku nie da się prosto analitycznie znaleźć estymaty punktowej MAP i trzeba używać algorytmów optymalizacji numerycznej. W ogólności można połączyć regularyzacje L1 i L2 tj. rozważać zregularyzowaną funkcję kosztu postaci

$$L(\mathcal{X}; \mathbf{w}) = \frac{1}{2} \sum_{i=1}^n [y_i - \phi(\mathbf{x}_i; \mathbf{w})]^2 + \lambda_1 \|\mathbf{w}\|_1 + \lambda_2 \|\mathbf{w}\|_2^2.$$

Zagadnienie minimalizacji takiej funkcji kosztu nazywamy ElasticNet i tak jak w przypadku LASSO musimy korzystać z algorytmów optymalizacji numerycznej. Często wykorzystuje się tutaj algorytmy bezgradientowe np. coordinate descent.

3.3 Regresja kwantylowa

Zastanówmy się wpierw na czym tak naprawdę polega modelowanie rozkładu warunkowego $y(\mathbf{x}) \mid \mathbf{w}$ za pomocą określonego rozkładu prawdopo-

dobieństwa. Na pierwszy rzut oka może się wydawać, iż takie podejście wprowadza bardzo silne założenia, a co za tym idzie ograniczenia w stosowaniu naszego modelu. Zauważmy jednak, iż w przypadku podejścia typu *likelihood* rozkład jest niejako wybierany w taki sposób, aby jego parametry były użyteczne. Istotnie w przypadku regresji zwykle nie ma jak zweryfikować rzeczywistego rozkładu $y(\mathbf{x})$, gdyż mamy tylko po jednej wartości y dla danego \mathbf{x} . Modelując $y(\mathbf{x})$ rozkładem normalnym chodzi nam zatem raczej o to, że w takim modelu chcemy znaleźć parametr (prostą) taką, że masa prawdopodobieństwa punktów po obu stronach prostej jest jednakowa i większość masy jest zgromadzona w bliskiej odległości od prostej (wynika to z kształtu rozkładu normalnego, który nie posiada tzw. ciężkich ogonów)

W takim ujęciu możemy zakładać różne inne rozkłady na $y(\mathbf{x})$ jeśli interesują nas proste, które inaczej mają rozdzielać masę prawdopodobieństwa między punkty lub też dopuszczać obserwacje leżące daleko od prostej regresji. Problemem w przypadku rozkładu normalnego jest jego czułość na wartości odstające, gdyż w rozkładzie normalnym ogony tego rozkładu mają stosunkowo niewielką masę prawdopodobieństwa. Chcielibyśmy zatem rozkład z ciężkimi ogonami (z ang. *heavy tails*). Dodatkowo chcielibyśmy mieć rozkład, który pozwala znaleźć prostą, która nie rozkłada masy po równo, ale np. tak, że 90% masy prawdopodobieństwa jest pod nią. Oba te problemy możemy rozwiązać modelując rozkład $y(\mathbf{x})$ przez tzw. *asymetryczny rozkład Laplace'a* (z ang. *Asymmetric Laplace Distribution, ALD*).

Definicja 3.1 (Asymetrycznego rozkładu Laplace'a). Mówimy, iż zmienna losowa rzeczywista X ma asymetryczny rozkład Laplace'a, tzn. $X \sim \text{ALD}(m, \lambda, q)$ jeśli jej gęstość wyraża się wzorem

$$p(x; m, \lambda, q) = \frac{q(1-q)}{\lambda} \begin{cases} e^{-\frac{q-1}{\lambda}(x-m)}, & x \leq m \\ e^{-\frac{q}{\lambda}(x-m)}, & x \geq m \end{cases}.$$

Zauważmy, że dystrybuanta rozkładu ALD ma postać

$$F(x; m, \lambda, q) = \begin{cases} qe^{\frac{1-q}{\lambda}(x-m)}, & x \leq m \\ 1 - (1-q)e^{-\frac{q}{\lambda}(x-m)}, & x \geq m \end{cases}$$

zatem parametr q określa rząd kwantyla m . W przypadku regresji możemy zatem modelować wartość $y(\mathbf{x})$ przez rozkład ALD dla ustalonego q postaci

$$y(\mathbf{x}) \mid \mathbf{w}, \lambda \sim \text{ALD}(\phi(\mathbf{x}; \mathbf{w}), \lambda, q),$$

gdzie λ pełni podobną rolę jak σ w przypadku rozkładu normalnego. Wiadomo wówczas, iż estymacja MLE \mathbf{w} daje prostą regresję taką, że ułamek $1 - q$ masy prawdopodobieństwa znajduje się pod prostą (estymujemy zatem warunkowy kwantyl rzędu q). Zanegowana logarytmiczna funkcja wiarygodności (inaczej funkcja kosztu) dla modelu ALD ma postać

$$L(\mathcal{X}; \mathbf{w}, \lambda) = \sum_{i=1}^n [(q-1)r_i\theta(-r_i) + qr_i\theta(r_i)] ,$$

gdzie

$$r_i := y_i - \phi(\mathbf{x}_i; \mathbf{w})$$

są tzw. *rezyduami*, a θ oznacza funkcję skokową Heaviside'a. Taką funkcję kosztu nazywamy *pinball loss*, a otrzymane zagadnienie optymalizacji – regresją kwantylową. Modelowanie $y(\mathbf{x})$ za pomocą rozkładu ALD pozwala nam w prosty i „robust” sposób znaleźć niepewność naszych predykcji punktowych, tj. dla danego zagadnienia dopasowujemy trzy modele oparte na ALD dla $q = 0.5$ (estymacja punktowa, mediana, odporna na outliery) oraz np. $q = 0.1$ i $q = 0.9$ będące oszacowaniem niepewności punktowej estymaty. Minimalizację powyższej funkcji kosztu można przedstawić jako zagadnienie programowania liniowego i efektywnie rozwiązać korzystając z odpowiednich solverów. Można również dodać człon regularyzujący w postaci normy L1 wektora estymowanych parametrów.

3.4 Procesy gaussowskie

Jak już wspomnieliśmy macierz kowariancji n -wymiarowej zmiennej losowej \mathbf{x} o wartości oczekiwanej $\boldsymbol{\mu}$ jest zdefiniowana jako

$$\boldsymbol{\Sigma} = \mathbb{E} [(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T] .$$

Wiemy również, iż macierz ta jest nieujemnie określona. Pokażemy teraz, iż dla każdej nieujemnie określonej macierzy symetrycznej \mathbf{K} wymiaru $n \times n$ istnieje n -wymiarowa zmienna losowa o wielowymiarowym rozkładzie normalnym, dla której \mathbf{K} jest macierzą kowariancji. Istotnie dla każdej nieujemnie określonej macierzy symetrycznej istnieje macierz \mathbf{L} taka, że

$$\mathbf{K} = \mathbf{L}\mathbf{L}^T ,$$

jest to tzw. *dekompozycja Choleskiego*. Niech $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{1})$, wówczas zmienna losowa $\mathbf{L}\mathbf{z}$ ma rozkład o zerowej wartości oczekiwanej i macierzy kowariancji

$$\mathbb{E} [(\mathbf{L}\mathbf{z})(\mathbf{L}\mathbf{z})^T] = \mathbb{E} [\mathbf{L}\mathbf{z}\mathbf{z}^T\mathbf{L}^T] = \mathbf{L}\mathbb{E}[\mathbf{z}\mathbf{z}^T]\mathbf{L}^T = \mathbf{L}\mathbf{L}^T = \mathbf{K} .$$

Powyższe własności wskazują, iż macierze kowariancji można w pewnym sensie utożsamiać z nieujemnie określonymi macierzami symetrycznymi.

Definicja 3.2 (Funkcji kowariancji). Funkcję $k : \mathbb{R}^n \times \mathbb{R}^n \mapsto \mathbb{R}$ taką, że $\forall m \in \mathbb{N} : \forall X = \{\mathbf{x}_1, \dots, \mathbf{x}_m\} \subset \mathbb{R}^n$ macierz

$$k(X, X) = \begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & k(\mathbf{x}_1, \mathbf{x}_2) & \cdots & k(\mathbf{x}_1, \mathbf{x}_m) \\ k(\mathbf{x}_2, \mathbf{x}_1) & k(\mathbf{x}_2, \mathbf{x}_2) & \cdots & k(\mathbf{x}_2, \mathbf{x}_m) \\ \vdots & \vdots & \ddots & \vdots \\ k(\mathbf{x}_m, \mathbf{x}_1) & k(\mathbf{x}_m, \mathbf{x}_2) & \cdots & k(\mathbf{x}_m, \mathbf{x}_m) \end{bmatrix}$$

jest dodatnio określoną macierzą symetryczną nazywamy funkcją kowariancji, jądrem dodatnio określonym (z ang. *positive definite kernel*) lub *jądrem Mercera*.

Dla dwóch zbiorów punktów $X = \{\mathbf{x}_1, \dots, \mathbf{x}_m\} \subset \mathbb{R}^n$ i $Y = \{\mathbf{y}_1, \dots, \mathbf{y}_s\} \subset \mathbb{R}^n$ i funkcji kowariancji k wprowadzimy oznaczenie

$$k(X, Y) := \begin{bmatrix} k(\mathbf{x}_1, \mathbf{y}_1) & k(\mathbf{x}_1, \mathbf{y}_2) & \cdots & k(\mathbf{x}_1, \mathbf{y}_s) \\ k(\mathbf{x}_2, \mathbf{y}_1) & k(\mathbf{x}_2, \mathbf{y}_2) & \cdots & k(\mathbf{x}_2, \mathbf{y}_s) \\ \vdots & \vdots & \ddots & \vdots \\ k(\mathbf{x}_m, \mathbf{y}_1) & k(\mathbf{x}_m, \mathbf{y}_2) & \cdots & k(\mathbf{x}_m, \mathbf{y}_s) \end{bmatrix}.$$

Poniżej podajemy kilka przykładów funkcji kowariancji

- *Gaussian kernel* dla normy $\|\cdot\|$ i hiper-parametrów a, l (amplituda i skala długości)

$$k(\mathbf{x}, \mathbf{y}) = a^2 \exp \left\{ -\frac{1}{2l^2} \|\mathbf{x} - \mathbf{y}\|^2 \right\}$$

- *Periodic kernel* dla normy $\|\cdot\|$ i hiper-parametrów a, l, p (amplituda, skala długości, okres zmienności)

$$k(\mathbf{x}, \mathbf{y}) = a^2 \exp \left\{ -\frac{2}{l^2} \sin^2 \left(\frac{\pi}{p} \|\mathbf{x} - \mathbf{y}\| \right) \right\}$$

- *White noise kernel* dla hiper-parametru σ

$$k(\mathbf{x}, \mathbf{y}) = \sigma^2 \delta_{\mathbf{x}, \mathbf{y}}$$

- *Matérn kernel* dla normy $\|\cdot\|$ i hiper-parametrów a, l, ν (amplituda, skala długości, regularność)

$$k(\mathbf{x}, \mathbf{y}) = a^2 \frac{2^{1-\nu}}{\Gamma(\nu)} \left(\frac{\sqrt{2\nu}}{l} \|\mathbf{x} - \mathbf{y}\| \right)^\nu K_\nu \left(\frac{\sqrt{2\nu}}{l} \|\mathbf{x} - \mathbf{y}\| \right),$$

gdzie $\Gamma(x)$ to funkcja gamma Eulera, a $K_\nu(x)$ to zmodyfikowana funkcja Bessela 2-go rodzaju rzędu ν .

Twierdzenie 3.1 (Własności funkcji kowariancji). Suma lub iloczyn dwóch funkcji kowariancji oraz złożenie funkcji kowariancji z wielomianem o nieujemnych współczynnikach jest również funkcją kowariancji.

