Spis treści

1	Wstęp 1						
	1.1	Notacja					
	1.2	Uczenie nadzorowane					
	1.3	Uczenie nienadzorowane					
	1.4	Praktyka uczenia maszynowego					
		1.4.1 Przygotowanie danych					
		1.4.2 Metody selekcji cech					
		1.4.3 Metryki do oceny regresji i klasyfikacji 6					
		1.4.4 Tuning hiperparametrów i walidacja skrośna					
	_	1.100.1					
2		babilistyka 10					
	2.1	Zmienne losowe					
	2.2	Ważne rozkłady jednowymiarowe					
	2.3	Wielowymiarowy rozkład normalny					
	2.4	Wnioskowanie statystyczne					
3	Podstawy statystycznego uczenia maszynowego 19						
	3.1	Podstawy wnioskowania bayesowskiego					
	3.2	Modele gaussowskie i liniowe modele gaussowskie					
	3.3	Regresja liniowa					
	3.4	Regularyzacja					
	3.5	Robust regression					
	3.6	Procesy gaussowskie					
	3.7	Klasyfikator najbliższych sąsiadów					
	3.8	Naiwny klasyfikator bayesowski					
	3.9	Estymator jądrowy gęstości					
		Wieloklasowa regresja logistyczna					
		Próbkowanie Monte Carlo łańcuchami Markowa					
	5.11	3.11.1 Algorytm Importance Sampling					
		3.11.2 Algorytm Metropolisa–Hastingsa 38					

1 Wstęp

Celem tych notatek jest zwięzłe 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. $\boldsymbol{x}, \boldsymbol{X}$. Wielkość \boldsymbol{x}^i będzie oznaczać dany element wektora (w tym przypadku ity element \boldsymbol{x}). Wielkość \boldsymbol{x}_μ będzie oznaczać pewien (w tym przypadku μ -ty) element pewnego zbioru wektorów. Macierze oraz wielowymiarowe tablice (zwane również niefortunnie tensorami) będziemy oznaczać (jedynie) wytłuszczonymi wielkimi literami np. $\boldsymbol{X}, \boldsymbol{\Phi}$. Analogicznie jak w przypadku wektorów przez $\boldsymbol{X}^{i_1i_2...i_k}$ będziemy oznaczać (i_1,i_2,\ldots,i_k) element k-wymiarowej tablicy \boldsymbol{X} , natomiast \boldsymbol{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 $\mathcal{X} = \{y_i(\boldsymbol{x}_i)\}_{i=1}^n$, gdzie $\boldsymbol{x} \in \mathbb{R}^m$ nazywamy wektorem cech a $y(\boldsymbol{x})$ 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 \boldsymbol{x} na y? W dalszych paragrafach poznamy różne metody, ale najczęściej modelem jest pewna rodzina funkcji postaci $\phi(\boldsymbol{x};\boldsymbol{w})$ parametryzowana skończoną liczbą parametrów, które możemy łącznie zapisać jako pewien wektor \boldsymbol{w} . Aby znaleźć parametry \boldsymbol{w} , dzięki którym dla konkretnego zagadnienia model będzie zadowalająco odwzorowywał cechy na predykcje (innymi słowy aby nauczyć model) wprowadzamy dodatkowo funkcjonał kosztu (z ang. loss function) $L(\mathcal{X};\boldsymbol{w})$, który kwantyfikuje odpowiedzi modelu ϕ w stosunku do znanych prawidłowych odpowiedzi y dla danych ze zbioru \mathcal{X} . Najczęściej ma on postać

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

gdzie $p(y(\boldsymbol{x}) \mid \boldsymbol{w})$ jest warunkową gęstością prawdopodobieństwa danej obserwacji $y(\boldsymbol{x})$ warunkowaną przez wartość parametrów \boldsymbol{w} . Do powyższego funkcjonału możemy również dodawać tzw. człony regularyzujące (z ang. regularizers). Trening modelu polega wówczas na znalezieniu parametrów \boldsymbol{w}^* , które minimalizują funkcjonał kosztu na zbiorze treningowym \mathcal{X}

$$\boldsymbol{w}^* = \arg\min_{\boldsymbol{w}} L(\mathcal{X}; \boldsymbol{w}). \tag{1.2.1}$$

Zauważmy, że takie podejście ma jedna zasadnicza wade – istotnie nie interesuje nas tak naprawde, 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 dogłębna 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 logarytmicznej, która poprawia stabilność numeryczną, gdy przewidywane wartości są dużymi dodatnimi liczbami rzeczywistymi, zmienia dziedzinę zmiennej objaśnianej z R+ na 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 skonstruowania 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 Metody selekcji cech

Zwykle dane w zbiorach wykorzystywanych w uczeniu maszynowym mają dużą wymiarowość (wiele cech). Nas jednak najbardziej interesuje tzw. *intrinsic dimensionality*, czyli podprzestrzeń, która jest realnie ważna dla problemu klasyfikacji lub regresji. Celem selekcji cech jest właśnie wybór tych cech, które są istotne. Typowo obliczamy ważność cech (z ang. *feature importance*) i usuwamy te najmniej wartościowe. Takie podejście może pomóc nam usunąć szum z danych i poprawić wyniki modeli (jest to szczególne ważne dla modeli, które używają cech wprost np. kNN, Naive Bayes). Najczęstszymi podejściami do selekcji cech są tzw. metody *filter* i *embedded*.

- <u>Metody filter</u>. Waga cechy jest obliczana na podstawie pewnej ogólnej miary jakości cech; są typowo oparte na pewnych statystykach (np. wariancja) i sprawdzają pojedynczą zmienną naraz (tzw. metody *univariate*).
 - Variance threshold. Przeskalowujemy cechy korzystajac z Min-Max, aby wszystkie miały ten sam zakres wartości. Obliczamy wariancję każdej cechy. Odrzucamy cechy o bardzo małej wariancji, oznacza to bowiem, iż są praktycznie stałe. Próg typowo jest ustalany na bardzo niski typu 0.01.
 - Korelacja cech. Obliczamy macierz korelacji i cech, po czym z każdej pary eliminujemy najmocniej skorelowane cechy, poprzez usunięcie tej o niższej wariancji lub tej, która ma średnio większe korelacje z innymi cechami.
 - Korelacja ze zmienną zależną. Obliczamy korelacje między cechami, a zmienną zależną i usuwamy te o wartości bezwzględnej bliskiej 0.
- <u>Metody embedded</u>. W metodach embedded trenujemy jakiś parametryczny model uczenia maszynowego na naszych danych i wagi tego

modelu interpretujemy jako wagi cech w zbiorze. Metody te wymagają dobrych modeli, zauważmy bowiem, że wagi modelu mówią tak naprawdę jak ważna jest dana cecha dla predykcji tego modelu, a nie ogólnie dla zagadnienia. Jako modele można wybrać proste modele liniowe, w których mamy prostą interpretację wag.

1.4.3 Metryki do oceny regresji i klasyfikacji

Wcześniej nie wspominaliśmy nic o metodach sprawdzenia jakości modelu po wytrenowaniu. 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

RMSE :=
$$\sqrt{\frac{1}{n} \sum_{i=1}^{n} (y_i - \hat{y}_i)^2}$$
, MAE := $\frac{1}{n} \sum_{i=1}^{n} |y_i - \hat{y}_i|$. (1.4.1)

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

MAPE :=
$$\frac{1}{n} \sum_{i=1}^{n} \left| \frac{y_i - \hat{y}_i}{y_i} \right|$$
, SMAPE := $\frac{1}{n} \sum_{i=1}^{n} \frac{|y_i - \hat{y}_i|}{|y_i| + |\hat{y}_i|}$. (1.4.2)

Obie te metryki mają zakres od 0 do 1, przy czym niższa wartość oznacza lepszy model. Metryki te mają jednak szereg problemów, z których najpoważniejsze 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*)

MASE :=
$$\frac{\sum_{i=1}^{n} |y_i - \hat{y}_i|}{\sum_{i=1}^{n} |y_i - \overline{y}|},$$
 (1.4.3)

gdzie $\overline{y} = \frac{1}{n} \sum_{i=1}^{n} y_i$. 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 klasa pozytywna i negatywna. 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

$$\text{Accuracy} := \frac{\text{TP} + \text{TN}}{n} \,, \quad \text{Precision} := \frac{\text{TP}}{\text{TP} + \text{FP}} \,, \quad \text{Recall} := \frac{\text{TP}}{\text{TP} + \text{FN}} \,. \tag{1.4.4}$$

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

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

$$F_{\beta} := (1 + \beta^2) \frac{\text{Precision} \cdot \text{Recall}}{\beta^2 \cdot \text{Precision} + \text{Recall}},$$
 (1.4.5)

		Predi		
	[Positive	Negative	
Actual Class	Positive	True Positive (TP)	False Negative (FN) Type II Error	Sensitivity $\frac{TP}{(TP+FN)}$
Actual Class	Negative	False Positive (FP) Type I Error	True Negative (TN)	Specificity $\frac{TN}{(TN+FP)}$
		$\frac{TP}{(TP+FP)}$	Negative Predictive Value $\frac{TN}{(TN+FN)}$	Accuracy $\frac{TP + TN}{(TP + TN + FP + FN)}$

Rysunek 1: Macierz pomyłek oraz możliwe metryki oceny jakości klasyfikatora

gdzie β określa ile razy bardziej ważny jest recall od precision. Typowo używa się $F_1{\rm -score}$

$$F_1 = 2 \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}.$$
 (1.4.6)

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 (FPR(threshold), TPR(threshold)) dla progów z zakresu [0;1], gdzie

$$TPR := \frac{TP}{TP + FN}, \quad FPR := \frac{FP}{FP + TN}. \tag{1.4.7}$$

Metryką niezależną od wybranego progu jest tzw. <u>AUROC</u> (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 nie-

zależnie od wartości progu klasyfikuje wszystkie przykłady poprawnie ma AUROC równy 1.

