

DATA1002 Week 10

Tutorial

Monday 13/10/25

Tutorial Outline

- Content revision (Machine Learning Intro) [15 min]
 - Content revision (Clustering) [30 min]
 - Content revision (Supervised Learning) [30 min]
 - Python Task & Research Task (for Assignment 2) [15 min]
 - Revision for Coding Test [30 min]



THE UNIVERSITY OF SYDNEY

Tutor: *Tommy Lu*

54

Housekeeping

Group Project Stage 2 (Presentation)

Due: 11:59 PM on Sunday at the end of week 12 (Nov 2nd)

Value: 8% of Total Mark

Note: Get started your project ASAP. Discuss with your tutors and make use of Ed to ask questions.

- Aggregate summaries
- Charts and visual representations
- Machine learning predictions
- Presentations



Data analysis 	Project stage 2 in-class group presentation on predictive model and evaluation of its success.	8%	Week 12
---	--	----	---------

Content

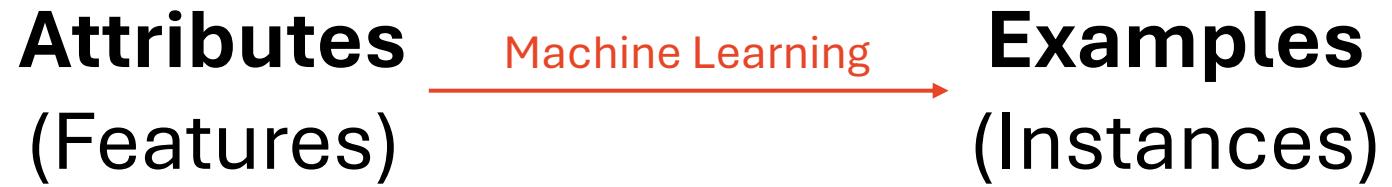
Introduction to Machine Learning

Data Science Process

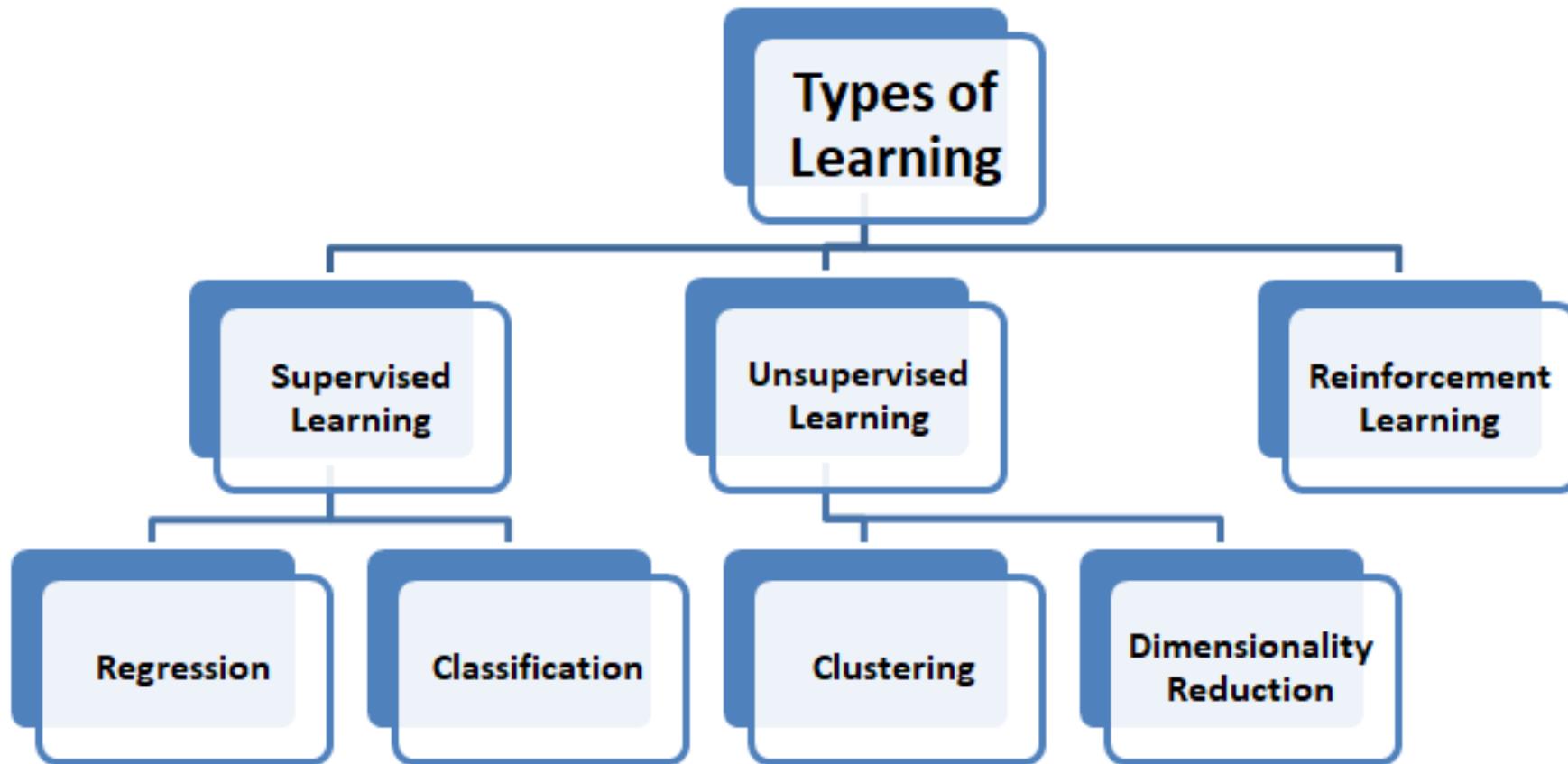
Attributes
(Features)