Definicja 3.3 (Procesu gaussowskiego). Procesem Gaussowskim (z ang. *Gaussian Process*) nazywamy rodzinę skalnych zmiennych losowych indeksowanych przez punkty $\mathbf{x} \in \mathbb{R}^n$

$$\mathcal{GP} = \{f_{\mathbf{x}} \mid \mathbf{x} \in \mathbb{R}^n\}$$

taką że każdy skończony podzbiór \mathcal{GP} ma łącznie wielowymiarowy rozkład normalny tj. dla dowolnego zbioru $X = \{\mathbf{x}_1, \dots, \mathbf{x}_m\} \subset \mathbb{R}^n$ zachodzi

$$\begin{bmatrix} f_{\mathbf{x}_1} \\ \vdots \\ f_{\mathbf{x}_m} \end{bmatrix} \sim \mathcal{N}(\boldsymbol{\mu}_X, \boldsymbol{\Sigma}_X).$$

Zauważmy, iż proces Gaussowski możemy jednoznacznie zdefiniować podając przepisy na parametry $\boldsymbol{\mu}_X$ i $\boldsymbol{\Sigma}_X$ dla dowolnego zbioru X . W praktyce często przyjmujemy $\boldsymbol{\mu}_X = \mathbf{0}$, natomiast przepisem na macierz kowariancji może być zdefiniowana wyżej funkcja kowariancji $k(X, X)$ tj.

$$\begin{bmatrix} f_{\mathbf{x}_1} \\ \vdots \\ f_{\mathbf{x}_m} \end{bmatrix} \sim \mathcal{N}(\mathbf{0}, k(X, X)).$$

Process Gaussowski daje nam w praktyce rozkład prawdopodobieństwa nad funkcjami $f : \mathbb{R}^n \mapsto \mathbb{R}$, których charakter jest określony przez jądro k (np.

funkcja gładka dla jądra Gaussowskiego, okresowa dla jądra periodycznego, itp.). Zauważmy, że nie wnioskujemy tu o parametrach konkretnej rodziny funkcji (jak w przypadku regresji liniowej); interesuje nas jedynie rozkład predykcyjny. Załóżmy, iż w dokładnie znanych przez nas punktach $X = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m\}$ zaobserwowaliśmy wartości pewnej funkcji, o których zakładamy, iż pochodzą z procesu Gaussowskiego zadanego jądrem k , które wyraża nasze założenia a priori co do charakteru badanej funkcji

$$\mathbf{f}_X = \begin{bmatrix} f_{\mathbf{x}_1} \\ \vdots \\ f_{\mathbf{x}_m} \end{bmatrix} \sim \mathcal{N}(\mathbf{0}, k(X, X)) .$$

Powiedzmy, iż chcemy znać wartości \mathbf{f}_Y tej funkcji w zadanych punktach $Y = \{\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_s\}$. Ponieważ założyliśmy, iż wartości funkcji pochodzą z procesu Gaussowskiego, więc rozkład łączny \mathbf{f}_X i \mathbf{f}_Y jest rozkładem normalnym

$$\begin{bmatrix} \mathbf{f}_X \\ \mathbf{f}_Y \end{bmatrix} \sim \mathcal{N}\left(\mathbf{0}, \begin{bmatrix} k(X, X) & k(X, Y) \\ k(Y, X) & k(Y, Y) \end{bmatrix}\right) .$$

Zauważmy, iż z twierdzenia o własnościach niezdegenerowanego rozkładu normalnego wnioskujemy, iż warunkowy $\mathbf{f}_Y \mid \mathbf{f}_X$ jest również rozkładem normalnym o parametrach

$$\begin{aligned} \boldsymbol{\mu} &= k(Y, X)k^{-1}(X, X)\mathbf{f}_X \\ \boldsymbol{\Sigma} &= k(Y, Y) - k(Y, X)k^{-1}(X, X)k(X, Y) \end{aligned} .$$

Dodatkową niepewność związaną z pomiarem wartości \mathbf{f}_X możemy uchwycić zmieniając postać jądra

$$k(\mathbf{x}, \mathbf{y}) \leftarrow k(\mathbf{x}, \mathbf{y}) + \mathcal{I}_X(\mathbf{x})\sigma^2\delta_{\mathbf{x}, \mathbf{y}} ,$$

gdzie σ jest hiper-parametrem określającym precyzję pomiaru. Oczywiście k jest dalej funkcją kowariancji, gdyż takie podstawienie powoduje jedynie dodanie dodatnich członów do pewnych elementów diagonalnych macierzy kowariancji, więc macierz ta jest nadal symetryczna i dodatnio określona. Wówczas rozkład predykcyjny ma parametry

$$\boxed{\begin{aligned} \boldsymbol{\mu} &= k(Y, X) [k(X, X) + \sigma^2 \mathbf{1}]^{-1} \mathbf{f}_X \\ \boldsymbol{\Sigma} &= k(Y, Y) - k(Y, X) [k(X, X) + \sigma^2 \mathbf{1}]^{-1} k(X, Y) \end{aligned}} .$$

3.5 Naiwny klasyfikator bayesowski

Przechodzimy teraz do zagadnienia klasyfikacji, tj. predykcji zmiennej t pochodzącej ze skończonego zbioru kategorii (inaczej klas) dla danego wektora cech \mathbf{x} . Będziemy rozważać jedynie modele probabilistyczne, których wyjściem jest zbiór prawdopodobieństw $p_i(\mathbf{x})$ dla $i = 1, \dots, k$, gdzie k jest liczbą klas, a $p_i(\mathbf{x})$ określa prawdopodobieństwo, iż prawidłową klasą dla wektora \mathbf{x} jest i -ta klasa. Oczywiście musi zachodzić

$$\sum_j p_j(\mathbf{x}) = 1$$

dla dowolnego \mathbf{x} . Przewidywaną klasą zostaje następnie ta z największym prawdopodobieństwem. W przypadku klasyfikacji binarnej są natomiast możliwe inne reguły, np. przewidywaną klasą jest 0 jeśli $p_0 \geq p_{\text{threshold}}$ dla pewnego ustalonego progu niekoniecznie równego 0.5.

Pierwszym modelem, który rozważymy będzie bardzo prosty model bayesowski. Korzystając z twierdzenia Bayesa możemy zapisać

$$p_i(\mathbf{x}) = p(c_i | \mathbf{x}) = \frac{p(\mathbf{x} | c_i)p(c_i)}{\sum_j p(\mathbf{x} | c_j)p(c_j)}.$$

Jeśli $\mathcal{X} = \{t_i(\mathbf{x}_i)\}_{i=1}^n$ oznacza nasz zbiór obserwacji i.i.d. to gęstości prawdopodobieństwa $p(\mathbf{x} | c_i)$ i $p(c_i)$ możemy oszacować z danych \mathcal{X} . W przypadku drugiego członu możemy go prosto oszacować jako

$$p(c_i) = \frac{|\{t | t = c_i, t \in \mathcal{X}\}|}{|\mathcal{X}|}.$$

W przypadku pierwszego członu sytuacja jest bardziej skomplikowana. Możemy założyć pewien rozkład parametryczny dla rozkładu warunkowego $p(\mathbf{x} | c_i)$ dla każdej klasy i następnie oszacować parametry tego rozkładu z danych. Możemy również wykorzystać KDE i estymować rozkłady w sposób nieparametryczny. Oba te podejścia mają jednak problem związany z tzw. *klątwą wymiaru* (z ang. *curse of dimensionality*). Istotnie przestrzenie o dużej liczbie wymiarów są bardzo „puste”, tj. aby wyczerpująco wypróbować taką przestrzeń potrzebujemy wykładniczo wielu punktów. Jedną z możliwości poradzenia sobie z klątwą jest wprowadzenie naiwnego założenia o warunkowej niezależności cech w wektorze $\mathbf{x} = [x_1, \dots, x_d]^T$. Wówczas gęstości prawdopodobieństw $p(\mathbf{x} | c_i)$ można zapisać jako

$$p(\mathbf{x} | c_i) = p(x_1, \dots, x_d | c_i) = \prod_{j=1}^d p(x_j | c_i).$$

Teraz estymacja rozkładów warunkowych $p(x_j | c_i)$ dla każdej z klasy jest o wiele prostsza i możemy wykorzystać jedno ze wspomnianych podejść, tj. albo estymację parametrów pewnego rozkładu, albo KDE. Przewidywaną klasą dla wektora cech zostaje więc

$$t(\mathbf{x}) = \arg \max_{c_i \in \{c_1, \dots, c_k\}} p_i(\mathbf{x}) = \arg \max_{c_i \in \{c_1, \dots, c_k\}} p(c_i) \prod_{j=1}^d p(x_j | c_i).$$

3.6 Klasyfikator najbliższych sąsiadów

Klasyfikator najbliższych sąsiadów (z ang. *k Nearest Neighbors*) jest jednym z najprostszych klasyfikatorów, w którym prawdopodobieństwo $p_i(\mathbf{x})$ jest dane wzorem

$$p_i(\mathbf{x}) = \frac{1}{\kappa} \sum_{\mathbf{x}' \in N(\mathbf{x}; \kappa, d)} [t(\mathbf{x}') = c_i],$$

gdzie $N(\mathbf{x}; \kappa, d)$ jest zbiorem κ najbliższych sąsiadów punktu \mathbf{x} ze zbioru obserwacji \mathcal{X} względem ustalonej metryki, półmetryki lub podobieństwa d , a $[\cdot]$ oznacza nawias Iwersona. Metodę najbliższych sąsiadów można również wykorzystać do regresji, gdzie zamiast sumowania wartości nawiasów Iwersona, sumujemy wartości zmiennej objaśnianej dla sąsiadów punktu \mathbf{x} ze zbioru \mathcal{X} (powoduje to, że model taki potrafi tylko interpolować wartości, więc nie jest dobrym modelem dla regresji).

Klasyfikator kNN jest bardzo elastycznym modelem z nieliniową granicą decyzyjną. Jakość klasyfikacji silnie zależy od lokalnej gęstości punktów w przestrzeni oraz wybranej wartości κ , będącej hiperparametrem tego modelu. Ogólnie małe κ powoduje, że kNN ma duży variance i dość „poszarpaną” granicę decyzyjną, natomiast duże κ powoduje, że kNN ma duży bias i „gładką” granicę decyzyjną. Typowo wykorzystywane funkcje d to metryka euklidesowa, półmetryka euklidesowa, metryka Manhattan i podobieństwo cosinusowe.

Jedną z modyfikacji klasyfikatora kNN, który może polepszyć wyniki w przypadku, gdy w naszej przestrzeni istnieją obszary, w których mamy małą gęstość punktów ze zbioru \mathcal{X} jest *ważenie sąsiadów*, które polega na tym, iż prawdopodobieństwa danej klasy wśród κ sąsiadów obliczamy teraz jako średnią ważoną, gdzie wagą jest odwrotność odległości danego sąsiada od nowego wektora cech $w(\mathbf{x}, \mathbf{x}') = 1/d(\mathbf{x}, \mathbf{x}')$, tj.

$$p_i(\mathbf{x}) = \frac{\sum_{\mathbf{x}' \in N(\mathbf{x}; \kappa, d)} w(\mathbf{x}, \mathbf{x}') [t(\mathbf{x}') = c_i]}{\sum_{\mathbf{x}' \in N(\mathbf{x}; \kappa, d)} w(\mathbf{x}, \mathbf{x}')}.$$

W przypadku naiwnego kNN podczas treningu zapamiętujemy jedynie zbiór \mathcal{X} natomiast naiwna implementacja predykcji ma złożoność czasową $O(\kappa nm)$, gdzie m jest wymiarem wektora cech, przy założeniu, że złożoność obliczenia funkcji d dla pary punktów ma złożoność $O(m)$. Jest to nieakceptowalna złożoność, gdyż zwykle chcemy używać klasyfikatora kNN do setek milionów punktów. Dwa podejścia, które stosuje się zwykle do rozwiązania tego problemu to:

- zbudowanie odpowiedniej struktury danych w fazie treningu, aby w czasie predykcji można było szybciej znajdować najbliższych sąsiadów (np. k-d tree, ball tree);
- wykorzystanie algorytmów aproksymacyjnych (z ang. *Approximate Nearest Neighbors*, *ANN*) do znajdowania sąsiadów, którzy niekoniecznie naprawdę są najbliżsi, ale aproksymacja jest wystarczająco dobra do praktycznych zastosowań.

