

POLITECHNIKA ŁÓDZKA

WYDZIAŁ FIZYKI TECHNICZNEJ, INFORMATYKI
I MATEMATYKI STOSOWANEJ

Kierunek: Matematyka Stosowana

Specjalność: Analiza Danych w Biznesie i Logistyce

Matematyczne modele wykorzystywane w systemach rekomendacji.

Anita Kudaj
Nr albumu: 220020

Praca magisterska napisana w Instytucie Matematyki
Politechniki Łódzkiej

Promotor: dr, mgr inż. Piotr Kowalski

ŁÓDŹ, 07.2019

Spis treści

1	Wstęp	3
2	Preliminaria	5
2.1	Oznaczenia używane w pracy	5
2.2	Elementy algebry liniowej	6
2.3	Elementy rachunku prawdopodobieństwa i statystyki	9
3	Elementy eksploracji danych wykorzystywane w systemach rekomendujących	11
3.1	Wstępne przetwarzanie danych	11
3.1.1	Miary podobieństwa	11
3.1.2	Redukcja wymiaru	13
3.2	Metody eksploracji danych	19
3.2.1	Algorytm k - najbliższych sąsiadów	20
3.2.2	Algorytm k - średnich	21
3.3	Szacowanie błędów obliczeń	23
3.3.1	Ocena dokładności metody	23
3.3.2	Ocena jakości modelu	25
4	Modele tworzenia rekomendacji	28
4.1	Systemy rekomendujące oparte na treści	31
4.1.1	Wygenerowanie profilu przedmiotu - algorytm TFIDF	32
4.1.2	Wygenerowanie profilu użytkownika	33
4.2	Filtrowanie kolaboratywne	36
4.2.1	Filtrowanie kolaboratywne oparte na użytkowniku	36
4.2.2	Filtrowanie kolaboratywne oparte na przedmiotach	39
4.2.3	Algorytmy eksploracji danych w filtrowaniu kolaboratywnym	40
4.3	Systemy rekomendujące kontekstowe	41
4.4	Dekompozycja macierzy ocen - SVD	42
4.5	Ocena jakości modeli w regułach rekomendujących	43

5	Eksperymenty / część praktyczne	45
5.1	Tworzenie rekomendacji z użyciem ALS i filtrowania kolaboratywnego .	45
5.1.1	Apache Spark	45
5.1.2	ALS i MLlib	46
5.1.3	Implementacja algorytmu	47
6	Podsumowanie	55

Rozdział 1

Wstęp

Znaczny wzrost informacji dostępnych w internecie wywołał zapotrzebowanie na systemy, które mogą pomóc użytkownikom znaleźć cenne dla nich dane. W tym kontekście filtrowanie informacji zapewnia konkretne narzędzia - reguły rekomendujące, których celem jest uszeregowanie danych według częściowo ujawnionych preferencji, otrzymanych z uprzednio przeprowadzonych badań, aby w efekcie zasugerować odbiorcy odpowiedni dla niego produkt.

Celem tej pracy jest przedstawienie matematycznej strony metod systemów rekomendujących poprzez analizę algorytmów ukrytych pod ich nazwami.

Niniejsza praca składa się z czterech części. Rozdział 2. zawiera wykaz oznaczeń, przypomnienie podstawowych definicji z zakresu algebry liniowej oraz rachunku prawdopodobieństwa i statystyki wykorzystanych w dalszych częściach pracy. W rozdziale 3. przedstawione zostały elementy eksploracji danych stanowiące nieodzowną część systemów rekomendujących. Zaczynamy od wstępnego przetwarzania danych oraz definicji i twierdzeń związanych z miarami podobieństwa i technikami redukcji wymiaru, tutaj szczególna uwaga została poświęcona rozkładowi według wartości osobliwych (ang. Singular Value Decomposition). Następnie przechodzimy do opisu dwóch metod eksploracji danych - algorytmu k-najbliższych sąsiadów i algorytmu k-średnich, aby zakończyć oceną dokładności i jakości modeli. Rozdział 4. opisuje modele tworzenia rekomendacji. Zostały tu zawarte sformułowania definicji związanych z systemami rekomendującymi, informacje dotyczące systemów opartych na treści, filtrowania kolaboratywnego oraz systemów kontekstowych. Rozdział został wzbogacony o algorytmy i przykłady, które pokazują zastosowanie poszczególnych metod. Ponadto znajdziemy tu przedstawienie dekompozycji macierzy ocen metodą SVD. Ostatni rozdział to praktyczne zastosowanie filtrowania kolaboratywnego. W rozdziale tym podstawę stanowią informacje na temat Apache Spark oraz algorytmu ALS (ang. Alternating Least Square) z biblioteki MLlib. Dodatkowo w ostatniej sekcji rozdziału została umieszczona przykładowa implementa-

cja algorytmu ALS wraz z komentarzami wyjaśniającymi poszczególne kroki.

Rozdział 2

Preliminaria

2.1 Oznaczenia używane w pracy

W niniejszej pracy zostały użyte następujące oznaczenia:

\mathbb{N} - zbiór liczb naturalnych,

\mathbb{R} - zbiór liczb rzeczywistych,

\mathbb{K} - ciało liczb rzeczywistych lub zespolonych,

X - (duże, pochylone litery) jako oznaczenia zbiorów,

\mathbf{x} - (małe, pogrubione litery) jako oznaczenia wektorów,

X - (duże litery) jako oznaczenia zmiennych losowych,

\mathbb{X} - (duże litery z wyłączeniem \mathbb{N} , \mathbb{R} , \mathbb{K}) jako oznaczenia macierzy,

$[a_{ij}]_{j=1,\dots,n}^{i=1,\dots,m}$ - macierz o m wierszach i n kolumnach,

$[a_{ij}]$ - macierz kwadratowa,

$\mathbb{M}_{m \times n}(\mathbb{K})$ - zbiór wszystkich macierzy o wymiarach m na n i elementach z ciała \mathbb{K} ,

\mathcal{V} - przestrzeń liniowa,

$\text{span}(X)$ - przestrzeń generowana przez zbiór X ,

(Ω, \mathcal{F}, P) - przestrzeń probabilistyczna,

Ω - zbiór zdarzeń elementarnych,

\mathcal{F} - rodzina podzbiorów zbioru Ω ,

P - funkcja prawdopodobieństwa,

$B(\mathbb{R}^n)$ - σ -ciało zbiorów borelowskich w \mathbb{R}^n ,

$E(X)$ - wartość oczekiwana,

$\text{Cov}(X; Y)$ - kowariancja zmiennych losowych X, Y ,

$\text{Var}(X)$ - wariancję zmiennej losowej X ,

$\sigma(X)$ - odchylenie standardowe zmiennej losowej X ,

$\rho(X, Y)$ - współczynnik korelacji zmiennych losowych X, Y ,

$d(x, y)$ - odległość punktów x i y ,

$d_e(x, y)$ - odległość euklidesowa punktów x i y ,
 $\text{sim}(X, Y)$ - współczynnik podobieństwa wektorów X i Y ,
 $\rho^p(X, Y)$ - współczynnik korelacji Pearsona,
 $J(A, B)$ - indeks Jaccarda,
 $\|\cdot\|_F$ - norma Frobeniusa.

2.2 Elementy algebry liniowej

W definicjach poniżej korzystamy z pojęcia ciała, którego wyjaśnienie odnajdziemy w książce Tadeusza Poredy i Jacka Jędrzejewskiego *Algebra liniowa z elementami geometrii analitycznej* [14, Sec 4.4].

Definicja 2.1 (Macierz [14, Sec 8.1 Def. 8.1]). *Niech $n, m \in \mathbb{N}$. Macierz o m wierszach, n kolumnach (o wymiarach m na n) i wyrazach w ciele \mathbb{K} nazywamy funkcję*

$$\mathbb{A} : \{1, 2, \dots, m\} \times \{1, 2, \dots, n\} \rightarrow \mathbb{K}.$$

Wartością funkcji dla argumentu (i, j) jest element a_{ij} należący do ciała \mathbb{K} . Macierz zapisujemy w postaci tabeli

$$\mathbb{A} = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{bmatrix}.$$

Przez $\mathbb{M}_{m \times n}(\mathbb{K})$ oznaczamy zbiór wszystkich macierzy o wymiarach m na n i elementach z ciała \mathbb{K} . Przykładowymi sposobami zapisu macierzy są:

$$[a_{ij}]_{j=1, \dots, n}^{i=1, \dots, m}, (a_{ij})_{j=1, \dots, n}^{i=1, \dots, m}, [a_{ij}]_{j \leq n}^{i=1 \leq m}, (a_{ij})_{j \leq n}^{i=1 \leq m}, [a_{ij}], (a_{ij}),$$

przy czym sposobów $[a_{ij}]$, (a_{ij}) używamy, gdy liczba kolumn i wierszy danej macierzy jest ustalona. W tej pracy używać będziemy zapisu $[a_{ij}]_{j=1, \dots, n}^{i=1, \dots, m}$ oraz zapisu $[a_{ij}]$ w przypadku macierzy kwadratowych.

Definicja 2.2 (Wyznacznik macierzy [14, Sec 10.1, Def. 10.1]). *Niech $\mathcal{M}(\mathbb{K}) = \bigcup_{n \in \mathbb{N}} \mathbb{M}_{n \times n}(\mathbb{K})$ oznacza zbiór wszystkich macierzy kwadratowych o wyrazach z \mathbb{K} .*

Funkcję:

$$\det : \mathcal{M}(\mathbb{K}) \rightarrow \mathbb{K}$$

określamy następująco:

- jeżeli $\mathbb{A} = [a_{11}]$, to $\det(\mathbb{A}) = a_{11}$,

$$\bullet \text{ jeżeli } \mathbb{A} = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nn} \end{bmatrix}, \text{ gdy } n > 1, \text{ to}$$

$$\det(\mathbb{A}) = \sum_{i=1}^n (-1)^{1+i} \cdot a_{i1} \cdot \det(\mathbb{A}_{i1}),$$

gdzie \mathbb{A}_{ij} jest macierzą powstałą z macierzy \mathbb{A} przez skreślenie i -tego wiersza i j -tej kolumny.

Funkcję \det nazywamy wyznacznikiem, natomiast wartość tej funkcji dla macierzy \mathbb{A} wyznacznikiem macierzy \mathbb{A} .

Definicja 2.3 (Macierz transponowana [14, Sec 8.1]). Niech $\mathbb{A} = [a_{ij}]_{j=1, \dots, n}^{i=1, \dots, m}$, będzie macierzą ze zbioru $\mathbb{M}_{m \times n}(\mathbb{K})$. Macierz $\mathbb{B} = [b_{ij}]_{j=1, \dots, m}^{i=1, \dots, n}$ nazywamy macierzą transponowaną macierzy \mathbb{A} , jeśli

$$b_{ji} = a_{ij}$$

dla każdego $i \in \{1, \dots, n\}$ oraz $j \in \{1, \dots, m\}$. Piszemy wtedy $\mathbb{B} = \mathbb{A}^T$.

Twierdzenie 2.4 (Własności transpozycji macierzy [1, Sec 5.1 Tw. 5.1]). Niech $\mathbb{A} \in \mathbb{M}_{n \times m}(\mathbb{K})$, $\mathbb{B} \in \mathbb{M}_{m \times n}(\mathbb{K})$, $\mathbb{C} \in \mathbb{M}_{m \times n}(\mathbb{K})$ oraz $\lambda \in \mathbb{K}$. Zachodzą następujące równości:

- $(\mathbb{A}^T)^T = \mathbb{A}$,
- $(\mathbb{A} + \mathbb{B})^T = \mathbb{A}^T + \mathbb{B}^T$,
- $(\lambda \mathbb{A})^T = \lambda \mathbb{A}^T$,
- $(\mathbb{A}\mathbb{C})^T = \mathbb{C}^T \mathbb{A}^T$.

Uwaga 2.5 (Rodzaje macierzy [14, Sec 8.1, Sec 10.4]). Poniżej zostały zdefiniowane niektóre rodzaje macierzy.

- Macierzą kwadratową nazywamy macierz, w której liczba wierszy i liczba kolumn są równe. Liczbę tę nazywamy stopniem macierzy kwadratowej.
- Macierzą diagonalną nazywamy macierz kwadratową $[a_{ij}]$, gdzie wszystkie elementy poza główną przekątną są równe 0. Macierz diagonalną oznaczamy $\text{diag}(a_{11}, a_{22}, \dots, a_{nn})$.
- Macierzą jednostkową stopnia n nazywamy macierz diagonalną, w której na głównej przekątnej wszystkie elementy są równe 1. Macierz jednostkową będziemy oznaczać \mathbb{I} .

- Macierz kwadratową \mathbb{C} , gdzie $\mathbb{C} = [c_{ij}]$, nazywamy macierzą ortogonalną, jeżeli spełniony jest warunek

$$\mathbb{C}^T \cdot \mathbb{C} = \mathbb{C} \cdot \mathbb{C}^T = \mathbb{I}.$$

- Macierzą nieosobliwą nazywamy macierz kwadratową, której wyznacznik jest różny od 0.
- Macierzą osobliwą nazywamy macierz kwadratową, której wyznacznik jest równy 0.

Uwaga 2.6. Podstawowymi działaniami wykonywanymi na macierzach, których będziemy używać w dalszych częściach tej pracy są:

- mnożenie macierzy [14, Sec 9.3 Def 9.13],
- suma macierzy [14, Sec 8.1],
- iloczyn macierzy przez element ciała [14, Sec 8.1].

Definicja 2.7 (Ślad macierzy). [1, Sec 6.4] Śladem macierzy $\mathbb{A} = [a_{ij}]$ nazywamy wielkość:

$$\text{tr}(\mathbb{A}) = \sum_{i=1}^n a_{ii} = a_{11} + a_{22} + \cdots + a_{nn}.$$

Twierdzenie 2.8 (Własności śladu macierzy [1, Sec 6.4]). Niech $\mathbb{A}, \mathbb{B}, \mathbb{C} \in \mathbb{M}_{n \times n}(\mathbb{K})$ oraz $\lambda \in \mathbb{K}$. Zachodzą następujące równości:

- $\text{tr}(\mathbb{A} + \mathbb{B}) = \text{tr}(\mathbb{A}) + \text{tr}(\mathbb{B})$,
- $\text{tr}(\lambda \mathbb{A}) = \lambda \text{tr}(\mathbb{A})$,
- $\text{tr}(\mathbb{A}) = \text{tr}(\mathbb{A}^T)$,
- $\text{tr}(\mathbb{A}\mathbb{B}) = \text{tr}(\mathbb{B}\mathbb{A})$,
- $\text{tr}(\mathbb{A}\mathbb{B}\mathbb{C}) = \text{tr}(\mathbb{C}\mathbb{A}\mathbb{B}) = \text{tr}(\mathbb{B}\mathbb{C}\mathbb{A})$.

W kolejnych definicjach korzystamy z pojęć przestrzeni liniowej, wymiaru oraz przekształcenia liniowego, których wyjaśnienia możemy odnaleźć w książce *Algebra liniowa z elementami geometrii analitycznej* odpowiednio w [14, Sec 7.1], [14, Sec 7.5] oraz [14, Sec 9.1]. Dodatkowo wykorzystamy także pojęcia rzędu kolumnowego i wierszowego macierzy, których definicje także odnajdziemy w powyższej publikacji, w [14, Sec 8.1].

Definicja 2.9 (Rząd macierzy [14, Sec 8.1]). Rzędem macierzy \mathbb{A} nazywamy wspólną wartość rzędu kolumnowego i wierszowego macierzy \mathbb{A} . Rząd macierzy oznaczamy symbolem $\text{rz}(\mathbb{A})$.

Definicja 2.10 (Przestrzeń generowana przez zbiór [14, Sec 7.1 Def 7.13]). *Niech X będzie dowolnym i niepustym podzbiorem przestrzeni liniowej \mathcal{V} . Podprzestrzenią generowaną przez zbiór X nazywamy zbiór wszystkich skończonych kombinacji liniowych wektorów ze zbioru X . Zbiór ten oznaczamy symbolem $\text{span}(X)$.*

Symbolicznie zapisujemy zbiór $\text{span}(X)$ jako:

$$\left\{ x \in \mathcal{V} : \exists_{n \in \mathbb{N}} \exists_{(\alpha_1, \dots, \alpha_n) \in \mathbb{K}^n} \exists_{(x_1, \dots, x_n) \in X^n} (x = \alpha_1 \cdot x_1 + \dots + \alpha_n \cdot x_n) \right\},$$

gdzie \mathcal{V} jest przestrzenią liniową nad ciałem liczb rzeczywistych lub ciałem liczb zespolonych \mathbb{K} .

Definicja 2.11 (Wartość własna macierzy kwadratowej [1, Sec 13.1]). *Liczbę $\lambda \in \mathbb{K}$ nazywamy wartością własną macierzy kwadratowej \mathbb{A} , jeżeli istnieje niezerowy wektor \mathbf{x} taki, że*

$$\mathbb{A}\mathbf{x} = \lambda\mathbf{x}.$$

Każdy niezerowy wektor \mathbf{x} spełniający powyższe równania nazywamy wektorem własnym macierzy \mathbb{A} odpowiadającym wartości własnej λ .

2.3 Elementy rachunku prawdopodobieństwa i statystyki

Definicja 2.12 (Przestrzeń probabilistyczna [7, Sec 1.2, Sec 1.4]). *Przestrzenią probabilistyczną nazywamy uporządkowaną trójkę (Ω, \mathcal{F}, P) , gdzie:*

- Ω to zbiór wszystkich zdarzeń elementarnych i $\Omega \neq \emptyset$,
- \mathcal{F} rodzina podzbiorów zbioru Ω taka, że:
 - $\emptyset \in \mathcal{F}$,
 - jeżeli $A \in \mathcal{F}$, to $\bar{A} = \Omega \setminus A \in \mathcal{F}$,
 - jeżeli $A_n \in \mathcal{F}$ dla $n = 1, 2, \dots$, to $\bigcup_{n=1}^{\infty} A_n \in \mathcal{F}$,
- $P : \mathcal{F} \rightarrow [0, 1]$ taka, że:
 - $\forall_{A \in \mathcal{F}} P(A) \geq 0$,
 - $P(\Omega) = 1$,
 - jeżeli $A_n \in \mathcal{F}$, $n = 1, 2, \dots$ są takie, że $A_i \cap A_j \neq \emptyset$ dla $i \neq j$ to

$$P\left(\bigcup_{n=1}^{\infty} A_n\right) = \sum_{n=1}^{\infty} P(A_n).$$

W poniższej definicji zostało użyte pojęcie σ -algebry, którego wyjaśnienie można odnaleźć w [7, Sec 1.2 Def. 1.2].

