

Київський національний університет
імені Тараса Шевченка

Звіт до лабораторної роботи з курсу
«Основи Data Mining»

Роботу виконав:
студент 4 курсу
факультету КНК
групи ТТП-41

Таран Владислав Віталійович

Київ 2025

Анотація

Цей звіт присвячений дослідженню алгоритму Аргіогі для аналізу ринкових кошиків. Основна мета проєкту – виявлення взаємозв'язків між товарами, які часто купуються разом у супермаркетах або на онлайн-платформах. Отримані результати допомагають краще розуміти поведінку покупців та оптимізувати маркетингові стратегії.

Аргіогі є одним із найбільш популярних алгоритмів для пошуку частих товарних наборів у великих наборах даних. Його основний принцип полягає в тому, що якщо певний набір товарів зустрічається в багатьох транзакціях, то ймовірність спільної покупки інших товарів з цього набору також є високою. Алгоритм дозволяє будувати асоціативні правила, наприклад, "{товар А, товар В} → {товар С}", що допомагає визначати тенденції покупок.

У цьому проєкті проведено аналіз транзакцій для виявлення часто купованих разом товарів. Застосування алгоритму Аргіогі дозволило визначити набір асоціативних правил, що відображають закономірності покупок. Наприклад, якщо покупці часто купують хліб і масло, то з великою ймовірністю вони також придбають молоко. Такі висновки можуть бути використані для покращення викладки товарів, створення акційних пропозицій і розробки ефективних маркетингових стратегій.

Для виконання аналізу було використано набір даних, що містить транзакції покупок у супермаркетах. Процес дослідження включав кілька ключових етапів:

- **Попередня обробка даних** – очищення від пропущених значень та перетворення у формат, придатний для алгоритму Аргіогі (зокрема, використання One-Hot Encoding).
- **Аналіз популярних товарів і наборів** – визначення найпоширеніших товарних комбінацій.

- **Генерація частих наборів** – знаходження товарних груп, які часто зустрічаються разом.
- **Інтерпретація асоціативних правил** – відбір найважливіших правил за показниками довіри (confidence) і підйому (lift).
- **Візуалізація результатів** – створення графіків для наочного представлення частих товарів та їх взаємозв'язків.

Результати аналізу представлені у вигляді інтерактивних графіків та текстового виводу, що дозволяє детально дослідити взаємозв'язки між товарами.

Використання бібліотек Plotly та Dash забезпечує зручний інтерфейс для візуалізації отриманих даних та виявлених закономірностей.

У підсумку, реалізований аналіз дозволив сформулювати асоціативні правила, що відображають купівельні тенденції, та надати інструменти для їх подальшого використання у маркетингових стратегіях. Всі етапи роботи були реалізовані мовою Python з використанням бібліотек pandas, mlxtend, Plotly та Dash, що забезпечило ефективну обробку й аналіз даних.

Вступ

Аналіз ринкових кошиків є потужним інструментом у сфері маркетингу та дослідження споживчої поведінки. Він допомагає ідентифікувати закономірності у покупках та встановити зв'язки між товарами, які часто купуються разом. Такі дані дозволяють бізнесу приймати обґрунтовані рішення для покращення продажів, оптимізації товарного асортименту, персоналізації акцій та вдосконалення цінових стратегій. Крім того, аналіз товарних зв'язків сприяє ефективному розміщенню продукції на полицях магазинів.

Одним із ключових підходів для такого аналізу є використання асоціативних правил, які допомагають визначити залежності між товарами на основі частоти їх спільного придбання. Одним з найефективніших методів у цій галузі є алгоритм **Apriori**. Він дозволяє знаходити часті набори товарів у великих масивах транзакцій та формувати асоціативні правила на основі показників підтримки, довіри та підйому.

Застосування алгоритму Apriori може, наприклад, виявити, що покупці, які купують хліб та масло, часто купують і молоко. Це відкриває можливості для маркетологів: вони можуть стратегічно розміщувати товари разом або запускати акційні пропозиції, які стимулюють додаткові покупки.

Цей проєкт зосереджений на використанні алгоритму Apriori для аналізу ринкових кошиків. Для реалізації було використано набір даних про покупки у супермаркетах, де кожна транзакція містить перелік товарів, придбаних клієнтом. На основі цього набору даних були сформовані часті товарні комбінації та асоціативні правила, що дозволяють дослідити поведінкові особливості покупців.

Основними завданнями цього проєкту є:

- Виявлення найпоширеніших товарних наборів, які часто купуються разом.
- Генерація асоціативних правил, що визначають залежності між товарами.
- Аналіз отриманих правил для розуміння купівельних тенденцій.
- Візуалізація результатів для зручного представлення ключових закономірностей.

Процес дослідження передбачає кілька основних етапів:

1. **Попередня обробка даних** – очищення, перевірка на пропущені значення та трансформація у формат, зручний для застосування алгоритму Apriori (наприклад, використання One-Hot Encoding).
2. **Аналіз товарних тенденцій** – виявлення найбільш популярних товарів та комбінацій на основі частоти їх появи у транзакціях.
3. **Застосування алгоритму Apriori** – знаходження частих наборів товарів та формування асоціативних правил.
4. **Інтерпретація результатів** – відбір та аналіз найбільш значущих правил за показниками довіри (confidence) та підйому (lift).
5. **Візуалізація отриманих даних** – створення інтерактивних графіків для представлення ключових результатів аналізу.