Obecnie to drugie podejście jest dominujące i wykorzystywane na przykład w przypadku dużych modeli językowych do wyszukiwania kontekstów dla danych zapytań (z ang. *Retrieval Augmented Generation*).

3.7 Regresja logistyczna

Rozpatrzmy teraz problem klasyfikacji binarnej z perspektywy estymacji MLE parametrów. Niech $\mathcal{X} = \{t_i(\mathbf{x}_i)\}_{i=1}^n$ będzie zbiorem obserwacji i.i.d., gdzie zakładamy, iż $t \in \{0, 1\}$. Jako model statystyczny przyjmujemy, iż klasa $t(\mathbf{x})$ pochodzi z rozkładu Bernoulliego z parametrem $\pi(\mathbf{x}; \mathbf{w})$ (prawdopodobieństwem klasy pozytywnej, zwanym również funkcją wiążącą, z ang. *link function*) zależnym od estymowanych parametrów \mathbf{w}

$$t(\mathbf{x}) \mid \mathbf{w} \sim \text{Ber}(\pi(\mathbf{x}; \mathbf{w})).$$

W przypadku regresji logistycznej jako funkcję $\pi(\mathbf{x}; \mathbf{w})$ przyjmujemy $\sigma(\phi(\mathbf{x}; \mathbf{w}))$, gdzie

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

to tzw. *funkcja logistyczna* (będąca szczególnym przypadkiem funkcji sigmoidalnej), natomiast $\phi(\mathbf{x}; \mathbf{w})$ może być dowolną funkcją wykorzystywaną do zagadnienia regresji. W przypadku regresji logistycznej przyjmujemy prosty model liniowy

$$\phi(\mathbf{x}; \mathbf{w}) = \mathbf{w}^T \mathbf{x}.$$

Ponownie zauważmy, iż taka postać uwzględnia człon stały poprzez podstawienie $\mathbf{x} \leftarrow [\mathbf{x}^T \ 1]^T$. Funkcja wiarygodności dla powyższego modelu statystycznego ma postać

$$\mathcal{L}(\mathcal{X}; \mathbf{w}) = \prod_{i=1}^n \pi(\mathbf{x}_i; \mathbf{w})^{t_i} (1 - \pi(\mathbf{x}_i; \mathbf{w}))^{1-t_i},$$

skąd funkcja kosztu

$$\begin{aligned} L(\mathcal{X}; \mathbf{w}) &= - \sum_{i=1}^n [t_i \log \pi(\mathbf{x}_i; \mathbf{w}) + (1 - t_i) \log(1 - \pi(\mathbf{x}_i; \mathbf{w}))] \\ &= - \sum_{i=1}^n [t_i \log \sigma(\phi(\mathbf{x}_i; \mathbf{w})) + (1 - t_i) \log(1 - \sigma(\phi(\mathbf{x}_i; \mathbf{w})))]. \end{aligned}$$

Taką funkcję kosztu nazywamy *entropią krzyżową* (z ang. *cross-entropy function*). Niestety dla takiej postaci funkcji kosztu nie można znaleźć minimum w postaci analitycznej, dlatego musimy wykorzystywać algorytmy optymalizacji numerycznej, które najczęściej wykorzystują pierwsze pochodne funkcji kosztu po parametrach. Wyznamy więc jeszcze pochodną funkcji L po parametrach \mathbf{w} . Wprowadzając macierz $\mathbf{X}^{ij} = \mathbf{x}_i^j$ oraz wektory $\mathbf{t}^i = t_i$, $\phi^i = \phi(\mathbf{x}_i; \mathbf{w}) = \sum_j \mathbf{X}^{ij} \mathbf{w}^j$, $\sigma^i = \sigma(\phi^i)$ możemy zapisać

$$L(\mathcal{X}; \mathbf{w}) = - \sum_{i=1}^n [\mathbf{t}^i \log \sigma^i + (1 - \mathbf{t}^i) \log(1 - \sigma^i)].$$

W takim razie mamy

$$\frac{\partial L}{\partial \phi^a} = \sum_b \frac{\partial L}{\partial \sigma^b} \frac{\partial \sigma^b}{\partial \phi^a},$$

gdzie

$$\frac{\partial L}{\partial \sigma^b} = - \sum_i \left[\frac{\mathbf{t}^i}{\sigma^i} \delta_{ib} - \frac{1 - \mathbf{t}^i}{1 - \sigma^i} \delta_{ib} \right] = \frac{1 - \mathbf{t}^b}{1 - \sigma^b} - \frac{\mathbf{t}^b}{\sigma^b}$$

oraz

$$\frac{\partial \sigma^b}{\partial \phi^a} = \sigma'(\phi^b) \delta_{ab}, \quad \sigma'(z) = \frac{e^{-z}}{(1 + e^{-z})^2} = \sigma(z)(1 - \sigma(z))$$

zatem

$$\frac{\partial L}{\partial \phi^a} = \left(\frac{1 - \mathbf{t}^a}{1 - \sigma^a} - \frac{\mathbf{t}^a}{\sigma^a} \right) \sigma^a (1 - \sigma^a) = \sigma^a - \mathbf{t}^a.$$

W takim razie, z powyższego mamy

$$\frac{\partial L}{\partial \mathbf{w}^a} = \sum_b \frac{\partial L}{\partial \phi^b} \frac{\partial \phi^b}{\partial \mathbf{w}^a} = \sum_b \mathbf{X}^{ba} (\sigma^b - t^b),$$

co możemy zapisać w zwartej postaci macierzowej

$$\frac{\partial L}{\partial \mathbf{w}} = \mathbf{X}^T (\boldsymbol{\sigma} - \mathbf{t}).$$

3.8 Regresja softmax

Regresję logistyczną możemy w prosty sposób uogólnić na przypadek klasyfikacji jednej z k klas. Zakładamy teraz, iż zbiór obserwacji i.i.d. ma postać $\mathcal{X} = \{t_i(\mathbf{x}_i)\}_{i=1}^n$ dla $t \in \{1, \dots, k\}$. Nasz model statystyczny uogólnimy do postaci rozkładu kategoryjnego (z ang. *categorical distribution* lub *multinomial distribution*) z parametrem $\boldsymbol{\pi}(\mathbf{x}; \mathbf{W})$ (wektorem prawdopodobieństw każdej klasy) zależnym od estymowanych parametrów \mathbf{W}

$$t(\mathbf{x}) \mid \mathbf{W} \sim \text{Cat}(\boldsymbol{\pi}(\mathbf{x}; \mathbf{W})).$$

W przypadku regresji softmax jako funkcję $\boldsymbol{\pi} : \mathbb{R}^m \mapsto [0; 1]^k$ przyjmujemy funkcję $\boldsymbol{\sigma}(\boldsymbol{\phi}(\mathbf{x}; \mathbf{W}))$, gdzie $\boldsymbol{\sigma} : \mathbb{R}^k \mapsto [0; 1]^k$

$$\sigma^i(\mathbf{z}) = \frac{\exp(\mathbf{z}^i)}{\sum_{j=1}^k \exp(\mathbf{z}^j)}$$

to tzw. *funkcja softmax*, natomiast $\boldsymbol{\phi} : \mathbb{R}^m \mapsto \mathbb{R}^k$, to dowolna funkcja. W przypadku regresji softmax przyjmujemy prosty model liniowy

$$\boldsymbol{\phi}(\mathbf{x}; \mathbf{W}) = \mathbf{W} \mathbf{x}.$$

Funkcja wiarygodności dla powyższego modelu ma postać

$$\mathcal{L}(\mathcal{X}; \mathbf{W}) = \prod_{i=1}^n \prod_{j=1}^k \pi^j(\mathbf{x}_i; \mathbf{W})^{[t_i=j]},$$

skąd

$$\begin{aligned} L(\mathcal{X}; \mathbf{W}) &= - \sum_{i=1}^n \sum_{j=1}^k [t_i = j] \log \pi^j(\mathbf{x}_i; \mathbf{W}) \\ &= - \sum_{i=1}^n \sum_{j=1}^k [t_i = j] \log \sigma^j(\boldsymbol{\phi}(\mathbf{x}_i; \mathbf{W})) \end{aligned}$$

Niestety dla takiej postaci funkcji kosztu nie można znaleźć minimum w postaci analitycznej, dlatego musimy wykorzystywać algorytmy optymalizacji numerycznej. Wyznamy więc jeszcze pierwsze pochodne funkcji L po parametrach \mathbf{W} . Wprowadźmy macierze

$$\begin{aligned} \mathbf{T}^{ij} &= [t_i = j] \\ \mathbf{X}^{ij} &= \mathbf{x}_i^j \\ \Phi^{ij} &= \phi^j(\mathbf{x}_i; \mathbf{W}) = \sum_l \mathbf{W}^{jl} \mathbf{X}^{il} \\ \mathbf{S}^{ij} &= \sigma^j(\phi(\mathbf{x}_i; \mathbf{W})) = \frac{\exp(\Phi^{ij})}{\sum_{l=1}^k \exp(\Phi^{il})} \end{aligned} \quad ,$$

gdzie każdy wiersz macierzy \mathbf{X} jest wektorem cech danego przykładu; każdy wiersz \mathbf{T} jest tzw. *wektorem one-hot* dla danego przykładu tj. wektorem binarnym, w którym dokładnie na jednej pozycji jest wartość 1 i pozycja ta odpowiada prawidłowej klasie dla danego przykładu; każdy wiersz macierzy Φ jest wektorem *mlogitów* tj. liczb rzeczywistych, które po zastosowaniu funkcji softmax dają wartości prawdopodobieństwa każdej klasy. Wówczas

$$L(\mathcal{X}; \mathbf{W}) = - \sum_{i=1}^n \sum_{j=1}^k \mathbf{T}^{ij} \log \mathbf{S}^{ij} \quad ,$$

skąd

$$\frac{\partial L}{\partial \Phi^{ab}} = \sum_{c,d} \frac{\partial L}{\partial \mathbf{S}^{cd}} \frac{\partial \mathbf{S}^{cd}}{\partial \Phi^{ab}} \quad ,$$

gdzie

$$\frac{\partial L}{\partial \mathbf{S}^{cd}} = - \sum_{i,j} \frac{\mathbf{T}^{ij}}{\mathbf{S}^{ij}} \delta_{ic} \delta_{jd} = - \frac{\mathbf{T}^{cd}}{\mathbf{S}^{cd}}$$

oraz

$$\begin{aligned} \frac{\partial \mathbf{S}^{cd}}{\partial \Phi^{ab}} &= \exp(\Phi^{cd}) \frac{\delta_{ac} \delta_{bd} [\sum_l \exp(\Phi^{cl})] - \delta_{ac} \exp(\Phi^{cb})}{[\sum_l \exp(\Phi^{cl})]^2} \\ &= \mathbf{S}^{cd} \delta_{ac} \delta_{bd} - \delta_{ac} \mathbf{S}^{cd} \mathbf{S}^{cb} \end{aligned} \quad .$$

Z powyższego

$$\frac{\partial L}{\partial \Phi^{ab}} = \sum_{c,d} \frac{\mathbf{T}^{cd}}{\mathbf{S}^{cd}} (\delta_{ac} \mathbf{S}^{cd} \mathbf{S}^{cb} - \mathbf{S}^{cd} \delta_{ac} \delta_{bd}) = \mathbf{S}^{ab} - \mathbf{T}^{ab} \quad ,$$

gdzie skorzystaliśmy z faktu, iż z konstrukcji dla dowolnego wiersza a macierzy \mathbf{T} zachodzi $\sum_d \mathbf{T}^{ad} = 1$. W takim razie

$$\frac{\partial L}{\partial \mathbf{W}^{ab}} = \sum_{c,d} \frac{\partial L}{\partial \Phi^{cd}} \frac{\partial \Phi^{cd}}{\partial \mathbf{W}^{ab}} = \sum_c (\mathbf{S}^{ca} - \mathbf{T}^{ca}) \mathbf{X}^{cb},$$

co możemy zapisać w zwartej postaci macierzowej

$$\frac{\partial L}{\partial \mathbf{W}} = (\mathbf{S} - \mathbf{T})^T \mathbf{X}.$$

4 Uczenie głębokie i sieci neuronowe

[*Introduction, Ch. 1* [1]] Artificial intelligence, or AI, is concerned with building systems that simulate intelligent behavior. It encompasses a wide range of approaches, including those based on logic, search, and probabilistic reasoning. Machine learning is a subset of AI that learns to make decisions by fitting mathematical models to observed data. This area has seen explosive growth and is now (incorrectly) almost synonymous with the term AI.