Definicja 2.13 ($B(\mathbb{R}^n)$ [7, Sec 1.12]). Niech \mathcal{I} oznacza klasę wszystkich zbiorów, składających się ze skończonych sum rozłącznych zbiorów $I = I_1 \times I_2 \times \dots \times I_k$, gdzie $I_k = [a_k, b_k)$. Najmniejszą σ -algebrę $\sigma(\mathcal{I})$ generowaną przez klasę zbiorów \mathcal{I} nazywa się σ -algebrą borelowską zbiorów w \mathbb{R}^n i oznacza się $B(\mathbb{R}^n)$.

Definicja 2.14 (Zmienna losowa [5, Sec 5.1 Def. 1]). Odwzorowanie $X : \Omega \rightarrow \mathbb{R}^n$ nazywamy zmienną losową o wartościach w \mathbb{R}^n , jeśli dla każdego $A \in B(\mathbb{R}^n)$ zbiór $X^{-1}(A) \in \mathcal{F}$.

Definicja 2.15 (Charakterystyki liczbowe zmiennej losowej [5], [7]).

- Wartością oczekiwaną zmiennej losowej X o wartościach w \mathbb{R} nazywamy liczbę:

$$E(X) = \int_{\Omega} X dP,$$

jeżeli X jest P -całkowalna, tzn. jeżeli zachodzi:

$$E(X) = \int_{\Omega} |X| dP < \infty.$$

- Kowariancją zmiennych losowych X, Y nazywamy liczbę:

$$\text{Cov}(X; Y) = E((X - E(X))(Y - E(Y))).$$

- Wariancją zmiennej losowej X nazywamy liczbę:

$$\text{Var}(X) = E((X - E(X))^2),$$

jeżeli wyznaczona wartość oczekiwana istnieje.

- Odchyleniem standardowym zmiennej losowej X nazywamy liczbę:

$$\sigma(X) = \sqrt{\text{Var}(X)}.$$

- Współczynnikiem korelacji nazywamy charakterystykę ilościową stopnia zależności dwóch zmiennych losowych X i Y zdefiniowaną następująco:

$$\rho(X, Y) = \frac{\text{Cov}(X; Y)}{\sigma(X) \sigma(Y)}.$$

Rozdział 3

Elementy eksploracji danych wykorzystywane w systemach rekomendujących

Większość systemów rekomendujących opiera się na algorytmach, które możemy rozumieć jako różne techniki eksploracji danych. Zazwyczaj proces eksploracji danych składa się z trzech kroków:

1. wstępnego przetwarzanie danych,
2. analizy danych,
3. interpretacji wyników.

W tym rozdziale zostaną przeanalizowane najważniejsze i najczęściej używane w regułach rekomendujących metody. Zaczniemy od miar podobieństw i technik redukcji wymiaru. W kolejnym etapie spojrzymy na metody klasyfikacji, grupowania i regresji, aby zakończyć interpretacją wyników i oceną błędów obliczeń.

3.1 Wstępne przetwarzanie danych

Przed przystąpieniem do kroku analizy dane wymagają przygotowania. W tej sekcji zostaną zawarte informacje na temat wstępnego przetwarzania danych, które spotykamy przy regułach rekomendujących.

3.1.1 Miary podobieństwa

W systemach rekomendujących bardzo częstym podejściem jest używanie metod klasyfikacji i grupowania. Metody te opierają się na obliczaniu podobieństw i odległości.

Najprostszym i jednocześnie najczęściej używanym kryterium jest odległość euklidesowa.

Definicja 3.1 (Odległość euklidesowa [15]). *Niech $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$, $n \in \mathbb{N}$. Odległością euklidesową między \mathbf{x} i \mathbf{y} nazywamy:*

$$d_e(x, y) = \sqrt{\sum_{k=1}^n (\mathbf{x}_k - \mathbf{y}_k)^2}.$$

Kolejnym podejściem, gdzie poszczególne elementy są postrzegane jako n - wymiarowe wektory, a podobieństwo między nimi jest obliczane na podstawie kąta, który tworzą jest odległość kosinusowa.

Definicja 3.2 (Odległość kosinusowa [15]). *Niech $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$, $n \in \mathbb{N}$. Odległością kosinusową nazywamy funkcję $d : \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$ opisaną wzorem:*

$$d(\mathbf{x}, \mathbf{y}) = 1 - \text{sim}(\mathbf{x}, \mathbf{y}),$$

gdzie $\text{sim} : \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$ wyznacza współczynnik podobieństwa wektorów \mathbf{x} i \mathbf{y} według formuły:

$$\text{sim}(\mathbf{x}, \mathbf{y}) = \frac{\sum_{k=1}^n x_k y_k}{\sqrt{\sum_{k=1}^n x_k^2} \sqrt{\sum_{k=1}^n y_k^2}}.$$

Innym podejściem, które pozwala na modelowanie podobieństwa między zmiennymi losowymi jest korelacja Pearsona, którą definiujemy następująco:

Definicja 3.3 (Współczynnik korelacji Pearsona [15]). *Niech X, Y będą zmiennymi losowymi o rozkładach ciągłych oraz niech $(x_1, \dots, x_n), (y_1, \dots, y_n)$ oznaczają losową próbę prostą. Przez \bar{x} i \bar{y} oznaczmy:*

$$\bar{x} = \frac{1}{n} \sum_{k=1}^n x_k, \quad \bar{y} = \frac{1}{n} \sum_{k=1}^n y_k.$$

Wówczas współczynnikiem korelacji Pearsona nazywamy:

$$\rho^p(X, Y) = \frac{\sum_{k=1}^n (x_k - \bar{x})(y_k - \bar{y})}{\sqrt{\sum_{k=1}^n (x_k - \bar{x})^2} \sqrt{\sum_{k=1}^n (y_k - \bar{y})^2}}.$$

Przy innych rodzajach danych do opisu podobieństwa można użyć Indeksu Jaccarda (współczynnik podobieństwa Jaccarda). Współczynnik ten mierzy podobieństwo między dwoma zbiorami, a definiujemy go niestępująco:

Definicja 3.4 (Indeks Jaccarda [4]). *Niech A i B oznaczają zbiory. Indeksem Jaccarda (podobieństwem Jaccarda) nazywamy funkcję:*

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|}.$$

3.1.2 Redukcja wymiaru

Zbyt duża ilość zmiennych, które opisują obserwacje, powoduje wzrost ryzyka, że zmienne te są ze sobą skorelowane, a informacje wnoszone przez część zmiennych są redundantne. W poniższym rozdziale przyjrzymy się najczęściej wybieranemu algorytmowi redukcji wymiarów w kontekście reguł rekomendujących. Jest to rozkład według wartości osobliwych (ang. Singular Value Decomposition (SVD)).

Definicja 3.5 (Rozkład według wartości osobliwych [13]). *Rozkładem według wartości osobliwych m na n - wymiarowej macierzy \mathbb{X} , gdzie $m \geq n$ nazywamy odszukanie takich macierzy $\mathbb{U}, \Sigma, \mathbb{V}$, że:*

$$\mathbb{X} = \mathbb{U}\Sigma\mathbb{V}^T,$$

gdzie:

- $\mathbb{U}^T\mathbb{U} = \mathbb{I}$, $\mathbb{V}^T\mathbb{V} = \mathbb{I}$, \mathbb{U} jest wymiaru m na m oraz \mathbb{V} wymiaru n na n ,
- Σ jest macierzą wymiaru m na n , gdzie $\sigma_{ij} = 0$, gdy $i \neq j$.

Uwaga 3.6. [13] Niezerowe wyrazy macierzy Σ nazywamy wartościami osobliwymi macierzy \mathbb{X} . Kolumny macierzy \mathbb{U} i \mathbb{V} nazywamy odpowiednio lewymi i prawymi wektorami szczególnymi macierzy \mathbb{X} .

Twierdzenie 3.7. [10] Niech $\mathbb{X} \in \mathbb{M}_{m \times n}$, gdzie $m \geq n$. Istnieje wtedy rozkład według wartości osobliwych:

$$\mathbb{X} = \mathbb{U}\Sigma\mathbb{V}^T.$$

Definicja 3.8 (Wartość osobliwa macierzy [17]). *Wartością osobliwą σ_k macierzy \mathbb{X} nazywamy*

$$\sigma_k = \sqrt{\lambda_k},$$

gdzie λ_k , $k \in \mathbb{N}$ jest nieujemną wartością własną macierzy $\mathbb{X}^T\mathbb{X}$.

Uwaga 3.9. Należy zauważyć, że macierz $\mathbb{X}^T\mathbb{X}$ jest macierzą symetryczną, zatem jej wartości własne należą do zbioru liczb $\mathbb{R}_+ \cup \{0\}$.

W przypadku gdy wyrazy macierzy \mathbb{X} są liczbami ze zbioru liczb zespolonych szukamy wartości własnych macierzy $\mathbb{X}^*\mathbb{X}$, gdzie \mathbb{X}^* jest sprzężeniem hermitowskim macierzy, tzn. złożeniem operacji transpozycji i sprzężenia zespolonego ($\mathbb{X}^* = \overline{\mathbb{X}^T}$).

Definicja 3.10 (Norma Frobeniusa[13]). Niech $\mathbb{A} \in \mathbb{M}_{m \times n}(\mathbb{R})$. Normą Frobeniusa nazywamy:

$$\|\mathbb{A}\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2} = \sqrt{\text{tr}(\mathbb{A}^T\mathbb{A})}.$$

Twierdzenie 3.11 (Warunki równoważne SVD [13]). *Niech rozkład według wartości osobliwych macierzy \mathbb{X} będzie dany wzorem*

$$\mathbb{X} = \mathbb{U}\Sigma\mathbb{V}^T$$

gdzie

$$\mathbb{U} = [\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_m], \mathbb{V} = [\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n], \Sigma = \begin{bmatrix} \sigma_1 & 0 & \dots & 0 \\ 0 & \sigma_2 & \dots & 0 \\ \dots & \dots & \ddots & \dots \\ 0 & 0 & \dots & \sigma_n \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 0 \end{bmatrix},$$

oraz $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > \sigma_{r+1} = \dots = \sigma_n = 0$. Przez $\mathbf{u}_1, \dots, \mathbf{u}_m$ i $\mathbf{v}_1, \dots, \mathbf{v}_n$ oznaczamy kolumny macierzy \mathbb{U} i \mathbb{V} . Wtedy:

1. $\text{rz}(\mathbb{X}) = r$,
2. $\mathbb{X} = \sum_{i=1}^r \mathbf{u}_i \cdot \sigma_i \cdot \mathbf{v}_i^T$,
3. $\|\mathbb{X}\|_F^2 = \sigma_1^2 + \dots + \sigma_r^2$.

Lemat 3.12. *Niech \mathbb{X} będzie jak w twierdzeniu 3.11.*

Wtedy:

$$\mathbb{X} = \sum_{i=1}^r \mathbf{u}_i \cdot \sigma_i \cdot \mathbf{v}_i^T.$$

Dowód. (opracowanie własne) Niech $\mathbb{X} \in \mathbb{M}_{m \times n}$, $m \geq n$ oraz $\mathbb{X} = \mathbb{U}\Sigma\mathbb{V}^T$, gdzie macierze $\mathbb{U}, \Sigma, \mathbb{V}$ mają postać:

$$\mathbb{U} = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1m} \\ a_{21} & a_{22} & \dots & a_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \dots & a_{mm} \end{bmatrix}, \mathbb{V} = \begin{bmatrix} b_{11} & b_{12} & \dots & b_{1n} \\ b_{21} & b_{22} & \dots & b_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ b_{n1} & b_{n2} & \dots & b_{nn} \end{bmatrix}, \Sigma = \begin{bmatrix} \sigma_{11} & 0 & \dots & 0 \\ 0 & \sigma_{22} & \dots & 0 \\ \vdots & \vdots & \ddots & \sigma_{nn} \\ 0 & 0 & \dots & 0 \end{bmatrix}.$$

Iloczyn $\mathbb{U}\Sigma$ jest równy

$$\mathbb{U}\Sigma = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1m} \\ a_{21} & a_{22} & \dots & a_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \dots & a_{mm} \end{bmatrix} \cdot \begin{bmatrix} \sigma_{11} & 0 & \dots & 0 \\ 0 & \sigma_{22} & \dots & 0 \\ \vdots & \vdots & \ddots & \sigma_{nn} \\ 0 & 0 & \dots & 0 \end{bmatrix} = \begin{bmatrix} a_{11}\sigma_{11} & a_{12}\sigma_{22} & \dots & a_{1n}\sigma_{nn} \\ a_{21}\sigma_{11} & a_{22}\sigma_{22} & \dots & a_{2n}\sigma_{nn} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1}\sigma_{11} & a_{m2}\sigma_{22} & \dots & a_{mn}\sigma_{nn} \end{bmatrix}.$$

Zatem :

$$\mathbb{U}\Sigma\mathbb{V}^T = \begin{bmatrix} a_{11}\sigma_{11} & a_{12}\sigma_{22} & \dots & a_{1n}\sigma_{nn} \\ a_{21}\sigma_{11} & a_{22}\sigma_{22} & \dots & a_{2n}\sigma_{nn} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1}\sigma_{11} & a_{m2}\sigma_{22} & \dots & a_{mn}\sigma_{nn} \end{bmatrix} \cdot \begin{bmatrix} b_{11} & b_{21} & \dots & b_{n1} \\ b_{12} & b_{22} & \dots & b_{n2} \\ \vdots & \vdots & \ddots & \vdots \\ b_{1n} & b_{2n} & \dots & b_{nn} \end{bmatrix} =$$

$$\begin{bmatrix} a_{11}\sigma_{11}b_{11} + a_{12}\sigma_{22}b_{12} + \dots + a_{1n}\sigma_{nn}b_{1n} & \dots & a_{11}\sigma_{11}b_{n1} + a_{12}\sigma_{22}b_{n2} + \dots + a_{1n}\sigma_{nn}b_{nn} \\ a_{21}\sigma_{11}b_{11} + a_{22}\sigma_{22}b_{12} + \dots + a_{2n}\sigma_{nn}b_{1n} & \dots & a_{21}\sigma_{11}b_{n1} + a_{22}\sigma_{22}b_{n2} + \dots + a_{2n}\sigma_{nn}b_{nn} \\ \vdots & \ddots & \vdots \\ a_{m1}\sigma_{11}b_{11} + a_{m2}\sigma_{22}b_{12} + \dots + a_{mn}\sigma_{nn}b_{1n} & \dots & a_{m1}\sigma_{11}b_{n1} + a_{m2}\sigma_{22}b_{n2} + \dots + a_{mn}\sigma_{nn}b_{nn} \end{bmatrix}.$$

Z drugiej strony pokażemy że $\mathbb{X} = \sum_{i=1}^r \mathbf{u}_i \cdot \sigma_{ii} \cdot \mathbf{v}_i^T$.

Niech $r = n$.

Mamy, że:

$$1. \begin{bmatrix} a_{11} \\ a_{21} \\ \vdots \\ a_{m1} \end{bmatrix} \sigma_{11} \begin{bmatrix} b_{11} \\ b_{21} \\ \vdots \\ b_{n1} \end{bmatrix}^T = \begin{bmatrix} a_{11}\sigma_{11} \\ a_{21}\sigma_{11} \\ \vdots \\ a_{m1}\sigma_{11} \end{bmatrix} [b_{11} \ b_{21} \ \dots \ b_{n1}]$$

$$= \begin{bmatrix} a_{11}\sigma_{11}b_{11} & a_{11}\sigma_{11}b_{21} & \dots & a_{11}\sigma_{11}b_{n1} \\ a_{21}\sigma_{11}b_{11} & a_{21}\sigma_{11}b_{21} & \dots & a_{21}\sigma_{11}b_{n1} \\ \vdots & \dots & \ddots & \vdots \\ a_{m1}\sigma_{11}b_{11} & a_{m1}\sigma_{11}b_{21} & \dots & a_{m1}\sigma_{11}b_{n1} \end{bmatrix}$$

,

$$2. \begin{bmatrix} a_{12} \\ a_{22} \\ \vdots \\ a_{m2} \end{bmatrix} \sigma_{22} \begin{bmatrix} b_{12} \\ b_{22} \\ \vdots \\ b_{n2} \end{bmatrix}^T = \begin{bmatrix} a_{12}\sigma_{22} \\ a_{22}\sigma_{22} \\ \vdots \\ a_{m2}\sigma_{22} \end{bmatrix} [b_{12} \ b_{22} \ \dots \ b_{n2}]$$

$$= \begin{bmatrix} a_{12}\sigma_{22}b_{12} & a_{12}\sigma_{22}b_{22} & \dots & a_{12}\sigma_{22}b_{n2} \\ a_{22}\sigma_{22}b_{12} & a_{22}\sigma_{22}b_{22} & \dots & a_{22}\sigma_{22}b_{n2} \\ \vdots & \dots & \ddots & \vdots \\ a_{m2}\sigma_{22}b_{12} & a_{m2}\sigma_{22}b_{22} & \dots & a_{m2}\sigma_{22}b_{n2} \end{bmatrix},$$

\vdots

$$\begin{aligned} \text{n. } \begin{bmatrix} a_{1n} \\ a_{2n} \\ \vdots \\ a_{mn} \end{bmatrix} \sigma_{nn} \begin{bmatrix} b_{1n} \\ b_{2n} \\ \vdots \\ b_{nn} \end{bmatrix}^T &= \begin{bmatrix} a_{1n}\sigma_{nn} \\ a_{2n}\sigma_{nn} \\ \vdots \\ a_{mn}\sigma_{nn} \end{bmatrix} [b_{1n} \ b_{2n} \ \dots \ b_{nn}] \\ &= \begin{bmatrix} a_{1n}\sigma_{nn}b_{1n} & a_{1n}\sigma_{nn}b_{2n} & \dots & a_{1n}\sigma_{nn}b_{nn} \\ a_{2n}\sigma_{nn}b_{1n} & a_{2n}\sigma_{nn}b_{2n} & \dots & a_{2n}\sigma_{nn}b_{nn} \\ \vdots & \dots & \ddots & \vdots \\ a_{mn}\sigma_{nn}b_{1n} & a_{mn}\sigma_{nn}b_{2n} & \dots & a_{mn}\sigma_{nn}b_{nn} \end{bmatrix}. \end{aligned}$$

Sumując macierze powstałe w n krokach otrzymujemy:

$$\begin{bmatrix} a_{11}\sigma_{11}b_{11} + a_{12}\sigma_{22}b_{12} + \dots + a_{1n}\sigma_{nn}b_{1n} & \dots & a_{11}\sigma_{11}b_{n1} + a_{12}\sigma_{22}b_{n2} + \dots + a_{1n}\sigma_{nn}b_{nn} \\ a_{21}\sigma_{11}b_{11} + a_{22}\sigma_{22}b_{12} + \dots + a_{2n}\sigma_{nn}b_{1n} & \dots & a_{21}\sigma_{11}b_{n1} + a_{22}\sigma_{22}b_{n2} + \dots + a_{2n}\sigma_{nn}b_{nn} \\ \vdots & \ddots & \vdots \\ a_{m1}\sigma_{11}b_{11} + a_{m2}\sigma_{22}b_{12} + \dots + a_{mn}\sigma_{nn}b_{1n} & \dots & a_{m1}\sigma_{11}b_{n1} + a_{m2}\sigma_{22}b_{n2} + \dots + a_{mn}\sigma_{nn}b_{nn} \end{bmatrix}.$$

Zatem, z dowolności \mathbb{U} , Σ i \mathbb{V} , jeżeli $\mathbb{X} = \mathbb{U}\Sigma\mathbb{V}^T$, to $\mathbb{X} = \sum_{i=1}^r \mathbf{u}_i \cdot \sigma_{ii} \cdot \mathbf{v}_i^T$. \square

Lemat 3.13. Niech \mathbb{X} będzie jak w twierdzeniu 3.11.

