

# Кластеризація. Модель k-means

<b>Актуальність</b>	Дев'яносто відсотків даних у світі сьогодні створено лише за останні два роки. Щодня створюється близько трьох мільярдів гігабайтів нової інформації. Збереження такого об'єму даних потребує немалої ресурсів та зусиль. Але найбільша проблема полягає не в збереженні даних, а в їх обробці – за прогнозами в 2020 році лише 35% даних будуть корисними. Тому інтелектуальний аналіз даних є актуальним проблемою сьогодення. Кластерний аналіз є потужним механізмом для спрощення даних для подальшого їх аналізу.
<b>Постановка проблеми</b>	Сьогодні використовуються різні модифікації даного алгоритму, але кожен із них має два основних недоліки: необхідність знати кількість кластерів перед виконанням алгоритму і оптимальність отриманих рішення не є гарантованою.
<b>Шляхи вирішення проблеми</b>	Тому наша модифікація буде направлена на прискорення роботи алгоритму за допомогою використання паралельних обчислень.  Розроблений алгоритм представляє собою ітераційну процедуру: <ul style="list-style-type: none"><li>- Крок 1 (CPU). Випадковим чином обрати центри кластерів із елементів вхідних даних, встановити номер ітерації <math>l = 0</math>.</li><li>- Крок 2 (CPU). Підключити GPU модуль, обрахувати конфігурацію запуску для функції-ядра:</li></ul>
<b>Результати та висновок</b>	В даному розділі проведено аналіз алгоритму k-means на предмет можливості розпаралелювання обчислень. Виявлено, що в цьому алгоритмі операції обчислення матриці приналежності та центрів кластерів можна виконати паралельно. Наведено модифіковані алгоритми кластеризації для масово паралельних обчислень на графічних процесорах Nvidia.  Математичний аналіз складності алгоритму показує велике теоретичне прискорення для розроблених модифікацій алгоритмів в порівнянні з виконанням на центральному процесорі, що може досягати більше ста для модифікованого алгоритму k-means.