A deep neural network is a type of machine learning model, and when it is fitted to data, this is referred to as deep learning.

4.1 Sieci w pełni połączone

Neuron modelujemy jako obiekt realizujący przekształcenie $f : \mathbb{R}^n \mapsto \mathbb{R}$

$$f(\mathbf{x}; \mathbf{w}, b) = a(\mathbf{w}^T \mathbf{x} + b),$$

gdzie $a : \mathbb{R} \mapsto \mathbb{R}$ jest pewnym nieliniowym przekształceniem. Model ten ma raczej niewiele wspólnego ze współczesnymi modelami matematycznymi biologicznych neuronów, ale jest o nie w luźny sposób oparty. W szczególności wyjście neuronu zależy od kombinacji liniowej wejść \mathbf{x} (tzw. *preaktywacji*) oraz odpowiedniej nieliniowej funkcji aktywacji (z ang. *activation function*). Jako funkcje a najczęściej przyjmujemy funkcje ReLU (z ang. *Rectified Linear Unit*) lub podobne – SiLU (z ang. *Sigmoid Linear Unit*), GELU (z ang. *Gaussian Error Linear Unit*)

$$\text{ReLU}(x) = \max(0, x) = x\theta(x), \quad \text{SiLU}(x) = x\sigma(x), \quad \text{GELU}(x) = x\Phi(x),$$

gdzie $\theta(x)$ jest funkcją skokową Heaviside'a, $\sigma(x)$ jest logistyczną funkcją sigmoidalną, a $\Phi(x)$ jest dystrybuantą standardowego rozkładu normalnego.

Neurony łączymy następnie w sieci, w których wyjście danego neuronu jest wejściem innego. Sieci takie realizują pewne przekształcenia $\phi : \mathbb{R}^{n_{\text{in}}} \mapsto \mathbb{R}^{n_{\text{out}}}$, które są parametryzowane przez *wagi* w na krawędziach łączących neurony oraz obciążenia b (z ang. *bias*). Sieć taka reprezentuje zatem pewien graf obliczeń.

Najprostszą siecią, którą najpierw opiszemy będzie sieć w pełni połączona (z ang. *fully connected network*, *FCN*) zwana również ze względów historycznych perceptronem wielowarstwowym (z ang. *Multilayer Perceptron*, *MLP*). W sieci takiej neurony ułożone są w warstwy, przy czym połączenia występują tylko między neuronami w sąsiednich warstwach i nie ma połączeń między neuronami w tej samej warstwie (warstwy takie nazywamy w pełni połączonym – *fully connected layer*, *FC*). Przykładową sieć w pełni połączoną przedstawiono na rysunku 2. Na rysunku krawędzie grafu reprezentują wagi sieci, a wierzchołki reprezentują wartości funkcji aktywacji. Wierzchołki w pierwszej i ostatniej warstwie zawierają wartości odpowiednio danych wejściowych i wyjściowych i zakładamy, iż dla nich funkcja aktywacji jest funkcją identycznościową. Warstwy te nazywamy *warstwami widocznymi*, natomiast pozostałe – *warstwami ukrytymi* (z ang. *hidden layers*).

Niech $\mathbf{L} : \mathbb{R}^n \mapsto \mathbb{R}^m$ oznacza przekształcenie afiniczne postaci

$$\mathbf{L}^\mu(\mathbf{x}; \mathbf{W}, \mathbf{b}) = \mathbf{b}^\mu + \sum_{\nu=1}^n \mathbf{W}^{\mu\nu} \mathbf{x}^\nu,$$

gdzie $\mathbf{W} \in \mathbb{R}^{m \times n}$, $\mathbf{b} \in \mathbb{R}^m$. Dodatkowo przez złożenie funkcji skalarnej $a : \mathbb{R} \mapsto \mathbb{R}$ z funkcją \mathbf{L} będziemy rozumieć funkcję $a \circ \mathbf{L} : \mathbb{R}^n \mapsto \mathbb{R}^m$

$$(a \circ \mathbf{L})^\mu(\mathbf{x}) = a(\mathbf{L}^\mu(\mathbf{x})) = a\left(\mathbf{b}^\mu + \sum_{\nu=1}^n \mathbf{W}^{\mu\nu} \mathbf{x}^\nu\right).$$

Sieć w pełni połączona złożona z $k+2$ warstw o rozmiarach $n_{\text{in}}, n_1, \dots, n_k, n_{\text{out}}$ realizuje zatem przekształcenie $\phi : \mathbb{R}^{n_{\text{in}}} \mapsto \mathbb{R}^{n_{\text{out}}}$

$$\phi\left(\mathbf{x}; \{(\mathbf{W}_i, \mathbf{b}_i)\}_{i=1}^{k+1}\right) = (\mathbf{L}_{k+1} \circ a \circ \mathbf{L}_k \dots a \circ \mathbf{L}_2 \circ a \circ \mathbf{L}_1)(\mathbf{x}),$$

gdzie $\mathbf{L}_i : \mathbb{R}^{n_{i-1}} \mapsto \mathbb{R}^{n_i}$ (przy czym $n_0 = n_{\text{in}}$ i $n_{k+1} = n_{\text{out}}$)

$$\mathbf{L}_i^\mu(\mathbf{x}; \mathbf{W}_i, \mathbf{b}_i) = \mathbf{b}_i^\mu + \sum_{\nu=1}^{n_{i-1}} \mathbf{W}_i^{\mu\nu} \mathbf{x}^\nu,$$



Rysunek 2: Przykładowa sieć w pełni połączona

$\mathbf{W}_i \in \mathbb{R}^{n_i \times n_{i-1}}, \mathbf{b}_i \in \mathbb{R}^{n_i}$ (dla uproszczenia zapisu w miejscach, gdzie nie prowadzi to do niejasności będziemy oznaczać zbiór parametrów sieci $\{(\mathbf{W}_i, \mathbf{b}_i)\}_{i=1}^{k+1}$ po prostu jako \mathbf{W}). Okazuje się, iż taka postać rodziny funkcji parametryzowanych przez wagi i obciążenia $\mathbf{W}_i, \mathbf{b}_i$ potrafi aproksymować niemal dowolną funkcję (*Universal Approximation Theorem*), zatem jest to niemal uniwersalny aproksymator, który możemy zastosować do wszelkich omawianych wcześniej zagadnień regresji i klasyfikacji. Zauważmy, iż już wcześniej rozwinęliśmy teorię w sposób dość ogólny wyprowadzając większość wyników dla pewnej dowolnej funkcji $\phi(\mathbf{x}; \mathbf{W})$ i jedynie na końcu znajdując dokładne rozwiązania dla prostego odwzorowania liniowego $\mathbf{W}\mathbf{x}$ (na marginesie takie proste odwzorowanie odpowiada płytkiej sieci neuronowej o $k = 0$ warstwach ukrytych).

W powyższej postaci sieć MLP przetwarza pojedynczo każdy wektor \mathbf{x} należący do zbioru danych $\mathcal{X} = \{y_i(\mathbf{x}_i)\}_{i=1}^n$, a najbardziej kosztowną operacją jest mnożenie wektora przez macierz. Jest to jednak dość nieefektywne, gdyż na współczesnym sprzęcie mamy bardzo dobrze zoptymalizowane implementacje mnożenia macierzy przez siebie. Możemy jednak rozszerzyć odwzorowania \mathbf{L} i ϕ w taki sposób, aby sieć przetwarzała od razu cały

batch (tj. pewien podzbiór $\mathcal{X} \supseteq \mathcal{B} = \{y_\beta(\mathbf{x}_\beta)\}$) przykładów uczących reprezentowany przez macierz $\mathbf{X}^{\beta\mu} = \mathbf{x}_\beta^\mu$ (stosujemy tutaj tzw. konwencję *batch first* zgodnie z którą pierwszy indeks odpowiada kolejnym przykładom z batcha). Możemy zatem uogólnić postać przekształcenia ϕ jakie realizuje sieć do postaci $\Phi : \mathbb{R}^{n_b \times n_{in}} \mapsto \mathbb{R}^{n_b \times n_{out}}$, gdzie $n_b = |\mathcal{B}|$

$$\Phi\left(\mathbf{X}; \{(\mathbf{W}_i, \mathbf{b}_i)\}_{i=1}^{k+1}\right) = (\mathbf{L}_{k+1} \circ a \circ \mathbf{L}_k \dots a \circ \mathbf{L}_2 \circ a \circ \mathbf{L}_1)(\mathbf{X}),$$

gdzie $\mathbf{L}_i : \mathbb{R}^{n_b \times n_{i-1}} \mapsto \mathbb{R}^{n_b \times n_i}$

$$\mathbf{L}_i^{\beta\mu}(\mathbf{X}; \mathbf{W}_i, \mathbf{b}_i) = \mathbf{b}_i^\mu + \sum_{\nu=1}^{n_{i-1}} \mathbf{W}_i^{\mu\nu} \mathbf{X}^{\beta\nu},$$

$\mathbf{W}_i \in \mathbb{R}^{n_i \times n_{i-1}}, \mathbf{b}_i \in \mathbb{R}^{n_i}$. Przez złożenie funkcji skalarnej $a : \mathbb{R} \mapsto \mathbb{R}$ z funkcją \mathbf{L} nadal rozumiemy

$$(a \circ \mathbf{L})^{\beta\mu}(\mathbf{X}) = a(\mathbf{L}^{\beta\mu}(\mathbf{X})).$$

Widzimy zatem, iż tak naprawdę w takim sformułowaniu sieć przetwarza n_b przykładów równoległe zamiast przetwarzać je sekwencyjnie pojedynczo, co nie zmienia nic koncepcyjnie, a jedynie sprawia, iż sieć przetwarza dane bardziej efektywnie.

[Notes, Ch. 4 [1]]

- **Deep learning:** It has long been understood that it is possible to build more complex functions by composing shallow neural networks or developing networks with more than one hidden layer. Indeed, the term “deep learning” was first used by Dechter (1986). However, interest was limited due to practical concerns; it was not possible to train such networks well. The modern era of deep learning was kick-started by startling improvements in image classification reported by [Alex] Krizhevsky et al. (2012) [AlexNet moment]. This sudden progress was arguably due to the confluence of four factors: larger training datasets, improved processing power for training, the use of the ReLU activation function, and the use of stochastic gradient descent.
- **Depth efficiency:** Several results show that there are functions that can be realized by deep networks but not

by any shallow network whose capacity is bounded above exponentially. In other words, it would take an exponentially larger number of units in a shallow network to describe these functions accurately. This is known as the depth efficiency of neural networks.

- **Width efficiency:** Lu et al. (2017) investigate whether there are wide shallow networks (i.e., shallow networks with lots of hidden units) that cannot be realized by narrow networks whose depth is not substantially larger. They show that there exist classes of wide, shallow networks that can only be expressed by narrow networks with polynomial depth. This is known as the width efficiency of neural networks. This polynomial lower bound on width is less restrictive than the exponential lower bound on depth, suggesting that depth is more important. Vardi et al. (2022) subsequently showed that the price for making the width small is only a linear increase in the network depth for networks with ReLU activations.