Inną analogiczną metryką jest $\underline{\text{AUPRC}}$, gdzie zamiast krzywej ROC stosujemy krzywą $\underline{\text{PRC}}$ (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 niezbalansowanej, 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 ktą klasę, FN_k to liczba przykładów z klasy k, które zaklasyfikowaliśmy jako inną klasę. Wówczas odpowiednie metryki mają postać

$$\begin{split} & \text{MicroPrecision} := \frac{\sum_{k} \text{TP}_{k}}{\sum_{k} \text{TP}_{k} + \sum_{k} \text{FP}_{k}} \,, \\ & \text{MacroPrecision} := \frac{1}{K} \sum_{k=1}^{K} \frac{\text{TP}_{k}}{\text{TP}_{k} + \text{FP}_{k}} \end{split} \tag{1.4.8}$$

oraz

$$\begin{aligned} & \text{MicroRecall} := \frac{\sum_{k} \text{TP}_{k}}{\sum_{k} \text{TP}_{k} + \sum_{k} \text{FN}_{k}}, \\ & \text{MacroRecall} := \frac{1}{K} \sum_{k=1}^{K} \frac{\text{TP}_{k}}{\text{TP}_{k} + \text{FN}_{k}}. \end{aligned} \tag{1.4.9}$$

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

1.4.4 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

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

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

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

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

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

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

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

Zauważmy, że w przypadku rozkładu ciągłego zachodzi

$$F(\boldsymbol{x}) = \int_{-\infty}^{\boldsymbol{x}^1} \cdots \int_{-\infty}^{\boldsymbol{x}^n} p(\boldsymbol{x}') \, \mathrm{d}^n \boldsymbol{x}' \ .$$

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

rem

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

odpowiednio dla rozkładu dyskretnego i ciągłego.

Definicja 2.4 (Macierzy kowariancji). Macierzą kowariancji funkcji zmiennej losowej f(X) nazywamy macierz

$$\mathbb{E}\left[(\boldsymbol{f}(\boldsymbol{x})-\boldsymbol{m_f})(\boldsymbol{f}(\boldsymbol{x})-\boldsymbol{m_f})^T\right]\,,$$

gdzie

$$m_f = \mathbb{E}\left[f(x)\right]$$
.

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 dystrybuantą 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 X ma rozkład ciągły o gęstości p_X i niech $Y^i = \varphi^i(X)$ dla $i = 1, \ldots, n$. Jeśli odwzorowanie φ jest różniczkowalne i odwracalne, przy czym odwzorowanie odwrotne $\psi = \varphi^{-1}$ jest różniczkowalne, to n-wymiarowa zmienna Y ma rozkład o gestości

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

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

Twierdzenie 2.2.

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

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

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

$$p(\boldsymbol{x} \mid \boldsymbol{y}) = \frac{p(\boldsymbol{y} \mid \boldsymbol{x}) p_{\boldsymbol{X}}(\boldsymbol{x})}{\int_{\mathbb{R}^k} p(\boldsymbol{y} \mid \boldsymbol{x}) p_{\boldsymbol{X}}(\boldsymbol{x}) \, \mathrm{d}^k \boldsymbol{x}} \quad \text{(Bayes theorem)}$$

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

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

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. $X\sim {\rm Pos}(\lambda).$ Jeśli $X\sim {\rm Pos}(\lambda),$ to

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

Twierdzenie 2.3 (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

$$\lim_{n \to \infty} n p_n = \lambda$$

dla pewnej liczby $\lambda > 0$. Wówczas

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

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 $X \sim \operatorname{Pos}(\lambda t)$. 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

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

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

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

2.3 Wielowymiarowy rozkład normalny

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

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

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

$$oldsymbol{X} = oldsymbol{A}oldsymbol{Z} + oldsymbol{\mu}$$
 .

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

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

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

Poziomice gęstości niezdegenerowanego wielowymiarowego rozkładu normalnego są elipsoidami, których półosie są skierowane wzdłuż wektorów własnych macierzy Σ i mają długości proporcjonalne do pierwiastka z wartości własnych.

Twierdzenie 2.4 (Własności niezdegenerowanego rozkładu normalnego). Niech $X \sim \mathcal{N}(\mu, \Sigma)$ dla dodatnio określonej macierzy Σ , wówczas

- 1. Wszystkie rozkłady brzegowe i warunkowe \boldsymbol{X} są rozkładami normalnymi.
- 2. Zmienne składowe X_1, \ldots, X_n są niezależne wtedy i tylko wtedy, gdy Σ jest macierzą diagonalną.

2.4 Wnioskowanie statystyczne

Niech zmienna losowa \boldsymbol{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 \boldsymbol{X} . Do podstawowych zagadnień wnioskowania statystycznego należą:

- oszacowanie wielkości charakteryzujących rozkład \boldsymbol{X} (np. wartości średniej albo wariancji);
- weryfikacja hipotez dotyczących rozkładu \boldsymbol{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 prawdopodobieństwa 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 X_1, \ldots, 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 $\hat{\theta}(\boldsymbol{X}_1,\ldots,\boldsymbol{X}_n)$ służącą do oszacowania wartości parametru θ . Liczbę $\hat{\theta}(\boldsymbol{x}_1,\ldots,\boldsymbol{x}_n)$ 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

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

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

Niech X_1,\dots,X_n będzie prostą próbą losową. Jeśli $p(\cdot\mid \pmb{\theta})$ opisuje rozkład warunkowy, z którego pochodzą obserwacje, to

$$\mathcal{L}(\boldsymbol{x}_1, \dots, \boldsymbol{x}_n; \boldsymbol{\theta}) = \prod_{i=1}^n p(\boldsymbol{x}_i \mid \boldsymbol{\theta}).$$
 (2.4.1)

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

$$L(\mathbf{x}; \boldsymbol{\theta}) = -\log \mathcal{L}(\mathbf{x}; \boldsymbol{\theta}). \tag{2.4.2}$$

Wówczas dla realizacji prostej próby losowej mamy

$$L(\boldsymbol{x}_1, \dots, \boldsymbol{x}_n; \boldsymbol{\theta}) = -\sum_{i=1}^n \log p(\boldsymbol{x}_i \mid \boldsymbol{\theta}).$$
 (2.4.3)

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) $\{\boldsymbol{x}_1,\ldots,\boldsymbol{x}_n\}$ maksymalizuje wartość funkcji wiarygodności lub, co równoważne, minimalizuje wartość zanegowanej logarytmicznej funkcji wiarygod-

ności tj.

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

Jeśli funkcja wiarygodności jest różniczkowalna względem θ dla dowolnych x^i , to MLE można czasem wyznaczyć analitycznie korzystając z warunku koniecznego optymalności, tzn. rozwiązując układ równań

$$\frac{\partial L(\mathbf{x}_1, \dots, \mathbf{x}_n; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = 0.$$
 (2.4.4)

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

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 Podstawy wnioskowania bayesowskiego

Niech $\mathcal{X} = \{x_1, \dots, x_n\}$ będzie realizacją prostej próby losowej, czyli inaczej zbiorem obserwacji i.i.d. Zakładamy, że obserwacje te pochodzą z pewnego parametrycznego modelu statystycznego \mathcal{P} z parametrami $\boldsymbol{\theta} \in \Theta$. Wcześniej $\boldsymbol{\theta}$ miał jedynie rangę parametru. We wnioskowaniu bayesowskim uznajemy, że parametry $\boldsymbol{\theta}$ są również zmiennymi losowymi, a model statystyczny modeluje warunkowy rozkład prawdopodobieństwa obserwacji pod warunkiem parametru. Mamy zatem rodzinę gęstości prawdopodobieństwa $p(\boldsymbol{x} \mid \boldsymbol{\theta})$ i chcemy wnioskować o parametrze $\boldsymbol{\theta}$ na podstawie obserwacji \mathcal{X} . Jeśli znamy rozkład a priori (inaczej <u>prior</u>) parametru $\boldsymbol{\theta}$ opisany przez $p(\boldsymbol{\theta})$, to z twierdzenia Bayesa rozkład a posteriori (posterior) jest dany przez

$$p(\boldsymbol{\theta} \mid \mathcal{X}) = \frac{p(\mathcal{X} \mid \boldsymbol{\theta})p(\boldsymbol{\theta})}{p(\mathcal{X})} = \frac{p(\mathcal{X} \mid \boldsymbol{\theta})p(\boldsymbol{\theta})}{\sum_{\boldsymbol{\theta}' \in \Theta} p(\mathcal{X} \mid \boldsymbol{\theta}')p(\boldsymbol{\theta}')}$$
(3.1.1)

Nie jesteśmy tutaj zbyt formalni z notacją, gdyż używamy p bardziej jakby był to rozkład prawdopodobieństwa (miara probabilistyczna), a w rzeczywistości jest to gęstość lub funkcja prawdopodobieństwa, więc powinniśmy używać indeksów dolnych określających zmienną losową, której rozkład jest opisany danym p, aby rozróżnić człony od siebie. Zapis taki jest jednak niezwykle wygodny i dość czytelny. Należy jedynie pamiętać, iż nazwa argumentu funkcji p określa teraz z jakiego wzoru powinniśmy skorzystać aby obliczyć jej wartość. Człon w mianowniku postaci

$$Z = \sum_{\boldsymbol{\theta}' \in \Theta} p(\mathcal{X} \mid \boldsymbol{\theta}') p(\boldsymbol{\theta}') \cong \int_{\Omega} p(\mathcal{X} \mid \boldsymbol{\theta}') p(\boldsymbol{\theta}') d^{n} \boldsymbol{\theta}'$$
(3.1.2)

jest tzw. czynnikiem normalizacyjnym (czyli po prostu liczbą, często oznaczaną przez Z), który zapewnia, iż $p(\theta \mid \mathcal{X})$ sumuje / całkuje się do 1.