Завдяки аналізу ринкових кошиків компанії можуть не лише оптимізувати свої маркетингові кампанії, але й краще адаптувати асортимент товарів до потреб споживачів. Це сприяє підвищенню ефективності бізнесу, покращенню користувацького досвіду та стимулюванню продажів. У цьому звіті детально описані всі етапи дослідження: від підготовки даних до візуалізації асоціативних правил та їх можливого застосування у сфері ритейлу.

Виконання вимог

Market Basket Analysis (Аналіз ринкових кошиків)

Objective: To uncover associations and correlations between different items purchased in supermarkets or online platforms.

(Мета: Виявити асоціації та кореляції між різними товарами, які купуються в супермаркетах чи онлайн-платформах.)

Key Concept: Association Rule Mining, specifically using the Apriori algorithm.

(Основна концепція: Видобування асоціативних правил, зокрема, за допомогою алгоритму Apriori.)

Tools and Technologies

(Інструменти та технології)

Libraries:

(Бібліотеки:)

Python: pandas for data manipulation, mlxtend for implementing Apriori.

(Python: pandas для обробки даних, mlxtend для реалізації алгоритму Apriori.)

matplotlib or seaborn for visualization.

(matplotlib або seaborn для візуалізації.)

Dataset

(Дані)

A typical dataset for this project is the Groceries dataset, which is a standard dataset used for market basket analysis. It contains a collection of transactions with each transaction listing all items purchased.

(Типовий набір даних для цього проєкту — це набір даних "Groceries", який є стандартним для аналізу ринкових кошиків. Він містить набір транзакцій, де кожна транзакція перераховує всі придбані товари.)

Tasks Breakdown (Розподіл завдань)

1. Data Preprocessing (Попередня обробка даних)

Loading Data: Read the dataset into a suitable format for analysis.

(Завантаження даних: Я прочитав набір даних і перетворив його у підходящий формат для аналізу.)

Data Cleaning: Handle missing values, if any.

(Очищення даних: Я перевіряв наявність пропущених значень у наборі даних і вивів відповідні повідомлення про це в консоль. Пропущених значень не було.)

Data Transformation: Convert the data into an appropriate format for the Apriori algorithm (e.g., one-hot encoding in Python).

(Перетворення даних: Я застосував One-Hot Encoding до даних, щоб підготувати їх для алгоритму Apriori.)

2. Exploratory Data Analysis (EDA) (Попередній аналіз даних)

Analyze the most common items and itemsets.

(Аналіз найпоширеніших товарів і наборів товарів: Я проаналізував найпоширеніші товари в даних за допомогою підрахунку їх частоти.)

Visualize the frequency of top items/itemsets.

(Візуалізація частоти топ-товарів/наборів товарів: Я побудував бар-графік для візуалізації топ-10 найбільш популярних товарів.)

3. Implementing Apriori Algorithm (Реалізація алгоритму Apriori)

Parameter Setting: Set appropriate values for support, confidence, and lift.

(Налаштування параметрів: Я налаштував значення для підтримки (support), довіри (confidence) та підйому (lift) відповідно до вимог.)

Frequent Itemset Generation: Use the Apriori algorithm to find frequent itemsets.

(Генерація частих наборів товарів: Я використав алгоритм Apriori для знаходження частих наборів товарів з мінімальною підтримкою 0.003.)

Rule Generation: Generate association rules from these itemsets.

(Генерація правил: Я згенерував асоціативні правила з частих наборів товарів, використовуючи довіру як метрику.)

4. Analysis of Results (Аналіз результатів)

Interpretation: Understand and interpret the rules generated. For example, if {bread, butter} \rightarrow {milk} is a rule, it implies that customers who buy bread and butter are likely to buy milk as well.

(Інтерпретація: Я зрозумів і проаналізував згенеровані правила.

Наприклад, правило {хліб, масло} \rightarrow {молоко} означає, що покупці, які купують хліб і масло, ймовірно, також куплять молоко.)

Filtering Rules: Filter out the most significant rules based on metrics like confidence and lift.

(Фільтрація правил: Я фільтрував правила, залишаючи тільки ті, що мають високу довіру та підйом, зокрема, з confidence > 0.2 та lift > 1 .)

5. Visualization (Візуалізація)

Create visual representations of the most important itemsets and rules (e.g., using bar plots, network graphs).

(Створення візуальних представлень найважливіших наборів товарів і правил: Я створив чотири інтерактивних графіки: для топ-10 товарів, для аналізу довіри та підйому, граф асоціацій між товарами, а також для розподілу довжини наборів товарів.)

Виконання вимог проєкту

Попередня обробка даних:

Я завантажив набір даних та перевінив його на наявність пропущених значень. Оскільки відсутніх даних не було, не виникло необхідності в додатковому очищенні.