4.2 Funkcje kosztu

Tak jak zauważyliśmy wyżej wyprowadzoną wcześniej teorię można bezpośrednio zaaplikować dla przypadku sieci neuronowych. Wartości parametrów \mathbf{W} znajdujemy zatem w sposób analogiczny jak dla modeli klasycznych (tj. nie głębokich) minimalizując odpowiednią funkcję kosztu $L(\mathcal{X}; \mathbf{W})$ zależną od zbioru treningowego \mathcal{X} i parametrów sieci \mathbf{W} , przy czym w czasie treningu \mathcal{X} jest ustalone, a dostosowujemy parametry \mathbf{W} , natomiast w czasie inferencji wagi \mathbf{W} są ustalone (wytrenowane), a do sieci podajemy nowe przykłady. Wymienione wcześniej funkcje kosztu dla zagadnień regresji, klasyfikacji binarnej oraz wieloklasowej możemy zatem bezpośrednio wykorzystać uwzględniając jednak, iż odwzorowanie $\phi(\mathbf{x}; \mathbf{W})$ definiujące parametry określonego rozkładu prawdopodobieństwa jest teraz odpowiednią siecią neuronową. Nie ma tutaj więc żadnej zmiany jeśli chodzi o tworzenie modeli matematycznych rozważanych zagadnień w postaci odpowiednich parametrycznych modeli statystycznych, to co się zmienia to większa ekspresywność takich modeli dzięki wykorzystaniu sieci neuronowej jako uniwersalnego aproksymatora dla parametrów modelowanych rozkładów.

Wcześniej wyprowadziliśmy funkcje kosztu (jako zanegowaną logarytmiczną funkcję wiarygodności) dla typowych zagadnień regresji i klasyfikacji – odpowiednio błąd średniokwadratowy oraz entropię krzyżową i

wieloklasową entropię krzyżową. Poniżej zamieszczamy tabelę zawierającą możliwe rozkłady prawdopodobieństwa, które można wykorzystać do skonstruowania funkcji kosztu dla mniej typowych zagadnień uczenia maszynowego.

| <i>Domain</i> | <i>Distribution</i> | <i>Use</i> |
|---------------------------------|---------------------------|---------------------------|
| $y \in \mathbb{R}$ | univariate normal | regression |
| $y \in \mathbb{R}$ | Laplace or t-distribution | robust regression |
| $y \in \mathbb{R}$ | mixture of gaussians | multimodal regression |
| $y \in \mathbb{R}_+$ | exponential or gamma | predicting magnitude |
| $y \in [0; 1]$ | beta | predicting proportion |
| $\mathbf{y} \in \mathbb{R}^d$ | multivariate normal | multivariate regression |
| $y \in (-\pi; +\pi]$ | von Mises | predicting direction |
| $y \in \{0, 1\}$ | Bernoulli | binary classification |
| $y \in \{1, 2, \dots, K\}$ | categorical | multiclass classification |
| $y \in \{0, 1, 2, 3, \dots\}$ | Poisson | predicting event counts |
| $\mathbf{y} \in \text{Perm}(n)$ | Plackett–Luce | ranking |

Tabela 1: *Figure 5.11, Ch. 5, [1]*

[*Multiple outputs, Ch. 5, [1]*] To make two or more prediction types simultaneously, we similarly assume the errors in each are independent. For example, to predict wind direction and strength, we might choose the von Mises distribution (defined on circular domains) for the direction and the exponential distribution (defined on positive real numbers) for the strength. The independence assumption implies that the joint likelihood of the two predictions is the product of individual likelihoods. These terms will become additive when we compute the negative log-likelihood.

[*Notes, Ch. 5, [1]*]

- **Class imbalance and focal loss:** Lin et al. (2017c) address data imbalance in classification problems. If the number of examples for some classes is much greater than for others, then the standard maximum likelihood loss does not work well; the model may concentrate on becoming more confident about well-classified examples from the dominant classes and classify less well-represented classes

poorly. Lin et al. (2017c) introduce focal loss, which adds a single extra parameter that down-weights the effect of well-classified examples to improve performance.

- **Learning to rank:** Cao et al. (2007), Xia et al. (2008), and Chen et al. (2009) all used the Plackett–Luce model in loss functions for learning to rank data. This is the listwise approach to learning to rank as the model ingests an entire list of objects to be ranked at once. Alternative approaches are the pointwise approach, in which the model ingests a single object, and the pairwise approach, where the model ingests pairs of objects.
- **Non-probabilistic approaches:** It is not strictly necessary to adopt the probabilistic approach discussed in this chapter, but this has become the default in recent years; any loss function that aims to reduce the distance between the model output and the training outputs will suffice, and distance can be defined in any way that seems sensible. There are several well-known non-probabilistic machine learning models for classification, including support vector machines, which use hinge loss and AdaBoost, which uses exponential loss.

4.3 Aspekty optymalizacyjne

Przedstawimy teraz algorytmy optymalizacji numerycznej najczęściej używane przy trenowaniu głębokich sieci neuronowych. Jest oczywiste, iż nie jest możliwe znalezienie wartości wag sieci neuronowej minimalizujących funkcję kosztu $L(\mathcal{X}; \mathbf{W})$ w sposób analityczny. Dodatkowo ze względu na ogromną liczbę parametrów współczesnych sieci (rzędu 10^9) jakiegokolwiek wyrafinowane algorytmy optymalizacji numerycznej (np. wykorzystujące drugie pochodne) również nie wchodzi w grę ze względu na zbyt duże wymagania pamięciowe. Biorąc pod uwagę wymienione ograniczenia skupimy się na prostych algorytmach optymalizacji pierwszego rzędu (z ang. *first order methods*).

Podstawowym algorytmem jest metoda spadku wzdłuż gradientu (z ang. *gradient descent*). Polega ona na iteracyjnym aktualizowaniu wartości parametrów \mathbf{W} zgodnie z kierunkiem najszybszego spadku wyznaczonego przez zanegowany gradient funkcji kosztu $-\frac{\partial L(\mathcal{X}; \mathbf{W})}{\partial \mathbf{W}}$. Wielkość kroków jakie wykonujemy w każdej iteracji jest zdeterminowana przez stałą $\lambda > 0$ nazywaną *stałą uczącą* (z ang. *learning rate*). Tutaj należy dodać iż w implementacji

wartość funkcji kosztu jest mnożona przez liczbę przykładów, w taki sposób, iż obliczamy tak naprawdę średni koszt na przykład, czyli $L \leftarrow \frac{1}{|\mathcal{X}|}L$. W przeciwnym wypadku wartość stałej uczącej należałoby dostosowywać do liczby przykładów w zbiorze.

```

1 # Gradient descent
2 # -----
3 # Let  $\mathcal{X}$  be the whole training set.
4 # Initialize heuristically parameters  $\mathbf{W}$  of the neural net.
5 for epoch in range(max_epochs):
6     # Step 1. Compute derivatives of the loss function with
7     # respect to the parameters of neural net  $\frac{\partial L(\mathcal{X}; \mathbf{W})}{\partial \mathbf{W}}$ .
8
9     # Step 2. Update parameters of the neural net according to
10     $\mathbf{W} \leftarrow \mathbf{W} - \lambda \frac{\partial L(\mathcal{X}; \mathbf{W})}{\partial \mathbf{W}}$ 
11    # where  $\lambda > 0$  is a positive constant called 'learning rate'.

```

[*Gradient descent variants*, [2]] There are three variants of gradient descent, which differ in how much data we use to compute the gradient of the objective function. Depending on the amount of data, we make a trade-off between the accuracy of the parameter update and the time it takes to perform an update.

Podstawowe modyfikacje powyższego algorytmu (zwanego również *full-batch gradient descent*) są związane z ograniczeniem zbioru treningowego dla którego obliczamy gradient. Podejście takie wynika z faktu, iż dla bardzo dużych zbiorów treningowych \mathcal{X} obliczane tensory mogą nie mieścić się w pamięci. Dodatkowo modyfikacja wprowadza losowość przy aktualizacji parametrów sieci, co pozwala na ucieczkę z lokalnych minimów funkcji kosztu. Aby rozwiązać ten problem dzielimy zbiór \mathcal{X} w każdej epoce na rozłączne *batche* \mathcal{B} o tej samej wielkości n_b składające się z losowych przykładów ze zbioru \mathcal{X} (losowane bez zwracania i każdorazowo przy rozpoczęciu nowej epoki). Algorytm taki nazywamy *batch gradient descent* lub *mini-batch gradient descent*.

```

1 # Mini-batch gradient descent
2 # -----
3 # Let  $\mathcal{X}$  be the whole training set.
4 # Initialize heuristically parameters  $\mathbf{W}$  of the neural net.
5 for epoch in range(max_epochs):
6     random_shuffle( $\mathcal{X}$ )
7     for  $\mathcal{B}$  in get_batches( $\mathcal{X}$ , batch_size):
8         # Step 1. Compute derivatives of the loss function with
9         # respect to the parameters of neural net  $\frac{\partial L(\mathcal{B}; \mathbf{W})}{\partial \mathbf{W}}$ .
10

```



```

11     # Step 2. Update parameters of the neural net according to
12      $\mathbf{W} \leftarrow \mathbf{W} - \lambda \frac{\partial L(\mathcal{B}; \mathbf{W})}{\partial \mathbf{W}}$ 
13     # where  $\lambda > 0$  is a positive constant called 'learning rate'.

```

W szczególności gdy $n_b = 1$, tj. każdy batch składa się z pojedynczego, losowo wybranego przykładu uczącego, to taki algorytm nazywamy *stochastycznym spadkiem wzdłuż gradientu* (z ang. *stochastic gradient descent*, *SGD*).

[*Stochastic gradient descent*, [2]] Gradient descent performs redundant computations for large datasets, as it recomputes gradients for similar examples before each parameter update. SGD does away with this redundancy by performing one update at a time. It is therefore usually much faster and can also be used to learn online. SGD performs frequent updates with a high variance that cause the objective function to fluctuate heavily.

While gradient descent converges to the minimum of the basin the parameters are placed in, SGD's fluctuation, on the one hand, enables it to jump to new and potentially better local minima. On the other hand, this ultimately complicates convergence to the exact minimum, as SGD will keep overshooting. However, it has been shown that when we slowly decrease the learning rate, SGD shows the same convergence behavior as batch gradient descent, almost certainly converging to a local or the global minimum for non-convex and convex optimization respectively.

[*Mini-batch gradient descent*, [2]] Mini-batch gradient descent finally takes the best of both worlds and performs an update for every mini-batch of training examples. This way, it a) reduces the variance of the parameter updates, which can lead to more stable convergence; b) can make use of highly optimized matrix optimizations common to state-of-the-art deep learning libraries that make computing the gradient w.r.t. a mini-batch very efficient.

Common mini-batch sizes range between 50 and 256, but can vary for different applications. Mini-batch gradient descent is typically the algorithm of choice when training a neural network and the term SGD usually is employed also when mini-batches are used.

4.3.1 Metoda pędu

W standardowym sformułowaniu SGD ma pewne problemy. Przede wszystkim dobór odpowiedniej wartości stałej uczącej do uzyskania dobrej zbieżności może być trudny. Dodatkowo istnieje poważny problem, który pojawia się przy minimalizacji silnie niewypukłych funkcji kosztu polegający na tym, iż algorytm utyka w lokalnym minimum lub w punkcie siodłowym funkcji. Standardowy algorytm SGD nie posiada żadnych mechanizmów pozwalających na ucieczkę z tych punktów. Podstawowym ulepszeniem algorytmu SGD jest wprowadzenie tzw. pędu (z ang. *momentum*) tj. członu postaci

$$\mathbf{V}_t = \gamma \mathbf{V}_{t-1} - \lambda \frac{\partial L(\mathcal{B}; \mathbf{W})}{\partial \mathbf{W}} \Big|_{\mathbf{W}}, \quad \mathbf{V}_0 = \mathbf{0}$$

i aktualizację parametrów sieci zgodnie z

$$\mathbf{W} \leftarrow \mathbf{W} + \mathbf{V}_t,$$

gdzie t jest numerem iteracji, a γ to tzw. *momentum term*, którego wartość typowo przyjmuje się równą 0.9.