Wtedy:

$$\|\mathbb{X}\|_F^2 = \sigma_1^2 + \dots + \sigma_r^2.$$

Dowód. (opracowanie własne) Niech $\mathbb{X} = \mathbb{U}\Sigma\mathbb{V}^T$. Mamy zatem

$$\|\mathbb{X}\|_F^2 = \|\mathbb{U}\Sigma\mathbb{V}^T\|_F^2.$$

Z definicji metryki Frobeniusa

$$\|\mathbb{U}\Sigma\mathbb{V}^T\|_F^2 = \text{tr} \left((\mathbb{U}\Sigma\mathbb{V}^T)^T (\mathbb{U}\Sigma\mathbb{V}^T) \right).$$

Opierając się na własności $(\mathbb{A}\mathbb{C})^T = \mathbb{C}^T\mathbb{A}^T$ transpozycji macierzy otrzymamy:

$$\|\mathbb{U}\Sigma\mathbb{V}^T\|_F^2 = \text{tr} \left((\mathbb{U}\Sigma\mathbb{V}^T)^T (\mathbb{U}\Sigma\mathbb{V}^T) \right) = \text{tr} \left(\mathbb{V}\Sigma^T\mathbb{U}^T\mathbb{U}\Sigma\mathbb{V}^T \right).$$

Skoro $\mathbb{U}^T\mathbb{U} = \mathbb{I}$, to

$$\text{tr} \left(\mathbb{V}\Sigma^T\mathbb{U}^T\mathbb{U}\Sigma\mathbb{V}^T \right) = \text{tr} \left(\mathbb{V}\Sigma^T\Sigma\mathbb{V}^T \right).$$

Następnie z własności śladu macierzy $\text{tr}(\mathbb{A}\mathbb{B}\mathbb{C}) = \text{tr}(\mathbb{C}\mathbb{A}\mathbb{B}) = \text{tr}(\mathbb{B}\mathbb{C}\mathbb{A})$ wynika, że:

$$\text{tr} \left(\mathbb{V}\Sigma^T\Sigma\mathbb{V}^T \right) = \text{tr} \left(\mathbb{V}^T\mathbb{V}\Sigma^T\Sigma \right).$$

Skoro $\mathbb{V}^T \mathbb{V} = \mathbb{I}$, to

$$\text{tr}(\mathbb{V}^T \mathbb{V} \Sigma^T \Sigma) = \text{tr}(\Sigma^T \Sigma).$$

Bezpośrednio możemy więc policzyć, że

$$\|\mathbb{X}\|_F^2 = \sigma_1^2 + \sigma_2^2 + \dots + \sigma_r^2.$$

□

Twierdzenie 3.14 (Twierdzenie Eckart - Younga [10, 13]). *Niech $\mathbb{X} \in \mathbb{M}_{m \times n}(\mathbb{R})$, $m \geq n$, $m, n \in \mathbb{N}$. $\mathbb{U}, \Sigma, \mathbb{V}$ niech będą takie, że $\mathbb{U} \in \mathbb{M}_{m \times m}(\mathbb{R})$, $\Sigma \in \mathbb{M}_{m \times n}(\mathbb{R})$, $\mathbb{V} \in \mathbb{M}_{n \times n}(\mathbb{R})$, $\mathbb{X} = \mathbb{U} \Sigma \mathbb{V}^T$, $\mathbb{U}^T \mathbb{U} = \mathbb{I}$, $\mathbb{V}^T \mathbb{V} = \mathbb{I}$ oraz:*

$$\mathbb{U} = [\mathbf{u}_1, \dots, \mathbf{u}_m], \Sigma = \begin{bmatrix} \sigma_1 & 0 & \dots & 0 \\ 0 & \sigma_2 & \dots & 0 \\ \dots & \dots & \ddots & \dots \\ 0 & 0 & \dots & \sigma_n \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 0 \end{bmatrix}, \mathbb{V} = [\mathbf{v}_1, \dots, \mathbf{v}_n],$$

gdzie przez $\mathbf{u}_1, \dots, \mathbf{u}_m$ i $\mathbf{v}_1, \dots, \mathbf{v}_n$ oznaczamy kolumny macierzy \mathbb{U} i \mathbb{V} . (Istnienie powyższych macierzy wynika z twierdzenia 3.7.)

Zdefiniujmy dla $k \in \{1, \dots, n\}$:

$$\mathbb{X}_k = \sum_{i=1}^k \mathbf{u}_i \cdot \sigma_i \cdot \mathbf{v}_i^T,$$

wtedy

$$\|\mathbb{X} - \mathbb{X}_k\|_F^2 = \min_{\text{rz}(\mathbb{B}) \leq k} \|\mathbb{X} - \mathbb{B}\|_F^2 = \sigma_{k+1}^2 + \dots + \sigma_n^2.$$

Aby udowodnić twierdzenie Eckart - Younga wprowadźmy twierdzenie Weylsa:

Twierdzenie 3.15 (Twierdzenie Weylsa [2]). *Niech $\mathbb{X}, \mathbb{Y} \in \mathbb{M}_{m \times n}(\mathbb{K})$. Dodatkowo niech odpowiednio $\sigma_1(\mathbb{X}) \geq \sigma_2(\mathbb{X}) \geq \dots \geq \sigma_r(\mathbb{X}) \geq 0$, $\sigma_1(\mathbb{Y}) \geq \sigma_2(\mathbb{Y}) \geq \dots \geq \sigma_r(\mathbb{Y}) \geq 0$ i $\sigma_1(\mathbb{Z}) \geq \sigma_2(\mathbb{Z}) \geq \dots \geq \sigma_r(\mathbb{Z}) \geq 0$ będą kolejnymi wartościami osobliwymi macierzy \mathbb{X}, \mathbb{Y} , $\mathbb{Z} = \mathbb{X} + \mathbb{Y}$. Wtedy:*

$$\sigma_{i+j-1}(\mathbb{Z}) \leq \sigma_i(\mathbb{X}) + \sigma_j(\mathbb{Y}),$$

gdzie $1 \leq i, j \leq r$, $i + j \leq r + 1$.

Lemat 3.16. *Niech \mathbb{X} będzie jak w twierdzeniu 3.11 oraz \mathbb{X}_k będą jak w twierdzeniu 3.14. Wtedy:*

$$\|\mathbb{X} - \mathbb{X}_k\|_F^2 = \sigma_{k+1}^2 + \sigma_{k+2}^2 + \dots + \sigma_r^2.$$

Dowód. (opracowanie własne) Niech $\mathbb{X} = \mathbb{U}\Sigma\mathbb{V}^T$ oraz $\mathbb{X}_k = \mathbb{U}\Sigma_k\mathbb{V}^T$. Mamy zatem

$$\|\mathbb{X} - \mathbb{X}_k\|_F^2 = \|\mathbb{U}\Sigma\mathbb{V}^T - \mathbb{U}\Sigma_k\mathbb{V}^T\|_F^2.$$

Stąd

$$\|\mathbb{U}\Sigma\mathbb{V}^T - \mathbb{U}\Sigma_k\mathbb{V}^T\|_F^2 = \|\mathbb{U}(\Sigma\mathbb{V}^T - \Sigma_k\mathbb{V}^T)\|_F^2 = \|\mathbb{U}(\Sigma - \Sigma_k)\mathbb{V}^T\|_F^2.$$

Z definicji metryki Frobeniusa

$$\|\mathbb{U}(\Sigma - \Sigma_k)\mathbb{V}^T\|_F^2 = \text{tr} \left((\mathbb{U}(\Sigma - \Sigma_k)\mathbb{V}^T)^T (\mathbb{U}(\Sigma - \Sigma_k)\mathbb{V}^T) \right).$$

Opierając się na własności $(\mathbb{A}\mathbb{C})^T = \mathbb{C}^T\mathbb{A}^T$ transpozycji macierzy otrzymamy:

$$\text{tr} \left((\mathbb{U}(\Sigma - \Sigma_k)\mathbb{V}^T)^T (\mathbb{U}(\Sigma - \Sigma_k)\mathbb{V}^T) \right) = \text{tr} \left(\mathbb{V}(\Sigma - \Sigma_k)^T \mathbb{U}^T \mathbb{U}(\Sigma - \Sigma_k)\mathbb{V}^T \right).$$

Skoro $\mathbb{U}^T\mathbb{U} = \mathbb{I}$, to

$$\text{tr} \left(\mathbb{V}(\Sigma - \Sigma_k)^T \mathbb{U}^T \mathbb{U}(\Sigma - \Sigma_k)\mathbb{V}^T \right) = \text{tr} \left(\mathbb{V}(\Sigma - \Sigma_k)^T (\Sigma - \Sigma_k)\mathbb{V}^T \right).$$

Następnie z własności śladu macierzy $\text{tr}(\mathbb{A}\mathbb{B}\mathbb{C}) = \text{tr}(\mathbb{C}\mathbb{A}\mathbb{B}) = \text{tr}(\mathbb{B}\mathbb{C}\mathbb{A})$ wynika, że:

$$\text{tr} \left(\mathbb{V}(\Sigma - \Sigma_k)^T (\Sigma - \Sigma_k)\mathbb{V}^T \right) = \text{tr} \left(\mathbb{V}^T \mathbb{V}(\Sigma - \Sigma_k)^T (\Sigma - \Sigma_k) \right).$$

Skoro $\mathbb{V}^T\mathbb{V} = \mathbb{I}$, to

$$\text{tr} \left(\mathbb{V}^T \mathbb{V}(\Sigma - \Sigma_k)^T (\Sigma - \Sigma_k) \right) = \text{tr} \left((\Sigma - \Sigma_k)^T (\Sigma - \Sigma_k) \right).$$

Bezpośrednio możemy więc sformułować, że

$$\|\mathbb{X} - \mathbb{X}_k\|_F^2 = \sigma_{k+1}^2 + \sigma_{k+2}^2 + \dots + \sigma_r^2.$$

□

Dowód. [Twierdzenia 3.14]

Niech $\mathbb{X} \in \mathbb{M}_{m \times n}(\mathbb{R})$ będzie macierzą o wartościach rzeczywistych, gdzie $m \geq n$. Załóżmy, że

$$\mathbb{X} = \mathbb{U}\Sigma\mathbb{V}^T$$

jest rozkładem według wartości osobliwych macierzy \mathbb{X} . Chcemy pokazać, że najlepszym przybliżeniem macierzy \mathbb{X} w normie Frobeniusa (oznaczamy $\|\cdot\|_F$) jest

$$\mathbb{X}_k = \sum_{i=1}^k \mathbf{u}_i \cdot \sigma_i \cdot \mathbf{v}_i^T,$$

gdzie \mathbf{u}_i i \mathbf{v}_i oznaczają odpowiednio i -te kolumny macierzy \mathbb{U} i \mathbb{V} .

Zauważmy, że z własności 2. twierdzenia 3.11 mamy

$$||\mathbb{X} - \mathbb{X}_k||_F^2 = || \sum_{i=k+1}^n \mathbf{u}_i \cdot \sigma_i \cdot \mathbf{v}_i^T ||_F^2.$$

Po obliczeniach wykonanych w lemacie 3.16 mamy, że

$$|| \sum_{i=k+1}^n \mathbf{u}_i \cdot \sigma_i \cdot \mathbf{v}_i^T ||_F^2 = \sum_{i=k+1}^n \sigma_i^2.$$

Niech $\mathbb{B}_k \in \mathbb{M}_{m \times n}(\mathbb{R})$, $\text{rz}(\mathbb{B}_k) = k$. Stąd należy udowodnić, że

$$||\mathbb{X} - \mathbb{X}_k||_F^2 = \sum_{i=k+1}^n \sigma_i^2 \leq ||\mathbb{X} - \mathbb{B}_k||_F^2.$$

Niech $\mathbb{X}' = \mathbb{X} - \mathbb{B}_k$, $\mathbb{X}'' = \mathbb{B}_k$. Wtedy $\mathbb{X} = \mathbb{X}' + \mathbb{X}''$.

Korzystając z Twierdzenia Weylsa otrzymujemy

$$\sigma_{i+j-1}(\mathbb{X}) \leq \sigma_i(\mathbb{X}') + \sigma_j(\mathbb{X}'').$$

Stąd

$$\sigma_{i+j-1}(\mathbb{X}) \leq \sigma_i(\mathbb{X} - \mathbb{B}_k) + \sigma_j(\mathbb{B}_k).$$

Niech $i \in \{1, \dots, r-k\}$ oraz $j = k+1$. Wtedy

$$\forall_{i \in \{1, \dots, r-k\}} \sigma_{i+k}(\mathbb{X}) \leq \sigma_i(\mathbb{X} - \mathbb{B}_k) + \sigma_{k+1}(\mathbb{B}_k).$$

Zauważmy, że

$$\sigma_{k+1}(\mathbb{B}_k) = 0,$$

wtedy

$$\forall_{i \in \{1, \dots, r-k\}} \sigma_i(\mathbb{X} - \mathbb{B}_k) \geq \sigma_{i+k}(\mathbb{X}).$$

Stąd

$$||\mathbb{X} - \mathbb{B}_k||_F^2 = \sum_{i=1}^r \sigma_i(\mathbb{X} - \mathbb{B}_k)^2 \geq \sum_{i=1}^{r-k} \sigma_i(\mathbb{X} - \mathbb{B}_k)^2 \geq \sum_{i=1}^{r-k} \sigma_{i+k}(\mathbb{X})^2 \geq \sum_{i=k+1}^r \sigma_i(\mathbb{X})^2 = ||\mathbb{X} - \mathbb{X}_k||_F^2.$$

Ostatecznie

$$||\mathbb{X} - \mathbb{B}_k||_F^2 \geq ||\mathbb{X} - \mathbb{X}_k||_F^2.$$

□

3.2 Metody eksploracji danych

Termin eksploracja danych jest często używany jako określenie procesu odkrywania wiedzy z danych. Często jednak terminem "proces odkrywania wiedzy" określamy cały proces pracy z danymi, natomiast termin "eksploracja danych" odnosi się do etapu odkrywania pewnego rodzaju reguł.

Jak podaje Tadeusz Morzy [11] w metodach eksploracji można wyróżnić:

- odkrywanie asocjacji,
- klastrowanie,
- odkrywanie wzorców sekwencji,
- odkrywanie klasyfikacji,
- odkrywanie podobieństw w przebiegach czasowych,
- wykrywanie zmian i odchyłeń.

W tej sekcji zostaną przedstawione te metody, które stosowane są najczęściej w regułach rekomendujących.

3.2.1 Algorytm k - najbliższych sąsiadów

Opis poniższego algorytmu oparty został na [15, Sec 2.3.1].

Algorytm k -najbliższych sąsiadów (k -NN) jest powszechnie używanym algorytmem klasyfikacji.

Przyporządkowanie nowych elementów zostaje przeprowadzone na podstawie porównania obserwacji z k najbardziej podobnymi jej obiektami ze zbioru treningowego. Podstawowe założenie algorytmu mówi, że jeżeli nowy rekord znajduje się w pewnym otoczeniu, to na podstawie k - najbliższych mu obserwacji zostanie przyporządkowana do niego klasa, której pojawienie się w rozważanym zbiorze jest najliczniejsze.

Niech q będzie punktem dla którego chcemy odnaleźć jego klasę l .

$X = \{\{x_1, l_1\}, \dots, \{x_n, l_n\}\}$ niech będzie zbiorem treningowym, gdzie x_j jest j -tym elementem zbioru obserwacji, natomiast l_j etykietką klasy do której obserwacja należy, $j \in \{1, \dots, n\}$.