Ponieważ założyliśmy, iż obserwacje ze zbioru \mathcal{X} są warunkowo niezależne względem parametru $\boldsymbol{\theta}$ oraz pochodzą z tego samego rozkładu opisanego przez $p(\boldsymbol{x}\mid\boldsymbol{\theta})$, więc człon $p(\mathcal{X}\mid\boldsymbol{\theta})$ zwany wiarygodnością możemy zapisać jako

$$p(\mathcal{X} \mid \boldsymbol{\theta}) = \prod_{i=1}^{n} p(\boldsymbol{x}_i \mid \boldsymbol{\theta}). \tag{3.1.3}$$

Całe wnioskowanie bayesowskie opiera się na wyznaczeniu rozkładu a posteriori dla danego zbioru obserwacji \mathcal{X} , który wyraża naszą wiedzę o estymowanym parametrze $\boldsymbol{\theta}$. Na podstawie tego rozkładu możemy wyznaczyć estymatę punktową MAP maksymalizującą gęstość prawdopodobieństwa a posteriori,

$$\hat{\boldsymbol{\theta}}_{\text{MAP}}(\mathcal{X}) = \arg \max_{\boldsymbol{\theta} \in \Theta} p(\boldsymbol{\theta} \mid \mathcal{X}),$$
 (3.1.4)

jak również niepewność związaną z wyznaczeniem tej estymaty np. poprzez wyznaczenie przedziału wiarygodności (nie należy mylić z przedziałem ufności). Możemy również skonstruować rozkład predykcyjny (z ang. posterior predictive distribution) określający prawdopodobieństwo uzyskania nowej obserwacji \boldsymbol{t}

$$p(\mathbf{t} \mid \mathcal{X}) = \int_{\Theta} p(\mathbf{t} \mid \boldsymbol{\theta}) p(\boldsymbol{\theta} \mid \mathcal{X}) d^{k} \boldsymbol{\theta} .$$
 (3.1.5)

Znając rozkład a posteriori estymowanego parametru θ możemy nie tylko wyznaczyć estymaty punktowe, wartości oczekiwane i przedziały wiarygodności, ale również znaleźć estymator Bayesa (z ang. Bayes estimator), który minimalizuje wartość oczekiwaną pewnej funkcji straty (z ang. loss

function) $L(\boldsymbol{\theta}, \hat{\boldsymbol{\theta}})$ po wszystkich estymatorach $\hat{\boldsymbol{\theta}}$

$$\hat{\boldsymbol{\theta}}_{\text{Bayes}}(\mathcal{X}) = \arg\min_{\hat{\boldsymbol{\theta}} \in \Theta} \int_{\Theta} L(\boldsymbol{\theta}, \hat{\boldsymbol{\theta}}) p(\boldsymbol{\theta} \mid \mathcal{X}) \, d^k \boldsymbol{\theta} . \tag{3.1.6}$$

Całkę w powyższym wzorze nazywa się również funkcją ryzyka (z ang. risk function) $R(\hat{\boldsymbol{\theta}})$, która określa oczekiwaną stratę spowodowaną wykorzystaniem danego estymatora parametru.

Są dwa zasadnicze problemy we wnioskowaniu bayesowskim: pierwszym jest potrzeba znania rozkładu a priori estymowanego parametru, drugim – problem z obliczeniem czynnika normalizującego, który może być skomplikowaną całką lub sumą po wykładniczo-wielu elementach. Oba te problemy można czasami rozwiązać wprowadzając tzw. prior sprzężony do wiarygodności, tzn. zakładamy taki rozkład a priori, aby dla danej wiarygodności rozkład a posteriori miał znaną postać (np. rozkładu normalnego, rozkładu beta), wówczas nie musimy obliczać czynnika normalizującego, gdyż jest on po prostu znany.

3.2 Modele gaussowskie i liniowe modele gaussowskie

Zajmiemy się teraz wnioskowaniem bayesowskim w modelach, w których potrafimy analitycznie znaleźć postać rozkład a posteriori. Jak już wspomnieliśmy (Tw. 2.4) gdy zmienna losowa X ma wielowymiarowy rozkład normalny z wartością oczekiwaną μ i dodatnio określoną macierzą kowariancji Σ to wszystkie rozkłady warunkowe i brzegowe są rozkładami normalnymi. Poniżej wyznaczymy parametry tych rozkładów.

Twierdzenie 3.1. Niech zmienne losowe $x \in \mathbb{R}^{n-k}$ i $y \in \mathbb{R}^k$ mają łącznie wielowymiarowy rozkład normalny

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

Wówczas

1. Zmienne losowe \boldsymbol{x} i \boldsymbol{y} mają odpowiednio rozkłady

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

2. Rozkład warunkowy $x \mid y$ jest rozkładem normalnym $x \mid y \sim \mathcal{N}(\mu_{x\mid y}, \Sigma_{x\mid y})$ o parametrach

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

Powyższe własności normalnych rozkładów łącznych pozwalają jawnie wnioskować w tzw. liniowych modelach gaussowskich (z ang. Linear Gaussian Models). Załóżmy, iż nasze obserwacje są modelowane zmienną losową \boldsymbol{y} o rozkładzie normalnym z estymowanym parametrem \boldsymbol{x} i znanymi parametrami $\boldsymbol{A}, \boldsymbol{b}, \boldsymbol{\Sigma}_{\boldsymbol{y}}$ tak, że wiarygodność ma postać

$$y \mid x \sim \mathcal{N}(Ax + b, \Sigma_y),$$
 (3.2.1)

gdzie \boldsymbol{A} jest w ogólności macierzą prostokątną. Jeśli jako prior na \boldsymbol{x} przyjmiemy również rozkład normalny

$$x \sim \mathcal{N}(\mu_x, \Sigma_x)$$
 (3.2.2)

to rozkład a posteriori jest również rozkładem normalnym. W szczególności załóżmy, że mamy zbiór obserwacji i.i.d. $\mathcal{X}=\{\boldsymbol{y}^1,\ldots,\boldsymbol{y}^n\}$. Wówczas wiarygodność ma postać

$$p(\mathcal{X} \mid \boldsymbol{x}) \propto \prod_{i=1}^{n} \exp\left(-\frac{1}{2}(\boldsymbol{y}^{i} - (\boldsymbol{A}\boldsymbol{x} + \boldsymbol{b}))^{T} \boldsymbol{\Sigma}_{\boldsymbol{y}}^{-1}(\boldsymbol{y}^{i} - (\boldsymbol{A}\boldsymbol{x} + \boldsymbol{b}))\right), \quad (3.2.3)$$

a rozkład a posteriori

$$p(\boldsymbol{x} \mid \mathcal{X}) \propto p(\mathcal{X} \mid \boldsymbol{x}) \exp\left(-\frac{1}{2}(\boldsymbol{x} - \boldsymbol{\mu_x})^T \boldsymbol{\Sigma_x}^{-1} (\boldsymbol{x} - \boldsymbol{\mu_x})\right).$$
 (3.2.4)

Rozpisując wszystkie czynniki i pomijając czynnik stałe otrzymujemy

$$\log p(\boldsymbol{x} \mid \mathcal{X}) \propto -\frac{1}{2} \boldsymbol{x}^{T} \left(\boldsymbol{\Sigma}_{\boldsymbol{x}}^{-1} + n \boldsymbol{A}^{T} \boldsymbol{\Sigma}_{\boldsymbol{y}}^{-1} \boldsymbol{A} \right) \boldsymbol{x}$$
$$+ \boldsymbol{x}^{T} \left(\boldsymbol{\Sigma}_{\boldsymbol{x}}^{-1} \boldsymbol{\mu}_{\boldsymbol{x}} + \boldsymbol{A}^{T} \boldsymbol{\Sigma}_{\boldsymbol{y}}^{-1} \left[\sum_{i=1}^{n} \boldsymbol{y}^{i} - n \boldsymbol{b} \right] \right)$$
(3.2.5)

skąd widzimy, iż rozkład $x \mid \mathcal{X}$ jest rozkładem normalnym o parametrach

$$\begin{bmatrix} \boldsymbol{\mu}_{\boldsymbol{x}|\mathcal{X}} = \boldsymbol{\Sigma}_{\boldsymbol{x}|\mathcal{X}} \left(\boldsymbol{\Sigma}_{\boldsymbol{x}}^{-1} \boldsymbol{\mu}_{\boldsymbol{x}} + \boldsymbol{A}^T \boldsymbol{\Sigma}_{\boldsymbol{y}}^{-1} \left[\sum_{i=1}^n \boldsymbol{y}^i - n \boldsymbol{b} \right] \right), \\ \boldsymbol{\Sigma}_{\boldsymbol{x}|\mathcal{X}} = \left(\boldsymbol{\Sigma}_{\boldsymbol{x}}^{-1} + n \boldsymbol{A}^T \boldsymbol{\Sigma}_{\boldsymbol{y}}^{-1} \boldsymbol{A} \right)^{-1}. \end{cases}$$
(3.2.6)