[*Momentum*, [2]] Essentially, when using momentum, we push a ball down a hill. The ball accumulates momentum as it rolls downhill, becoming faster and faster on the way (until it reaches its terminal velocity, if there is air resistance, i.e. $\gamma < 1$). The same thing happens to our parameter updates: The momentum term increases for dimensions whose gradients point in the same directions and reduces updates for dimensions whose gradients change directions. As a result, we gain faster convergence and reduced oscillation.

4.3.2 Metoda Nesterova

[*Nesterov accelerated gradient*, [2]] However, a ball that rolls down a hill, blindly following the slope, is highly unsatisfactory. We would like to have a smarter ball, a ball that has a notion of where it is going so that it knows to slow down before the hill slopes up again.

Nesterov accelerated gradient is a way to give our momentum term this kind of prescience. We know that we will use our momentum term $\gamma \mathbf{V}_{t-1}$ to move the parameters \mathbf{W} . Computing $\mathbf{W} + \gamma \mathbf{V}_{t-1}$ thus gives us an approximation of the next position

of the parameters (the gradient is missing for the full update), a rough idea where our parameters are going to be. We can now effectively look ahead by calculating the gradient not at our current parameters \mathbf{W} but at the approximate future position of our parameters:

$$\mathbf{V}_t = \gamma \mathbf{V}_{t-1} - \lambda \frac{\partial L(\mathcal{B}; \mathbf{W})}{\partial \mathbf{W}} \Big|_{\mathbf{W} + \gamma \mathbf{V}_{t-1}}, \quad \mathbf{V}_0 = \mathbf{0},$$

$$\mathbf{W} \leftarrow \mathbf{W} + \mathbf{V}_t.$$

4.3.3 Adam

[*Adam*, Ch.6, [1]] Gradient descent with a fixed step size has the following undesirable property: it makes large adjustments to parameters associated with large gradients (where perhaps we should be more cautious) and small adjustments to parameters associated with small gradients (where perhaps we should explore further). When the gradient of the loss surface is much steeper in one direction than another, it is difficult to choose a learning rate that (i) makes good progress in both directions and (ii) is stable.

A straightforward approach is to normalize the gradients so that we move a fixed distance (governed by the learning rate) in each direction.

$$\mathbf{G} = \frac{\partial L(\mathcal{B}; \mathbf{W})}{\partial \mathbf{W}} \Big|_{\mathbf{W}},$$

$$\mathbf{W} \leftarrow \mathbf{W} - \lambda \frac{\mathbf{G}}{\sqrt{\mathbf{G}^2 + \epsilon}}.$$

where the square root and division are both elementwise, λ is the learning rate, and ϵ is a small constant that prevents division by zero when the gradient magnitude is zero.

The result is that the algorithm moves a fixed distance λ along each coordinate, where the direction is determined by whichever way is downhill. This simple algorithm makes good progress in both directions but will not converge unless it happens to land exactly at the minimum. Instead, it will bounce back and forth around the minimum.

Adaptive moment estimation, or *Adam*, takes this idea and adds momentum to both the estimate of the gradient and the squared

gradient:

$$\begin{aligned}
\mathbf{G} &= \left. \frac{\partial L(\mathcal{B}; \mathbf{W})}{\partial \mathbf{W}} \right|_{\mathbf{w}}, \\
\mathbf{M}_t &= \beta_1 \mathbf{M}_{t-1} + (1 - \beta_1) \mathbf{G}, \quad \tilde{\mathbf{M}}_t = \frac{\mathbf{M}_t}{1 - \beta_1^t}, \\
\mathbf{V}_t &= \beta_2 \mathbf{V}_{t-1} + (1 - \beta_2) \mathbf{G}^2, \quad \tilde{\mathbf{V}}_t = \frac{\mathbf{V}_t}{1 - \beta_2^t}, \\
\mathbf{W} &\leftarrow \mathbf{W} - \lambda \frac{\tilde{\mathbf{M}}_t}{\sqrt{\tilde{\mathbf{V}}_t} + \epsilon}.
\end{aligned}$$

This algorithm can converge to the overall minimum and makes good progress in every direction in the parameter space.

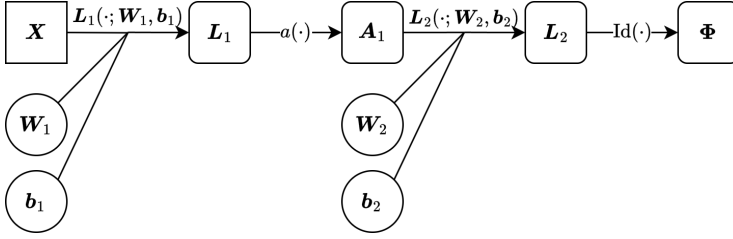
The choices of learning algorithm, batch size, learning rate schedule, and momentum coefficients are all considered hyperparameters of the training algorithm; these directly affect the final model performance but are distinct from the model parameters. Choosing these can be more art than science, and it is common to train many models with different hyperparameters and choose the best one. This is known as *hyperparameter search*.

[Adam, [2]] The authors propose default values of 0.9 for β_1 , 0.999 for β_2 , and 10^{-8} for ϵ . They show empirically that Adam works well in practice and compares favorably to other adaptive learning-method algorithms.

4.4 Algorytm wstecznej propagacji błędu

W poprzednim paragrafie opisaliśmy podstawowe algorytmy wykorzystywane do minimalizacji funkcji kosztu, a co za tym idzie trenowania sieci neuronowych. Nie opisaliśmy jednak dwóch dość istotnych szczegółów: 1) jak właściwie obliczać gradient funkcji kosztu po parametrach sieci $\frac{\partial L(\mathcal{B}; \mathbf{W})}{\partial \mathbf{W}}$; 2) w jaki sposób zainicjalizować wartości parametrów sieci. W tym paragrafie zajmujemy się pierwszym z wymienionych problemów, natomiast w następnym poruszamy kwestie regularyzacji oraz inicjalizacji głębokich sieci neuronowych.

Sieć neuronowa (przetwarzająca od razu cały batch przykładów reprezentowany przez tablicę \mathbf{X}) jest tak naprawdę pewnym skierowanym, acyklicznym grafem (z ang. *Directed Acyclic Graph*, *DAG*), zwanym *grafem obliczeń*. Przykładowy graf obliczeń reprezentujący prostą sieć MLP z jedną warstwą ukrytą przedstawiono na rysunku 3.



Rysunek 3: Graf obliczeń reprezentujący prostą sieć MLP z jedną warstwą ukrytą

Chcemy w sposób automatyczny z dokładnością do epsilon maszynowego obliczyć wartości pochodnych $\frac{\partial L(\mathcal{X}; \mathbf{W})}{\partial \mathbf{W}_i}$, $\frac{\partial L(\mathcal{X}; \mathbf{W})}{\partial \mathbf{b}_i}$. Zauważmy, iż zgodnie z regułą łańcuchową oznaczając

$$\Delta_i^{\mu_1 \mu_2} = \frac{\partial L}{\partial L_i^{\mu_1 \mu_2}}$$

mamy

$$\frac{\partial L}{\partial \mathbf{W}_i^{\mu_1 \mu_2}} = \sum_{\nu_1, \nu_2} \Delta_i^{\nu_1 \nu_2} \frac{\partial L_i^{\nu_1 \nu_2}}{\partial \mathbf{W}_i^{\mu_1 \mu_2}} = \sum_{\nu_1} \Delta_i^{\nu_1 \mu_1} A_{i-1}^{\nu_1 \mu_2}$$

oraz

$$\frac{\partial L}{\partial \mathbf{b}_i^\mu} = \sum_{\nu} \Delta_i^{\nu \mu},$$

gdzie oznaczamy $\mathbf{A}_0 = \mathbf{X}$. Jednocześnie zauważmy, że

$$\Delta_i^{\mu_1 \mu_2} = \sum_{\nu_1, \nu_2} \frac{\partial L}{\partial A_i^{\nu_1 \nu_2}} \frac{\partial A_i^{\nu_1 \nu_2}}{\partial L_i^{\mu_1 \mu_2}}$$

oraz

$$\frac{\partial L}{\partial A_i^{\nu_1 \nu_2}} = \sum_{\gamma_1, \gamma_2} \frac{\partial L}{\partial L_{i+1}^{\gamma_1 \gamma_2}} \frac{\partial L_{i+1}^{\gamma_1 \gamma_2}}{\partial A_i^{\nu_1 \nu_2}} = \sum_{\gamma} \Delta_{i+1}^{\nu_1 \gamma} \mathbf{W}_{i+1}^{\gamma \nu_2}$$

i

$$\frac{\partial A_i^{\nu_1 \nu_2}}{\partial L_i^{\mu_1 \mu_2}} = a'(L_i^{\nu_1 \nu_2}) \delta_{\mu_1 \nu_1} \delta_{\mu_2 \nu_2},$$

skąd

$$\Delta_i^{\mu_1 \mu_2} = a'(L_i^{\mu_1 \mu_2}) \sum_{\gamma} \Delta_{i+1}^{\mu_1 \gamma} \mathbf{W}_{i+1}^{\gamma \mu_2}.$$

Powyższe wyniki możemy zapisać w zwartych postaciach macierzowych

$$\boxed{\frac{\partial L}{\partial \mathbf{W}_i} = \Delta_i^T \mathbf{A}_{i-1}, \quad \frac{\partial L}{\partial \mathbf{b}_i} = \Delta_i^T \mathbf{1}, \quad \Delta_i = a'(\mathbf{L}_i) \odot (\Delta_{i+1} \mathbf{W}_{i+1}),}$$

gdzie $\mathbf{A}_0 = \mathbf{X}$ oraz $\Delta_{k+1} = \frac{\partial L}{\partial \Phi}$ (zauważmy tutaj, iż człon ten wyznaczyliśmy już powyżej dla typowych funkcji kosztu). Powyższy przykład to wyprowadzenie algorytmu wstecznej propagacji błędu dla typowej sieci w pełni połączonej (reprezentowanej przez graf obliczeń 3 z powielonymi warstwami ukrytymi – na rysunku przedstawiono sieć z jedną warstwą ukrytą). Algorytm obliczania pochodnych możemy zatem zapisać następująco.

```

1 # Backpropagation through MLP
2 # -----
3 # Let k be the number of hidden layers of the MLP network.
4
5 # Step 1. For given batch  $\mathbf{X}$  perform a forward pass
6 # by computing and memorizing arrays  $\mathbf{L}_i, \mathbf{A}_i$ .
7
8  $\mathbf{A}_0 = \mathbf{X}$ 
9
10 for i in range(1, k+2):
11     # Compute and memorize
12      $\mathbf{L}_i = \mathbf{A}_{i-1} \mathbf{W}_i^T + \mathbf{b}_i$ 
13      $\mathbf{A}_i = a(\mathbf{L}_i)$ 
14
15 # Compute loss  $L(\mathcal{B}; \mathbf{W})$  for the output  $\Phi = \mathbf{L}_{k+1}$ .
16
17 # Step 2. Perform backpropagation
18
19 # Initialize seed
20  $\Delta_{k+1} = \frac{\partial L(\mathcal{B}; \mathbf{W})}{\partial \Phi}$ 
21
22 for i in range(k+1, 0, -1):
23     # Compute gradients
24      $\frac{\partial L}{\partial \mathbf{W}_i} = \Delta_i^T \mathbf{A}_{i-1}, \quad \frac{\partial L}{\partial \mathbf{b}_i} = \Delta_i^T \mathbf{1}$ 
25
26     # Update seed
27     if i > 1:
28          $\Delta_{i-1} = a'(\mathbf{L}_{i-1}) \odot (\Delta_i \mathbf{W}_i)$ 

```

Powyższy algorytm jest specjalnie dostosowany do sieci w pełni połączonej. Możemy go jednak prosto uogólnić do dowolnego acyklicznego i skierowanego grafu obliczeń. Załóżmy, iż rozpatrujemy warstwę (wierzchołek w grafie zaznaczony zaokrąglonym kwadratem) \mathbf{F} zależną od pewnych parametrów \mathbf{W} (w szczególności warstwa może nie mieć żadnych parametrów). Dla dowolnej warstwy \mathbf{F} przez $\dot{\mathbf{F}}$ oznaczmy pochodną $\dot{\mathbf{F}} = \frac{\partial L}{\partial \mathbf{F}}$.