Przeprowadzając algorytm k -NN zaczynamy od wyboru podzbioru

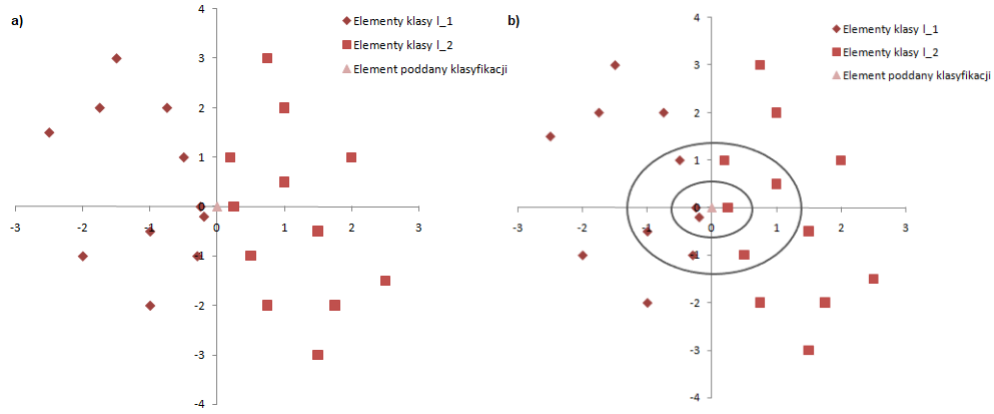
$$Y = \{\{y_1, l_1\}, \dots, \{y_k, l_k\}\},$$

$k \in \{1, \dots, n\}$ takiego, że $Y \in X$ oraz

$$\sum_{i=1}^k d(q, y_i)$$

jest minimalna. Y zawiera więc k punktów z X , które leżą najbliżej rozważanego punktu q . Następnie do punktu q zostaje przyporządkowana klasa taka, że

$$l = f(\{l_1, \dots, l_k\}).$$



Rysunek 3.1: Metoda k - najbliższych sąsiadów. (źródło własne)

Na powyższych wykresach widzimy przykładowe zastosowanie algorytmu k -NN. Na pierwszym rysunku przedstawiony został zbiór treningowy z podziałem na dwie klasy (rąby, koła) oraz punkt (trójkąt), który będziemy chcieli przyporządkować do jednej z nich. Na drugim przedstawiono natomiast dwa okręgi, jedno prezentujące najbliższe sąsiedztwo dla $k = 3$, drugie dla $k = 9$. W obu przypadkach nowy punkt zostanie przyporządkowany do klasy l_1 . Warto jednak zauważyć, że znajduje się on na granicy dwóch klastrów, przez co przy innym wyborze k może zostać przyporządkowany do klasy l_2 .

Opisana metoda przyporządkowuje wybranemu rekordowi najbardziej mu podobne. Wykorzystuje do tego pewne miary odległości.

Najtrudniejszym zadaniem przy przeprowadzaniu algorytmu k -NN jest często wybór k . Jeżeli k będzie zbyt małe - klasyfikator stanie się bardzo wrażliwy, jeżeli jednak k będzie zbyt duże sąsiedztwo może zawierać zbyt dużo punktów z innych klas. Rozważając przypadek z rysunku 3.1 łatwo zauważyć, że nawet mała zmiana w obserwacjach zbioru treningowego może doprowadzić do zmiany wyniku.

3.2.2 Algorytm k - średnich

Opis poniższego algorytmu oparty został na [12, Sec 2.3.1].

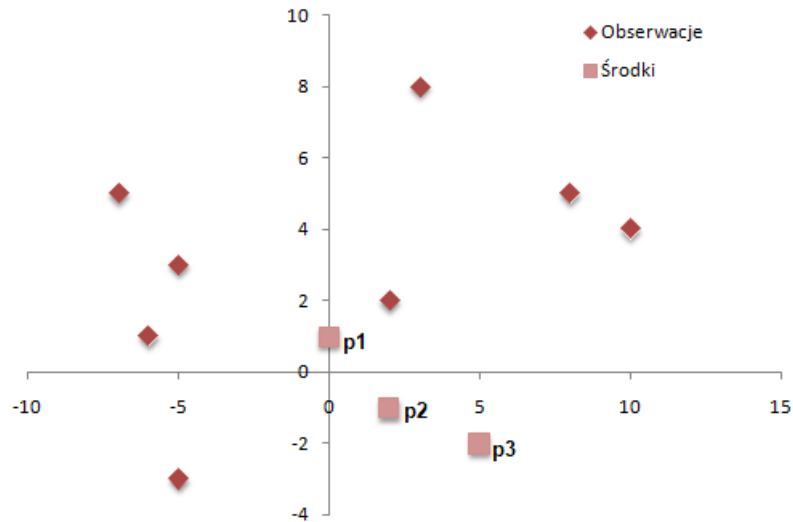
Algorytm k -średnich jest prostym i zarazem efektywnym algorytmem grupowania.

Głównym celem algorytmu jest podział pewnego zbioru X :

$$X = \{\mathbf{x}_i = (\mathbf{x}_{i1}, \dots, \mathbf{x}_{id}) : i \in \{1, \dots, N\}\},$$

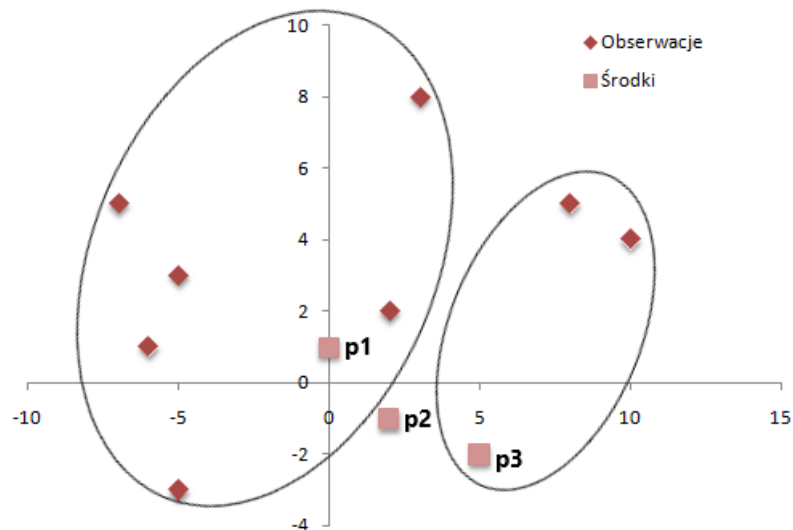
gdzie \mathbf{x}_i jest d - wymiarowym wektorem cech opisującym obiekt na podzbiory.

W wyniku grupowania n - elementowego zbioru X na k podgrup otrzymujemy macierz podziału \mathbb{A} o wymiarach $k \times n$. Każdy z elementów tej macierzy a_{ik} oznacza stopień w jakim wektor \mathbf{x}_k przynależy do grupy. Na wstępie algorytmu ustalamy wartość parametru k jako liczbę grup, które chcemy wyodrębnić. Wybieramy k reprezentantów, które stanowią prototypy grup.



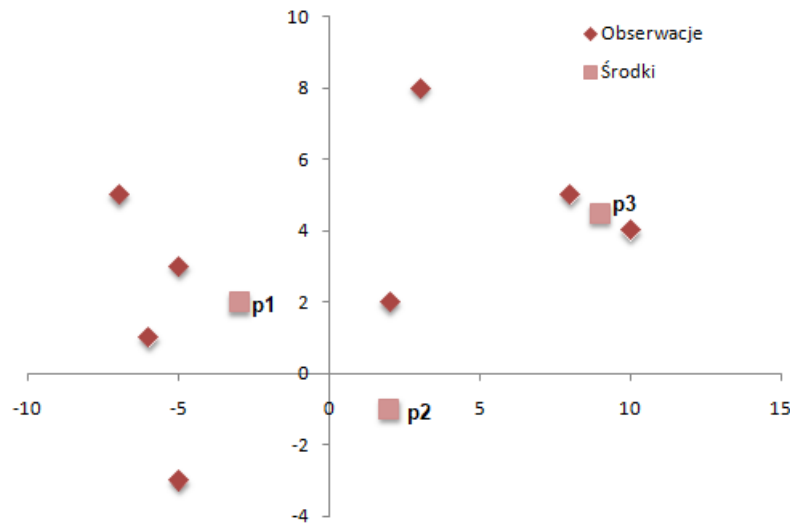
Rysunek 3.2: Metoda k -średnich. Wybór początkowych środków. (źródło własne)

W powyższym przykładzie (rysunek 3.2) wybranymi środkami są punkty $p1$, $p2$, $p3$. Kolejnym krokiem jest przypisanie każdego z elementów do najbliższej mu grupy.



Rysunek 3.3: Metoda k -średnich. Przypisanie elementów do grup. (źródło własne)

Dla każdej z tak ustalonych grup obliczamy średnią arytmetyczną współrzędnych, które staną się kolejnymi środkami.



Rysunek 3.4: Metoda k -średnich. Wybór nowych środków. (źródło własne)

Kroki te są wykonywane do momentu występowania migracji między obiektami. W algorytmie k -średnich liczba grup pozostaje więc niezmienną, zmienna jest tylko przynależność do grup. W metodzie tej poszukiwanie optymalnego podziału odpowiada wyznaczaniu takich grup, które minimalizują następującą funkcję kryterialną:

$$J(\mathbf{p}, \mathbb{A}) = \sum_{i=1}^k \sum_{k=1}^N a_{ki} d(\mathbf{p}_i, \mathbf{x}_k)^2,$$

gdzie

- $d(\mathbf{p}_i, \mathbf{x}_k)$ oznacza odległość elementu reprezentowanego przez wektor \mathbf{x} od grupy wyznaczonej przez środek \mathbf{p} ,
- N to liczebność zbioru X ,
- \mathbb{A} oznacza macierz podziału.

3.3 Szacowanie błędów obliczeń

3.3.1 Ocena dokładności metody

Najczęściej używanymi miarami dokładności modelu są:

- błąd średni (Mean Error),
- średni błąd bezwzględny (Mean Absolute Error),
- średni błąd kwadratowy (Mean Squared Error).

Niech dla elementu i ze zbioru $P = \{x_1, \dots, x_n\}$ będzie dostarczona predykcja \hat{r}_i . Aby ocenić jakość jej wyniku należy porównać ją ze znaną wartością r_i .

Definicja 3.17 (Błąd średni [15, Sec 4.1.1]). *Średnim błędem nazywamy wartość wyrażenia:*

$$ME = \frac{1}{|P|} \sum_{x_i \in P} (\hat{r}_i - r_i).$$

Definicja 3.18 (Średni błąd bezwzględny [15, Sec 4.1.1]). *Średnim błędem bezwzględnym nazywamy wartość wyrażenia:*

$$MAE = \frac{1}{|P|} \sum_{x_i \in P} |\hat{r}_i - r_i|.$$

Definicja 3.19 (Średni błąd kwadratowy [15, Sec 4.1.1]). *Średnim błędem kwadratowym nazywamy wartość wyrażenia:*

$$MSE = \frac{1}{|P|} \sum_{x_i \in P} (\hat{r}_i - r_i)^2.$$

Uwaga 3.20. [15, Sec 4.1.1] *Funkcja kwadratowa jest funkcją wypukłą co pozwala na dość częste zastępowanie średniego błędu kwadratowego przez średnią kwadratową błędów (Root Mean Squared Error (RMSE)):*

$$RMSE = \sqrt{MSE}$$

Normalized RMSE (NRMSE) oraz Normalized MAE (NMAE) są znormalizowanymi, przez użycie zakresu wartości $r_{max} - r_{min}$, wersjami błędów RMSE i MAE.

Kolejnym rodzajem powszechnie używanego błędu, który pozwala na użycie sum ważonych, jest średni błąd RMSE (Average RMSE).

Definicja 3.21 (Średni błąd RMSE [15, Sec 4.1.1]). *Niech $w_i > 0$ będzie wagą dla elementu i oraz niech $\sum w_i = 1$.*

Średnim błędem RMSE nazywamy wartość wyrażenia:

$$ARMSE = \sqrt{\sum_{x_i \in P} w_i (\hat{r}_i - r_i)^2}.$$

3.3.2 Ocena jakości modelu

Ocenę jakości modelu przeprowadza się na zbiorze testowym. Dla każdego z rekordów jest znana jego etykieta. Rekordy te są poddawane działaniu modelu, a następnie etykiety przypisane rekordom przez model są porównywalne z rzeczywistymi wartościami etykiet.

W następnym kroku zliczana jest liczba rekordów poprawnie i niepoprawnie zaklasyfikowanych przez model, a wynik testu zostaje przedstawiony w postaci macierzy pomyłek.

Definicja 3.22 (Macierz pomyłek [11, Sec 4.8.1]). *Macierz pomyłek nazywamy macierz kwadratową $m \times m$ (m oznacza liczbę etykiet), gdzie wiersze opisują etykiety praktyczne, natomiast kolumny etykiety przyporządkowane rekordom przez model. Element macierzy - $f_{i,j}$ oznacza liczbę rekordów z etykietą E_i , którym została przypisana etykieta E_j .*

\mathbb{F}	E_1	E_2
E_1	f_{11}	f_{12}
E_2	f_{21}	f_{22}

Tabela 3.1: Macierz pomyłek.

Uwaga 3.23. [11, Sec 4.8.1] Często elementy macierzy pomyłek dla problemów klasyfikacji binarnej oznacza się symbolami : TP , TN , FN , FP . Oznaczenia te symbolizują cztery możliwe przypadki występujące w klasyfikacji binarnej. Załóżmy, że wyróżniamy klasę pozytywną (+) i negatywną (-). Wtedy :

- TP (ang. *true - positive*) - liczba pozytywnych rekordów testowych zaklasyfikowanych do klasy pozytywnej,
- FN (ang. *false - negative*) - liczba pozytywnych rekordów testowych zaklasyfikowanych do klasy negatywnej,
- FP (ang. *false - positive*) - liczba negatywnych rekordów testowych zaklasyfikowanych do klasy pozytywnej,
- TN (ang. *true - negative*) - liczba negatywnych rekordów testowych zaklasyfikowanych do klasy negatywnej.

Macierz pomyłek przyjmuje wtedy postać:

\mathbb{F}	+	-
+	TP	FN
-	FP	TN

Tabela 3.2: Macierz pomyłek - przypadek klasyfikacji binarnej.

Poprzez analizę macierzy pomyłek bez problemu obliczymy łączną liczbę rekordów zaklasyfikowanych poprawnie oraz rekordów przypisanych błędnie przez klasyfikator.

Analizę zawartości macierzy można rozszerzyć o dodatkową informację - koszt błędnej klasyfikacji (ang. misclassification cost).

Definicja 3.24 (Koszt błędnej klasyfikacji [11, Sec 4.8.1]). *Oznaczmy przez e_{ij} koszt błędnego zaklasyfikowania do klasy E_j rekordu, który w rzeczywistości należy do klasy E_i . Koszt poprawnej klasyfikacji oznaczmy przez e_{ii} oraz załóżmy, że $\forall_i e_{ii} = 0$. Dodatkowo niech f_t oznacza liczbę wszystkich przykładów testowych, f_p liczbę poprawnie zaklasyfikowanych rekordów testowych oraz $f_p = \sum_{i=1}^m f_{ii}$, f_b niech natomiast oznacza liczbę błędnych klasyfikacji i $f_b = f_t - f_p$.*

Koszt błędnej klasyfikacji $E(f_b)$ nazywamy sumę:

$$E(f_b) = \sum_{i=1}^m \sum_{j=1}^m f_{ij} \cdot e_{ij}.$$

W przypadkach, gdy błędne zaklasyfikowania rekordów nie różnią się kosztami, do oceny jakości klasyfikatora można wykorzystać miary takie jak trafność klasyfikacji (ang. accuracy) oraz błąd klasyfikacji (ang. error rat).

Definicja 3.25 (Trafność klasyfikacji [11, Sec 4.8.1]). *Trafnością klasyfikacji nazywamy stosunek liczby poprawnie zaklasyfikowanych rekordów testowych do łącznej liczby rekordów testowych:*

$$TR = \frac{f_p}{f_t} = \frac{\sum_{i=1}^m f_{ii}}{f_t}.$$

Definicja 3.26 (Błąd klasyfikacji [11, Sec 4.8.1]). *Błędem klasyfikacji nazywamy stosunek liczby błędnie zaklasyfikowanych rekordów testowych do łącznej liczby rekordów testowych:*

$$BK = \frac{f_b}{f_t} = \frac{\sum_{i=1}^m \sum_{j=1}^m f_{ij}}{f_t} = 1 - TR.$$

Uwaga 3.27. [11, Sec 4.8.1] Innymi miarami, które można wywnioskować bezpośrednio z macierzy pomyłek dla klasyfikacji binarnej (tabla 3.2) są:

- współczynnik TP (czułość):

$$WTP = \frac{TP}{TP + FN},$$

- współczynnik FP :

$$WFP = \frac{FP}{FP + TN},$$

- współczynnik TN (specyficzność):

$$WTN = \frac{TN}{FP + TN},$$

- precyzja:

$$precyzja = \frac{TP}{TP + FP},$$

- zwrot:

$$zwrot = \frac{TP}{TP + FN}.$$

Rozdział 4

Modele tworzenia rekomendacji

W niniejszym rozdziale zajmiemy się formalnym zdefiniowaniem zadania, które ukrywa się pod nazwą tworzenia rekomendacji. Do jego poprawnego określenia będą przydatne następujące pojęcia.

Definicja 4.1 (Przedmiot [6, Sec 1.3]). *Przedmiotem nazwiemy klasę obiektów tego samego typu, nierozróżnialnych dla obserwatora i reprezentowanych przez co najmniej jeden element. W dalszej części pracy zbiór przedmiotów będziemy oznaczać przez P .*

Przedmioty stanowią podstawową grupę elementów w rozważaniach systemach rekomendujących.

Definicja 4.2 (Użytkownik [6, Sec 1.3]). *Użytkownikiem nazywamy osobę zdolną do przedstawienia własnej oceny wybranego przedmiotu. W dalszej części pracy zbiór użytkowników będziemy oznaczać przez U .*

W pracy [6] użyty jest zawsze ten sam zbiór ocen, jednak łatwo możemy pokusić się o jego uogólnioną definicję.

Definicja 4.3 (Zbiór ocen [6, Sec 1.3]). *Podzbiór skończony $\{0, 1, \dots, n\}$, $n \in \mathbb{N}$ nazywamy zbiorem ocen dla przedmiotów. W dalszej części pracy zbiór ocen będziemy oznaczać przez O .*

Definicja 4.4 (Macierz preferencji [6, Sec 1.3]). *Rozważmy zbiór przedmiotów o liczności n oraz grupę użytkowników o liczności m . Macierzą preferencji nazywamy macierz $\mathbb{O} \in \mathbb{M}_{n \times m}(O)$.*

Macierz preferencji \mathbb{O} możemy utożsamiać z tabelą, która jako nazwy kolumn przyjmuje poszczególnych użytkowników, natomiast jako nazwy wierszy przyjmuje przedmioty. Wypełnienie tabeli stanowią oceny, które użytkownicy wystawili przedmiotom. Ideę pojęcia macierzy preferencji przedstawia poniższa tabela:

	Użytkownicy
P	Oceny
r	Oceny
z	Oceny
e	Oceny
d	Oceny
m	Oceny
i	Oceny
o	Oceny
t	Oceny
y	Oceny

zastosowanie tego pojęcia możemy odnaleźć w przykładzie 4.1.

