

# РЕКОМЕНДАЦІЙНА СИСТЕМА ДЛЯ ВИБОРУ ФІЛЬМІВ






Автор: Яковець Михайло Вікторович

Науковий керівник: Дідковська Марина Віталіївна

# Актуальність

- Стрімке зростання кількості даних в мережі Інтернет
- Зростання кількості користувачів
- Переміщення комерційної сфери в Інтернет-простір

Рекомендации на основе ваших просмотров

 <p>Електрочайник PHILIPS HD4646/00 Белый</p> <p>649 грн</p> <p>★★★★★ 152 отзыва</p>	 <p>Мультиварка-скороварка ROTEX REPC75-B</p> <p>1 739 грн</p> <p>★★★★★ 168 отзывов</p>	 <p>Haibike Big Curve 9.10 29 40 Black-Blue (4153024540)</p> <p>11 875 грн</p> <p>10 094 грн</p> <p>★★★★★ 2 отзыва</p>	 <p>Optima F-1 Disc 26 S 2015 Black/Red (SKDCH-OP-26-025-1)</p> <p>Заканчивается</p> <p>5 534 грн</p> <p>★★★★★</p>	 <p>УМБ Xiaomi Mi Power Bank 5000 mAh Silver</p> <p>279 грн</p> <p>★★★★★ 199 отзывов</p>
--	--	--	--	--

## ■ Мета роботи:

1. Дослідження та розробка алгоритмів формування рекомендацій на основі структурованих та неструктурованих даних.
2. Опис принципів роботи рекомендаційної системи для вибору фільмів.
3. Практична реалізація рекомендаційної системи для вибору фільмів на основі розробленого алгоритму.

■ **Об'єкт дослідження:** рекомендаційні системи

■ **Предмет дослідження:** методи та алгоритми формування списку рекомендацій фільму для користувача

# Постановка задачі

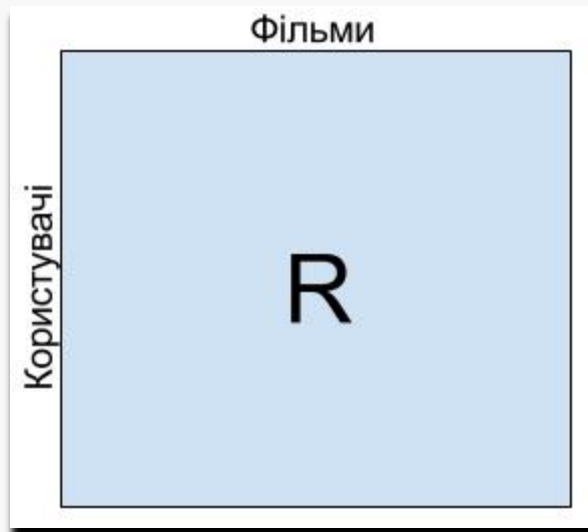
- На основі
  - *вподобань конкретного користувача;*
  - *вподобань інших користувачів;*
  - *інформації про фільм*видати релевантну рекомендацію.
- Перевірити якість рекомендацій.

# Задача формування рекомендацій як практична системна задача

- Принципові невизначеності;
- наявність ризиків;
- суперечності численних цілей;
- наявність елементів різної природи та інформації різного характеру (кількісного і якісного), що циркулює між цими елементами.