Wówczas z reguły łańcuchowej

$$\frac{\partial L}{\partial \mathbf{W}^{\mu_1 \dots \mu_k}} = \sum_{\nu_1 \dots \nu_l} \dot{\mathbf{F}}^{\nu_1 \dots \nu_l} \frac{\partial \mathbf{F}^{\nu_1 \dots \nu_l}}{\partial \mathbf{W}^{\mu_1 \dots \mu_k}}.$$

Niech $\text{succ}(\mathbf{F})$ oznacza zbiór następników wierzchołka \mathbf{F} , tj. warstw (wierzchołków w grafie zaznaczonych zaokrąglonym kwadratem) do których istnieje skierowana krawędź od warstwy \mathbf{F} . Wówczas

$$\dot{\mathbf{F}}^{\mu_1 \dots \mu_k} = \frac{\partial L}{\partial \mathbf{F}^{\mu_1 \dots \mu_k}} = \sum_{\mathbf{H} \in \text{succ}(\mathbf{F})} \sum_{\nu_1 \dots \nu_l} \dot{\mathbf{H}}^{\nu_1 \dots \nu_l} \frac{\partial \mathbf{H}^{\nu_1 \dots \nu_l}}{\partial \mathbf{F}^{\mu_1 \dots \mu_k}}.$$

Z powyższego zatem widzimy, iż możemy obliczać wielkości $\dot{\mathbf{F}}$ w sposób iteracyjny przechodząc po warstwach w kolejności odwrotnej do ich porządku topologicznego (z ang. *reversed topological order*) i zapamiętywać wartości $\dot{\mathbf{F}}$. Wówczas dla każdej warstwy pochodne $\dot{\mathbf{H}}$ jej następników będą już obliczone i będziemy mogli obliczyć $\dot{\mathbf{F}}$. Naiwna implementacja powyższego schematu będzie wymagać bardzo dużo pamięci, gdyż dla każdej warstwy musimy przechowywać oprócz jej wartości \mathbf{F} również wartość pochodnej $\dot{\mathbf{F}}$. Zauważmy, iż w przypadku algorytmu wstecznej propagacji błędu dla sieci MLP nie musimy przechowywać pochodnych funkcji kosztu po wartościach warstw, a jedynie dynamicznie obliczamy i przekazujemy wstecz tzw. *ziarno* (z ang. *seed*) Δ . Wynika to oczywiście z faktu, iż w przypadku sieci MLP warstwy tworzą skierowaną ścieżkę, czyli najprostszy DAG. W przypadku innych architektur sieci neuronowych tak nie musi być i wówczas każdą grupę warstw w obrębie której warstwy nie tworzą ścieżki należy oddzielnie przeanalizować pod kątem wykorzystania pamięci.

[Notes, Ch.7, [1]]

- **Reducing memory requirements:** Training neural networks is memory intensive. We must store both the model parameters and the pre-activations at the hidden units for every member of the batch during the forward pass. Two methods that decrease memory requirements are gradient checkpointing (Chen et al., 2016a) and micro-batching (Huang et al., 2019). In gradient checkpointing, the activations are only stored every N layers during the forward pass. During the backward pass, the intermediate missing activations are recalculated from the nearest checkpoint. In this manner, we can drastically reduce the memory requirements at the computational cost of performing the

forward pass twice. In micro-batching, the batch is subdivided into smaller parts, and the gradient updates are aggregated from each sub-batch before being applied to the network.

4.5 Inicjalizacja

Jak zaznaczyliśmy wyżej ten paragraf poświęcimy zagadnieniu inicjalizacji parametrów sieci neuronowej, który jest wymaganym elementem opisanych algorytmów minimalizacji numerycznej.

[*Parameter initialization, Ch.7, [1]*] The backpropagation algorithm computes the derivatives that are used by stochastic gradient descent and Adam to train the model. We now address how to initialize the parameters before we start training. To see why this is crucial, consider that during the forward pass, each set of pre-activations \mathbf{L}_i is computed as:

$$\mathbf{L}_i = \mathbf{b}_i + a(\mathbf{L}_{i-1})\mathbf{W}_i^T$$

Imagine that we initialize all the biases to zero and the elements of \mathbf{W} according to a normal distribution with mean zero and variance σ^2 . Consider two scenarios:

- If the variance is very small (e.g., 10^{-5}), then each element of \mathbf{L}_i will be a weighted sum of \mathbf{A}_{i-1} where the weights are very small; the result will likely have a smaller magnitude than the input. In addition, the ReLU function clips values less than zero, so the range of \mathbf{A}_i will be half that of \mathbf{L}_i . Consequently, the magnitudes of the pre-activations at the hidden layers will get smaller and smaller as we progress through the network.
- If the variance is very large (e.g., 10^5), then each element of \mathbf{L}_i will be a weighted sum of \mathbf{A}_{i-1} where the weights are very large; the result is likely to have a much larger magnitude than the input. The ReLU function halves the range of the inputs, but if variance is large enough, the magnitudes of the pre-activations will still get larger as we progress through the network.

In these two situations, the values at the pre-activations can become so small or so large that they cannot be represented with finite precision floating point arithmetic.

Even if the forward pass is tractable, the same logic applies to the backward pass. Each gradient update consists of multiplying by \mathbf{W}_i . If the values of \mathbf{W}_i are not initialized sensibly, then the gradient magnitudes may decrease or increase uncontrollably during the backward pass. These cases are known as the *vanishing gradient problem* and the *exploding gradient problem*, respectively. In the former case, updates to the model become vanishingly small. In the latter case, they become unstable.

Najprostszym pomysłem inicjalizacji parametrów sieci neuronowych jest przyjęcie wartości 0 dla obciążeń \mathbf{b}_i i wylosowanie elementów macierzy wag \mathbf{W}_i niezależnie z rozkładu normalnego $\mathcal{N}(0, \sigma^2)$. Powstaje pytanie jaka powinna być wartość wariancji σ^2 . Pewną odpowiedź na to pytanie sugeruje analiza wariancji elementów aktywacji w kolejnych warstwach

$$\mathbf{L}_{i+1}^{\beta\mu} = \mathbf{b}_i^\mu + \sum_{\nu=1}^{n_{i-1}} \mathbf{A}_i^{\beta\nu} \mathbf{W}_i^{\mu\nu}.$$

W przypadku inicjalizacji mamy $\mathbf{b}_i^\mu = 0$ oraz $\mathbf{W}_i^{\mu\nu} \sim \mathcal{N}(0, \sigma_{W_i}^2)$, zatem

$$\mathbb{E} \left[\mathbf{L}_{i+1}^{\beta\mu} \right] = \sum_{\nu=1}^{n_{i-1}} \mathbb{E} \left[\mathbf{A}_i^{\beta\nu} \right] \cdot \mathbb{E} \left[\mathbf{W}_i^{\mu\nu} \right] = 0$$

oraz

$$\mathbb{V} \left[\mathbf{L}_{i+1}^{\beta\mu} \right] = \mathbb{E} \left[(\mathbf{L}_{i+1}^{\beta\mu})^2 \right] - \mathbb{E} \left[\mathbf{L}_{i+1}^{\beta\mu} \right]^2 = \sum_{\nu=1}^{n_{i-1}} \mathbb{V} \left[\mathbf{A}_i^{\beta\nu} \right] \cdot \sigma_{W_i}^2.$$

Zakładając, iż wariancje składowych aktywacji są równe $\sigma_{A_i}^2$ oraz pochodzą z rozkładu symetrycznego wokół 0 mamy (dla funkcji aktywacji ReLU)

$$\sigma_{A_{i+1}}^2 = \sum_{\nu=1}^{n_{i-1}} \frac{1}{2} \sigma_{W_i}^2 \sigma_{A_i}^2 = \frac{n_{i-1}}{2} \sigma_{W_i}^2 \sigma_{A_i}^2.$$

Widzimy zatem, iż jeśli chcemy aby rozrzut wartości kolejnych aktywacji był podobny do poprzednich musi zachodzić

$$\sigma_{W_i}^2 = \frac{2}{n_{i-1}}.$$

Otrzymujemy zatem następującą heurystykę inicjalizacji

$$\mathbf{b}_i = \mathbf{0}, \quad \mathbf{W}_i^{\mu\nu} \sim \mathcal{N} \left(0, \frac{2}{n_{i-1}} \right)$$

zwaną *inicjalizacją He*. Analogiczna analiza dla propagacji wstecznej, gdzie aktualizujemy Δ sugeruje z kolei przyjęcie

$$\sigma_{W_i}^2 = \frac{2}{n_i}.$$

Jeśli $n_i \neq n_{i-1}$, tj. macierz \mathbf{W}_i nie jest kwadratowa, to jedną z możliwości jest przyjęcie średniej arytmetycznej $\bar{n} = \frac{n_{i-1} + n_i}{2}$

$$\sigma_{W_i}^2 = \frac{2}{\bar{n}} = \frac{4}{n_{i-1} + n_i}.$$

4.6 Regularyzacja

Jak wspomnieliśmy we wstępie ważnym pojęciem w uczeniu maszynowym jest pojemność modeli. Model o dużej pojemności (czyli na przykład właśnie głęboka sieć neuronowa) może nastroić poważnych problemów w trakcie uczenia przez zjawisko overfittingu. Model o dużej pojemności może perfekcyjnie dopasować się do zbioru trenującego. Równocześnie może on zignorować strukturę, której wykrycie było w istocie celem uczenia. Przyczyną przeuczenia jest fakt, że każdy zbiór trenujący jest skończony: 1) jest skończoną próbą obserwacji z zagadnienia, które jest celem uczenia; 2) każdy zbiór trenujący opisuje więc problem uczenia w sposób niepełny; 3) model przeucza, gdy jego pojemność pozwala mu dopasować się do konkretnego zbioru trenującego w sposób uniemożliwiający uogólnianie na nowe przypadki. Problem ten jest tym istotniejszy, im mniejszy jest zbiór trenujący. Algorytmy i techniki służące do kontroli pojemności modelu nazywamy metodami regularyzacji (z ang. *regularization methods*).

- **Weight decay.** Analogicznie jak w przypadku prostych modeli liniowych jedną z możliwości regularyzacji jest dodanie do funkcji kosztu L czynnika regularyzującego zawierającego odpowiednią sumę norm poszczególnych parametrów sieci

$$L^*(\mathcal{X}; \mathbf{W}) = L(\mathcal{X}; \mathbf{W}) + \frac{\gamma_2}{2} \sum_{i=1}^{k+1} \sum_{\mu, \nu} (\mathbf{W}_i^{\mu\nu})^2 + \gamma_1 \sum_{i=1}^{k+1} \sum_{\mu, \nu} |\mathbf{W}_i^{\mu\nu}|.$$

Różnica między regularyzacją L^2 i L^1 jest taka, że norma L^2 „zachęca” sieć do wybierania małych wag, ale nie dokładnie będących 0, natomiast L^1 zmniejsza wagi, ale w szczególności prowadzi do wag wynoszących dokładnie 0. W związku z tym regularyzacja L^1 przeprowadza selekcję cech - „wyłączy” te, które są mało ważne. Zwykle

człony regularyzujące wprowadza się jedynie dla wag sieci, a nie dla obciążeń, skąd nazwa tej metody regularyzacji – *weight decay*.