Examples
(Instances)

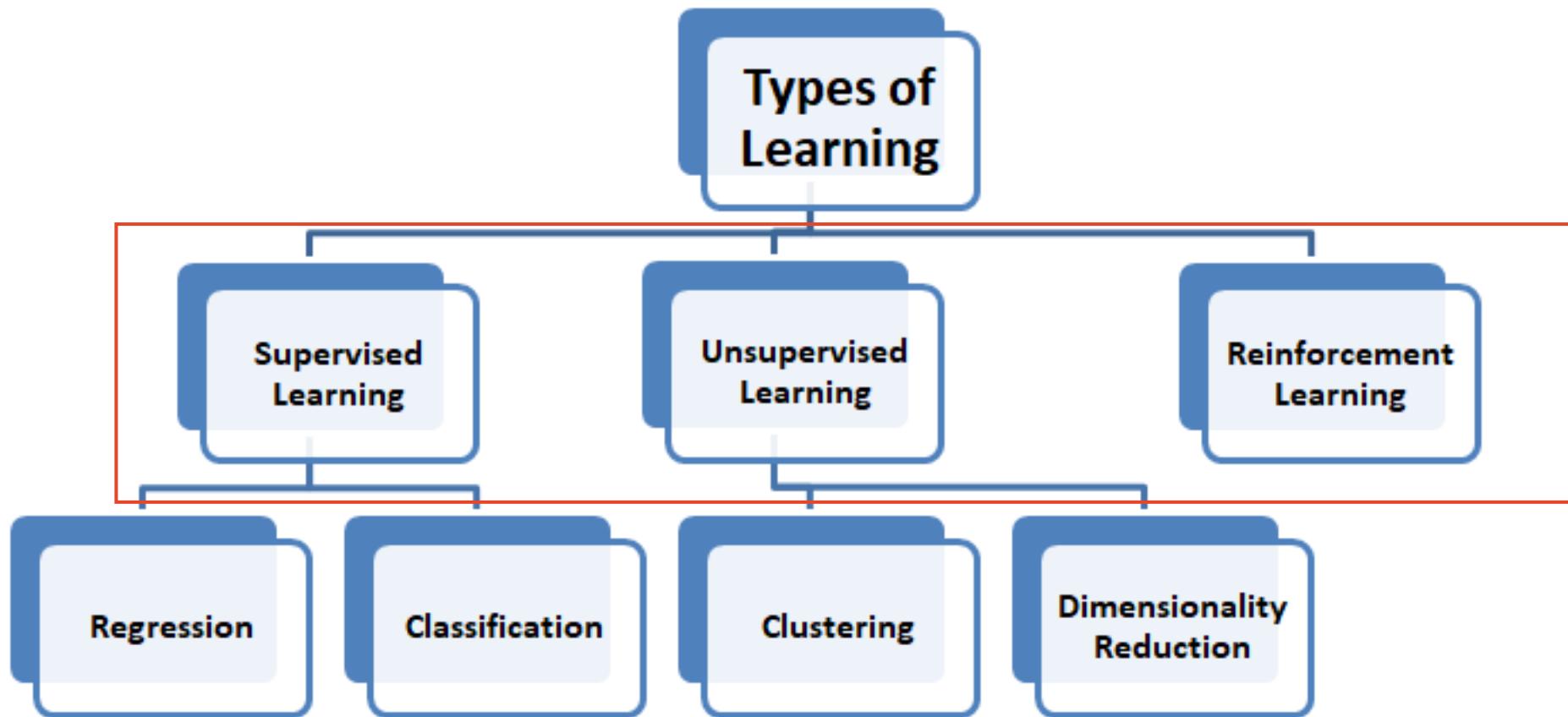
Data Science Process



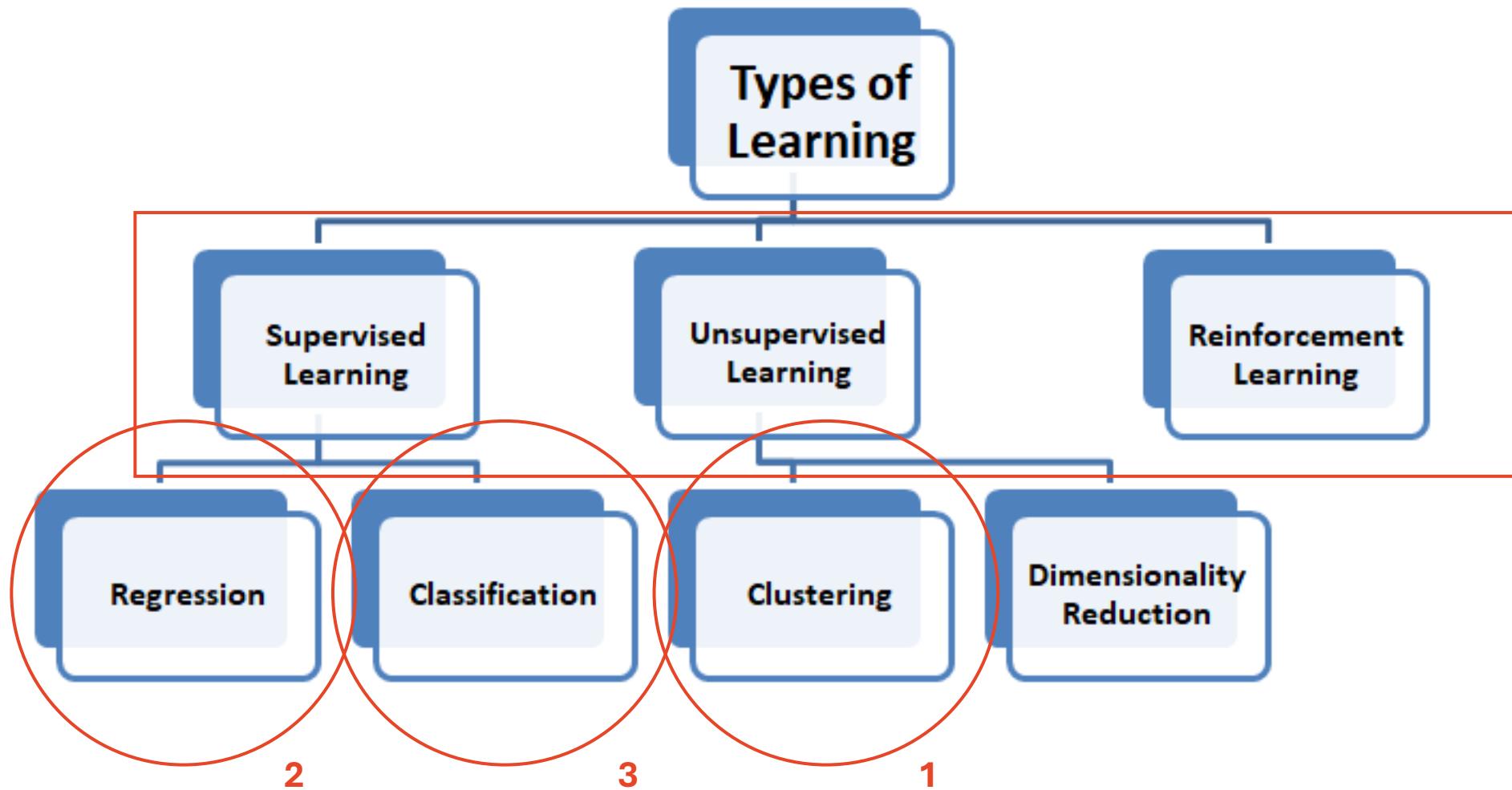
Machine Learning



Machine Learning