3.3 Regresja liniowa

Załóżmy, iż modelujemy obserwacje postaci (y, \boldsymbol{x}) , gdzie y to skalar zwany zmienną objaśnianą, którego wartość obserwujemy, a \boldsymbol{x} to wektor zmiennych objaśniających, co do którego zakładamy, iż dla danego pomiaru y jest

on znany dokładnie. Dodatkowo załóżmy liniowy model $\hat{y}(x; w)$ postaci

$$\hat{y}(\boldsymbol{x}; \boldsymbol{w}) = \boldsymbol{x}^T \boldsymbol{w}, \qquad (3.3.1)$$

w którym mamy dodatkowo <u>błąd losowy</u> $\epsilon \sim \mathcal{N}(0, \sigma^2)$ z nieznanym σ . Możemy zatem napisać model statystyczny postaci

$$y(\mathbf{x}) \mid \mathbf{w}, \sigma \sim \mathcal{N}(\hat{y}(\mathbf{x}; \mathbf{w}), \sigma^2),$$
 (3.3.2)

gdzie \boldsymbol{w} , σ to estymowane parametry. Powiedzmy, iż mamy zbiór obserwacji i.i.d. $\mathcal{X} = \{y_1(\boldsymbol{x}_1), \dots, y_n(\boldsymbol{x}_n)\}$. Wiarygodność ma zatem postać

$$\mathcal{L}(\mathcal{X}; \boldsymbol{w}, \sigma) = \frac{1}{(2\pi\sigma^2)^{n/2}} \prod_{i=1}^n \exp\left(-\frac{1}{2\sigma^2} \left[y_i - \hat{y}(\boldsymbol{x}_i; \boldsymbol{w})\right]^2\right).$$
(3.3.3)

W przypadku regresji liniowej zamiast pełnego wnioskowania bayesowskiego często stosuje się prostsze podejście polegające na ograniczeniu się do znalezienia estymaty punktowej MLE. Zanegowana logarytmiczna funkcja wiarygodności, którą będziemy również nazywać funkcją kosztu ma postać (pomijamy człony stałe, gdyż nie są one istotne przy dalszej minimalizacji)

$$L(\mathcal{X}; \boldsymbol{w}, \sigma) = n \log \sigma + \frac{1}{2\sigma^2} \sum_{i=1}^{n} [y_i - \hat{y}(\boldsymbol{x}_i; \boldsymbol{w})]^2 + \text{const.}$$
 (3.3.4)

Minimalizując funkcję L względem \boldsymbol{w} i σ otrzymamy estymaty MLE tych parametrów. Dla ustalonego stałego σ otrzymana funkcja L ma postać formy kwadratowej i otrzymany przy takim uproszczeniu problem optymalizacyjny nazywamy metodą najmniejszych kwadratów (z ang. Ordinary Least Squares, OLS). W przypadku modelu liniowego estymatory można znaleźć analitycznie rozwiązując układ równań

$$\frac{\partial L}{\partial \boldsymbol{w}^{j}} = \frac{1}{2\sigma^{2}} \frac{\partial}{\partial \boldsymbol{w}^{j}} \sum_{i=1}^{n} \left[y_{i} - \boldsymbol{w}^{T} \boldsymbol{x}_{i} \right]^{2} = -\frac{1}{\sigma^{2}} \sum_{i=1}^{n} (y_{i} - \boldsymbol{w}^{T} \boldsymbol{x}_{i}) \boldsymbol{x}_{i}^{j} = 0,
\frac{\partial L}{\partial \sigma} = \frac{n}{\sigma} - \frac{1}{\sigma^{3}} \sum_{i=1}^{n} \left[y_{i} - \hat{y}(\boldsymbol{x}_{i}; \boldsymbol{w}) \right]^{2} = 0.$$
(3.3.5)

Z powyższego

$$\sum_{i=1}^{n} y_i \boldsymbol{x}_i^j - \boldsymbol{w}_{\text{MLE}}^T \sum_{i=1}^{n} \boldsymbol{x}_i \boldsymbol{x}_i^j = 0,$$

$$\sigma_{\text{MLE}}^2 = \frac{1}{n} \sum_{i=1}^{n} \left[y_i - \hat{y}(\boldsymbol{x}_i; \boldsymbol{w}_{\text{MLE}}) \right]^2.$$
(3.3.6)

Wprowadzając wektor $\boldsymbol{y}^i := y_i$ oraz macier
z $\boldsymbol{X}^{ij} := \boldsymbol{x}_i^j$ możemy zapisać pierwsze równanie jako

$$-\boldsymbol{y}^T \boldsymbol{X} + \boldsymbol{w}_{\text{MLE}}^T \boldsymbol{X}^T \boldsymbol{X} = 0, \qquad (3.3.7)$$

skad

$$\begin{bmatrix}
\boldsymbol{w}_{\text{MLE}} = (\boldsymbol{X}^T \boldsymbol{X})^{-1} \boldsymbol{X}^T \boldsymbol{y} = \boldsymbol{X}^+ \boldsymbol{y}, \\
\sigma_{\text{MLE}}^2 = \frac{1}{n} (\boldsymbol{y} - \boldsymbol{X} \boldsymbol{w}_{\text{MLE}})^T (\boldsymbol{y} - \boldsymbol{X} \boldsymbol{w}_{\text{MLE}}),
\end{bmatrix} (3.3.8)$$

gdzie X^+ oznacza pseudoodwrotność Moore'a–Penrose'a, którą można efektywnie obliczyć korzystając z rozkładu SVD macierzy X.

3.4 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 $\gamma \| \boldsymbol{w} \|$ dla pewnej normy $\| \cdot \|$ i hiperparametru γ . W kontekście bayesowskim regularyzację można również rozumieć jako pewną "niechęć" (tłumienie, zachowawczość) modelu do zmiany rozkładu a priori estymowanego parametru.

W przypadku regresji liniowej jeśli zamiast poszukiwania estymaty MLE będziemy poszukiwać estymaty MAP (z ang. Maximum a Posteriori estimate) z rozkładem a priori na parametr \boldsymbol{w} danym przez $\boldsymbol{w} \sim \mathcal{N}(\boldsymbol{0}, \tau^2 \boldsymbol{1})$, to logarytm gęstości rozkładu a posteriori (który również będziemy nazywać zregularyzowaną funkcją kosztu) ma postać (tutaj zakładamy, że σ jest znaną stałą)

$$L(\mathcal{X}; \boldsymbol{w}) = \frac{1}{2\sigma^2} \sum_{i=1}^{n} \left[y_i - \hat{y}(\boldsymbol{x}_i; \boldsymbol{w}) \right]^2 + \frac{1}{2\tau^2} \boldsymbol{w}^T \boldsymbol{w} + \text{const.}$$
 (3.4.1)

skąd możemy bez problemy wyznaczyć estymatę punktową MAP parametru ${m w}$

$$\boldsymbol{w}_{\text{MAP}} = \left(\gamma \mathbf{1} + \boldsymbol{X}^T \boldsymbol{X}\right)^{-1} \boldsymbol{X}^T \boldsymbol{y}, \qquad (3.4.2)$$

gdzie $\gamma=\frac{\sigma^2}{\tau^2}$ nazywamy siłą regularyzacji. Zagadnienie minimalizacji funkcji kosztu będącej formą kwadratową z dodanym członem regularyzującym w postaci sumy kwadratów współrzędnych wektora \boldsymbol{w} (normy L2 wektora)

nazywamy regresją grzbietową (z ang. $ridge\ regression$), natomiast taką postać członu regularyzującego – regularyzacją L2. Zauważmy, że im większa jest wartość γ (mniejsza niepewność związana z rozkładem a priori) tym drugi człon w nawiasie staje się mniej istotny.

Innym przykładem regularyzacji jest tzw. regularyzacja L1, która polega na dodaniu do funkcji kosztu członu postaci $\gamma \sum_{j=1}^d |\boldsymbol{w}^j|$ tj. normy L1 wektora wag. Zagadnie 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 (tutaj parametr σ nie występuje, tj. ta funkcja kosztu nie ma bezpośredniej interpretacji probabilistycznej jako funkcja wiarygodności)

$$L(\mathcal{X}; \boldsymbol{w}) = \frac{1}{2} \sum_{i=1}^{n} [y_i - \hat{y}(\boldsymbol{x}_i; \boldsymbol{w})]^2 + \frac{\gamma_1}{2} \|\boldsymbol{w}\|_1 + \frac{\gamma_2}{2} \|\boldsymbol{w}\|_2^2.$$
 (3.4.3)

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.

3.5 Robust regression

Zastanówmy się wpierw na czym tak naprawdę polega modelowanie rozkładu warunkowego $y(x) \mid w$ za pomocą określonego rozkładu prawdopodobieństwa. Można by było pomyśleć, 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(x), gdyż mamy tylko po jednej wartości y dla danego x. Modelując y(x) rozkładem normalnym chodzi nam zatem raczej o to, że w takim modelu chcemy znaleźć parametr (prostą) taką, że punkty w zbiorze są po odpowiednich stronach owej prostej zgodnie z gęstością prawdopodobieństwa, która jest symetryczna i ma właśnie kształt dzwonu, tj. dużo masy prawdopodobieństwa jest zebrane w niewielkiej odległości od estymowanej prostej.