- **Metoda wczesnego stopu.** Najprostszym, a zarazem powszechnie stosowanym sposobem regularyzacji sieci neuronowych jest metoda wczesnego stopu (z ang. *early stopping*). Polega ona na podziale zbioru *treningowego* na mniejszy zbiór treningowy i tzw. *zbiór walidacyjny*. W trakcie treningu monitorujemy błąd popełniany przez model na zbiorze walidacyjnym (ale do treningu używamy tylko zbioru treningowego). Przerywamy uczenie, gdy błąd na zbiorze walidacyjnym zaczyna rosnąć lub nie maleje przez wybraną liczbę epok. Dlaczego postęp uczenia monitorowany jest na zbiorze walidacyjnym (wydzielonym ze zbioru trenującego) a nie na zbiorze testowym? Monitorując przebieg uczenia na zbiorze testowym dobieralibyśmy parametry uczenia (np. ilość epok) tak, by dawały dobry wynik na zbiorze testowym. Wówczas wynik na zbiorze testowym mógłby być wyraźnie lepszy niż faktyczna skuteczność sieci w docelowym zastosowaniu. Nie możemy dobierać parametrów modelu wprost pod zbiór, na którym będziemy oceniać jego docelową skuteczność. Celem zbioru testowego jest ocena skuteczności modelu dla nowych obserwacji (nieдоступnych w czasie uczenia).

Takie podejście wymaga dość dużego zbioru danych, z których musimy wydzielić trzy zbiory: zbiór treningowy (na którym uczymy model), zbiór walidacyjny (na którym ewaluujemy model przy procedurze *early stoppingu*) oraz zbiór testowy (na którym ewaluujemy model po treningu). Okres, przez który czekamy na uzyskanie lepszego wyniku, to cierpliwość (z ang. *patience*). Im mniejsze, tym mocniejszy jest ten rodzaj regularyzacji, ale trzeba z tym uważać, bo łatwo jest przesadzić i zbyt szybko przerywać trening. Niektóre implementacje uwzględniają tzw. *grace period*, czyli gwarantowaną minimalną liczbę epok, przez którą będziemy trenować sieć, niezależnie od wybranej cierpliwości.

- **Dropout.** Dobrym sposobem regularyzacji modeli uczenia maszynowego jest uśrednianie odpowiedzi wielu modeli popełniających różne błędy. Chcemy więc uczyć wiele modeli (rozwiązujących ten sam problem) w taki sposób, by modele te różniły się między sobą. Jest to tzw. uczenie zespołowe (z ang. *ensemble learning*). Jak zapewnić, by wyuczone modele się między sobą różniły? Jednym ze sposobów budowy modelu zespołowego jest *bagging* (od bootstrap aggregating) polegający na trenowaniu modeli na losowo wybranych (ze zwraca-

niem) podzbiorach zbioru treningowego, odpowiedzią modelu będzie średnia z odpowiedzi wszystkich wyuczonych w taki sposób modeli. Oczywiście w przypadku obliczeniowo kosztownych modeli jak sieci neuronowe zaimplementowanie takiego schematu wprost jest nieefektywne. Metodą, która efektywnie realizuje właśnie uczenie zespołowe na pojedynczej sieci neuronowej jest dropout. Polega on na tym, iż w procesie treningu ze z góry zadany prawdopodobieństwem „wyłączamy” wyjście neuronu (tj. dla każdej aktualizacji wartości parametrów sieci, w algorytmie wstecznej propagacji błędu mnożymy elementwise wyjście \mathbf{A}_i przez tensor \mathbf{D}_i , którego każdy element może być równy 0 z prawdopodobieństwem p_i lub 1 z prawdopodobieństwem $1 - p_i$). Mechanizmu dropout używamy tylko w procesie treningu. W przypadku inferencji wyjście warstwy \mathbf{A}_i jest z kolei mnożone przez $1 - p_i$, czyli prawdopodobieństwo pozostania neuronu w sieci (*weight scaling inference rule*). Typowo usuwa się dość sporo neuronów w warstwach ukrytych, np. 50%. Dzięki dropoutowi można używać dość dużych stałych uczących oraz współczynników pędu, a także trenować znacznie większe i głębsze sieci. Okazuje się, iż w praktyce takie duże, mocno trenowane sieci z regularyzacją uczą się lepiej niż mniejsze, które by tej regularyzacji nie potrzebowały. Jeśli sieć nie przeucza na zbiorze trenującym to najpewniej lepsze wyniki uzyskamy ucząc większą sieć (która normalnie mogłaby przeuczyć) i stosując dropout. Obecnie dropout jest powszechnie stosowanym sposobem regularyzacji warstw w pełni połączonych (np. ostatnich warstw w niektórych sieciach konwolucyjnych).

- **Transfer learning.**

[*Transfer learning and multi-task learning, Ch.9, [1]*] When training data are limited, other datasets can be exploited to improve performance. In transfer learning, the network is pre-trained to perform a related secondary task for which data are more plentiful. The resulting model is then adapted to the original task. This is typically done by removing the last layer and adding one or more layers that produce a suitable output. The main model may be fixed, and the new layers trained for the original task, or we may fine-tune the entire model. The principle is that the network will build a good internal representation of the data from the secondary task, which can subsequently be exploited for the original task. Equivalently, transfer learning can be

viewed as initializing most of the parameters of the final network in a sensible part of the space that is likely to produce a good solution.

- **Augmentation.** Metoda data augmentation polega na dodaniu do zbioru treningowego replik przykładów z tego zbioru poddanych pewnym transformacjom względem których oczekujemy, iż odpowiedzi sieci będą niezmienione.

[Notes, Ch. 9, [1]]

- **Finding wider minima:** It is thought that wider minima generalize better. Here, the exact values of the weights are less important, so performance should be robust to errors in their estimates. One of the reasons that applying noise to parts of the network during training is effective is that it encourages the network to be indifferent to their exact values. Keskar et al. (2017) showed that SGD finds wider minima as the batch size is reduced. This may be because of the batch variance term that results from implicit regularization by SGD. Ishida et al. (2020) use a technique named flooding, in which they intentionally prevent the training loss from becoming zero. This encourages the solution to perform a random walk over the loss landscape and drift into a flatter area with better generalization.

4.7 Hiperparametry

4.7.1 Dobór hiperparametrów

[*Choosing hyperparameters*, Ch.8, [1]] For a deep network, the model capacity depends on the numbers of hidden layers and hidden units per layer as well as other aspects of architecture that we have yet to introduce. Furthermore, the choice of learning algorithm and any associated parameters (learning rate, etc.) also affects the test performance. These elements are collectively termed hyperparameters. The process of finding the best hyperparameters is termed hyperparameter search or (when focused on network structure) neural architecture search.

Hyperparameters are typically chosen empirically; we train many models with different hyperparameters on the same training set, measure their performance, and retain the best model. However,

we do not measure their performance on the test set; this would admit the possibility that these hyperparameters just happen to work well for the test set but don't generalize to further data. Instead, we introduce a third dataset known as a validation set. For every choice of hyperparameters, we train the associated model using the training set and evaluate performance on the validation set. Finally, we select the model that worked best on the validation set and measure its performance on the test set. In principle, this should give a reasonable estimate of the true performance.

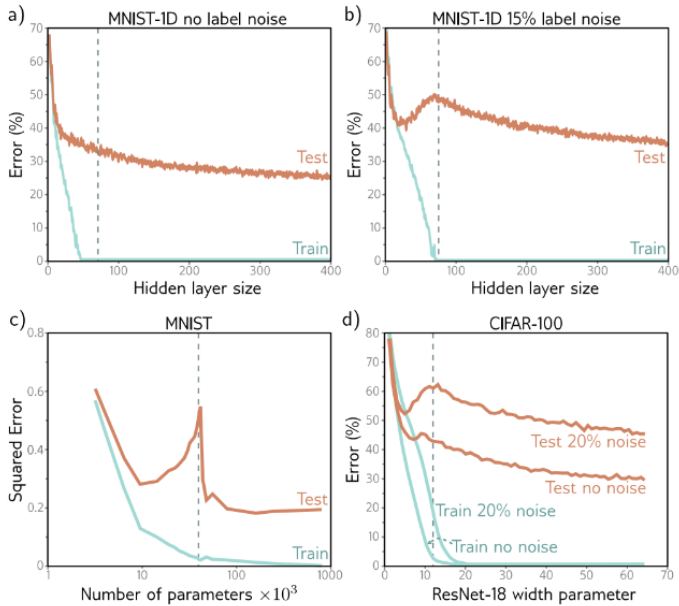
The hyperparameter space is generally smaller than the parameter space but still too large to try every combination exhaustively. Unfortunately, many hyperparameters are discrete (e.g., the number of hidden layers), and others may be conditional on one another (e.g., we only need to specify the number of hidden units in the tenth hidden layer if there are ten or more layers). Hence, we cannot rely on gradient descent methods as we did for learning the model parameters. Hyperparameter optimization algorithms intelligently sample the space of hyperparameters, contingent on previous results. This procedure is computationally expensive since we must train an entire model and measure the validation performance for each combination of hyperparameters.

4.7.2 Zjawisko double descent

Zjawisko double descent polega na nieoczekiwanym zachowaniu wartości błędów popełnianych przez model uczenia maszynowego przy zwiększaniu liczby jego parametrów. Zgodnie z bias-variance tradeoff oczekivalibyśmy bowiem iż w miarę zwiększania liczby parametrów będziemy coraz bardziej przeuczać model (dopasowywać go wprost do danych treningowych), a przez to zmniejszać jego zdolność generalizacji, co przekładałoby się na gorszy wynik na zbiorze testowym. Okazuje się jednak, iż krzywa błędów na zbiorze testowym w zależności od liczby parametrów początkowo spada, następnie (zgodnie z oczekiwaniami) gdy liczba parametrów jest rzędu liczebności zbioru treningowego błąd rośnie, ale (co niespodziewane) gdy liczba parametrów jest znacznie większa od liczebności zbioru treningowego, błąd testowy znowu spada.

[*Double descent*, Ch.8, [1]] This phenomenon is known as double descent. For some datasets like MNIST, it is present with the

original data. For others, like MNIST-1D and CIFAR-100, it emerges or becomes more prominent when we add noise to the labels. The first part of the curve is referred to as the *classical* or *under-parameterized regime*, and the second part as the *modern* or *over-parameterized regime*. The central part where the error increases is termed the critical regime.



Rysunek 4: Przykładowa krzywa wartości błędu na zbiorze testowym w zależności od liczby parametrów modelu wykazująca efekt double descent

[*Double descent*, Ch. 8, [1]] To understand why performance continues to improve as we add more parameters, note that once the model has enough capacity to drive the training loss to near zero, the model fits the training data almost perfectly. This implies that further capacity cannot help the model fit the training data any better; any change must occur between the training points. The tendency of a model to prioritize one solution over another as it extrapolates between data points is known as its inductive bias.

The putative explanation for double descent is that as we add capacity to the model, it interpolates between the nearest data points increasingly smoothly. In the absence of information about what happens between the training points, assuming smoothness is sensible and will probably generalize reasonably to new data.

However, this does not explain why over-parameterized models should produce smooth functions. The answer to this question is uncertain, but there are two likely possibilities. First, the network initialization may encourage smoothness, and the model never departs from the sub-domain of smooth function during the training process. Second, the training algorithm may somehow “prefer” to converge to smooth functions. Any factor that biases a solution toward a subset of equivalent solutions is known as a regularizer, so one possibility is that the training algorithm acts as an implicit regularizer.

4.8 Sieci konwolucyjne

4.9 Sieci rezydualne

4.10 Sieci transformer

Literatura

- [1] Simon J.D. Prince. *Understanding Deep Learning*. The MIT Press, 2023. Available at <http://udlbook.com>.
- [2] Sebastian Ruder. An overview of gradient descent optimization algorithms, 2017. arXiv:1609.04747, available at <https://arxiv.org/abs/1609.04747>.