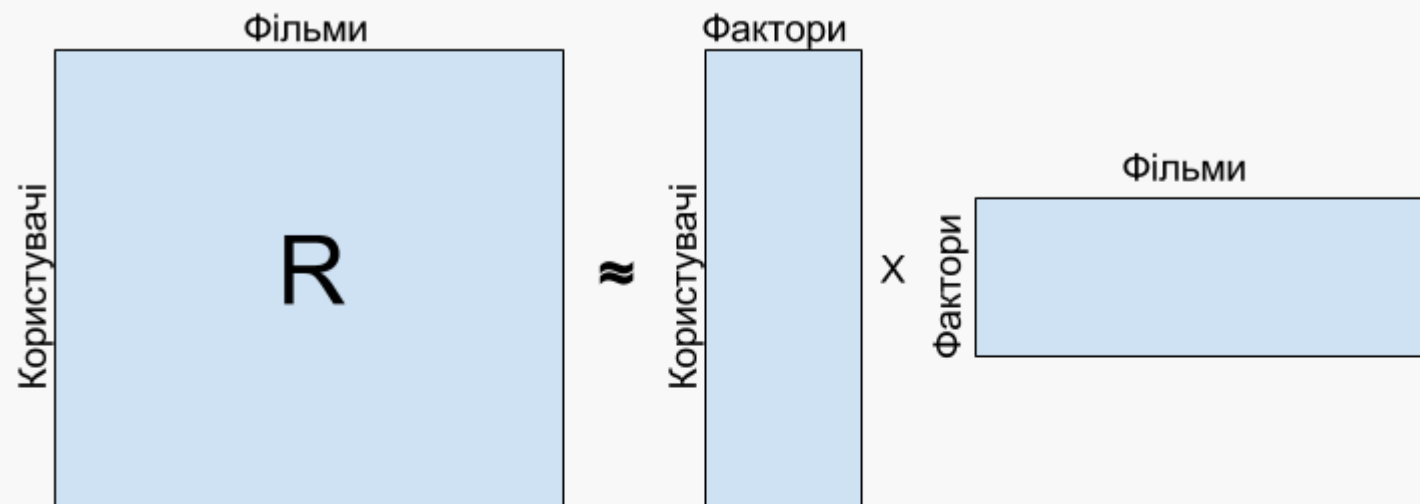
# Декомпозиція предметної області

- $U$  – множина користувачів, які оцінили фільми;
- $I$  – множина фільмів, яка була оцінена;
- $r_{u,i}$  – рейтинг (від 1 до 10), який користувачи поставив фільмуі;



# Існуючі підходи до розв'язку та їхні недоліки

- Фільтрація на основі вмісту (Content filtering)
- Колаборативна фільтрація (Collaborative filtering)
  - *User-based, item-based*
  - *Виявлення прихованих факторів*



# Пошук прихованих факторів

- Задача: за відомими оцінками користувачів передбачити невідомі оцінки.
- Модель представлення оцінок:

$$r_{ui} = \mu + b_i + b_u + \mathbf{q}_i^* \mathbf{p}_u$$

де  $r_{u,i}$  - оцінка з матриці рейтингів  $R$ ,  $\mu$  - середній рейтинг по вибірці,  $b_i$  - базовий рейтинг користувача,  $b_u$  - базовий рейтинг фільму,  $q_i$  - вектор факторів фільму,  $p_u$  - вектор факторів користувача.



# Пошук прихованих факторів

- За відомим значенням оцінок знайдемо такі значення  $\mu$ ,  $b_i$ ,  $b_u$ ,  $q_i$ ,  $p_u$  для усіх фільмів та користувачів, які найкраще наближають значення  $r_{u,i}$  до реальної оцінки, з точки зору квадрату похибки:

$$\underbrace{\sum_{(u,i) \in \mathcal{D}} (\hat{r}_{ui}(\Theta) - r_{ui})^2}_{\text{Якість на навчальній вибірці}} + \underbrace{\lambda \sum_{\theta \in \Theta} \theta^2}_{\text{регуляризація}} \rightarrow \min_{\Theta} \quad (\lambda - \text{параметр регуляризації})$$

- Для пошуку невідомих параметрів використано метод градієнтного спуску.

$$\begin{aligned} b_i &= b_i + \gamma(e_{i,a} - \lambda b_i), \\ b_a &= b_a + \gamma(e_{i,a} - \lambda b_a), \\ q_{a,j} &= q_{a,j} + \gamma(e_{i,a} p_{i,j} - \lambda q_{a,j}) \\ p_{i,j} &= p_{i,j} + \gamma(e_{i,a} p_{i,j} - \lambda p_{i,j}) \end{aligned}$$