Z uwagi na to, że przedmioty jako wytwory świata rzeczywistego są niemożliwe do opisanego za pomocą skończonej liczby cech rozważa się ich skończoną reprezentację nazywaną wektorem własności.

Definicja 4.5 (Własność [6, Sec 1.3]). *Własnością nazwiemy cechę wyrażoną za pomocą wartości liczbowej lub pewnej zmiennej kategorycznej, która reprezentuje cechę przedmiotu istotną dla użytkownika w procesie tworzenia oceny. Zbiór wszystkich własności w rozważanym modelu oznaczamy W . Dla każdej $w \in W$ poprzez V_w rozumiemy zbiór wszystkich dopuszczalnych wartości własności w .*

Definicja 4.6 (Kontekst [4, Sec 3]). *Kontekstem nazywamy wektor współlistniejących własności, które odzwierciedlają chwilowy stan użytkownika oraz wpływają na wartości jego preferencji.*

Pod pojęciem kontekstu możemy uwzględnić czas, dzień tygodnia, pogodę, wiek użytkownika, stanowisko pracy i wiele innych.

Definicja 4.7 (Funkcja anotująca [6, Sec 1.3]). *Funkcją anotującą własność $w \in W$ nazwiemy funkcję*

$$a_w: P \rightarrow V_w.$$

Mając na uwadze, że zbiór W jest skończony (jak również zbiór P) można utożsamiać funkcję anotującą z wektorem o długości $|W|$ nazywanym wektorem własności.

Problem 4.1 (Problem tworzenia rekomendacji [6, Sec 1.3]). *Rozważmy pewien zbiór przedmiotów P , pewien zbiór użytkowników U oraz pewien zbiór ocen O . Niech ponadto R będzie funkcją taką, że:*

$$R: P \times U \rightarrow O.$$

Załóżmy, że dla funkcji R znane są wartości dla pewnych par przedmiotów i użytkowników. Naszym zadaniem jest zaproponowanie sposobu predykcji brakujących wartości funkcji R w sposób minimalizujący wybrany funkcjonal błędu zdefiniowanego przez kwadrat normy Frobeniusa.

Problem tworzenia rekomendacji możemy utożsamić z wypełnieniem macierzy preferencji \mathbb{O} co zostało dokładniej opisane w sekcji 4.4.

Przyjrzyjmy się następującemu przykładowi, który ilustruje istotę problemu 4.1.

Przykład 4.1. Niech zbiór przedmiotów będzie w tym przypadku zbiorem sześciu książek, $P = \{p_1, p_2, p_3, p_4, p_5, p_6\}$, gdzie p_i dla $i \in \{1, 2, 3, 4, 5, 6\}$ oznacza i -tą książkę. Zbiór użytkowników niech będzie zbiorem czytelników, $U = \{u_1, u_2, u_3, u_4, u_5, u_6\}$, gdzie u_i dla $i \in \{1, 2, 3, 4, 5, 6\}$ oznacza i -tego czytelnika.

Poniższa tabela to macierz preferencji dla ustalonego zbioru P i ustalonego zbioru U . Znak '?' oznacza brakujące wartości funkcji R , zatem czytelnik danej książki nie czytał lub czytał lecz jego ocena jest nieznana.

<i>Czytelnicy</i>		u_1	u_2	u_3	u_4	u_5	u_6
<i>Książki</i>	p_1	6	3	?	6	4	?
	p_2	?	6	6	5	6	?
	p_3	7	7	8	7	8	9
	p_4	8	10	10	7	6	8
	p_5	9	6	6	6	6	?
	p_6	5	7	7	5	4	2

Zadaniem jest przewidzieć brakujące wartości funkcji R , czyli oceny, które mogą zostać nadane przez użytkowników w sposób minimalizujący błąd popełniany przez model.

Uwaga 4.8 (Podział systemów rekomendujących). Podziału systemów rekomendujących dokonujemy ze względu na zakres wykorzystywanych informacji. Wyróżniamy:

- systemy rekomendujące oparte na treści,
- filtrowanie kolaboratywne,
- systemy rekomendujące kontekstowe.

W przypadku systemów rekomendujących opartych na treści predykcja jest dokonywana na podstawie ocen wystawionych przedmiotom przez użytkowników oraz wektorów własności rozważanych przedmiotów.

W filtrowaniu kolaboratywnym wektory cech zostają pominięte, a predykcja dokonywana jest na podstawie ocen. Wyróżniamy dwa typy filtrowania kolaboratywnego:

- *filtrowanie oparte na użytkownikach*

Niech $\mathbf{u}_1 = [o_{1,1}, \dots, o_{1,n}]$ oraz $\mathbf{u}_2 = [o_{2,1}, \dots, o_{2,n}]$ będą wektorami opisującymi odpowiednio użytkownika u_1 i u_2 . Elementami wektorów są wartości funkcji R w przypadkach, gdzie zarówno użytkownik u_1 , jak i u_2 ocenili ten sam przedmiot. Główne założenie filtrowania kolaboratywnego opartego na użytkownikach mówi, że jeżeli odległość między wektorami \mathbf{u}_1 i \mathbf{u}_2 jest mała oraz użytkownik u_1 ocenił pewien przedmiot, dla którego użytkownik u_2 jeszcze nie wystawił oceny, to prawdopodobnie ocena użytkownika u_2 będzie podobna do oceny użytkownika u_1 .

- *filtrowanie oparte na przedmiotach*

Niech $\mathbf{p}_1 = [o_{1,1}, \dots, o_{1,n}]$ oraz $\mathbf{p}_2 = [o_{2,1}, \dots, o_{2,n}]$ będą wektorami opisującymi odpowiednio przedmiot p_1 i p_2 . Elementami wektorów są wartości funkcji R w przypadkach, gdzie przedmiot p_1 , jak i p_2 został oceniony przez tego samego użytkownika. Główne założenie filtrowania kolaboratywnego opartego na przedmiotach mówi, że jeżeli odległość między wektorami \mathbf{p}_1 i \mathbf{p}_2 jest mała oraz użytkownik ocenił w pewien sposób przedmiot p_1 w przeszłości to będzie skłonny w podobny sposób ocenić przedmiot p_2 .

Systemy rekomendujące kontekstowe są natomiast systemami rekomendującymi opartymi na treści w których zostaje uwzględniony dodatkowy wymiar - kontekst.

4.1 Systemy rekomendujące oparte na treści

Rozważania zawarte w tej sekcji stanowią formalizację treści zawartych w książce Gorakala S. K.: *Building Recommendation Engines* [4, Sec 3] oraz opierają się na treści kursu autorstwa Andrew Ng: *Recommender Systems* [8], którego treść została również załączona do pracy na płycie CD.

Systemy oparte na treści wyróżnia ukierunkowanie na spersonalizowany poziom użytkownika oraz treść produktu. Metoda ta opiera się na obliczaniu podobieństw oraz wykorzystuje techniki uczenia maszynowego, takie jak klasyfikacja.

Algorytm 4.9. *W metodzie tej celem jest stworzenie rekomendacji i wygenerowania listy przedmiotów, które mogą być odpowiednie użytkownikowi. Opieramy się na treści rozważanych przedmiotów. Algorytm tego rodzaju rekomendacji możemy przedstawić w następujących krokach:*

1. *stworzenie wektora własności $\mathbf{w} = [w_1, \dots, w_n]$, $n \in \mathbb{N}$, gdzie $\forall_{i \in \{1, \dots, n\}} w_i \in W$,*

2. wygenerowanie profilów produktów - przedstawienie wektorów wartości \mathbf{w}_{p_i} gdzie elementy wektora określają wartości funkcji anotującej poszczególnych własności w_1, \dots, w_n dla przedmiotu p_i dla $i \in \mathbb{N}$,
3. wygenerowanie profilów użytkowników - stworzenie wektorów własności \mathbf{w}_{u_j} przypisanych użytkownikom, gdzie poszczególne elementy wektora określają preferencje użytkownika u_j dla $j \in \mathbb{N}$ w stosunku do elementów wektora własności \mathbf{w} określonego w kroku 1.,
4. obliczmy ocenę $\hat{o}_{j,i}$ jaką użytkownik j wystawiłby dla przedmiotu i , którego wcześniej nie oceniał, za pomocą funkcji

$$\hat{o}_{j,i} = \mathbf{w}_{u_j}^T \mathbf{w}_{p_i},$$

5. porównując otrzymane w kroku 3. oceny, dokonujemy rekomendacji nowego przedmiotu.

Algorytm w ogólny sposób wyjaśnia istotę problemu generowania profili przedmiotów i użytkowników. Istnieje wiele metod, które pozwalają na precyzyjne wyznaczanie wektorów \mathbf{w}_{p_i} , \mathbf{w}_{u_j} . Szczegółowe opisy wybranych metod zostaną zawarte w kolejnych sekcjach rozdziału.

4.1.1 Wygenerowanie profilu przedmiotu - algorytm TFIDF

Rozważania na temat algorytmu TFIDF zawarte w tym rozdziale zostały przeprowadzone na podstawie książki *Recommender Systems Handbook* [15, Sec 3.3.1.1].

W większość systemów rekomendacji opartych na treści używamy gotowych modeli wyszukiwających. W przypadku rozważań przeprowadzanych na dokumentach tekstowych jednym z najbardziej popularnych jest model przestrzeni wektorowej (*ang. Vector Space Model*) z algorytmem TFIDF.

Warto na początku zaznaczyć, że przedmiotem p_i jest tu opisany przez dokument tekstowy (artykuł, książka), natomiast własnościami, które go charakteryzują są słowa. Mając na uwadze te założenia możemy zdefiniować kolejne etapy algorytmu *TFIDF*.

Niech $n \in \mathbb{N}$, $P = \{p_1, p_2, \dots, p_n\}$ będzie zestawem analizowanych przedmiotów. $W = \{w_1, w_2, \dots, w_n\}$, $n \in \mathbb{N}$ niech będzie zbiorem rozważanych własności.

Definicja 4.10 (Model przestrzeni wektorowej [15, Sec 3.3.1.1]). *Modelem przestrzeni wektorowej nazywamy formę reprezentacji przedmiotów, w której przedmiot p_i jest reprezentowany przez wektor z przestrzeni n -wymiarowej, a każdy z n wymiarów reprezentuje jedną z rozważanych własności przedmiotu.*

Definicja 4.11 (Liczność [15, Sec 3.3.1.1]). *Licznością $f_{k,j}$ nazywamy liczbę wystąpień własności w_k w przedmiocie p_j .*

Definicja 4.12 (TF [15, Sec 3.3.1.1]). *TF (ang. term frequency) nazywamy macierz przedstawiającą zależność własności w_k od przedmiotu p_j :*

$$TF(w_k, p_j) = \frac{f_{k,j}}{\max_z f_{z,j}},$$

gdzie:

- $\max_z f_{z,i}$ - maksymalna w odniesieniu do wszystkich wartości $w_z \in W$, $z \in \{1, \dots, n\}$, które pojawiły się w przedmiocie p_i , licznosc wystąpień własności.

Definicja 4.13 (IDF [15, Sec 3.3.1.1]). *IDF (ang. inverse document frequency) nazywamy funkcję:*

$$IDF(w_k) = \log \frac{N}{n_k},$$

gdzie:

- N - całkowita liczba przedmiotów w zbiorze P ,
- n_k - liczba przedmiotów, w których własność w_k , $k \in \{1, \dots, n\}$ wystąpiła przynajmniej raz.

Definicja 4.14 (TFIDF [15, Sec 3.3.1.1]). *TFIDF (ang. TF – term frequency, IDF – inverse document frequency) nazywamy funkcję:*

$$TFIDF(w_k, p_i) = TF(w_k, p_i) \cdot IDF(w_k).$$

Definicja 4.15 (Waga własności w przedmiocie). *Wagą własności w_k w przedmiocie p_i nazywamy wartość:*

$$s_{k,i} = \frac{TFIDF(w_k, p_i)}{\sqrt{\sum_{j=1}^{|W|} TFIDF(w_j, p_i)^2}}.$$

Uwaga 4.16. *Każdy z dokumentów p_i , $i \in \{1, \dots, n\}$ przedstawiamy jako wektor wag własności (słowa) w_k w przedmiocie p_i . Zatem $p_i = [s_{1,i}, s_{2,i}, \dots, s_{n,i}]$.*

4.1.2 Wygenerowanie profilu użytkownika

Rozważania zawarte w tej sekcji zostały oparte na: *Building Recommendation Engines* [4, Sec 3].

W tym kroku tworzymy profil użytkownika pasujący do profilu produktu rozważając własności wspólne z treścią produktu jako, że stwarza to możliwość porównywanie profili użytkowników i przedmiotów.

Algorytm 4.17 ([4, Sec 3]). *Generowanie profilu użytkownika odbywa się w dwu krokach:*

1. *Stworzenie macierzy $\mathbb{A} \in \mathbb{M}_{|P| \times |U|}$, gdzie wyrazy macierzy przyjmują wartości ze zbioru $\{0, 1\}$:*
 - $a_{ij} = 0$, *gdy użytkownik j nie ocenił przedmiotu i ,*
 - $a_{ij} = 1$, *gdy użytkownik j ocenił przedmiotu i .*
2. *Wygenerowanie macierzy profili użytkowników \mathbb{B}*

$$\mathbb{B} = TFIDF(w_k, p_i)^T \cdot \mathbb{A}$$

Przykład 4.2. *Niech wektor własności będzie określony następująco $\mathbf{w} = [w_1, w_2]$, a każdy z elementów wektora \mathbf{w} niech reprezentuje inny gatunek. Dodatkowo zakładamy, że istnieje gatunek w_0 , którego cechy reprezentują wszystkie książki oraz dla każdej z książek $w_0 = 1$. Poniższa tabela określa wartości funkcji anotującej dla poszczególnych książek:*

<i>Gatunki</i>		\mathbf{w}_1	\mathbf{w}_1	\mathbf{w}_2
<i>Książki</i>	\mathbf{p}_1	1	0.8	0.01
	\mathbf{p}_2	1	1	0
	\mathbf{p}_3	1	0.88	0.02
	\mathbf{p}_4	1	0.2	0.9
	\mathbf{p}_5	1	0	1
	\mathbf{p}_6	1	0.7	0.2

W celu wygenerowanie powyższych wartości macierzy możemy się posłużyć algorytmem generującym profil produktu z sekcji 4.1.1. Zakładamy, że macierz TF będzie zawierać informacje o odniesieniu poszczególnych gatunków w_0, w_1, w_2 do książek. Przy obliczaniu wartości funkcji IDF liczba dokumentów N to liczba wszystkich książek w zbiorze P , natomiast n_k to liczba książek, które reprezentują cechy gatunku w_k .

Zatem wektory wartości odpowiadające poszczególnym książkom mają postać

$$\mathbf{w}_{p_1} = [1, 0.8, 0.01], \mathbf{w}_{p_2} = [1, 1, 0], \mathbf{w}_{p_3} = [1, 0.88, 0.02],$$

$$\mathbf{w}_{p_4} = [1, 0.2, 0.9], \mathbf{w}_{p_5} = [1, 0, 1], \mathbf{w}_{p_6} = [1, 0.7, 0.2].$$

Dla każdego użytkownika j wyznaczamy wektor parametrów $\mathbf{w}_{u_j} \in \mathbb{R}^3$, który przedstawia przynależność opinii użytkownika do elementów wektora własności. Preferencje czytelników zostaną więc opisane za pomocą wektorów:

$$\mathbf{w}_{u_1}, \mathbf{w}_{u_2}, \mathbf{w}_{u_3}, \mathbf{w}_{u_4}, \mathbf{w}_{u_5}, \mathbf{w}_{u_6}.$$

Aby wygenerować profile użytkowników posłużmy się algorytmem z sekcji 4.1.2.

Założmy, że macierz \mathbb{A} przyjmuje postać:

<i>Czytelnicy</i>		\mathbf{u}_1	\mathbf{u}_2	\mathbf{u}_3	\mathbf{u}_4	\mathbf{u}_5	\mathbf{u}_6
<i>Książki</i>	\mathbf{p}_1	1	1	1	1	0	1
	\mathbf{p}_2	1	1	1	1	0	1
	\mathbf{p}_3	0	1	0	1	1	1
	\mathbf{p}_4	1	0	0	1	1	0
	\mathbf{p}_5	1	0	0	1	1	0
	\mathbf{p}_6	1	1	1	1	1	1

oraz, że $TFIDF(w_k, p_i)$ jest jak w początkowych założeniach:

<i>Gatunki</i>		\mathbf{w}_1	\mathbf{w}_1	\mathbf{w}_2
<i>Książki</i>	\mathbf{p}_1	1	0.8	0.01
	\mathbf{p}_2	1	1	0
	\mathbf{p}_3	1	0.88	0.02
	\mathbf{p}_4	1	0.2	0.9
	\mathbf{p}_5	1	0	1
	\mathbf{p}_6	1	0.7	0.2

Zatem:

$$\mathbb{B} = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 & 1 \\ 0.8 & 1 & 0.88 & 0.2 & 0 & 0.7 \\ 0.01 & 0 & 0.02 & 0.9 & 1 & 0.2 \end{bmatrix} \cdot \begin{bmatrix} 1 & 1 & 1 & 1 & 0 & 1 \\ 1 & 1 & 1 & 1 & 0 & 1 \\ 0 & 1 & 0 & 1 & 1 & 1 \\ 1 & 0 & 0 & 1 & 1 & 0 \\ 1 & 0 & 0 & 1 & 1 & 0 \\ 1 & 1 & 1 & 1 & 1 & 1 \end{bmatrix}$$

$$\mathbb{B} = \begin{bmatrix} 5 & 4 & 3 & 6 & 4 & 4 \\ 1.1 & 3.4 & 2.5 & 3.6 & 1.8 & 3.4 \\ 2.1 & 0.2 & 0.2 & 2.1 & 2.1 & 0.2 \end{bmatrix}$$

Obliczmy następnie ocenę jaką książce p_3 wystawiłby użytkownik u_1 wiedząc, że wektor preferencji użytkownika u_1 jest postaci $\mathbf{w}_{u_1} = [5, 1.1, 2.1]^T$. Użytkownik ten preferuje więc książki gatunku w_0 , gdy książki gatunków w_1 i w_2 są dla niego mniej atrakcyjne. Zatem:

$$\hat{o}_{1,3} = \mathbf{w}_{u_1}^T \mathbf{w}_{p_3} = [5, 1.1, 2.1] \cdot [1, 0.88, 0.02]^T = 5 \cdot 1 + 1.1 \cdot 0.88 + 2.1 \cdot 0.02 = 6.01.$$

Przewidywaną oceną jest zatem 6.01.