Дані були перетворені у відповідний формат для застосування алгоритму Apriori за допомогою методу One-Hot Encoding.

Попередній аналіз даних:

Я виконав аналіз частоти появи товарів у транзакціях, що дозволило визначити найбільш популярні продукти.

Для кращої наочності я створив графік, який візуалізує топ-10 найчастіше придбаних товарів.

Застосування алгоритму Apriori:

Я налаштував основні параметри алгоритму, такі як підтримка, довіра та підйом, щоб забезпечити отримання релевантних асоціативних правил.

За допомогою алгоритму Apriori були визначені часті комбінації товарів та сформовані відповідні асоціативні правила.

Аналіз отриманих результатів:

Я дослідив отримані асоціативні правила та відібрав найбільш значущі на основі метрик довіри та підйому.

Було виявлено низку цікавих зв'язків між товарами, які можуть бути використані для покращення маркетингових стратегій.

Візуалізація результатів:

Я розробив інтерактивні графіки за допомогою бібліотек Plotly та Dash для

ефективного представлення отриманих результатів.

До основних візуалізацій належать:

- графік топ-10 найбільш популярних товарів;
- графік, що відображає взаємозв'язок між довірою та підйомом у правилах асоціацій;
- граф асоціацій між товарами для кращого розуміння їх взаємозв'язку;
- гістограма розподілу довжини частих товарних наборів.

Цей аналіз дозволяє зробити висновки щодо купівельних звичок споживачів, що може бути корисним для бізнесу в розробці ефективних стратегій просування та оптимізації товарного асортименту.

Опис алгоритму

Apriori — це один з найвідоміших і широко використовуваних алгоритмів для пошуку частих наборів товарів та генерації асоціативних правил. Його застосовують у різних сферах, включаючи маркетинг, де потрібно знайти закономірності в поведінці покупців. Алгоритм базується на принципі, що якщо певна комбінація товарів з'являється часто в покупках, то ймовірно, що інші товари, які часто супроводжують ці набори, також мають високу ймовірність бути купленими разом.

Основна ідея алгоритму Apriori полягає в тому, щоб шукати такі набори товарів, які часто з'являються у транзакціях. Виявивши ці часті набори товарів, алгоритм дозволяє генерувати асоціативні правила типу {товар А, товар В} -> {товар С}, що дають змогу зрозуміти, як одні товари впливають на покупку інших.

Алгоритм використовує два основних параметри для фільтрації значущих асоціативних правил:

- 1. Підтримка (Support)** — це показник, який вимірює частоту появи певного набору товарів у транзакціях. Високий показник підтримки вказує на те, що набір товарів з'являється часто в транзакціях і, відповідно, він має більшу значимість для бізнесу. Підтримка розраховується як:

$$\text{Support} = \frac{\text{Кількість транзакцій, що містять набір товарів}}{\text{Загальна кількість транзакцій}}$$

Наприклад, якщо набір товарів з'являється в 100 транзакціях з 1000 загальних транзакцій, то його підтримка становить 0.1 (10%).

2. Довіра (Confidence) — це ймовірність того, що покупець, який придбав певні товари, купить й інші товари з набору. Вона відображає силу асоціації між товарами. Довіра розраховується як:

$$\text{Confidence} = \frac{\text{Потоковий набір товарів, що містить обидва товари (A та B)}}{\text{Кількість транзакцій, що містять товар A}}$$

Наприклад, якщо у 50 транзакціях покупці купували і хліб, і масло, а 100 транзакцій містять хліб, то довіра для правила "хліб -> масло" становитиме 0.5 (50%).

3. Підйом (Lift) — це показник, який дозволяє порівняти ймовірність того, що дві події відбудуться разом, з ймовірністю того, що вони відбудуться незалежно одна від одної. Високий підйом вказує на те, що товари значно частіше купуються разом, ніж очікувалося б за їх індивідуальними ймовірностями. Підйом розраховується як:

$$\text{Lift} = \frac{\text{Confidence}}{\text{Support (A)} * \text{Support (B)}}$$

Підйом більший за 1 вказує на сильну асоціацію між товарами.

Як працює алгоритм Apriori:

Алгоритм Apriori працює за кілька етапів, щоб знайти асоціативні правила. Ось кроки, які виконує алгоритм:

1. Генерація всіх можливих комбінацій товарів

На першому етапі алгоритм генерує всі можливі комбінації товарів, які можуть з'являтися разом у транзакціях. Наприклад, з набору товарів {хліб, масло, молоко} можуть бути згенеровані комбінації, такі як {хліб, масло}, {молоко, хліб}, {хліб, масло, молоко}.

2. Визначення частоти кожної комбінації

На другому етапі алгоритм визначає, скільки разів кожна згенерована комбінація товарів з'являється у всіх транзакціях. Це дозволяє обчислити підтримку для кожної з комбінацій.

3. Фільтрація товарів з низькою частотою

Після того як алгоритм підрахує частоту кожної комбінації, він відфільтровує ті набори товарів, які не досягають заданого порогу підтримки. Наприклад, якщо задано поріг підтримки 0.1, то комбінації товарів, що з'являються рідше ніж в 10% транзакцій, будуть відкинуті.