# Перевірка якості моделі

- Середнє квадратичне відхилення:

$$RMSE = \sqrt{\frac{1}{|\mathcal{D}|} \sum_{(u,i) \in \mathcal{D}} (\hat{r}_{ui} - r_{ui})^2}$$

- Середня абсолютна похибка

$$MAE = \frac{1}{|\mathcal{T}|} \sum_{(u,i) \in \mathcal{T}} |\hat{r}_{ui} - r_{ui}|$$

# Наївний баєсів класифікатор

- Задача: класифікувати відповіді до фільму: позитивні, негативні, нейтральні.
- Теорема Баєса:

$$p(C|F_1, \dots, F_n) = \frac{p(C) p(F_1, \dots, F_n|C)}{p(F_1, \dots, F_n)}.$$

- Модель наївного баєсового класифікатора:

$$p(C|F_1, \dots, F_n) = \frac{1}{Z} p(C) \prod_{i=1}^n p(F_i|C)$$

Де  $Z$  - масштабний множник, що залежить від  $F_i$  (const);

$$P(C) = \frac{\text{(кількість коментарів з класу } C)}{\text{(кількість усіх коментарів)}};$$

$$P(F_i|C) = \frac{\text{(кількість входжень признаку } F_i \text{ до класу } C)}{\text{(кількість усіх признаков в класі } C)}.$$