Po przeprowadzeniu podobnych obliczeń dla wszystkich wcześniej nieznanymi ocen możemy zarekomendować naszemu użytkownikowi nową lekturę.

4.2 Filtrowanie kolaboratywne

Rozważania na temat filtrowania kolaboratywnego zostały przeprowadzone na podstawie książki Gorakala S. K.: *Building Recommendation Engines* [4, Sec 3].

Podejście kolaboratywne omija niektóre ograniczenia występujące w metodach opartych na treści. Dzięki temu systemowi możemy dokonywać rekomendacji z pominięciem wektorów preferencji.

4.2.1 Filtrowanie kolaboratywne oparte na użytkowniku

Algorytm 4.18. Stworzenie rekomendacji filtrowania kolaboratywnego opartej na użytkownikach wykonamy w następujących krokach:

1. wybór użytkowników $u_j, u_k \in U$ dla $j, k \in \mathbb{N}$, między którymi chcemy obliczyć podobieństwo,
2. wybór przedmiotów $p_i \in P$ dla $i \in \mathbb{N}$, dla których znane wartości funkcji $R(p_i, u_j)$ i $R(p_i, u_k)$,
3. stworzenie wektorów ocen $o_{j,k}^{(j)}$ i $o_{j,k}^{(k)}$ dla użytkowników u_j i u_k wybranych w kroku 1., których elementy stanowią wartości $R(p_i, u_j)$ oraz $R(p_i, u_k)$, gdzie p_i to przedmioty wybrane w kroku 2.,
4. wyznaczenie odległości między czytelnikami u_j i u_k - najczęstszymi stosowanymi podejściami do obliczania odległości są metryka euklidesowa i współczynnik korelacji Pearsona,
5. wyznaczenie macierzy odległości \mathbb{U}_1 między wszystkimi czytelnikami ze zbioru U ,

6. wyznaczenie macierzy odległości \mathbb{U}_2 między czytelnikami poprzez normalizację danych w celu uzyskania wartości z przedziału $[0, 1]$, wyrazy macierzy przyjmują wartości:

$$u_{ij}^{(2)} = \frac{u_{ij}^{(1)}}{\max_{o_i} \{o_i : o_i \in O, i \in \mathbb{N}\} - \min_{o_i} \{o_i : o_i \in O, i \in \mathbb{N}\}}, i, j \in \mathbb{N},$$

gdzie $u_{ij}^{(1)}$ i $u_{ij}^{(2)}$ są odpowiednio elementami macierzy \mathbb{U}_1 i \mathbb{U}_2 ,

7. wyznaczenie macierzy podobieństwa \mathbb{U}_3 między użytkownikami - zakładając, że największa wartość prawdopodobieństwa to 1, macierz podobieństwa przyjmuje wartości:

$$u_{ij}^{(3)} = 1 - u_{ij}^{(2)}, i, j \in \mathbb{N},$$

gdzie $u_{ij}^{(2)}$ i $u_{ij}^{(3)}$ są odpowiednio elementami macierzy \mathbb{U}_2 i \mathbb{U}_3 ,

8. wyestymowanie nieznanych wartości funkcji R dla $u_j \in U, j \in \mathbb{N}$ oraz $p_i \in P, i \in \mathbb{N}$ - niech u_j będzie konkretnie ustalonym użytkownikiem, w celu obliczenia brakujących wartości funkcji R dla użytkownika u_j obliczmy średnią ważoną wykorzystując oceny i przyjmując wartości podobieństwa między u_j i innymi użytkownikami jako wagi.

W celu dokładniejszego zrozumienia rozważmy ponownie przykład 4.1.

Przykład 4.3. Chcąc obliczyć podobieństwo między użytkownikiem u_2 i u_3 wybierzmy książki, które zostały przeczytane przez obu użytkowników. W tym przypadku są to: p_2, p_3, p_4, p_5, p_6 . Wektorami ocen uwzględniającymi książki ocenione przez obu użytkowników są więc odpowiednio: dla użytkownika u_2 wektor $o_{2,3}^{(2)} = [6, 7, 10, 6, 7]^T$ oraz dla użytkownika u_3 wektor $o_{2,3}^{(3)} = [6, 8, 10, 6, 7]^T$.

Obliczamy odległość euklidesową między użytkownikami u_2 i u_3 :

$$d_e(o_{2,3}^{(2)}, o_{2,3}^{(3)}) = \sqrt{(6-6)^2 + (7-8)^2 + (10-10)^2 + (6-6)^2 + (7-7)^2} = \sqrt{1} = 1.$$

Postępując w podobny sposób dla każdej z par użytkowników otrzymamy następującą macierz odległości \mathbb{U}_1 :

\mathbb{U}_1	\mathbf{u}_1	\mathbf{u}_2	\mathbf{u}_3	\mathbf{u}_4	\mathbf{u}_5	\mathbf{u}_6
\mathbf{u}_1	0	5,099	4,243	3	4,359	3,606
\mathbf{u}_2	5,099	0	1	4,796	5,196	5,745
\mathbf{u}_3	4,243	1	0	3,873	5	5,477
\mathbf{u}_4	3	4,796	3,873	0	2,828	3,742
\mathbf{u}_5	4,359	5,196	5	2,828	0	3
\mathbf{u}_6	3,606	5,745	5,477	3,742	3	0

W procesie normalizacji danych dzielimy elementy macierzy przez

$$(\max\{0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10\} - \min\{0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10\}) = 10$$

i otrzymujemy macierz \mathbb{U}_2 postaci:

\mathbb{U}_2	\mathbf{u}_1	\mathbf{u}_2	\mathbf{u}_3	\mathbf{u}_4	\mathbf{u}_5	\mathbf{u}_6
\mathbf{u}_1	0	0,5099	0,4243	0,3	0,4359	0,3606
\mathbf{u}_2	0,5099	0	0,1	0,4796	0,5196	0,5745
\mathbf{u}_3	0,4243	0,1	0	0,3873	0,5	0,5477
\mathbf{u}_4	0,3	0,4796	0,3873	0	0,2828	0,3742
\mathbf{u}_5	0,4359	0,5196	0,5	0,2828	0	0,3
\mathbf{u}_6	0,3606	0,5745	0,5477	0,3742	0,3	0

Macierz podobieństwa \mathbb{U}_3 przyjmuje więc wartości:

\mathbb{U}_3	\mathbf{u}_1	\mathbf{u}_2	\mathbf{u}_3	\mathbf{u}_4	\mathbf{u}_5	\mathbf{u}_6
\mathbf{u}_1	1	0,4901	0,5757	0,7	0,5641	0,6394
\mathbf{u}_2	0,4901	1	0,9	0,5204	0,4804	0,4255
\mathbf{u}_3	0,5757	0,9	1	0,6127	0,5	0,4523
\mathbf{u}_4	0,7	0,5204	0,6127	1	0,7172	0,6258
\mathbf{u}_5	0,5641	0,4804	0,5	0,7172	1	0,7
\mathbf{u}_6	0,6394	0,4255	0,4523	0,6258	0,7	1

Obliczamy ocenę jaką użytkownik u_1 zaproponuje dla książki p_2 :

$$\frac{0,4901 \cdot 6 + 0,5757 \cdot 6 + 0,7 \cdot 5 + 0,5641 \cdot 6}{0,4901 + 0,5757 + 0,7 + 0,5641} = 5,7$$

Na podstawie metody filtrowania kolaboratywnego opartej na użytkownikach wnioskujemy, że użytkownik u_1 wystawiłby książce p_2 ocenę 5,7.

Postępując w analogiczny sposób przewidzimy wszystkie brakujące oceny :

Czytelnicy		\mathbf{u}_1	\mathbf{u}_2	\mathbf{u}_3	\mathbf{u}_4	\mathbf{u}_5	\mathbf{u}_6
Książki	\mathbf{p}_1	6	3	4,57	6	4	4,88
	\mathbf{p}_2	5,7	6	6	5	6	5,72
	\mathbf{p}_3	7	7	8	7	8	9
	\mathbf{p}_4	8	10	10	7	6	8
	\mathbf{p}_5	9	6	6	6	6	6,67
	\mathbf{p}_6	5	7	7	5	4	2

Możemy więc wnioskować, że w tym przypadku dla użytkownika u_3 książka p_1 prawdopodobnie nie będzie zbyt atrakcyjna. Użytkownik u_6 , natomiast, z chęcią przeczyta książkę p_5 .

4.2.2 Filtrowanie kolaboratywne oparte na przedmiotach

W przypadku filtrowania kolaboratywnego opartego na przedmiotach wartości podobieństwa między użytkownikami zostaje zastąpiona przez wartości podobieństwa między przedmiotami.

Algorytm 4.19. W tym rodzaju rekomendacji należy wykonać następujące kroki:

1. wybór przedmiotów $p_i, p_k \in P$, $i, k \in \mathbb{N}$, dla których znamy wartość funkcji $R(p_i, u_j)$ oraz wartość funkcji $R(p_k, u_j)$, gdzie u_j jest użytkownikiem, który wystawia ocenę w obu przypadkach,
2. stworzenie wektorów ocen $\overline{o(p_i)}$ i $\overline{o(p_k)}$ dla przedmiotów p_i i p_k wybranych w kroku 2., których elementy stanowią wartości funkcji R dla użytkownika u_j ,
3. wyznaczenie odległości między przedmiotami p_i i p_k - najczęstszymi stosowanymi podejściami do obliczania odległości jest podobieństwo kosinusów,
4. wyznaczenie macierzy podobieństwa \mathbb{P} między wszystkimi przedmiotami P ,
5. wyestymowanie nieznanych wartości funkcji R dla $u_j \in U$ dla $j \in \mathbb{N}$ oraz $p_i \in P$ dla $i \in \mathbb{N}$ - niech p_i będzie konkretnie ustalonym przedmiotem oraz u_j będzie konkretnie ustalonym użytkownikiem. W celu obliczenia brakujących wartości funkcji R dla u_j i p_i obliczymy średnią ważoną wykorzystując oceny oraz przyjmując wartości podobieństwa między p_i i innymi przedmiotami ocenionymi przez użytkownika jako wagi.

Przykład 4.4. Aby obliczyć podobieństwo między książkami p_1 i p_2 wyznaczmy wektory ocen w których uwzględnimy przypadki, gdzie jeden użytkownik ocenił obie pozycje.

Zatem: $\overline{o(p_1)} = [3, 6, 4]^T$, $\overline{o(p_2)} = [6, 5, 6]^T$.

Następnie używając wzoru na podobieństwo kosinusowe obliczamy podobieństwo między wybranymi książkami

$$\text{sim}(p_1, p_2) = \frac{\overline{o(p_1)} \cdot \overline{o(p_2)}}{|\overline{o(p_1)}| |\overline{o(p_2)}|} = \frac{3 \cdot 6 + 6 \cdot 5 + 4 \cdot 6}{\sqrt{6^2 + 3^2 + 6^2} \sqrt{6^2 + 5^2 + 6^2}} = 0,6339.$$

Postępując w analogiczny sposób otrzymamy macierz podobieństwa:

\mathbb{P}	$\mathbf{p_1}$	$\mathbf{p_2}$	$\mathbf{p_3}$	$\mathbf{p_4}$	$\mathbf{p_5}$	$\mathbf{p_6}$
$\mathbf{p_1}$	1	0,6339	0,7372	0,7195	0,8935	0,7599
$\mathbf{p_2}$	0,6339	1	0,7951	0,8150	0,7977	0,8898
$\mathbf{p_3}$	0,7372	0,7951	1	0,9780	0,8586	0,9200
$\mathbf{p_4}$	0,7195	0,8150	0,9780	1	0,8860	0,9681
$\mathbf{p_5}$	0,8935	0,7977	0,8586	0,8860	1	0,9413
$\mathbf{p_6}$	0,7599	0,8898	0,9200	0,9681	0,9413	1

Wystymujemy teraz ocenę jaką użytkownik u_6 zaproponuje dla książki p_2 . Ponownie obliczymy średnią ważoną ocen, tym razem, wykorzystując wartość podobieństwa między książką p_1 , a książkami ocenionymi wcześniej przez użytkownika oraz oceny jakie nadał on tym pozycjom:

$$\frac{(0,7951 \cdot 9 + 0,8150 \cdot 8 + 0,8898 \cdot 2)}{(0,7951 + 0,8150 + 0,8898)} = 6,16.$$

Na podstawie przeprowadzonych obliczeń zakładamy, że ocena jaką wystawiłby po przeczytaniu użytkownik u_6 książce p_2 to 6,16.

Powtarzając powyższe obliczenia dla każdej z pozycji wcześniej nieocenionej przez wybranego użytkownika otrzymamy wszystkie brakujące opinie. Następnie bazując na zdobytych danych z łatwością odnajdziemy pozycję najbardziej odpowiednią do zarekomendowania użytkownikowi.

4.2.3 Algorytmy eksploracji danych w filtrowaniu kolaboratywnym

Rozważania w tej sekcji zostały oparte na publikacji *A comprehensive survey of neighborhood-based recommendation methods* [3, Sec 2.1].

Metodami rekomendującymi opartymi na użytkownikach dokonujemy predykcji oceny o_{u_j, p_i} , $j, i \in \mathbb{N}$, użytkownika u_j dla przedmiotu p_i . Uwzględniamy oceny, jakie zostały nadane przedmiotowi p_i przez innych użytkowników oraz podobieństwo między nimi i użytkownikiem u_j . Przewidując oceny dla konkretnego użytkownika, zbiór rozważanych użytkowników możemy zawęzić uwzględniając w obliczeniach najbardziej podobnych rozważanemu. Podejście to stosujemy, gdy wektory opisujące użytkowników są wektorami o dużym wymiarze.

Załóżmy, że dla każdej pary użytkowników u_j i u_k znamy wartość podobieństwa między nimi, którą będziemy tu oznaczać w_{u_j, u_k} . Metodą k - najbliższych sąsiadów wyznaczamy użytkowników wykazujących największe podobieństwo do użytkownika u_j , ich zbiór oznaczmy jako $\mathcal{N}(u_j)$. Zauważmy jednak, że do estymacji oceny o_{u_j, p_i} mogą zostać uwzględnieni tylko użytkownicy, którzy ocenili przedmiot p_i . Dlatego też będziemy rozważać użytkowników najbardziej podobnych u_j , którzy dokonali oceny przedmiotu p_i . Zbiór takich użytkowników oznaczmy jako $\mathcal{N}_{p_i}(u_j)$.

Ocena o_{u_j, p_i} może zostać wyliczona jako średnia ocen użytkowników ze zbioru $\mathcal{N}_{p_i}(u_j)$:

$$o_{u_j, p_i} = \frac{1}{|\mathcal{N}_{p_i}(u_j)|} \sum_{u_i \in \mathcal{N}_{p_i}(u_j)} o_{u_i, p_i}.$$

Podejście to nie uwzględnia jednak wartości podobieństwa występujących między użytkownikami, co jest w tym przypadku bardzo istotne.

Możemy więc wykorzystać średnią ważoną, której zastosowanie zostało pokazane w sekcji 4.2.1.

W konsekwencji, estymacja nowej oceny odbywa się następująco

$$o_{u_j, p_i} = \frac{\sum_{u_i \in \mathcal{N}_{p_i}(u_j)} w_{u_j, u_i} o_{u_i, p_i}}{\sum_{u_i \in \mathcal{N}_{p_i}(u_j)} |w_{u_j, u_i}|}.$$

W przypadku, gdy rozważamy przypadek filtrowania opartego na przedmiotach zakładamy, że zbiór $\mathcal{N}_{u_j}(p_i)$ oznacza zbiór przedmiotów ocenionych przez użytkownika u_j , które są najbardziej podobnymi do przedmiotu p_i . Estymacja oceny użytkownika u_j dla przedmiotu p_i jest uzyskana, podobnie jak w przypadku opartym na użytkownikach, poprzez obliczenie średniej ważonej ocen nadanych przez użytkownika u_j przedmiotom ze zbioru $\mathcal{N}_{u_j}(p_i)$:

$$o_{u_j, p_i} = \frac{\sum_{p_j \in \mathcal{N}_{u_j}(p_i)} w_{p_j, p_i} o_{u_i, p_j}}{\sum_{p_j \in \mathcal{N}_{u_j}(p_i)} |w_{p_j, p_i}|}.$$

Autorzy artykułu: *Scalable Collaborative Filtering Using Cluster-based Smoothing* [18] przedstawiają wykorzystanie algorytmu k-średnich jako etapu wstępnego przetwarzania, który pomaga nam uzyskania informacji na temat klastrów. Odległość użytkownika do środków klastrów służy jako krok wstępnego wyboru sąsiadów.

Innym zastosowaniem opartym na metodzie klastrowania jest zastąpienie brakujących wartości dla użytkowników w klastrze przez wartości środków klastra.

4.3 Systemy rekomendujące kontekstowe

Rozważania w tej sekcji zostały również przeprowadzone na podstawie książki Gorakala S. K.: *Building Recommendation Engines* [4, Sec 3].

Upřednio opisane metody opierały się głównie na rozważaniu problemów dwuwymiarowych. W tym podejściu, przez dodanie nowego wymiaru, jakim jest kontekst (K), zaczynamy rozważać problemy trójwymiarowe:

$$R : U \times P \times K \rightarrow O$$

Algorytm 4.20. W modelu kontekstowym rekomendacje są generowane w następujący sposób:

1. za pomocą algorytmu systemów rekomendujących opartych na treści zostają wygenerowana lista rekomendacji bazująca na preferencjach użytkownika,
2. odfiltrowanie rekomendacji, które odpowiadają przyjętemu kontekstowi - wyróżniamy tutaj dwa podejścia:

- *filtrowanie jako etap wstępny (ang. Pre-Filtering)* - informacje kontekstowe są tu używane do odfiltrowania najbardziej istotnych informacji i skonstruowania dwuwymiarowego zbioru danych,
- *filtrowanie jako etap końcowy (ang. Post-Filtering)* - informacje o kontekście są ignorowane w wejściowych danych, rekomendacja dokonywana jest na całym zbiorze, a w następnym kroku lista rekomendacji stworzona dla użytkownika jest zawężana przez uwzględnienie kontekstu.