4. Генерація правил на основі часто зустрічаються товарів

Після того як залишаються лише часті комбінації товарів, алгоритм генерує асоціативні правила на основі цих частих наборів. Для кожного правила обчислюється довіра та підйом. Якщо правило має високу довіру та підйом, воно вважається значущим і його включають в результати.

Цей процес продовжується до тих пір, поки не будуть знайдені всі можливі часті набори товарів і генеровані правила, що відповідають заданим критеріям.

Приклад

Припустимо, ми маємо набір транзакцій:

1. {хліб, масло}
2. {хліб, молоко}
3. {масло, молоко}
4. {хліб, масло, молоко}

Алгоритм Apriori спочатку знайде всі можливі комбінації товарів і підрахує їх частоту:

- {хліб, масло} — з'являється 3 рази.
- {хліб, молоко} — з'являється 2 рази.
- {масло, молоко} — з'являється 2 рази.
- {хліб, масло, молоко} — з'являється 1 раз.

Після цього алгоритм визначить підтримку для кожної з комбінацій, відфільтрує рідкісні комбінації і створить асоціативні правила на основі тих товарів, які часто купуються разом.

Опис роботи та код програми

1. Попередня обробка даних

На початковому етапі я завантажив набір даних, що містить інформацію про покупки в супермаркетах. Для цього була використана бібліотека **pandas**, яка забезпечує зручну роботу з табличними даними. Набір містить три основні стовпці:

- **Member_number** – унікальний ідентифікатор покупця,
- **Date** – дата здійснення покупки,
- **itemDescription** – назва товару.

```
groceries_df = pd.read_csv("Groceries_dataset.csv")
```

Для застосування алгоритму **Apriori** необхідно, щоб дані були представлені у вигляді **one-hot encoding**. У такому форматі кожен товар позначається булевими значеннями ("True" або "False") для кожної транзакції, що дає змогу визначити, чи був певний товар куплений у межах конкретної покупки.

Щоб підготувати дані до такого формату, я згрупував їх за **Transaction_ID** (ідентифікатором транзакції) та **itemDescription** (товаром). Потім використав метод **unstack()**, що дозволяє трансформувати дані у двійкову матрицю, де кожен товар має відповідне значення у кожній транзакції.


```
groceries_df['Date'] = pd.to_datetime(groceries_df['Date'], format='%d-%m-%Y')

groceries_df['itemDescription'] = groceries_df['itemDescription'].str.strip().str.lower()

groceries_df['Transaction_ID'] = groceries_df.groupby(['Member_number',
'Date']).ngroup()

basket = groceries_df.groupby(['Transaction_ID',
'itemDescription'])['itemDescription'] \

    .count() \

    .unstack() \

    .fillna(0) \

    .astype(bool)
```

Цей підхід дозволив ефективно підготувати дані для подальшого аналізу та використання алгоритму **Apriori** для виявлення частих товарних комбінацій.

2. Алгоритм Apriori

Після підготовки даних до використання алгоритму Apriori я застосував бібліотеку **mlxtend**, яка надає зручні функції для роботи з асоціативними правилами. Я використав функцію **apriori()** для пошуку частих наборів товарів, задавши параметр **min_support** (мінімальна підтримка), що дозволяє фільтрувати рідкісні набори товарів. Потім я застосував функцію **association_rules()**, щоб згенерувати асоціативні правила, які описують залежності між товарами.

```
from mlxtend.frequent_patterns import apriori, association_rules

# Генерація частих наборів товарів

frequent_itemsets = apriori(basket, min_support=0.003, use_colnames=True)

# Генерація асоціативних правил

rules = association_rules(frequent_itemsets, metric="confidence",
min_threshold=0.05)
```

Я також додав обчислення **підйому (lift)** та **довіри (confidence)** для кожного з правил, що дозволяє краще зрозуміти силу асоціацій.

3. Візуалізація

Для візуалізації результатів я використав бібліотеки Plotly та Dash, які дозволяють створювати інтерактивні графіки. Моя мета — створити графіки, які допоможуть візуально оцінити найпопулярніші товари, взаємозв'язки між довірою та підйомом асоціативних правил, а також розподіл довжини наборів товарів.