W takim ujęciu możemy zakładać różne inne rozkłady na y(x) jeśli interesują nas proste, które inaczej mają rozdzielać masę prawdopodobieństwa

między punkty. 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 tzw. 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(x) przez tzw. asymetryczny rozkład Laplace'a (z ang. Asymmetric Laplace Distribution, \overline{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

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

Zauważmy, że dystrybuanta rozkładu ALD ma postać

$$F(x; m, \lambda, q) = \begin{cases} q e^{\frac{1-q}{\lambda}(x-m)}, & x \le m \\ 1 - (1-q)e^{-\frac{q}{\lambda}(x-m)}, & x \ge m \end{cases}$$
(3.5.1)

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

$$y(\boldsymbol{x}) \mid \boldsymbol{w}, \lambda \sim \text{ALD}(\hat{y}(\boldsymbol{x}; \boldsymbol{w}), \lambda, q),$$
 (3.5.2)

gdzie λ pełni podobną rolę jak σ w przypadku rozkładu normalnego. Widzimy wówczas, iż estymacja MLE \boldsymbol{w} daje najlepsze \hat{y} takie, 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

$$L(\mathcal{X}; \boldsymbol{w}, \lambda) = n \log \lambda + \frac{1}{\lambda} \sum_{i=1}^{n} \left[(q-1)z_i \theta(-z_i) + q z_i \theta(z_i) \right], \qquad (3.5.3)$$

gdzie

$$z_i := y_i - \hat{y}(\boldsymbol{x}_i; \boldsymbol{w}), \qquad (3.5.4)$$

a θ oznacza funkcję skokową Heaviside'a. W przypadku ustalonego, stałego λ taką funkcję kosztu nazywamy pinball loss. Modelowanie y(x) za

pomocą rozkładu ALD pozwala nam w prosty i "robust" sposób znaleźć również 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 (tzw. "widełki") będące oszacowaniem niepewności punktowej estymaty.

3.6 Procesy gaussowskie

Jak już wspomnieliśmy macierz kowariancji n-wymiarowej zmiennej losowej x o wartości oczekiwanej μ jest zdefiniowana jako

$$\Sigma = \mathbb{E}\left[(\boldsymbol{x} - \boldsymbol{\mu})(\boldsymbol{x} - \boldsymbol{\mu})^T \right]. \tag{3.6.1}$$

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

$$K = LL^T, (3.6.2)$$

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

$$\mathbb{E}\left[(\boldsymbol{L}\boldsymbol{z})(\boldsymbol{L}\boldsymbol{z})^{T}\right] = \mathbb{E}\left[\boldsymbol{L}\boldsymbol{z}\boldsymbol{z}^{T}\boldsymbol{L}^{T}\right] = \boldsymbol{L}\mathbb{E}[\boldsymbol{z}\boldsymbol{z}^{T}]\boldsymbol{L}^{T} = \boldsymbol{L}\boldsymbol{1}\boldsymbol{L}^{T} = \boldsymbol{K}. \quad (3.6.3)$$

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 $\forall m \in \mathbb{N} : \forall X = \{x_1, \dots, x_m\} \subset \mathbb{R}^n$ macierz

$$k(X,X) = \begin{bmatrix} k(x_1, x_1) & k(x_1, x_2) & \cdots & k(x_1, x_m) \\ k(x_2, x_1) & k(x_2, x_2) & \cdots & k(x_2, x_m) \\ \vdots & \vdots & \ddots & \vdots \\ k(x_m, x_1) & k(x_m, x_2) & \cdots & k(x_m, x_m) \end{bmatrix}$$

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

Dla dwóch zbiorów punktów $X = \{x_1, \dots, x_m\} \subset \mathbb{R}^n$ i $Y = \{y_1, \dots, y_s\} \subset \mathbb{R}^n$ i funkcji kowariancji k wprowadzimy oznaczenie

$$k(X,Y) := \begin{bmatrix} k(\boldsymbol{x}_1, \boldsymbol{y}_1) & k(\boldsymbol{x}_1, \boldsymbol{y}_2) & \cdots & k(\boldsymbol{x}_1, \boldsymbol{y}_s) \\ k(\boldsymbol{x}_2, \boldsymbol{y}_1) & k(\boldsymbol{x}_2, \boldsymbol{y}_2) & \cdots & k(\boldsymbol{x}_2, \boldsymbol{y}_s) \\ \vdots & \vdots & \ddots & \vdots \\ k(\boldsymbol{x}_m, \boldsymbol{y}_1) & k(\boldsymbol{x}_m, \boldsymbol{y}_2) & \cdots & k(\boldsymbol{x}_m, \boldsymbol{y}_s) \end{bmatrix}.$$
(3.6.4)

Poniżej podajemy kilka przykładów funkcji kowariancji

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

$$k(\mathbf{x}, \mathbf{y}) = a^2 \exp\left\{-\frac{1}{2l^2} \|\mathbf{x} - \mathbf{y}\|^2\right\}$$
 (3.6.5)

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

$$k(\boldsymbol{x}, \boldsymbol{y}) = a^2 \exp\left\{-\frac{2}{l^2} \sin^2\left(\frac{\pi}{p} \|\boldsymbol{x} - \boldsymbol{y}\|\right)\right\}$$
(3.6.6)

• White noise kernel dla hiper-parametru σ

$$k(\boldsymbol{x}, \boldsymbol{y}) = \sigma^2 \delta_{\boldsymbol{x}, \boldsymbol{y}} \tag{3.6.7}$$

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

$$k(\boldsymbol{x}, \boldsymbol{y}) = a^2 \frac{2^{1-\nu}}{\Gamma(\nu)} \left(\frac{\sqrt{2\nu}}{l} \|\boldsymbol{x} - \boldsymbol{y}\| \right)^{\nu} K_{\nu} \left(\frac{\sqrt{2\nu}}{l} \|\boldsymbol{x} - \boldsymbol{y}\| \right), \quad (3.6.8)$$

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

Twierdzenie 3.2 (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 skalarnych zmiennych losowych indeksowanych przez punkty $\boldsymbol{x} \in \mathbb{R}^n$

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

taką że każdy skończony podzbiór \mathcal{GP} ma łącznie wielowymiarowy rozkład normalny tj. dla dowolnego zbioru $X = \{x_1, \dots, x_m\} \subset \mathbb{R}^n$ zachodzi

 $egin{bmatrix} f_{oldsymbol{x}_m} \ \vdots \ f_{oldsymbol{x}_m} \end{bmatrix} \sim \mathcal{N}(oldsymbol{\mu}_X, oldsymbol{\Sigma}_X) \,.$

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

$$\begin{bmatrix} f_{\boldsymbol{x}_1} \\ \vdots \\ f_{\boldsymbol{x}_m} \end{bmatrix} \sim \mathcal{N}(\boldsymbol{0}, k(X, X)). \tag{3.6.9}$$

Process Gaussowski daje nam w praktyce rozkład prawdopodobieństwa nad funkcjami $f: \mathbb{R}^n \to \mathbb{R}$, 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 $X = \{x_1, x_2, \dots, x_m\}$ 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

$$\mathbf{f}_{X} = \begin{bmatrix} f_{\mathbf{x}_{1}} \\ \vdots \\ f_{\mathbf{x}_{m}} \end{bmatrix} \sim \mathcal{N}(\mathbf{0}, k(X, X)). \tag{3.6.10}$$

Powiedzmy, iż chcemy znać wartości f_Y tej funkcji w zadanych punktach $Y = \{y_1, y_2, \dots, y_s\}$. Ponieważ założyliśmy, iż wartości funkcji pochodzą z procesu Gaussowskiego, więc rozkład łączny f_X i f_Y jest rozkładem nor-

malnym

$$\begin{bmatrix} \mathbf{f}_X \\ \mathbf{f}_Y \end{bmatrix} \sim \mathcal{N} \left(\mathbf{0}, \begin{bmatrix} k(X, X) & k(X, Y) \\ k(Y, X) & k(Y, Y) \end{bmatrix} \right). \tag{3.6.11}$$

Zauważmy, iż jest to instancja modelu gaussowskiego (Tw. 3.1), więc rozkład warunkowy $f_Y \mid f_X$ jest również rozkładem normalnym o parametrach

$$\mu = k(Y, X)k^{-1}(X, X)\mathbf{f}_{X} \Sigma = k(Y, Y) - k(Y, X)k^{-1}(X, X)k(X, Y)$$
(3.6.12)

Dodatkową niepewność związaną z pomiarem wartości f_X możemy uchwycić zmieniając postać jądra

$$k(\boldsymbol{x}, \boldsymbol{y}) \leftarrow k(\boldsymbol{x}, \boldsymbol{y}) + \mathcal{I}_X(\boldsymbol{x}) \sigma^2 \delta_{\boldsymbol{x}, \boldsymbol{y}},$$
 (3.6.13)

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

$$\begin{bmatrix}
\boldsymbol{\mu} = k(Y, X) \left[k(X, X) + \sigma^2 \mathbf{1} \right]^{-1} \boldsymbol{f}_X \\
\boldsymbol{\Sigma} = k(Y, Y) - k(Y, X) \left[k(X, X) + \sigma^2 \mathbf{1} \right]^{-1} k(X, Y)
\end{bmatrix} (3.6.14)$$

3.7 Klasyfikator najbliższych sąsiadów

Do tej pory zajmowaliśmy się problemami regresji ciągłej zmiennej skalarnej. Przechodzimy teraz do metod uczenia maszynowego wykorzystywanych dla problemów klasyfikacji. Jako pierwszy rozważmy jeden z najprostszych (ale bardzo mocnych) modeli – klasyfikator k najbliższych sasiadów (z ang. k Nearest Neighbors, kNN). Załóżmy, iż mamy zbiór obserwacji i.i.d. postaci $\mathcal{X} = \{y_i(\boldsymbol{x}_i)\}_{i=1}^n$ przy czym y może być zarówno skalarną zmienną ciągłą jak i elementem skończonego zbioru klas. Zakładamy, że wektory cech $x \in \mathbb{R}^d$, a d: $\mathbb{R}^d \times \mathbb{R}^d \mapsto I \subseteq \mathbb{R}$ jest metryka, półmetryka lub pewna miarą podobieństwa między punktami \mathbb{R}^d . Reguła decyzyjna klasyfikatora k najbliższych sąsiadów polega na znalezieniu dla nowego wektora cech t, k najbliższych względem funkcji d punktów ze zbioru $\{x_1,\ldots,x_n\}$ i zwróceniu najczęstszej klasy dla tych sąsiadów. Metodę najbliższych sąsiadów można również wykorzystać do regresji, gdzie zwracamy wartość średnia arytmetyczną wartości y dla znalezionych k najbliższych sąsiadów (powoduje to, że model taki potrafi tylko interpolować wartości, więc nie jest dobrym modelem dla regresji). W przypadku klasyfikacji możemy natomiast zwracać również rozkład prawdopodobieństwa klas dla nowego wektora cech \boldsymbol{t} przez podanie stosunków występowania danej klasy wśród k najbliższych sąsiadów do k.

