



НАЦІОНАЛЬНИЙ ТЕХНІЧНИЙ УНІВЕРСИТЕТ УКРАЇНИ
«КИЇВСЬКИЙ ПОЛІТЕХНІЧНИЙ ІНСТИТУТ
імені ІГОРЯ СІКОРСЬКОГО»

ДИПЛОМНА РОБОТА НА ТЕМУ:

*Методи кластеризації даних в пакеті Orange
та їх узагальнення*

Виконав:

студент гр. КА-44

Буханенко Андрій Олександрович

Науковий керівник:

професор, д.ф.м.н. Данилов Валерій Якович

Об'єкт дослідження: кластерний аналіз в пакеті Orange.

Предмет дослідження: узагальнення методів кластеризації.

Мета дипломної роботи: розробити порівняльну характеристику методів кластеризації, створити методичне забезпечення для програмного забезпечення Orange.

Програмні засоби: програмне забезпечення Orange.

Результати: розроблено порівняльну характеристику методів кластеризації, створено методичне забезпечення для програмного забезпечення Orange.

Напрямок подальших досліджень: збільшити кількість доступних методів кластеризації в пакеті Orange, вдосконалити на розширити методичне забезпечення для Orange.

Методи кластерного аналізу дозволяють вирішувати такі завдання:

- проведення класифікації об'єктів з урахуванням ознак, що відображають сутність, природу об'єктів. Рішення такого завдання, як правило, призводить до поглиблення знань про сукупність об'єктів, що класифікуються;
- перевірка висунутих припущень про наявність певної структури у досліджуваній сукупності об'єктів, тобто пошук існуючої структури;
побудова нових класифікацій для недостатньо вивчених явищ, коли необхідно встановити наявність зв'язків усередині сукупності і спробувати принести в неї структуру.

Обґрунтуванням актуальності

- вирішення задачі кластеризації є необхідність усунення усереднення при ідентифікації невідомих залежностей, оскільки в цьому випадку прогнозні оцінки є зміщеними, і їх не можна вважати адекватними. Тому ідентифікацію необхідно здійснювати в класі однотипних об'єктів або процесів, передумовою чого є попередня кластеризація.

Міра відстаней

1. Евклідова відстань

$$\rho(x, x') = \sqrt{\sum_i^n (x_i - x'_i)^2}$$

2. Квадрат евклідової відстані

$$\rho(x, x') = \sum_i^n (x_i - x'_i)^2$$

*3. Відстань міських кварталів
(Манхеттенська відстань)*

$$\rho(x, x') = \sum_i^n |x_i - x'_i|$$

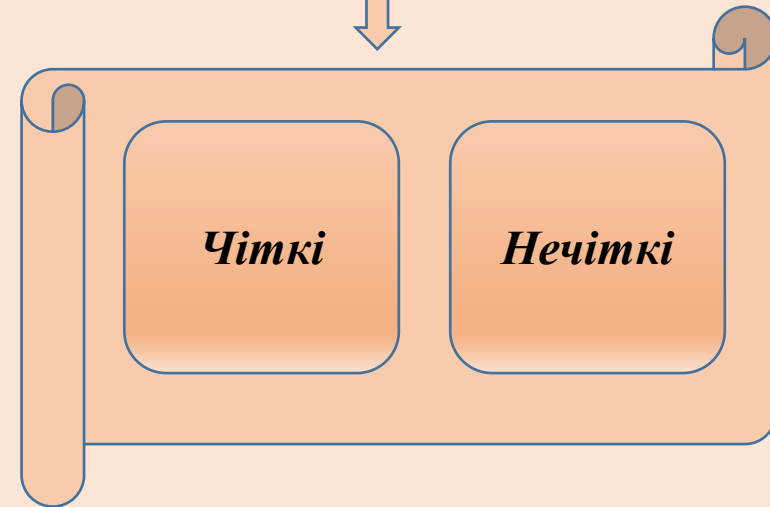
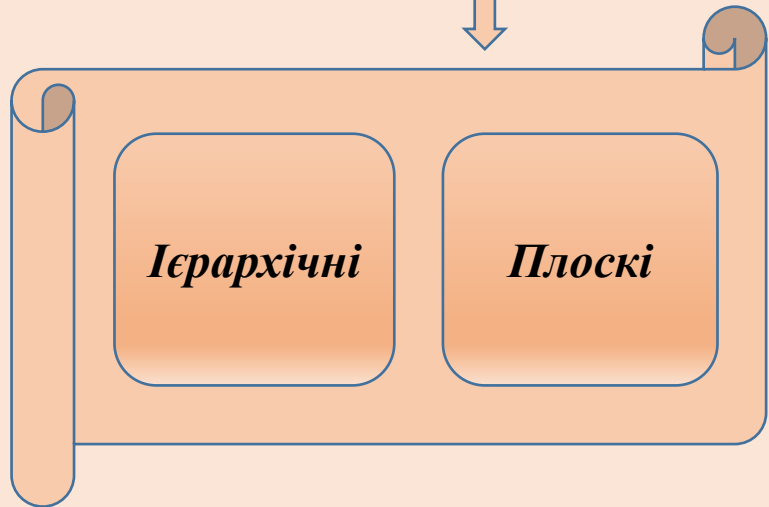
4. Відстань Чебишева