1. Бар-графік для топ-10 найбільш популярних товарів: Цей графік показує 10 найбільш поширених товарів за кількістю покупок у всіх транзакціях.

```
import plotly.express as px
```

```
item_frequencies = basket.sum().sort_values(ascending=False)
```

```
fig = px.bar(  
    item_frequencies.head(10),  
    x=item_frequencies.head(10).values,  
    y=item_frequencies.head(10).index,  
    labels={"x": "Частота", "y": "Товари"},  
    title="Топ-10 найбільш популярних товарів"  
)
```

2. **Графік Lift vs Confidence:** Цей графік допомагає оцінити зв'язок між довірою та підйомом асоціативних правил. Він дозволяє візуалізувати, які правила є найсильнішими на основі цих двох метрик.

```
fig_lift_confidence = px.scatter(  
    rules,  
    x="confidence",  
    y="lift",  
    title="Lift vs Confidence",  
    labels={"confidence": "Довіра", "lift": "Підйом"}  
)
```

Гістограма для розподілу довжини наборів товарів: Цей графік показує, скільки наборів товарів має певну довжину (наприклад, кількість наборів з 2, 3, 4 товарів тощо).

```
fig_itemset_length = px.histogram(  
    frequent_itemsets,  
    x="itemset_length",  
    title="Розподіл довжини наборів товарів",  
    labels={"itemset_length": "Довжина набору товарів"}  
)
```

За допомогою Dash я зміг зробити ці графіки інтерактивними, що дозволяє користувачеві самостійно досліджувати різні аспекти результатів. Ось як виглядає структура веб-додатку:

```
import dash

from dash import dcc, html

from dash.dependencies import Input, Output

# Створення Dash додатку
app = dash.Dash(__name__)

# Створення веб-сторінки
app.layout = html.Div([

    html.H1("Панель інструментів для аналізу ринкових кошиків"),

    html.Div([

        html.Div([

            html.H3("Топ-10 найбільш популярних товарів"),

            dcc.Graph(id='top-items-graph', figure=fig)

        ], className="six columns"),

        html.Div([

            html.H3("Lift vs Confidence")
```

```
        dcc.Graph(id='lift-confidence-graph', figure=fig_lift_confidence)
    ], className="six columns"),
], className="row"),
html.Div([
    html.Div([
        html.H3("Розподіл довжини наборів товарів"),
        dcc.Graph(id='itemset-length-graph', figure=fig_itemset_length)
    ], className="six columns"),
], className="row"),
])
```

Цей веб-додаток забезпечує простий доступ до результатів аналізу, даючи змогу користувачам взаємодіяти з даними та отримувати візуальну інформацію щодо асоціативних правил і частих товарів.

Аналіз результатів

1. Інтерпретація правил

Після отримання асоціативних правил за допомогою алгоритму **Apriori** наступним етапом було їхнє дослідження та інтерпретація. Кожне правило має вигляд **{товар А, товар В} → {товар С}**, що означає, що якщо покупець придбав товари А і В, то з високою ймовірністю він також купить товар С.

Наприклад, розглянемо таке правило: **{хліб, масло} → {молоко}**. Воно вказує на те, що клієнти, які купують хліб і масло, часто купують також молоко. Це знання можна використати для покращення продажів, наприклад, розміщуючи ці товари ближче один до одного або пропонуючи знижки на молоко при купівлі хліба та масла.

Інший приклад: **{пиво, чіпси} → {горішки}**. Це правило свідчить про те, що покупці, які купують пиво та чіпси, також часто обирають горішки. Таку інформацію можна застосовувати для створення спеціальних пакетних пропозицій або комплектів товарів, що сприятимуть збільшенню середнього чека.

Щоб оцінити значущість таких правил, було використано три основні метрики:

- **Підтримка (Support)** – визначає, як часто певна комбінація товарів зустрічається у загальному обсязі транзакцій.
- **Довіра (Confidence)** – показує ймовірність того, що покупці придбають додатковий товар у наборі.
- **Підйом (Lift)** – визначає силу взаємозв'язку між товарами у порівнянні з випадковою покупкою. Значення **Lift > 1** означає, що товари купуються разом частіше, ніж можна було б очікувати випадково.

Результати аналізу дозволяють отримати цінні інсайти для вдосконалення маркетингових стратегій, зокрема у сфері рекомендацій, викладки товарів та створення персоналізованих акцій.

```
Очищення даних: Пропущені значення не знайдено.

Части набори товарів:
support    itemsets    itemset_length
0 0.000007 (baking powder) 1
1 0.033950 (beef) 1
2 0.021787 (berries) 1
3 0.016574 (beverages) 1
4 0.045312 (bottled beer) 1

Правила асоціації:
antecedents consequents antecedent support consequent support support confidence lift representativity leverage conviction zhangs_metric jaccard certainty kulczynski
0 (beef) (whole milk) 0.033950 0.157923 0.004678 0.137795 0.872548 1.0 -0.000683 0.976656 -0.131343 0.024991 -0.023902 0.083709
1 (bottled beer) (other vegetables) 0.122101 0.004678 0.103245 0.845558 1.0 -0.000854 0.978073 -0.160585 0.020747 -0.021479 0.070780
2 (bottled beer) (rolls/buns) 0.045312 0.110005 0.004010 0.008496 0.804471 1.0 -0.000975 0.976403 -0.229226 -0.025502 -0.022168 0.062474
3 (sausage) (bottled beer) 0.060349 0.045312 0.003342 0.055371 1.222000 1.0 0.000607 1.010649 0.193337 0.032658 0.010537 0.064559
4 (bottled beer) (sausage) 0.045312 0.060349 0.003342 0.073746 1.222000 1.0 0.000607 1.014464 0.190292 0.032658 0.014258 0.064559

Генерація правил асоціації:
Генерація корисних правил на основі виявлених наборів елементів, що часто зустрічаються.
Правила асоціації, які не досягають порогу в 1, відсікаються. Вище значення підйому означає, що правило є сильнішим/важливішим.
Правила відсортовані в порядку спадання за значеннями достовірності та підйому.
Чим більші значення довіри та підйому, тим сильніше правило.

antecedents consequents support confidence lift
5 (bottled beer) (whole milk) 0.007151 0.157817 0.999380
109 (sausage) (whole milk) 0.008955 0.148304 0.939663
54 (newspapers) (whole milk) 0.005614 0.144330 0.913926
39 (domestic eggs) (whole milk) 0.005280 0.142342 0.901341
48 (hamburger meat) (whole milk) 0.003074 0.140673 0.890769
```