4.4 Dekompozycja macierzy ocen - SVD

Poniższe rozważania zostały przeprowadzone na podstawie publikacji Desrosiers K. i Karypis G. *A comprehensive survey of neighborhood-based recommendation methods* [3, Sec 4.1.1].

Idea aproksymacji macierzy preferencji \mathbb{O} w normie Frobeniusa jest aproksymacja macierzy preferencji \mathbb{O} o wymiarach $|P| \times |U|$ i $\text{rz}(\mathbb{O}) = n$ przez macierz $\hat{\mathbb{O}}$ taką, że $\text{rz}(\hat{\mathbb{O}}) = k$, $k < n$.

Zatem

$$\hat{\mathbb{O}} = \mathbb{P}\mathbb{Q}^T,$$

gdzie:

- \mathbb{P} jest macierzą zawierającą koordynaty użytkowników,
- \mathbb{Q} jest macierzą zawierającą koordynaty przedmiotów.

Intuicyjnie, u -ty rząd macierzy \mathbb{P} , $\mathbf{p}_u \in \mathbb{R}^k$, reprezentuje współrzędne użytkownika u rzutowane w k -wymiarowej przestrzeni. Podobnie i -ty wiersz macierzy \mathbb{Q} , $\mathbf{q}_i \in \mathbb{R}^k$, reprezentuje współrzędne przedmiotu i w tej przestrzeni.

Macierze \mathbb{P} i \mathbb{Q} są odnajdowane przez minimalizowanie błędu przybliżenia zdefiniowanego przez kwadrat normy Frobeniusa:

$$E(\mathbb{P}, \mathbb{Q}) = \|\mathbb{O} - \mathbb{P}\mathbb{Q}^T\|_F^2 = \sum_{u,i} (o_{u,i} - \mathbf{p}_u \mathbf{q}_i^T)^2.$$

Podejście to jest równoważne z wyprowadzeniu SVD macierzy \mathbb{O} :

$$\mathbb{O} = \mathbb{U}\Sigma\mathbb{V}^T,$$

gdzie:

- \mathbb{U} jest lewą macierzą wektorów szczególnych,
- \mathbb{V} jest prawą macierzą wektorów szczególnych,

- Σ jest $(m \times n)$ -wymiarową macierzą wartości osobliwych.

Przez $\Sigma_k, \mathbb{U}_k, \mathbb{V}_k$ oznaczmy macierze uzyskane w wyniku wyboru k największych wartości osobliwych oraz ich odpowiednich wektorów.

Macierze \mathbb{P} i \mathbb{Q} odpowiadają więc postaciom:

$$\mathbb{P} = \mathbb{U}_k \sqrt{\Sigma_k},$$

oraz

$$\mathbb{Q} = \mathbb{V}_k \sqrt{\Sigma_k}.$$

Predykcja oceny zostaje dokonana na podstawie równania:

$$o_{u,i} = \mathbf{p}_u \mathbf{q}_i^T.$$

Główny problem jaki pojawia się przy implementacji metody SVD jest brak dużej liczby wyrazów macierzy \mathbb{O} . Rozwiązaniem tego problemu jest odnalezienie brakujących elementów macierzy \mathbb{P} i \mathbb{Q} używając znanych ocen oraz minimalizując równanie:

$$E(\mathbb{P}, \mathbb{Q}) = \sum_{o_{u,i} \in \mathbb{O}} (o_{u,i} - \mathbf{p}_u \mathbf{q}_i^T)^2 + \lambda(\|\mathbf{p}_u\|^2 + \|\mathbf{q}_i\|^2),$$

gdzie λ jest parametrem kontroli poziomu regularyzacji.

To samo podejście może zostać zastosowane w przypadku obliczania podobieństwa między użytkownikami lub przedmiotami w metodzie filtrowania opartego na treści.

Minimalizujemy tu problem postaci:

$$E(\mathbb{P}, \mathbb{Q}) = \sum_{z_{u,i} \in \mathbb{O}} (o_{u,i} - \mathbf{p}_u \mathbf{q}_i^T)^2,$$

gdzie:

- $\forall_{u \in U} \|\mathbf{p}_u\| = 1,$
- $\forall_{i \in P} \|\mathbf{q}_i\| = 1,$
- z_{ui} jest średnią ocen o_{ui} znormalizowaną do zakresu $[-1, 1]$.

4.5 Ocena jakości modeli w regułach rekomendujących

W wielu przypadkach celem reguł rekomendującym nie jest przedstawienie przewidywanych ocen użytkowników, ale przewidywanie podjęcia działania. W przypadku sklepów

internetowych, pierwszym krokiem jest nakłonienie użytkownika do zakupu przedmiotu, a w kolejnym kroku zarekomendowanie produktów, które najbardziej odpowiadają jego preferencjom.

W ocenie takiego systemu wybieramy użytkownika testowego ukrywając kilka jego wyborów. Następnie próbujemy przewidzieć zestaw elementów, które użytkownika zainteresują lub nie. Możliwe wyniki takiego doświadczenia odzwierciedla tabela:

	Zarekomendowane	Niezarekomendowane
Wybrane przez użytkownika	TP	FN
Niewybrane przez użytkownika	FP	TN

Problem, który często pojawia się w ocenie takiego modelu to odpowiednie wykorzystanie zbioru danych w kontekście jego konstrukcji oraz estymacji oceny jego jakości.

Bardzo często używaną i wartą wspomnienia techniką, która pozwala na rozwiązanie tego problemu jest k -krotna krosvalidacja (ang. k -fold cross-validation).

Algorytm 4.21 (k -krotna krosvalidacja[11, Sec 4.8.4]). *Kroki w technice k -krotnej krosvalidacji są następujące:*

1. *Podzielenie zbioru danych D na k równych i rozłącznych podzbiorów D_1, D_2, \dots, D_k .*
2. *Wybór $k - 1$ z podzbiorów i wytrenowanie na nich modelu.*
3. *Testowanie modelu przez użycie k -tego zbioru - pominiętego w poprzednim kroku.*
4. *Powtarzanie kroków 1. i 2. do momentu, gdy każdy z k zbiorów zostanie przyjęty jako zbiór testowy.*
5. *Dokonanie krosvalidacyjnego oszacowania poprawności klasyfikatora - zsumowanie liczby poprawnych klasyfikacji dla wszystkich k klastrów i podzielenie przez licznosc oryginalnego zbioru danych.*

Rozdział 5

Eksperymenty / część praktyczne

Rozważany problem rekomendacji może być sformułowany jako problem uczenia maszynowego, w którym znane są oceny jakie użytkownicy wystawili pewnym przedmiotom, i którego zadaniem jest predykcja ocen użytkowników dla przedmiotów przez nich nieocenionych.

Założmy, że mamy n użytkowników i m przedmiotów. Otrzymujemy $n \times m$ - wymiarową macierz \mathbb{O} , w której wyrazy $o_{i,j}$ są wartościami funkcji $R(u_i, p_j)$ wystawioną przez użytkownika u_i przedmiotowi p_j , $j \in \{1, \dots, m\}$, $i \in \{1, \dots, n\}$. Naszym celem jest wypełnić macierz \mathbb{O} brakującymi ocenami.

5.1 Tworzenie rekomendacji z użyciem ALS i filtrowania kolaboratywnego

5.1.1 Apache Spark

Apache Spark to ciesząca się ostatnio dużą popularnością platforma obliczeniowa stworzona w celu przetwarzania dużych zbiorów danych (BigData). Powstała ona w odpowiedzi na MapReduce wykorzystywaną przez Apache Hadoop. Wspomniany MapReduce przetwarza dane w trybie wsadowym co oznacza, że podczas każdej operacji są one wczytywane i zapisywane na dysku (HDFS), przez co spada znacznie jego wydajność przy algorytmach iteracyjnych. W przypadku Apache Spark i głównej jego idei jaką jest Resilient Distributed Dataset zbiory danych są wczytywane do pamięci i dzięki temu są wykorzystywane przez kolejne kroki algorytmu bez konieczności ponownego wczytywania ich na dysk. Zwiększa to znacznie wydajność i szybkość wykonywania operacji.

Jedną z głównych bibliotek Apache Spark jest biblioteka MLlib. Jest to biblioteka uczenia maszynowego, której celem jest uczynić je łatwym i skalowalnym. MLlib

zapewnia narzędzia do obsługi algorytmów klasyfikacji, regresji, klastrowania, redukcji wymiaru, narzędzia algebry liniowej, statystyki i wiele innych. Biblioteka ta wspiera również narzędzia do obsługi reguł rekomendujących, a w szczególności filtrowania kolaboratywnego.

5.1.2 ALS i MLlib

Alternating Least Square (ALS) jest algorytmem faktoryzacji macierzy, który został zaimplementowany bibliotece uczenia maszynowego MLlib należącej do Apache Spark. Algorytm ten został opracowany z myślą o rozwiązywaniu problemów filtrowania na dużą skalę. Jest prosty, a zarazem dobrze skalowalny w stosunku do dużych zbiorów danych.

Niech macierz $\mathbb{O} = [o_{i,j}]_{m \times n}$, gdzie $o_{i,j}$ oznacza ocenę nadaną przedmiotowi p_j przez użytkownika u_i , m to liczba użytkowników, natomiast n jest liczbą przedmiotów.

Zauważmy, że niektóre wartości macierzy \mathbb{O} są nieznane. W procesie faktoryzacji macierzy \mathbb{O} z liczbą n_f faktorów każdy użytkownik u_i jest utożsamiany z wektorem $\mathbf{u}_i \in \mathbb{R}^{n_f}$ oraz każdy z przedmiotów p_j z wektorem $\mathbf{p}_j \in \mathbb{R}^{n_f}$. Elementy wektora \mathbf{p}_j mierzą stopień w jakim przedmiot p_j posiada czynnik lub cechę, elementy wektora \mathbf{u}_i natomiast powinowactwo użytkownika u_i do każdej z cech lub czynników.

Iloczyn skalarny $\mathbf{u}_i^T \mathbf{p}_j$ przedstawia interakcję między użytkownikiem u_i i przedmiotem p_j przybliżając ocenę użytkownika u_i dla przedmiotu p_j : $o_{i,j} \approx \mathbf{u}_i^T \mathbf{p}_j$.

Oznaczając $\mathbb{U} = [\mathbf{u}_i] \in \mathbb{M}_{n_f, m}(\mathbb{R})$ jako macierz cech użytkowników oraz $\mathbb{P} = [\mathbf{p}_j] \in \mathbb{M}_{n_f, n}(\mathbb{R})$ jako macierz cech przedmiotów staramy się zminimalizować błąd najmniejszych kwadratów postaci [9]:

$$\min_{\mathbb{U}, \mathbb{P}} \sum_{o_{i,j}} (o_{i,j} - \mathbf{u}_i^T \mathbf{p}_j)^2 + \lambda (\sum_{u_i} \|\mathbf{u}_i\|^2 + \sum_{p_j} \|\mathbf{p}_j\|^2),$$

gdzie

- $o_{i,j}$ są znanymi elementami macierzy \mathbb{O} wystawionymi przez użytkowników,
- $\lambda (\sum_{u_i} \|\mathbf{u}_i\|^2 + \sum_{p_j} \|\mathbf{p}_j\|^2)$ jest czynnikiem powszechnie stosowanym w funkcji straty zapobiegającym przeuczeniu modelu .

Powyższa funkcja nie jest funkcją wypukłą (ze względu na obiekt $u_i^T p_j$). Ustalając jednak jedną z macierzy \mathbb{U} lub \mathbb{P} , otrzymujemy postać kwadratową, którą można rozwiązać. Rozwiązanie zmodyfikowanego problemu gwarantuje monotoniczne obniżenie ogólnej funkcji kosztów. Stosując ten krok naprzemiennie do macierzy \mathbb{U} i \mathbb{P} , możemy iteracyjnie poprawiać dopasowanie modelu. Podejście to określamy jako algorytm ALS (Alternating Least Squares).

Algorytm 5.1 (ALS [9]). Zainicjowanie \mathbb{P} z losowymi wartościami.

Powtarzamy:

- Wykonujemy:

$$\mathbf{u}_i = \left(\sum_{o_{i,j} \in o_{i,*}} \mathbf{p}_j \mathbf{p}_j^T + \lambda \mathbb{I}_k \right)^{-1} \sum_{o_{i,j} \in o_{i,*}} o_{i,j} \mathbf{p}_j$$

dla $i = 1, \dots, n$,

- wykonujemy:

$$\mathbf{p}_j = \left(\sum_{o_{i,j} \in o_{*,j}} \mathbf{u}_i \mathbf{u}_i^T + \lambda \mathbb{I}_k \right)^{-1} \sum_{o_{i,j} \in o_{*,j}} o_{i,j} \mathbf{u}_i$$

dla $j = 1, \dots, m$.

do momentu spełnienia kryteriów końcowych.

Otrzymujemy macierze \mathbb{U} i \mathbb{P} .

5.1.3 Implementacja algorytmu

Poniższa implementacja została stworzona na podstawie dokumentacji Apache Spark [19]. Dane użyte w implementacji zostały zaczerpnięte ze strony internetowej groupLens.org [20]. Wszystkie elementy graficzne umieszczone w rozdziale pochodzą ze źródła własnego i stanowią części kodu, którego całość została załączona do pracy w plikach *rekomendacja_podejscie_1.py* i *rekomendacja_podejscie_2.py* (płyta CD).

1. Biblioteki:

Ładowanie bibliotek pyspark i MLlib.

- pyspark - język do przeprowadzania eksploracyjnej analizy danych, budowania algorytmów uczenia maszynowego i tworzenia narzędzi ETL,
- MLlib - jest częścią frameworku Apache Spark. Pozwala aplikować uczenie maszynowe na dużych zbiorach danych bez problemów ze skalowalnością. Dysponuje dużą liczbą algorytmów uczenia maszynowego, które można zastosować w zależności od przypadku biznesowego.

```
%python
from pyspark import *
from pyspark.mllib.recommendation import ALS, Rating
```

2. Dane:

Zbiór danych składa się z następujących plików:

- *ratings.dat* w formacie: UserID::MovieID::Rating::Timestamp,

- *users.dat* w formacie: UserID::Gender::Age::Occupation::Zip-code,
- *movies.dat* w formacie: MovieID::Title::Genres.

Podczas tworzenia algorytmu zostaną wykorzystane dwa z trzech plików *ratings.dat* i *movies.dat*. Ze względu na wygodę przetwarzania i czytelność separator w plikach został zamieniony z "::" na ",". Plik *ratings.dat* składa się kolejno z id_użytkownika w zakresie od 1 do 6040, id_filmu w zakresie od 1 do 3952, oceny wystawionej przez użytkownika dla filmu w pięciostopniowej skali i znacznika czasu. Istotne jest, że plik nie zawiera oceny każdego użytkownika dla każdego filmu, ale każdy z użytkowników wystawił co najmniej dwadzieścia ocen. Plik *movies.dat* zawiera id_filmu w zakresie od 1 do 3952, tytuł filmu oraz gatunek.

3. Przygotowanie danych:

Wczytanie pliku *ratings.dat* oraz sprawdzenie ilości obserwacji.

```
%python
oceny = sc.textFile("/FileStore/tables/ratings.dat")

print ('Plik liczy ' + str(oceny.count()) + ' obserwacji.')
```

► (1) Spark Jobs

Plik liczy 1000209 obserwacji.

Po wczytaniu łańcuch danych zostaje podzielony na wieloelementową listę. Proces ten zostaje przeprowadzony przy użyciu dwóch funkcji:

- `map()` - pobiera jako parametry funkcję oraz listę, a zwraca nową listę, która powstaje w wyniku wywołania funkcji przekazanej w pierwszym parametrze dla każdego elementu listy przekazanej w drugim parametrze,
- `split()` - dzieli łańcuch znaków na wieloelementową listę, jako argument funkcji podaje się separator. W tym przypadku jest to ','. Domyślnym parametrem funkcja `split()` jest biały znak.

Z pliku *ratings.dat* wybieramy tylko te dane, które są interesujące pod względem tworzenia modelu. Pominęto tutaj ostatnie pole, które jest znacznikiem czasu. W procesie rekomendacji nie wnosi ono bowiem żadnej wartości, a jest jedynie zbędną informacją, która w efekcie wydłuża czas przetwarzania. Dodatkowo, w celu wypisania elementów nowej listy użyta zostaje funkcja `collect()`, która zwraca tablicę wszystkich elementów w RDD (ang. Resilient Distributed Dataset). RDD to rdzeń Sparka, zbiór elementów podzielonych na węzły klastra, które mogą być obsługiwane równolegle.

```
%python
oceny2 = oceny.map(lambda l: l.split(',')).map(lambda l: (int(l[0]), int(l[1]), float(l[2])))
for wiersz in oceny2.collect():
    print(wiersz)
```

► (1) Spark Jobs

```
(1, 1193, 5.0)
(1, 661, 3.0)
(1, 914, 3.0)
(1, 3408, 4.0)
(1, 2355, 5.0)
(1, 1197, 3.0)
(1, 1287, 5.0)
(1, 2804, 5.0)
(1, 594, 4.0)
(1, 919, 4.0)
(1, 595, 5.0)
(1, 938, 4.0)
(1, 2398, 4.0)
```

Dane zostają przekształcone w obiekty typu Rating, gdzie jako argumenty podajemy id_użytkownika, id_produktu oraz ocenę: Rating(int user, int product, double rating). Obiekt pochodzi z biblioteki MLlib, a jego składowe zawierają dane postaci użytkownik-produkt-wartość.

```
%python
oceny3 = oceny.map(lambda l: l.split(',')).map(lambda l: Rating(int(l[0]), int(l[1]), float(l[2])))
for wiersz in oceny3.collect():
    print(wiersz)
```

► (1) Spark Jobs

```
Rating(user=1, product=1193, rating=5.0)
Rating(user=1, product=661, rating=3.0)
Rating(user=1, product=914, rating=3.0)
Rating(user=1, product=3408, rating=4.0)
Rating(user=1, product=2355, rating=5.0)
Rating(user=1, product=1197, rating=3.0)
Rating(user=1, product=1287, rating=5.0)
Rating(user=1, product=2804, rating=5.0)
Rating(user=1, product=594, rating=4.0)
Rating(user=1, product=919, rating=4.0)
Rating(user=1, product=595, rating=5.0)
Rating(user=1, product=938, rating=4.0)
```

Wczytanie pliku *movies.dat*. Zostają tu wykorzystane informacje z pierwszej i drugiej kolumny. Pomijamy parametry dotyczące gatunków filmów, gdyż tak jak dane o czasie są one bezwartościowe w modelu filtrowania kolaboratywnego.

```
%python
filmy = sc.textFile("/FileStore/tables/movies.dat")
filmy2 = filmy.map(lambda l: l.split(','))
filmy3 = filmy2.map(lambda l: (int(l[0]), l[1]))
for wiersz in filmy3.collect():
    print(wiersz)
```

► (1) Spark Jobs

```
(1, 'Toy Story (1995)')
(2, 'Jumanji (1995)')
(3, 'Grumpier Old Men (1995)')
(4, 'Waiting to Exhale (1995)')
(5, 'Father of the Bride Part II (1995)')
(6, 'Heat (1995)')
(7, 'Sabrina (1995)')
(8, 'Tom and Huck (1995)')
(9, 'Sudden Death (1995)')
(10, 'GoldenEye (1995)')
(11, 'American President')
(12, 'Dracula: Dead and Loving It (1995)')
(13, 'Balto (1995)')
(14, 'Nixon (1995)')
(15, 'Cutthroat Island (1995)')
(16, 'Casino (1995)')
(17, 'Sense and Sensibility (1995)')
```

Po przygotowaniu danych przechodzimy do budowy algorytmu.