$$\rho(x, x') = \max(|x_i - x'_i|)$$

5. Відстань степені

$$\rho(x, x') = \sqrt[p]{\sum_i^n (x_i - x'_i)^p}$$

Класифікація алгоритмів



Розглянуті методи кластеризації

- 1. Алгоритми ієрархічної кластеризації*
- 2. Алгоритми квадратичної помилки*
- 3. Нечіткі алгоритми*
- 4. Алгоритми, засновані на теорії графів*
- 5. Алгоритм виділення зв'язкових компонент*
- 6. Алгоритм мінімального покривного дерева*
- 7. Пошарова кластеризація*

Обчислювальна складність алгоритмів

Алгоритм	Обчислювальна складність
Ієрархічний O (n ²)	O (n ²)
k-середніх	
c-середніх	O (nkl), де k - число кластерів, l - число ітерацій
Виділення зв'язкових компонент	залежить від алгоритму
Мінімальна покривне дерево	O (n ² log n)
Пошарова кластеризація	O (max (n, m)), де m < n (n-1) / 2

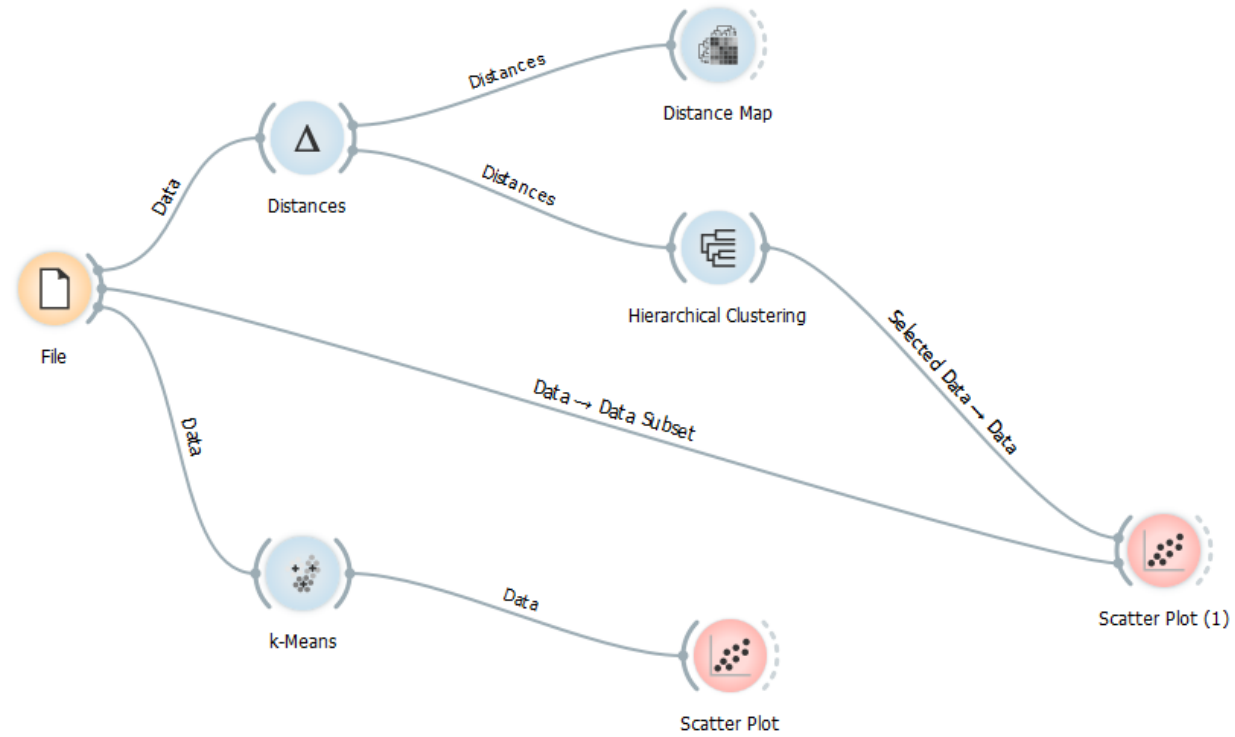
Порівняльна таблиця алгоритмів

Алгоритм	Форма кластерів	Вхідні дані	Результати
Ієрархічний	Довільна	Число кластерів або поріг відстані для усічення ієрархії	Бінарне дерево кластерів
k-середніх	Гіперсфера	Число кластерів	Центри кластерів
c-середніх	Гіперсфера	Число кластерів, ступінь нечіткості	Центри кластерів, матриця приналежності
Виділення зв'язкових	Зв'язкових компонент	Поріг відстані R	Деревовидна структура кластерів
Мінімальна покриває дерево	Довільна	Число кластерів або поріг відстані для видалення ребер	Деревовидна структура кластерів
Пошарова кластеризація	Довільна	Послідовність порогів відстані	відстані Деревовидна структура кластерів з різними рівнями ієрархії