Klasyfikator kNN jest bardzo elastycznym modelem z nieliniową granicą decyzyjną. Jakość klasyfikacji silnie zależy od lokalnej gęstości punktów w przestrzeni \mathbb{R}^d oraz wybranej wartości k, będącej hiperparametrem tego modelu. Ogólnie niskie k powoduje, że kNN ma duży variance i dość "poszarpaną" granicę decyzyjną, natomiast wysokie k powoduje, że kNN ma duży bias i "gładką" granicę decyzyjną. Typowo wykorzystywane funkcje d to

- metryka euklidesowa $\mathsf{d}(\boldsymbol{x},\boldsymbol{y}) = \sqrt{(\boldsymbol{x}-\boldsymbol{y})^T(\boldsymbol{x}-\boldsymbol{y})};$
- pół-metryka euklidesowa $\mathsf{d}(\boldsymbol{x},\boldsymbol{y}) = (\boldsymbol{x}-\boldsymbol{y})^T(\boldsymbol{x}-\boldsymbol{y});$
- metryka Manhattan $d(x, y) = \sum_{i=1}^{d} |x^i y^i|$;
- podobieństwo cosinusowe $\mathsf{d}(\boldsymbol{x},\boldsymbol{y}) = \frac{\boldsymbol{x}^T\boldsymbol{y}}{\|\boldsymbol{x}\|_2\|\boldsymbol{y}\|_2}$

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 k sąsiadów obliczamy teraz jako średnią ważoną, gdzie wagą jest odwrotność odległości danego sąsiada od nowego wektora cech $w(t, x) = 1/\mathsf{d}(t, x)$.

W przypadku naiwnego kNN podczas treningu zapamiętujemy jedynie zbiór \mathcal{X} natomiast naiwna implementacja predykcji ma złożoność czasową O(knd), przy założeniu, że złożoność obliczenia funkcji d dla pary punktów ma złożoność O(d). 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 znajdywać 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).

3.8 Naiwny klasyfikator bayesowski

Rozważamy dalej problem klasyfikacji, w którym mamy zbiór przykładów i.i.d. postaci $\mathcal{X} = \{y_i(\boldsymbol{x}_i)\}_{i=1}^n$, gdzie $y \in \{c_1, \dots, c_K\}$ oraz $\boldsymbol{x} \in \mathbb{R}^d$. W ogólności mamy rozkład warunkowy danej klasy pod warunkiem wektora cech, który jest dany przez gęstość prawdopodobieństwa $p(c_k \mid \boldsymbol{x})$. Korzystając z twierdzenia Bayesa możemy zapisać

$$p(c_k \mid \boldsymbol{x}) = \frac{p(\boldsymbol{x} \mid c_k)p(c_k)}{p(\boldsymbol{x})}.$$
 (3.8.1)

Jeśli umiemy obliczyć licznik wyrażenia po prawej stronie, to korzystając z reguły decyzyjnej MAP klasę dla nowego wektora cech wybieramy jako

$$\arg \max_{c_k \in \{c_1, \dots, c_K\}} p(\mathbf{x} \mid c_k) p(c_k). \tag{3.8.2}$$

Powstaje pytanie jak obliczyć to wyrażenie przy tak luźnych założeniach. Wprowadzamy naiwne założenie warunkowej niezależności cech względem danej klasy tj.

$$p(\mathbf{x} \mid c_k) = p(\mathbf{x}^1, \dots, \mathbf{x}^d \mid c_k) = \prod_{j=1}^d p(\mathbf{x}^j \mid c_k).$$
 (3.8.3)

Teraz jednowymiarowe rozkłady warunkowe $p(\mathbf{x}^j \mid c_k)$ możemy estymować z danych \mathcal{X} np. korzystając z jądrowego estymatora gęstości lub zakładając konkretny model parametryczny (np. jednowymiarowy rozkład normalny, rozkład dwumianowy) i estymując jego parametry dla każdej z klas osobno. Naiwne założenie warunkowej niezależności pozwala złagodzić problemy wynikające z przekleństwa wymiarowości (z ang. curse of dimensionality), takie jak potrzeba zbiorów danych skalujących się wykładniczo wraz z liczbą cech d. Człon $p(c_k)$ można natomiast prosto oszacować jako stosunek przykładów danej klasy w zbiorze \mathcal{X} .

3.9 Estymator jądrowy gęstości

Załóżmy, że mamy zbiór obserwacji i.i.d. $\{x_1, \ldots, x_n\}$ taki, że $x \sim \mathcal{D}$ dla pewnego d—wymiarowego ciągłego rozkładu prawdopodobieństwa \mathcal{D} z nieznaną gęstością prawdopodobieństwa p(x). Chcemy znaleźć estymator $\hat{p}(x)$ tej funkcji. Estymatorem jądrowym gęstości funkcji p (z ang. kernel density estimator) nazywamy funkcję

$$\hat{p}(\boldsymbol{x}) := \frac{1}{mh^d} \sum_{i=1}^n K\left(\frac{\boldsymbol{x} - \boldsymbol{x}_i}{h}\right), \qquad (3.9.1)$$

gdzie $h \in \mathbb{R}$ jest pewnym hiperparametrem zwanym bandwidth, a $K : \mathbb{R}^d \mapsto [0; \infty)$ to tzw. funkcja jądrowa będąca parzystą funkcją posiadającą w 0 maksimum globalne oraz spełniającą warunek unormowania

$$\int_{\mathbb{D}^d} K(\boldsymbol{x}) \, \mathrm{d}^d \boldsymbol{x} = 1. \tag{3.9.2}$$

Ze statystycznego punktu widzenia, postać jądra nie ma istotnego znaczenia i wybór funkcji K może być arbitralny, uwzględniający przede wszystkim pożądane własności otrzymanego estymatora, na przykład klasę jego regularności (ciągłość, różniczkowalność itp.). W przypadku jednowymiarowym jako funkcję K przyjmuje się klasyczne postacie gęstości rozkładów probabilistycznych, na przykład gęstość rozkładu normalnego. W przypadku wielowymiarowym stosuje się tzw. jądro radialne tj. dla jądra jednowymiarowego K wielowymiarowe jądro radialne definiujemy jako

$$K(\boldsymbol{x}) = K(\|\boldsymbol{x}\|) \tag{3.9.3}$$

dla pewnej normy (typowo normy euklidesowej) $\|\cdot\|$.

KDE w praktycznych zastosowaniach często przyspiesza się za pomocą struktur danych analogicznych jak w przypadku kNN, tj. zamiast sumować przyczynki od wszystkich punktów \boldsymbol{x}_i dla danego \boldsymbol{x} , znajdujemy k najbliższych sąsiadów \boldsymbol{x} ze zbioru $\{\boldsymbol{x}_i\}_{i=1}^n$ stosując np. ANN i obliczamy przyczynki do $\hat{p}(\boldsymbol{x})$ tylko od nich.

Jedynym większym problemem w przypadku KDE jest wybór odpowiedniej wartości parametru h. Jeśli używamy jądra gaussowskiego to <u>reguła Silvermana</u>, która podaje prosty przepis na szerokość kernela, przy założeniu, iż rozkład jest unimodalny

$$h_{\text{Silverman}} = \left(\frac{4}{d+2}\right) n^{\frac{-1}{d+4}} \sigma, \qquad (3.9.4)$$

gdzie σ to odchylenie standardowe, n to liczba próbek, a d – liczba wymiarów.

3.10 Wieloklasowa regresja logistyczna

Powróćmy do zagadnienia klasyfikacji, w którym mamy zbiór obserwacji i.i.d. $\mathcal{X} = \{y_i(\boldsymbol{x}_i)\}_{i=1}^n$, gdzie $y \in \{c_1, \dots, c_K\}$ i $\boldsymbol{x} \in \mathbb{R}^d$. Pokażemy teraz podejście oparte na estymacji MLE podobnie jak w przypadku regresji. Założymy, że obserwacje $y(\boldsymbol{x})$ pochodzą z rozkładu kategorialnego (wielopunktowego) zależnego od parametrów \boldsymbol{W} jako

$$y(\mathbf{x}) \mid \mathbf{W} \sim \operatorname{Cat} \left(\boldsymbol{\pi}^{1}(\boldsymbol{\phi}(\mathbf{x}; \mathbf{W})), \dots, \boldsymbol{\pi}^{K}(\boldsymbol{\phi}(\mathbf{x}; \mathbf{W})) \right),$$
 (3.10.1)

gdzie funkcja $\boldsymbol{\pi}: \mathbb{R}^K \mapsto [0;1]^K$ musi spełniać warunek unormowania

$$orall oldsymbol{z} \in \mathbb{R}^K : \sum_{j=1}^K oldsymbol{\pi}^j(oldsymbol{z}) = 1$$
 ,

aby jej składowe były odpowiadały prawdopodobieństwom danej klasy, Funkcja $\phi: \mathbb{R}^d \mapsto \mathbb{R}^K$ jest natomiast dowolną funkcją przekształcającą wektor cech na rzeczywisty wektor wymiaru K. W szczególności przyjmiemy następującą postać funkcji π zwaną funkcją softmax

$$\pi^{j}(z) = \frac{\exp(z^{j})}{\sum_{k=1}^{K} \exp(z^{k})}.$$
(3.10.2)