4. Algorytm ALS:

Podzielmy zbiór danych *oceny3* na dwa zbiory : *testowy* i *treningowy*. Zbiór *testowy* będzie stanowił 20% całego zbioru, natomiast zbiór *treningowy* 80%.

```
%python
(treningowy, testowy) = oceny3.randomSplit([0.8, 0.2])
```

Sprawdźmy liczbę rekordów w każdym ze zbiorów.

```
%python

print("Cały zbiór oceny3 liczy " + str(oceny3.count()) + " rekordów, zbiór treningowy to: " +
str(treningowy.count()) + " rekordów, a zbiór testowy: " + str(testowy.count())+" rekordów.")
```

► (3) Spark Jobs

```
Cały zbiór oceny3 liczy 1000209 rekordów, zbiór treningowy to: 799876 rekordów, a zbiór testowy: 200333 r
ekordów.
```

Matematyczny opis algorytmu został przedstawiony we wcześniejszej sekcji pracy. Aby zaimplementować algorytm ALS przyjrzyjmy się metodzie

```
train(ratings, rank, iterations=5, lambda=0.01, blocks=-1, nonnegative=False,
      seed=None)
```

z klasy `pyspark.mllib.recommendation.ALS`. Według dokumentacji metoda przyjmuje następujące parametry:

- `ratings` – RDD ocen lub (`id_użytkownika`, `id_produktu`, `ocena`) krotka,

- rank - liczba wyszukiwanych cech,
- iterations – liczba iteracji algorytmu ALS - domyślnie przyjmuje wartość 5,
- lambda – parametr regulujący algorytmu ALS - domyślnie przyjmuje wartość 0.01,
- blocks – liczba bloków używanych do przeprowadzania równoległych obliczeń - domyślnie przyjmuje wartość -1, co powoduje użycie automatycznie skonfigurowanej liczby bloków,
- nonnegative – określa, czy należy stosować nieujemne ograniczenia - domyślnie wartość false,
- seed – losowy parametr dla zainicjalizowania modelu faktoryzacji macierzy - domyślnie przyjmuje wartość None, co powoduje użycie czasu systemowego jako wartości parametru.

Przeprowadzimy uczenie modelu faktoryzującego macierzy na zbiorze RDD ocen nadanych przez użytkowników określone podzbiorowi produktów. Macierz ocen jest aproksymowana jako iloczyn dwóch macierzy niższego rzędu. Aby rozwiązać rozważany problem algorytm ALS jest uruchamiany iteracyjnie z konfigurowalnym poziomem równoległości.

Algorytm został przeprowadzony z wyborem różnych parametrów aby następnie porównać wyniki i wybrać najlepszy z nich. Poniżej został przedstawiony model dla którego otrzymano najmniejszy błąd średni kwadratowy. Pozostałe modele wykazujące przebieg prób zostały załączone w postaci załącznika do pracy *rekomendacja_podejscie_1.py*.

W próbie tej przyjęto następujące parametry rank = 5 i iterations = 20. Założmy również, że seed = 2019. Zbiorem danych (ratings) jest przygotowany wcześniej zbiór treningowy.

```
%python
rank = 5
iterations = 20
model_2 = ALS.train(treningowy, rank, iterations, seed = 2019)
```

Aby ocenić jakość zaproponowanego model_2 przygotowano zbiór dane_testowe. Zbiór ten stanowić będą wszystkie rekordy ze zbioru *testowego* z pominięciem trzeciej wartości jaką jest ocena.

```
%python
dane_testowe = testowy.map(lambda l: (l[0], l[1]))
for rekord in dane_testowe.collect():
    print(rekord)
```

► (1) Spark Jobs

```
(1, 1193)
(1, 661)
(1, 595)
(1, 1035)
(1, 2687)
(1, 720)
```

Używając modelu_2 dokonano predykcji wcześniej usuniętej wartości. Predykcja zostaje przeprowadzona za pomocą metody predictAll(user,product). Metoda ta zwraca listę przewidywanych ocen dla par użytkowników i produktów podanych jako parametr.

```
%python
predykcja_2 = model_2.predictAll(dane_testowe).map(lambda l: ((l[0], l[1]), l[2]))
```

Dokonajmy teraz zestawienia danych pokazując dla każdej z par (id_użytkownika, id_przedmiotu) ocenę pierwotną i ocenę wygenerowaną przez model_2.

```
%python
oceny3_vs_predykcja_2 = testowy.map(lambda l: ((l[0], l[1]), l[2])).join(predykcja_2)
```

Miarą oceny jakości modelu niech będzie średni błąd kwadratowy (MSE).

```
%python
MSE2 = oceny3_vs_predykcja_2.map(lambda l: (l[1][0] - l[1][1])**2).mean()
print("Średni błąd kwadratowy model_2 (MSE) = " + str(MSE2))
```

► (1) Spark Jobs

```
Średni błąd kwadratowy model_2 (MSE) = 0.7681980311146712
```

W jednej z prób (wszystkie są dostępne w postaci załącznika *rekomendacja_podejscie_2.py*) przyjęto większe parametry: rank = 50 oraz iterations = 100. Warto zauważyć, że większa liczba iteracji jednocześnie zapewnia dokładniejszy wynik. Parametry seed oraz ratings zostają niezmienione. Schemat postępowania zostaje zachowany. Zbiorem treningowym jest jednak teraz cały zbiór *oceny3*.

```
%python
rank = 50
iterations = 100
model_4 = ALS.train(oceny3, rank, iterations, seed = 2019)
```

Dane testowe stanowi teraz cały zbiór *oceny3* z usunięciem trzeciej wartości jaką jest ocena.

```
%python
dane_testowe2 = oceny3.map(lambda l: (l[0], l[1]))
for rekord in dane_testowe.collect():
    print(rekord)
```

► (1) Spark Jobs

```
(1, 1193)
(1, 661)
(1, 595)
(1, 1035)
(1, 2687)
(1, 720)
(1, 2340)
(1, 1566)
(1, 1836)
(1, 1028)
(2, 648)
```

5. Wybór modelu:

Z wszystkich modeli wybieramy model_4, gdzie średni błąd kwadratowy model_4 (MSE)=0.27165466938970756 . Mimo, że jest to nietypowe podejście oceny algorytmu użyjemy go do stworzenia rekomendacji. Wybierzmy zatem użytkownika, któremu chcemy zaproponować nowy film. Niech w tym przypadku będzie to użytkownik o $id_uzytkownika = 3$.

```
%python
id_uzytkownika = 3
```

Poniższa funkcja wypisuje wszystkie oceny wystawione przez wybranego użytkownika dla poszczególnych filmów (na zdjęciu można zobaczyć tylko kilka przykładowych).

```
%python
oceny_uzytkownika = oceny3.filter(lambda l: l[0] == id_uzytkownika)

print("Użytkownik: " + str(id_uzytkownika) + " zaproponował dla książki:")
for i in oceny_uzytkownika.collect():
    id_filmu = i[1]
    tytul = filmy3.filter(lambda l: l[0] == id_filmu)
    for j in tytul.collect():
        print("'" + j[1] + "', ocenę: " + str(i[2]) + ".")
```

► (52) Spark Jobs

```
Użytkownik: 3 zaproponował dla książki:
'Animal House (1978)', ocenę: 4.0.
'Full Monty', ocenę: 2.0.
'Mission: Impossible (1996)', ocenę: 3.0.
'Raising Arizona (1987)', ocenę: 4.0.
'28 Days (2000)', ocenę: 3.0.
```

6. Rekomendacja:

Rekomendacja zostanie dokonana za pomocą metody `recommendProducts(int user,int num)`, gdzie:

- user - id_użytkownika dla którego będzie dokonywana rekomendacja,
- num - liczba rekomendacji, którą chcemy otrzymać. Warto zauważyć, że w niektórych przypadkach liczba elementów może być mniejsza.

Funkcja zwraca obiekty, z których każdy zawiera id_użytkownika, id_produktu i wynik w polu oceny. Każdy z obiektów reprezentuje jeden zalecany produkt, są one sortowane malejąco według pola z przewidzianą oceną. Pierwszy ze zwróconych obiektów jest tym, który najprawdopodobniej będzie najbardziej odpowiadał użytkownikowi. Warto zauważyć, że wynik w polu oceny nie jest z zakresu 1-5. Jest to bowiem nie ocena, a wskaźnik mówiący jak bardzo produkt jest zalecany.

```
%python
print("Najlepsze rekomendacje dla wybranego użytkownika:")
rekomendacje = model_4.recommendProducts(id_użytkownika,3)
print(rekomendacje)
```

► (1) Spark Jobs

```
Najlepsze rekomendacje dla wybranego użytkownika:
[Rating(user=3, product=1627, rating=7.523104223641156), Rating(user=3, product=1649, rating=6.752193042078868), Rating(user=3, product=1046, rating=6.6399033990142176)]
```

Mając oceny przewidziane dla filmów zastąpiono ich id tytułami:

```
%python

print("Najlepszymi rekomendacjami dla użytkownika: " + str(id_użytkownika) + " są filmy:")
for i in rekomendacje:
    id_filmu = i[1]
    tytul = filmy3.filter(lambda l: l[0] == id_filmu)
    for j in tytul.collect():
        print("'" + j[1] + "', przewidywana ocena: " + str(round(i[2],1)))
```

► (3) Spark Jobs

```
Najlepszymi rekomendacjami dla użytkownika: 3 są filmy:
'U Turn (1997)', przewidywana ocena: 7.5
'Fast', przewidywana ocena: 6.8
'Beautiful Thing (1996)', przewidywana ocena: 6.6
```

Pozostało tylko zaproponować użytkownikowi nowy film.

Rozdział 6

Podsumowanie

Rosnącą popularność systemów rekomendujących przypisujemy szeregowi funkcji, które pełnią. Pozwalają one nie tylko na efektywne wyszukanie konkretnych przedmiotów, ale także na tworzenie propozycji dla użytkowników. Interesującymi są fakty, że dwie trzecie z filmów oglądanych na portalu Netflix zostało użytkownikom zarekomendowane przez stronę, a 35% sprzedaży w sklepie Amazon wynika z polecanych produktów. Do znanych portali, wykorzystujących silniki rekomendujące, należy również Google News. Modele rekomendacji zyskały w ostatnim czasie bardzo dużą popularność, zarówno w świecie internetowym, jak i środowiskach naukowych. Ciągłe jednak zauważa się brak źródeł, które opisują zagadnienia rekomendacji w ujęciu algorytmów matematycznych. Celem niniejszej pracy było więc przedstawienie matematycznych modeli reguł rekomendujących, ich charakterystyka oraz wskazanie praktycznych zastosowań.

W rozdziale 3. zaprezentowano wybrane elementy eksploracji danych spotykane w systemach rekomendujących. Na podstawie analizy różnych metod nasunął się wniosek, iż ważnym aspektem wykorzystywanym w rekomendacjach jest faktoryzacja macierzy. Pozwala ona na zmniejszenie wymiarowości problemu do najbardziej znaczących cech. Istotnym wkładem własnym, który został zaprezentowany w tym rozdziale było przeanalizowanie modelu SVD, wyprowadzenie warunków równoważnych temu rozkładowi w postaci lematów, ich udowodnienie oraz analiza twierdzenia Eckard-Younga i uzupełnienie jego dowodu.

W rozdziale 4. zaprezentowano trzy podstawowe metody systemów rekomendujących. Pierwszym były systemy rekomendujące oparte na treści, ukierunkowane na spersonalizowane potrzeby użytkownika oraz treść produktu. Reguły te znajdują zastosowanie w analizie tekstu, czego przykładem był przedstawiony w pracy algorytm TFIDF. Kolejnym rozważanym modelem w rozdziale było filtrowanie kolaboratywne, które pomijało ograniczenia wektorów preferencji, a sam model pozwalał na zwięzłe i intuicyjne wyjaśnianie obliczeń prognostycznych. Kierowaliśmy się tu ocenami użyt-

kowników wyszukując rekomendacji za pomocą dwóch podejść - filtrowania opartego na użytkownikach i opartego na przedmiotach. Trzecim wspomnianym modelem były systemy rekomendujące kontekstowe, gdzie rozważany problem był trójwymiarowy, z dodatkowym elementem jakim był kontekst. Po szczegółowej analizie każdego z systemów rekomendujących, zaproponowaliśmy algorytm systematyzujący metodę rozwiązania problemu, który uwzględnił elementy obliczeń matematycznych kryjące się pod modelem. Przedstawiliśmy również przykłady przedstawiające zastosowania algorytmów.

W rozdziale 5. zaprezentowaliśmy implementację jednego z rozważanych modeli w kontekście uczenia maszynowego. Problem, którego rozwiązanie chcieliśmy uzyskać - na podstawie zbiorów użytkowników, filmów i ocen zaproponować użytkownikowi kolejną pozycję ze zbioru, która będzie dla niego interesująca. Próbą rozwiązania tego problemu było przygotowanie algorytmu filtrowania kolaboratywnego przy zastosowaniu języka Python i technologii Big Data - Spark. Schemat wyprowadzony w rozdziale 4. znacznie ułatwił zrozumienie kolejnych kroków i metod w kodzie. Ostatecznie, efekty uzyskane przy użyciu utworzonego programu zostały uznane za zadowalające. Program pozwala na znalezienie ustalonej liczby pozycji, które są odpowiednie do zaproponowania wybranemu użytkownikowi.

Niniejsza praca nie zawiera oczywiście informacji na temat wszystkich istniejących systemów rekomendujących. Warto uwagi są na przykład systemy oparte na modelach, czy hybrydowe systemy rekomendujące łączące w sobie wiele podejść jednocześnie. Rozważania tutaj zawarte są jednak podstawą do zrozumienia zagadnienia rekomendacji, a sam algorytm przedstawia zarys praktycznego podejścia do analizowanej kwestii.

Bibliografia

- [1] Banaszak G., Gajda W.: *Elementy algebry liniowej część I i II*. Wydawnictwo Naukowe - Techniczne, 2002.
- [2] Bystrov V., COMISEF: <http://comisef.wikidot.com/concept:eigenvalue-and-singular-value-inequalities> (dostęp 07.05.2019).
- [3] Desrosiers K., Karypis G.: A comprehensive survey of neighborhood-based recommendation methods. Springer, 2011.
- [4] Gorakala S. K.: *Building Recommendation Engines*. Packt, 2016.
- [5] Jakubowski J., Sztencel R.: *Wstęp do teorii prawdopodobieństwa*. SCRIPT, 2001.
- [6] Kidziński Ł.: Statistical foundation of recommender systems. Master's thesis, University of Warsaw, 2011.
- [7] Krzyśko M.: *Wykład z teorii prawdopodobieństwa*. Wydawnictwo Naukowe - Techniczne, 2000.
- [8] *Kurs Recommender Systems, Andrew Ng* :
www.youtube.com/watch?v=giIXNoiq0_U&list=PL-6SiIrhTAi6x40q28s7yy94ubLzVXabj
(dostęp 19.05.2019).
- [9] Lublin M. Pere Y. Haoming L., Bangzheng H.: Matrix completion via alternating least square (ALS). Stanford University, 2015.
- [10] Mazeika M.: *The singular value decomposition and low rank approximation*. University of Chicago, 2016.
- [11] Morzy T.: *Eksploracja danych. Metody i algorytmy*. Wydawnictwo Naukowe PWN, 2013.
- [12] Nowak Brzezińska A.: Analiza skupień. konspekt do zajęć: Statystyczne metody analizy danych. Uniwersytet Śląski, 2012.

- [13] O'Brien G.W. Berry M.W., Dumais S.T.: Using linear algebra for intelligent information retrieval. SIAM, 1995.
- [14] Poreda T., Jędrzejewski J.: *Algebra liniowa z elementami geometrii analitycznej*. Politechnika Łódzka, 2011.
- [15] Ricci F. and Rokach L. and Shapira B. and Kantor P.: *Recommender Systems Handbook*. Springer, 2011.
- [16] Walesiak M., Gantar E.: *Statystyczna analiza danych z wykorzystaniem programu R*. Wydawnictwo Naukowe PWN, 2009.
- [17] Paszke W.: *Wartosci osobliwe i norma H_∞* . Uniwersytet Zielonogórski Instytut Sterowania i Systemów Informatycznych.
- [18] Xue G., Lin C., Yang Q, Xi W., Zeng H. Yul Y., Chen Z.: *Scalable Collaborative Filtering Using Cluster-based Smoothing*. ACM SIGIR Conference, Brazil, 2005.
- [19] *Dokumentacja Apache Spark*
Przewodnik biblioteki MLlib:
<https://spark.apache.org/docs/latest/ml-guide.html>
Informacje o filtrowaniu kolaboratywnym:
<https://spark.apache.org/docs/2.1.0>
<https://spark.apache.org/docs/2.2.0>
 (dostęp 24.05.2019).
- [20] *Zestaw danych o filmach, użytkownikach i ocenach ze strony internetowej grouplens.org:*
<https://grouplens.org/datasets/movielens/> (dostęp 25.05.2019).