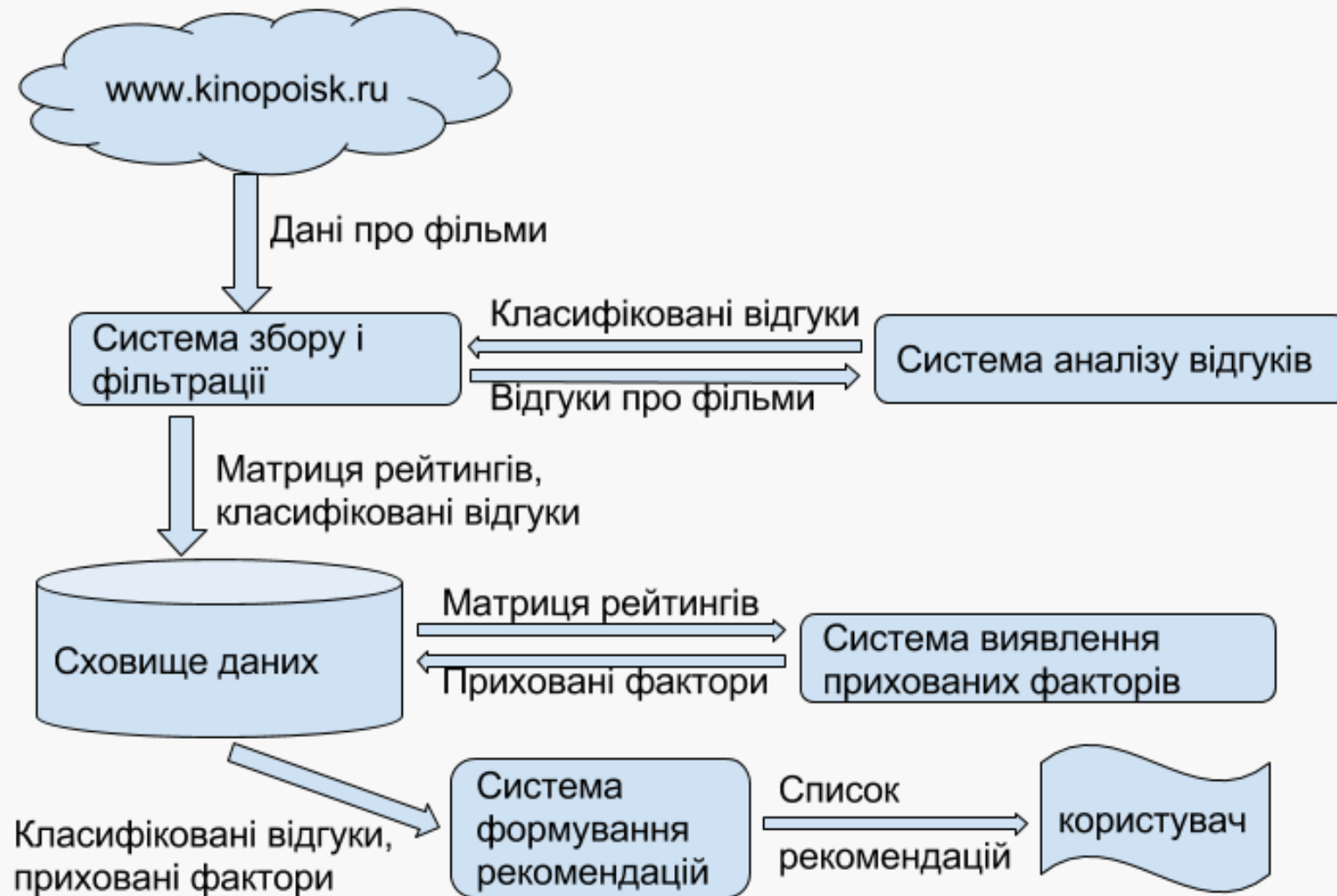
# Перевірка якості моделі

- Для перевірки якості моделі використовується крос-валідація:
  1. Вибірка ділиться  $N$  пособами на навчальну і перевірочну.
  2. Класифікатор навчається на тестовій вибірці і випробовується на перевірочній.
  3. Для кожного експерименту рахується помилка, як відношення вдалих визначень класу, до кількості елементів в тестовій вибірці.
  4. Загальна помилка рахується як середнє арифметичне помилок на кожному з розбиттів.