W przypadku wieloklasowej regresji logistycznej (zwanej również regresją softmax) jako funkcję ϕ przyjmiemy proste przekształcenie liniowe

$$\phi(x; W) = Wx, \qquad (3.10.3)$$

gdzie W jest macierzą estymowanych parametrów wymiaru $K \times d$. Wyprowadzimy teraz wzór na funkcję kosztu przy takim modelu statystycznym. Zauważmy wpierw, iż wiarygodność pojedynczego przykładu ze zbioru $\mathcal X$ możemy zapisać jako

$$p(y_i(\mathbf{x}_i) \mid \mathbf{W}) = \prod_{j=1}^K \pi^j (\phi(\mathbf{x}_i; \mathbf{W}))^{[y_i = c_j]},$$
 (3.10.4)

gdzie $[\ldots]$ oznacza nawias Iversona. W takim razie wiarygodność całego zbioru $\mathcal X$ ma postać

$$\mathcal{L}(\mathcal{X}; \boldsymbol{W}) = \prod_{i=1}^{n} \prod_{j=1}^{K} \boldsymbol{\pi}^{j} \left(\boldsymbol{\phi}(\boldsymbol{x}_{i}; \boldsymbol{W}) \right)^{[y_{i} = c_{j}]}, \qquad (3.10.5)$$

skąd funkcja kosztu ma postać

$$L(\mathcal{X}; \boldsymbol{W}) = -\sum_{i=1}^{n} \sum_{j=1}^{K} [y_i = c_j] \log \left[\boldsymbol{\pi}^j \left(\boldsymbol{\phi}(\boldsymbol{x}_i; \boldsymbol{W}) \right) \right].$$
 (3.10.6)

Wprowadzając macierze

$$egin{aligned} oldsymbol{X}^{ij} &= oldsymbol{x}_j^i, \quad oldsymbol{X} &= egin{bmatrix} oldsymbol{x}_1 & \dots & oldsymbol{x}_n \end{bmatrix} \ oldsymbol{T}^{ij} &= [y_j = c_i] \ oldsymbol{\Phi}^{ij}(oldsymbol{X}; oldsymbol{W}) &= oldsymbol{\sum}_{k=1}^d oldsymbol{W}^{ik} oldsymbol{X}^{kj} \,, \quad oldsymbol{\Phi}(oldsymbol{X}; oldsymbol{W}) &= oldsymbol{W} oldsymbol{X} \ oldsymbol{\Pi}^{ij}(oldsymbol{\Phi}) &= \pi^i \left(oldsymbol{\phi}(oldsymbol{x}_j; oldsymbol{W}) \right) = rac{\exp oldsymbol{\Phi}^{ij}}{\sum_{k=1}^K \exp oldsymbol{\Phi}^{kj}} \end{aligned}$$

gdzie każda kolumna macierzy \boldsymbol{X} jest wektorem cech danego przykładu; każda kolumna \boldsymbol{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; każda kolumna macierzy $\boldsymbol{\Phi}$ jest wektorem mlogitów tj. liczb rzeczywistych, które po zastosowaniu funkcji softmax dają wartości prawdopodobieństwa każdej klasy; możemy zapisać

$$L(\mathcal{X}; \boldsymbol{W}) = -\sum_{j=1}^{n} \sum_{i=1}^{K} \boldsymbol{T}^{ij} \log \left[\boldsymbol{\Pi}^{ij} (\boldsymbol{\Phi}(\boldsymbol{X}; \boldsymbol{W})) \right].$$
(3.10.7)

Niestety dla tak zdefiniowanej funkcji kosztu nie można znaleźć wzoru na minimum w postaci analitycznej (jak w przypadku regresji liniowej), dlatego do znalezienia estymaty MLE wykorzystujemy algorytmy optymalizacji numerycznej, w tym przypadku zwykle algorytmy gradientowe. Wyprowadzimy więc jeszcze wzór na pochodną funkcji kosztu po parametrach \boldsymbol{W} . Obliczmy najpierw pochodną L po $\boldsymbol{\Phi}^{ij}$ (dla przejrzystości zapisu nie piszemy granic sumowania – wynikają one naturalnie z wymiarów macierzy)

$$\frac{\partial L}{\partial \mathbf{\Phi}^{pq}} = -\sum_{i,j,r,s} \frac{\mathbf{T}^{ij}}{\mathbf{\Pi}^{rs}} \delta_{ir} \delta_{js} \frac{\partial \mathbf{\Pi}^{rs}}{\partial \mathbf{\Phi}^{pq}} = -\sum_{i,j} \frac{\mathbf{T}^{ij}}{\mathbf{\Pi}^{ij}} \frac{\partial \mathbf{\Pi}^{ij}}{\partial \mathbf{\Phi}^{pq}}.$$
 (3.10.8)

Jednocześnie

$$\frac{\partial \mathbf{\Pi}^{ij}}{\partial \mathbf{\Phi}^{pq}} = \frac{\delta_{pi}\delta_{qj} \exp \mathbf{\Phi}^{ij} \left[\sum_{k} \exp \mathbf{\Phi}^{kj} \right] - \delta_{qj} \exp \mathbf{\Phi}^{ij} \exp \mathbf{\Phi}^{pj}}{\left[\sum_{k} \exp \mathbf{\Phi}^{kj} \right]^{2}}, \quad (3.10.9)$$

skad

$$\frac{1}{\mathbf{\Pi}^{ij}} \frac{\partial \mathbf{\Pi}^{ij}}{\partial \mathbf{\Phi}^{pq}} = \delta_{pi} \delta_{qj} - \delta_{qj} \mathbf{\Pi}^{pj} \,. \tag{3.10.10}$$

Z powyższego zatem

$$\frac{\partial L}{\partial \mathbf{\Phi}^{pq}} = -\sum_{i,j} \mathbf{T}^{ij} \delta_{pi} \delta_{qj} + \sum_{i,j} \mathbf{T}^{ij} \delta_{qj} \mathbf{\Pi}^{pj} = \mathbf{\Pi}^{pq} - \mathbf{T}^{pq}, \qquad (3.10.11)$$

gdzie skorzystaliśmy z faktu, iż z konstrukcji dla dowolnej kolumny q macierzy T zachodzi $\sum_i T^{iq}=1$. Możemy zapisać powyższy wzór w eleganckiej postaci macierzowej

$$\boxed{\frac{\partial L}{\partial \mathbf{\Phi}} = \mathbf{\Pi} - \mathbf{T}.} \tag{3.10.12}$$

Pochodną funkcji kosztu po parametrach W możemy zatem obliczyć jako

$$\frac{\partial L}{\partial \mathbf{W}^{rs}} = \sum_{p,q} \frac{\partial L}{\partial \mathbf{\Phi}^{pq}} \frac{\partial \mathbf{\Phi}^{pq}}{\partial \mathbf{W}^{rs}}, \qquad (3.10.13)$$

gdzie

$$\frac{\partial \mathbf{\Phi}^{pq}}{\partial \mathbf{W}^{rs}} = \sum_{t} \delta_{rp} \delta_{st} \mathbf{X}^{tq} = \delta_{rp} \mathbf{X}^{sq}, \qquad (3.10.14)$$

skąd

$$\frac{\partial L}{\partial \mathbf{W}^{rs}} = \sum_{q} (\mathbf{\Pi}^{rq} - \mathbf{T}^{rq}) \mathbf{X}^{sq}, \qquad (3.10.15)$$

co również możemy zapisać w zwartej postaci macierzowej

$$\frac{\partial L}{\partial \boldsymbol{W}} = (\boldsymbol{\Pi} - \boldsymbol{T}) \boldsymbol{X}^T.$$
 (3.10.16)

3.11 Próbkowanie Monte Carlo łańcuchami Markowa

Okazuje się, iż do generowania próbek ze skomplikowanego rozkładu p(x) wystarcza znajomość tego rozkładu z dokładnością do stałej normalizującej, a zatem wystarczy znać rozkład łączny $\tilde{p}(x) = Z_p p(x)$. Generowanie próbek z kolei wystarcza natomiast, na mocy silnego prawa wielkich liczb, do szacowania wartości średnich dowolnych funkcji estymowanego parametru θ . Przypomnijmy, iż na mocy silnego prawa wielkich liczb ciąg średnich częściowych (\overline{X}_n) ciągu zmiennych losowych (X_n) i.i.d. z rozkładu $X \sim \mathcal{D}$ jest zbieżny z prawdopodobieństwem 1 do wartości oczekiwanej $\mathbb{E}[X]$ tj.

$$\Pr\left(\lim_{n\to\infty}\overline{X}_n = \mathbb{E}[X]\right) = 1. \tag{3.11.1}$$

Wartość oczekiwaną $\mathbb{E}[X]$ możemy zatem przybliżyć średnią \overline{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).