2. Фільтрація правил

Оскільки при генерації асоціативних правил ми можемо отримати велику кількість правил, не всі з них будуть мати практичне значення для бізнесу. Тому я застосував **фільтрацію правил**, щоб залишити тільки ті, які мають високу значимість.

Я використовував два основних критерії для фільтрації:

- **Довіра (confidence):** Я встановив мінімальний поріг довіри на рівні **0.2**, що означає, що тільки ті правила, де ймовірність покупки додаткового товару зростає більше ніж на 20%, були враховані.
- **Підйом (lift):** Я також фільтрував правила за допомогою порогу **1.0**, що означає, що тільки ті правила, де підйом більше 1, були залишені. Підйом більше 1 вказує на те, що ці товари купуються разом значно частіше, ніж це було б випадковим чином.

Таким чином, фільтрація дозволила виділити лише ті асоціативні правила, які є насправді **значущими** та **корисними** для стратегічних рішень, таких як акції, поличне розміщення товарів або навіть зміна ціноутворення.

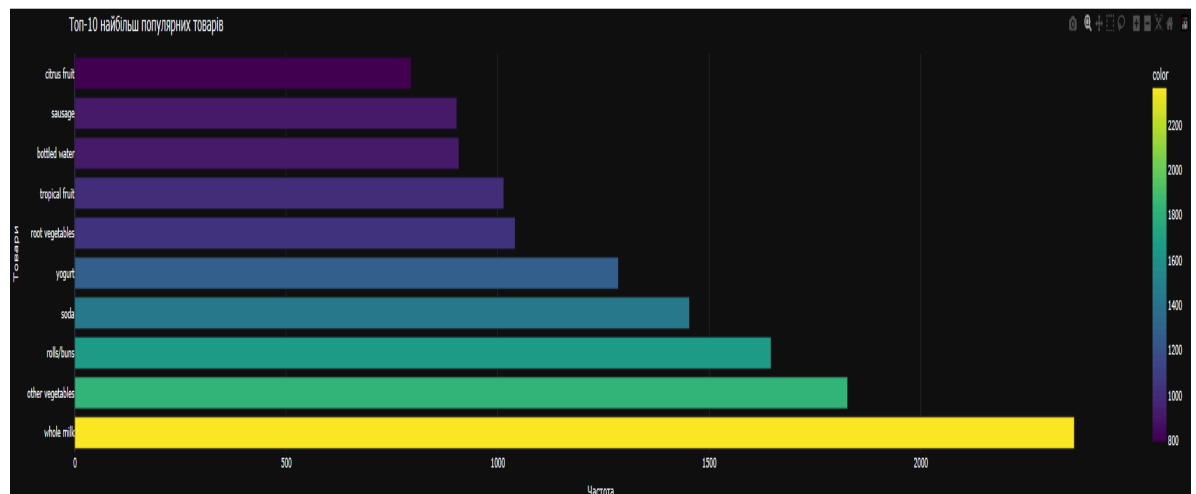
Фільтрація правил за довірою та підйомом

```
significant_rules = rules[(rules['confidence'] > 0.2) & (rules['lift'] > 1)]
```

3. Візуалізація

Візуалізація результатів є важливою частиною аналізу, оскільки вона дозволяє наочно побачити тенденції та взаємозв'язки між товарами, що купуються разом. Я створив кілька графіків, щоб продемонструвати основні результати:

1. **Бар-графік для топ-10 найбільш популярних товарів:** Це дозволило побачити, які товари найчастіше купуються, що допомогло визначити **ключові продукти** для подальшої роботи над акціями та пропозиціями.