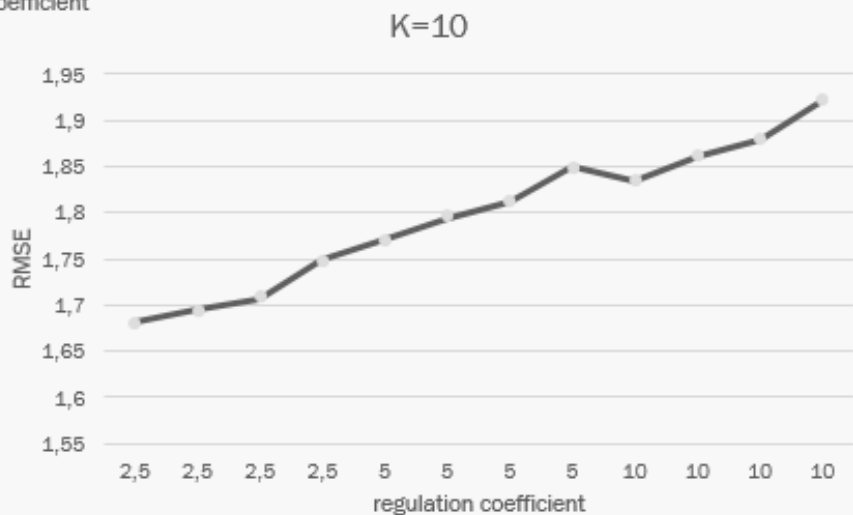
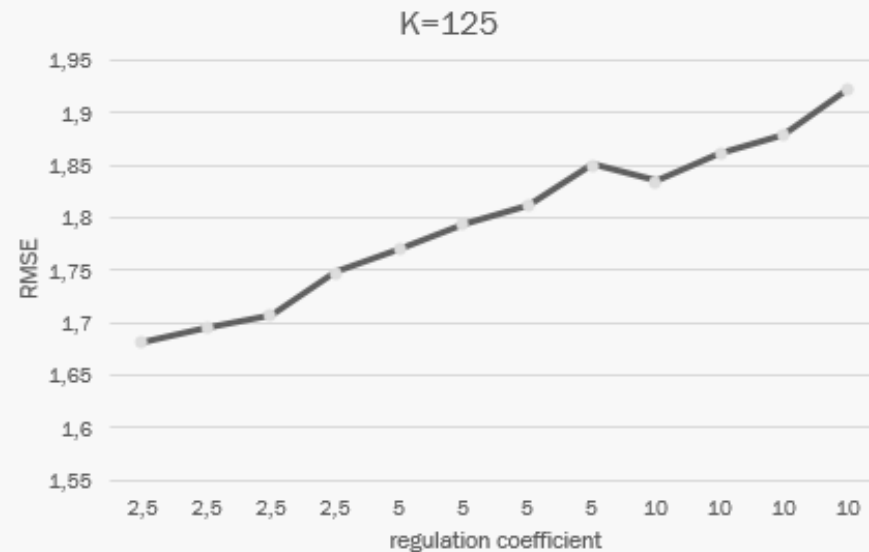
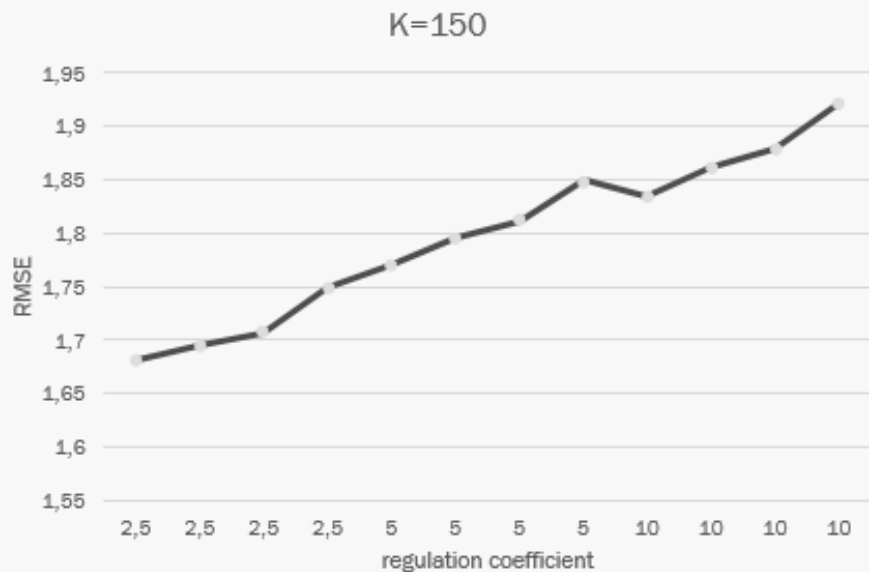
# Двохетапна рекомендаційна система на основі оцінок користувачів і відгуків про фільм

- Система збору і фільтрації даних
- Система виявлення прихованих факторів на основі оцінок користувачів
- Система формування рейтингу фільмів на основі класифікації відгуків про фільми
- Система формування списку рекомендацій