3.11.1 Algorytm Importance Sampling

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

$$p(\mathbf{x}) = \frac{1}{Z_p} \tilde{p}(\mathbf{x}) \tag{3.11.2}$$

tj. szukamy

$$\mathbb{E}_p[f(\boldsymbol{x})] = \int f(\boldsymbol{x})p(\boldsymbol{x})\mathrm{d}^n \boldsymbol{x}. \qquad (3.11.3)$$

Jeśli umiemy generować próbki \boldsymbol{x} z innego (prostszego) rozkładu $q(\boldsymbol{x})$, który nazywamy rozkładem proponującym kandydatów (z ang. proposal distribution) to możemy zapisać

$$\mathbb{E}_{p}[f(\boldsymbol{x})] = \int_{\mathbb{R}^{n}} f(\boldsymbol{x})p(\boldsymbol{x})d^{n}\boldsymbol{x} = \int_{\mathbb{R}^{n}} f(\boldsymbol{x})\frac{p(\boldsymbol{x})}{q(\boldsymbol{x})}q(\boldsymbol{x})d^{n}\boldsymbol{x}$$

$$= \mathbb{E}_{q}\left[f(\boldsymbol{x})\frac{p(\boldsymbol{x})}{q(\boldsymbol{x})}\right] = \frac{Z_{q}}{Z_{p}}\mathbb{E}_{q}\left[f(\boldsymbol{x})\frac{\tilde{p}(\boldsymbol{x})}{\tilde{q}(\boldsymbol{x})}\right].$$
(3.11.4)

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

$$Z_p = \int_{\mathbb{R}^n} \tilde{p}(\boldsymbol{x}) d^n \boldsymbol{x} = Z_q \int_{\mathbb{R}^n} \frac{\tilde{p}(\boldsymbol{x})}{\tilde{q}(\boldsymbol{x})} q(\boldsymbol{x}) d^n \boldsymbol{x} = Z_q \mathbb{E}_q \left[\frac{\tilde{p}(\boldsymbol{x})}{\tilde{q}(\boldsymbol{x})} \right], \qquad (3.11.5)$$

skąd ostatecznie

$$\mathbb{E}_{p}[f(\boldsymbol{x})] = \frac{\mathbb{E}_{q}\left[f(\boldsymbol{x})\frac{\tilde{p}(\boldsymbol{x})}{\tilde{q}(\boldsymbol{x})}\right]}{\mathbb{E}_{q}\left[\frac{\tilde{p}(\boldsymbol{x})}{\tilde{q}(\boldsymbol{x})}\right]}.$$
(3.11.6)

Jeśli z rozkładu q wygenerowaliśmy próbki $X = \{x_1, \dots, x_m\}$ to na mocy silnego prawa wielkich liczb mamy

$$\mathbb{E}_{p}[f(\boldsymbol{x})] \approx \frac{\sum_{i=1}^{m} f(\boldsymbol{x}_{i}) \frac{\tilde{p}(\boldsymbol{x}_{i})}{\tilde{q}(\boldsymbol{x}_{i})}}{\sum_{j=1}^{m} \frac{\tilde{p}(\boldsymbol{x}_{j})}{\tilde{q}(\boldsymbol{x}_{j})}} = \sum_{i=1}^{m} \lambda_{i} f(\boldsymbol{x}_{i}),$$
(3.11.7)

gdzie

$$\lambda_i = \frac{\tilde{p}(\boldsymbol{x}_i)/\tilde{q}(\boldsymbol{x}_i)}{\sum_{j=1}^m \tilde{p}(\boldsymbol{x}_j)/\tilde{q}(\boldsymbol{x}_j)}.$$
 (3.11.8)

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

$$\overline{f}_m(\boldsymbol{x}_1, \dots, \boldsymbol{x}_m) = \frac{1}{m} \sum_{i=1}^m \frac{f(\boldsymbol{x}_i) p(\boldsymbol{x}_i)}{q(\boldsymbol{x}_i)}$$
(3.11.9)

dla $x_i \sim q$ mamy

$$\mathbb{V}[\overline{f}_m] = \frac{1}{m} \mathbb{V}_q \left[f(\boldsymbol{x}) \frac{p(\boldsymbol{x})}{q(\boldsymbol{x})} \right] = \frac{1}{m} \int_{\mathbb{R}^n} \frac{(f(\boldsymbol{x})p(\boldsymbol{x}) - \mu_f q(\boldsymbol{x}))^2}{q(\boldsymbol{x})} d^n \boldsymbol{x} .$$
(3.11.10)

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 $f(\boldsymbol{x})p(\boldsymbol{x})$, co może być trudne do praktycznego zrealizowania.

3.11.2 Algorytm Metropolisa–Hastingsa

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

Definicja 3.4 (Łańcucha Markowa). Łańcuchem Markowa nazwiemy ciąg zmiennych losowych (X_t) o wartościach w \mathbb{R}^n taki,

że spełnione jest kryterium Markowa

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

Elementy ciągu nazywamy stanami łańcucha.

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

$$\int_{\mathbb{R}^n} \pi(\boldsymbol{y} \mid \boldsymbol{x}) \, \mathrm{d}^n \boldsymbol{y} , \qquad (3.11.11)$$

istotnie prawdopodobieństwo przejścia gdziekolwiek ze stanu \boldsymbol{x} jest równe 1. Będziemy zakładać dodatkowo, iż $\forall \boldsymbol{x}, \boldsymbol{y} \in \mathbb{R}^n : \pi(\boldsymbol{y} \mid \boldsymbol{x}) > 0$. Rozkład $p(\boldsymbol{x})$ 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

$$p(\mathbf{y}) = \int_{\mathbb{R}^n} \pi(\mathbf{y} \mid \mathbf{x}) p(\mathbf{x}) \, \mathrm{d}^n \mathbf{x} . \tag{3.11.12}$$

Rozkład stacjonarny danego łańcucha oznaczymy przez $p^*(x)$. Zauważmy, iż jeśli stan początkowy łańcucha X_0 pochodzi z rozkładu stacjonarnego p^* to każdy kolejny stan 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

$$p_t(\boldsymbol{y}) = \int_{\mathbb{R}^n} \pi(\boldsymbol{y} \mid \boldsymbol{x}) p_{t-1}(\boldsymbol{x}) d^n \boldsymbol{x} , \quad \text{dla } t > 1.$$
 (3.11.13)

Rozkładem granicznym łańcucha Markowa nazwiemy granicę w sensie zbieżności punktowej

$$\lim_{t \to \infty} p_t(\boldsymbol{x}). \tag{3.11.14}$$

Przy podanych wyżej założeniach istnieje twierdzenie, które mówi iż taki łańcuch Markowa posiada jednoznaczny rozkład stacjonarny tożsamy z roz-

kładem granicznym. Ponadto warunkiem wystarczającym, aby dany rozkład p(x) był rozkładem stacjonarnym łańcucha Markowa jest

$$\forall \boldsymbol{x}, \boldsymbol{y} \in \mathbb{R}^n : \pi(\boldsymbol{y} \mid \boldsymbol{x}) p(\boldsymbol{x}) = \pi(\boldsymbol{x} \mid \boldsymbol{y}) p(\boldsymbol{y}), \qquad (3.11.15)$$

co wynika z scałkowania powyższego równania

$$\int_{\mathbb{R}^n} \pi(\boldsymbol{y} \mid \boldsymbol{x}) p(\boldsymbol{x}) d^n \boldsymbol{x} = \int_{\mathbb{R}^n} \pi(\boldsymbol{x} \mid \boldsymbol{y}) p(\boldsymbol{y}) d^n \boldsymbol{x} = p(\boldsymbol{y}) \int_{\mathbb{R}^n} \pi(\boldsymbol{x} \mid \boldsymbol{y}) d^n \boldsymbol{x} = p(\boldsymbol{y}).$$
(3.11.16)

Kryterium to nazywamy <u>kryterium lokalnego balansu</u> (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 \boldsymbol{X}^0 po wykonaniu dużej liczby kroków (etap ten nazywamy okresem przejściowym z ang. burn-in period) stan \boldsymbol{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 $\pi(\boldsymbol{y}\mid\boldsymbol{x})$ dla danego rozkładu granicznego $p(\boldsymbol{x})$. Podstawową konstrukcję podaje algorytm Metropolisa–Hastingsa.

Algorytm Metropolisa–Hastingsa

- 1. Jako stan początkowy przyjmij dowolną dopuszczalną wartość $\boldsymbol{x}_0.$
- 2. Będąc w aktualnym stanie \boldsymbol{x} z prostego rozkładu proponującego kandydatów $q(\boldsymbol{y}\mid \boldsymbol{x})$ wylosuj kandydata \boldsymbol{y} na wartość łańcucha w kolejnym stanie.
- 3. Z prawdopodobieństwem

$$r(\boldsymbol{y} \mid \boldsymbol{x}) = \min \left\{ 1, \frac{p(\boldsymbol{y})q(\boldsymbol{x} \mid \boldsymbol{y})}{p(\boldsymbol{x})q(\boldsymbol{y} \mid \boldsymbol{x})} \right\}$$

zaakceptuj kandydata jako nowy stan i przejdź do stanu \boldsymbol{y} . W przeciwnym razie pozostań w stanie \boldsymbol{x} .

4. GOTO 2.

Funkcja przejścia ma zatem postać

$$\pi_{\mathrm{MH}}(\boldsymbol{y} \mid \boldsymbol{x}) = q(\boldsymbol{y} \mid \boldsymbol{x})r(\boldsymbol{y} \mid \boldsymbol{x}). \tag{3.11.17}$$

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

$$\pi_{\mathrm{MH}}(\boldsymbol{y} \mid \boldsymbol{x})p(\boldsymbol{x}) = \min \left\{ q(\boldsymbol{y} \mid \boldsymbol{x})p(\boldsymbol{x}), q(\boldsymbol{x} \mid \boldsymbol{y})p(\boldsymbol{y}) \right\} \pi_{\mathrm{MH}}(\boldsymbol{x} \mid \boldsymbol{y})p(\boldsymbol{y}) = \min \left\{ q(\boldsymbol{x} \mid \boldsymbol{y})p(\boldsymbol{y}), q(\boldsymbol{y} \mid \boldsymbol{x})p(\boldsymbol{x}) \right\}$$
(3.11.18)

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

$$\frac{p(\mathbf{y})}{p(\mathbf{x})} = \frac{\tilde{p}(\mathbf{y})/Z_p}{\tilde{p}(\mathbf{x})/Z_p} = \frac{\tilde{p}(\mathbf{y})}{\tilde{p}(\mathbf{x})}.$$
 (3.11.19)

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).