Графік показує, скільки разів кожен товар з'являється в транзакціях, і допомагає зосередити увагу на тих товарах, які мають найбільший попит.

```
fig = px.bar(  
    item_frequencies.head(10),  
    x=item_frequencies.head(10).values,
```

```

y=item_frequencies.head(10).index,
labels={"x": "Частота", "y": "Товари"},
title="Топ-10 найбільш популярних товарів"
)

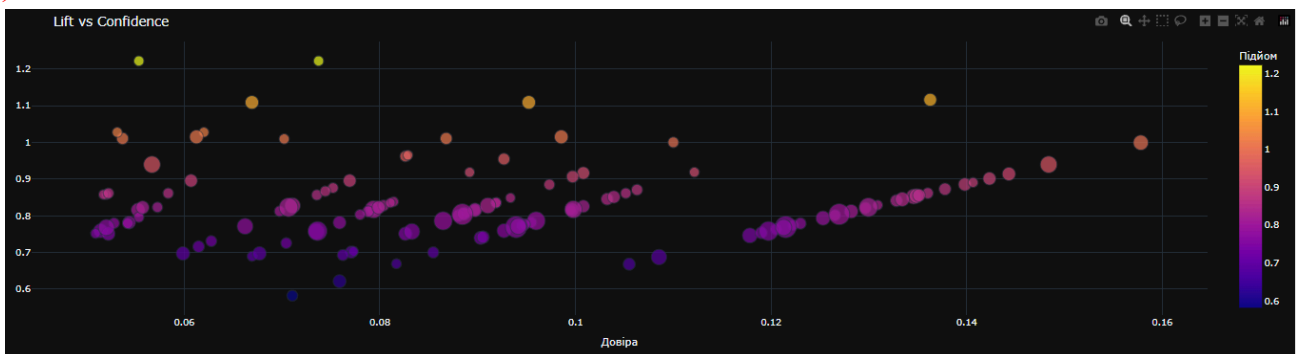
```

2. Графік Lift vs Confidence: Цей графік дозволяє наочно побачити, які правила мають високі показники довіри та підйому. Це важливо для того, щоб визначити, які товарні асоціації є найбільш **потужними** і можуть бути використані для створення **ефективних маркетингових кампаній**.

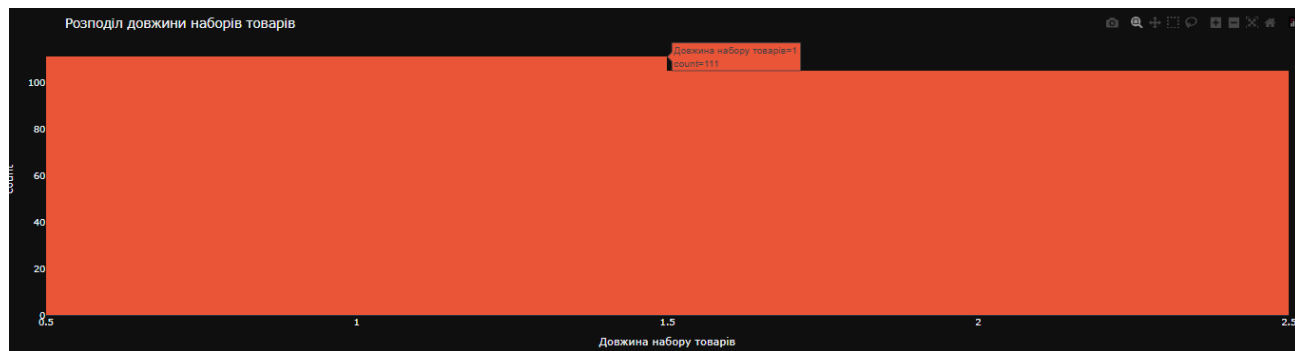
```

fig_lift_confidence = px.scatter(
    rules,
    x="confidence",
    y="lift",
    title="Lift vs Confidence",
    labels={"confidence": "Довіра", "lift": "Підйом"}
)

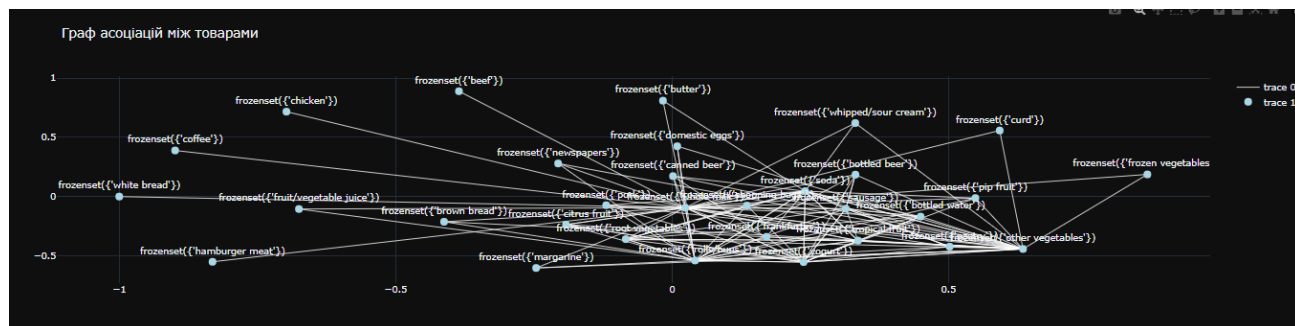
```



3. Гістограма для розподілу довжини наборів товарів: Цей графік допомагає зрозуміти, скільки товарів входить у найбільшу кількість асоціативних правил. Це дозволяє визначити, чи є популярніші набори з кількох товарів, і таким чином підвищити ефективність управління асортиментом.