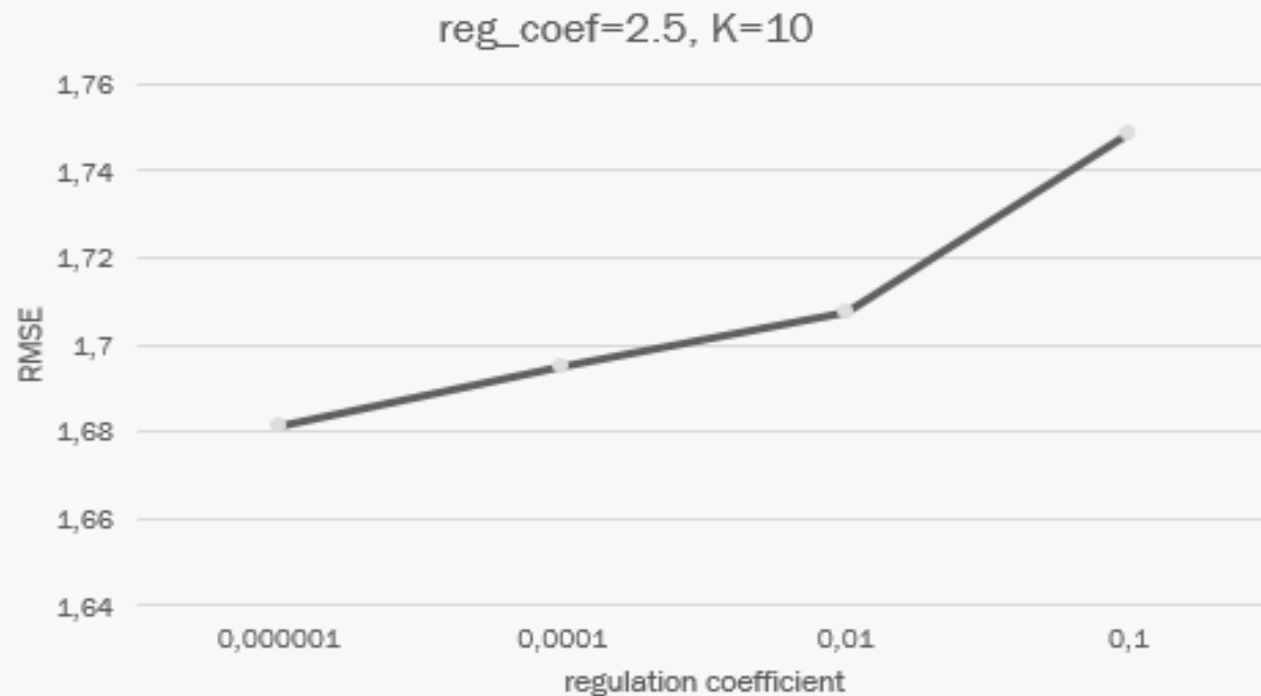
# Схема рекомендаційної системи



# Аналіз практичних результатів



# Аналіз практичних результатів



- Для зменшення похибки RMSE, потрібно зменшувати коефіцієнт регуляризації і збільшення кількості ітерацій навчання.



# Аналіз практичних результатів

Признаки	Точність “Погані – Інші”	Точність “Хороші – Нейтральні”
Уніграми	81.5%	74.6%
Бігами	79.5%	74.8%
Змішані	79.8%	76.6%

- В ролі признаков належності до класу було обрано змішаний тим, тобто множину біграм і уніграм.
- Ієрархічний класифікатор є більш точним при класифікації на початковому рівні (Погані, інші відгуки)

# Демо роботи продукту

**Movie Recommender System**

Load it!

Choose file

Create and save it!

K

Max Users

Max Movies

Movie Title

Movie score

Choose file

*Movie Preferences*

Люди Икс 2 7  
Волчьи дети Амэ и Юки 8  
Матрица: Перезагрузка 9  
Матрица 6  
Паранойд парк 8

*Recommendation*

Лиззи Магуайр Latent:8 Comments:0.9  
Как я теперь люблю Latent:8 Comments:0.8081395  
Вне/себя Latent:7 Comments:0.83203125  
Полкино Latent:7 Comments:0.75

# Висновки

- Розроблено підсистему збору і фільтрації даних про фільми і користувачів
- Запропоновано і розроблено підсистему формування рекомендацій на основі оцінок користувачів та неструктурованих даних про фільми у вигляді відгуків про фільми
- Налаштовано параметри для алгоритму пошуку рекомендацій на основі прихованих факторів, щоб збільшити точність прогнозу за RMSE і MAE
- Налаштовано ієрархічний наївний баєсів класифікатор для аналізу відгуків

# Шляхи подальшого розвитку

- Для формування рейтингу на основі даних про кількість позитивних, негативних, нейтральних відгуків, застосовується середня оцінка фільму, тобто:

$$\frac{(\text{поз.відгуки} + 0.75 * \text{нейтр.відгуки})}{(\text{кількість усіх відгуків})}$$

- Для збереження даних про фільми використовувалась файлова система. Щоб підвищити швидкість роботи з даними, варто перенести систему на реляційну БД.