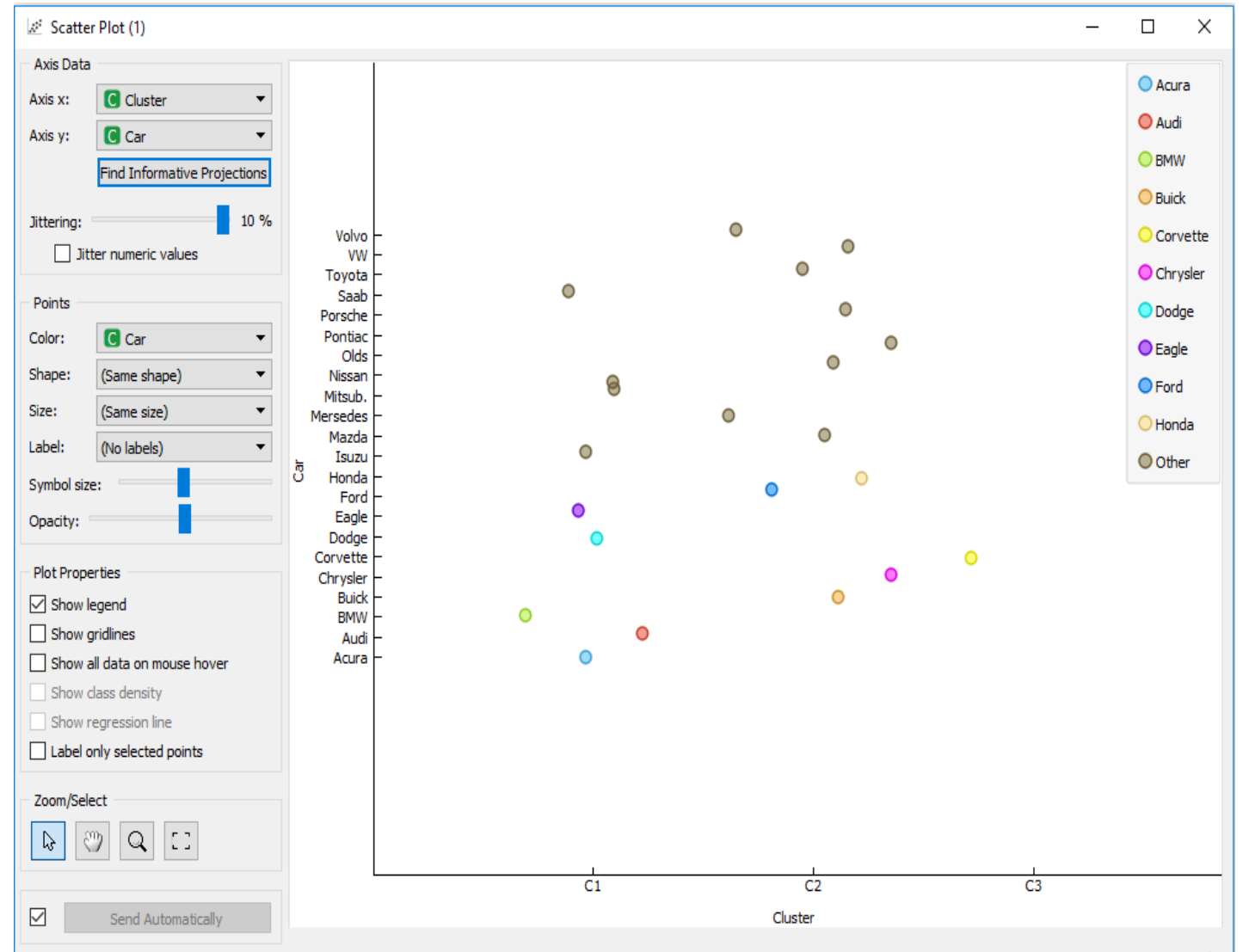
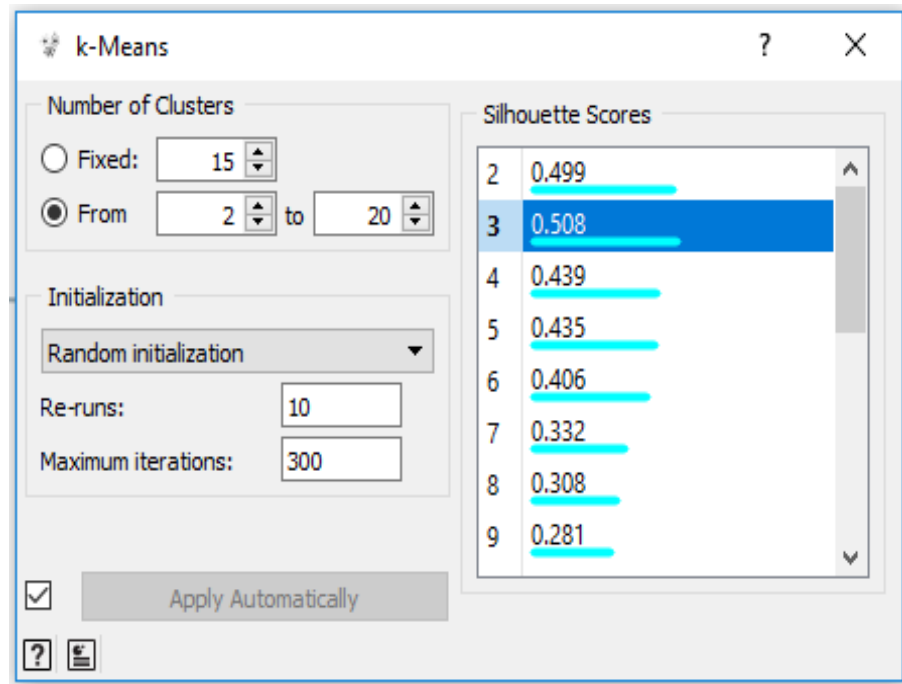
Panel of widgets categorized into Data, Visualize, Model, Evaluate, and Unsupervised. The Evaluate section is highlighted in light blue.

- Data**
- Visualize**
- Model**
- Evaluate**
 - Test & Score
 - Predictions
 - Confusion Matrix
 - ROC Analysis
 - Lift Curve
 - Calibration Plot
- Unsupervised**