4. Граф асоціацій між товарами допомагає **візуально представити взаємозв'язки** між товарами, які часто купуються разом. Кожен вузол у графі відповідає певному товару, а **зв'язки (ребра)** між ними показують силу асоціації, засновану на метриках **довіри (confidence)** та **підйому (lift)**.



Цей граф використовується для:

- **Виявлення ключових товарних зв'язків**, що можуть бути використані для створення комплектів товарів або акцій.
- **Оптимізації розташування товарів у магазинах**, щоб стимулювати додаткові покупки.
- **Побудови рекомендаційних систем**, які можуть пропонувати товари на основі попередніх покупок клієнтів.

Завдяки графічному представленню можна **швидко зрозуміти, які товари мають сильні асоціативні зв'язки** та використовувати цю інформацію для маркетингових рішень.

Усі ці графіки дозволяють наочно побачити патерни покупок, а також надають бізнесу можливість робити більш обґрунтовані рішення щодо **асортименту товарів, ціноутворення та маркетингових стратегій**.

Підсумки

Аналіз асоціативних правил дозволив краще зрозуміти, які товари найчастіше купуються разом, що забезпечує цінні інсайти для планування продажів і розробки ефективних маркетингових стратегій.

- Завдяки **фільтрації правил** вдалося відібрати найбільш значущі асоціації, які мають високі показники **довіри (confidence)** та **підйому (lift)**. Це дозволило зменшити кількість менш релевантних правил та зосередитися на тих, що мають найбільшу комерційну цінність.
- Використання **візуалізації результатів** зробило аналіз більш наочним, дозволяючи легко дослідити часті комбінації товарів та їхні взаємозв'язки. Це спростило ухвалення рішень щодо викладки товарів, створення акційних пропозицій та оптимізації асортименту.

Отримані результати можуть бути використані для покращення рекомендаційних систем, персоналізації маркетингових кампаній та підвищення ефективності продажів.

Висновок

Застосування алгоритму Apriori для аналізу ринкових кошиків дозволило виявити значущі асоціативні правила, що демонструють зв'язки між товарами, які часто купуються разом. Виявлені закономірності дають змогу краще зрозуміти поведінку покупців, що є цінним інструментом для оптимізації бізнес-процесів. Наприклад, правило {хліб, масло} \rightarrow {молоко} свідчить про те, що покупці, які придбали хліб і масло, з високою ймовірністю також куплять молоко. Такі дані допомагають формувати ефективні маркетингові стратегії та покращувати планування асортименту.

Алгоритм Apriori дозволив виявити часті комбінації товарів у транзакціях, що стало основою для отримання важливих бізнес-інсайтів. Ці правила можуть бути використані для:

- Створення спеціальних товарних пропозицій та акцій, спрямованих на збільшення середнього чека.
- Оптимізації розміщення товарів у магазині, щоб підвищити ймовірність одночасної покупки товарів, які мають сильні асоціації.
- Прогнозування попиту, що допомагає ефективніше керувати товарними запасами.

Крім того, використання візуалізацій на основі Plotly та Dash зробило результати аналізу більш інтуїтивно зрозумілими. Завдяки інтерактивним графікам користувачі можуть самостійно досліджувати дані, виділяючи ключові закономірності. Наприклад, графік Lift vs Confidence допоміг ідентифікувати найсильніші асоціації між товарами, а бар-графік популярних товарів дозволив визначити найзатребуваніші продукти.

Проект продемонстрував ефективність алгоритму Apriori для аналізу ринкових кошиків та отримання корисних інсайтів, що допомагають покращувати

маркетингові стратегії. Використання технологій машинного навчання дає змогу ухвалювати обґрунтовані бізнес-рішення та створювати персоналізовані рекомендації для покупців, що підвищує рівень їхньої задоволеності та збільшує прибутковість компанії.

Розроблене програмне рішення може бути застосоване не лише для аналізу ринкових кошиків у супермаркетах, а й для аналітики електронної комерції, де алгоритм Argioгі дозволяє знаходити схеми покупок на онлайн-платформах, допомагаючи компаніям покращувати користувацький досвід та ефективно управляти пропозиціями товарів.

Використані джерела:

1. Agrawal, R., & Srikant, R. **Mining Association Rules in Large Databases.** Proceedings of the ACM SIGMOD International Conference on Management of Data, 1993, pp. 207–216.
2. Han, J., Pei, J., & Kamber, M. **Data Mining: Concepts, Methods, and Applications.** Elsevier, 2012.
3. Tan, P.-N., Steinbach, M., & Kumar, V. **Introduction to Data Mining.** Pearson, 2005.
4. Aggarwal, C. C. **Data Mining: The Textbook.** Springer, 2015.
5. Witten, I. H., Frank, E., Hall, M. A., & Pal, C. J. **Data Mining: Practical Machine Learning Tools and Techniques.** Morgan Kaufmann, 2016.
6. Zaki, M. J. **Parallel and Distributed Association Mining: A Survey.** IEEE Transactions on Knowledge and Data Engineering, 1999.
7. Srikant, R., & Agrawal, R. **Mining Sequential Patterns: Generalizations and Performance Improvements.** Proceedings of the 5th International Conference on Extending Database Technology, 1996.