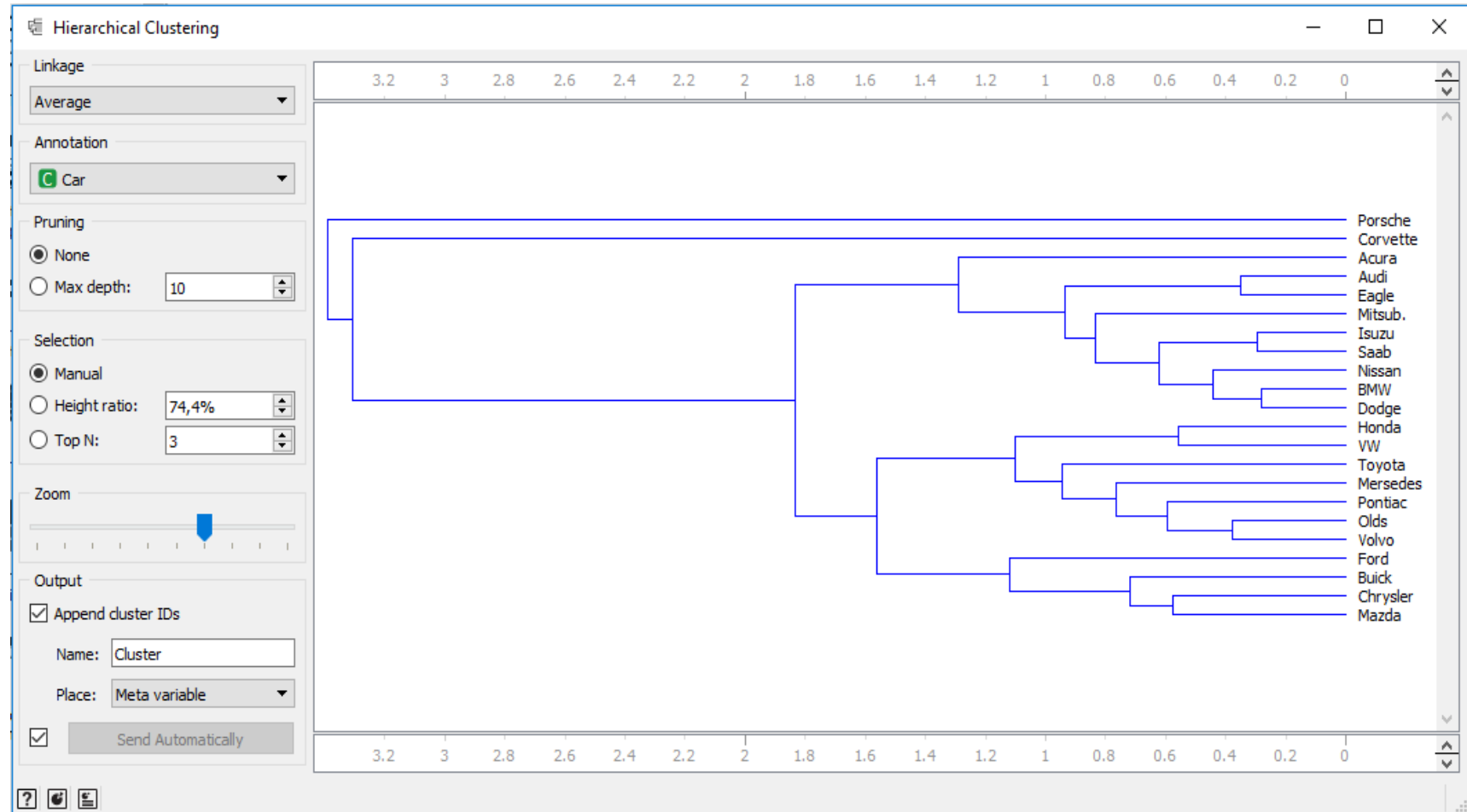
Bottom toolbar: i # T ↘ || ?



«k-Means»



«Hirarchical Clustering»



Акт впровадження

Дипломній комісії
"Інституту прикладного
системного аналізу"

Доповізь, записки

Класіфічний аналіз є невід'ємною частиною
обробки даних в нашій компанії.

Саме тому, працівник, Буханенко Андрій
Олександрович, постійно використовує програмне
забезпечення Orange для класифікації наших
даних. Orange дозволяє працювати з великою
кількістю даних як із базою даних, так із таблицями
Excel, робить необхідні звіти з великою
візуалізацією результатів.

05.06.2018

Керівник: Британ В.В.



Методичне забезпечення

РОЗДІЛ 1	ШВИДКИЙ СТАРТ
РОЗДІЛ 2	ДЕРЕВА РІШЕНЬ
РОЗДІЛ 3	МЕТОД ОПОРНИХ ВЕКТОРІВ
РОЗДІЛ 4	ШТУЧНІ НЕЙРОННІ МЕРЕЖІ
РОЗДІЛ 5	РЕГРЕСІЙНИЙ АНАЛІЗ
РОЗДІЛ 6	КЛАСТЕРНИЙ АНАЛІЗ

Висновки

- Розглянуто основні методи кластеризації даних, розроблена їхня порівняльна характеристика.
- Опрацьовані доступні методи кластеризації в програмному забезпеченні Orange.
- Розроблене методичне забезпечення для програмного забезпечення Orange.

Дякую за увагу!