Machine Learning



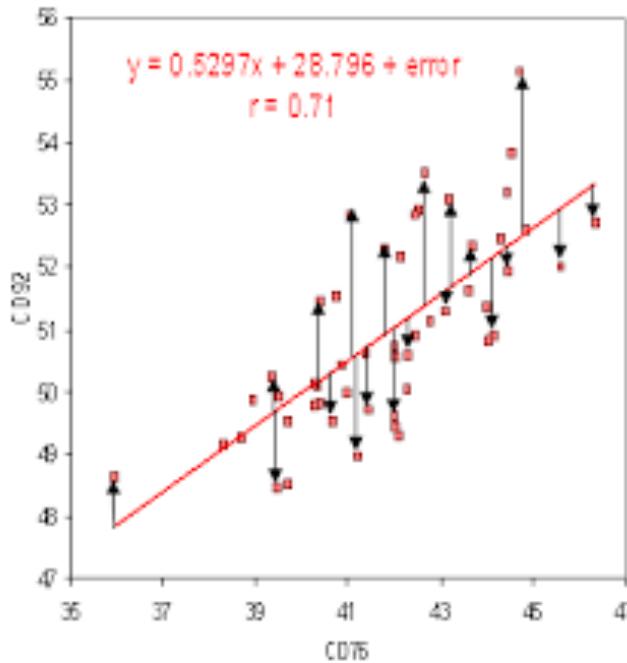
Model Training

Data

Loss Function

Optimization

Model training uses data, a loss function, and an optimization method (e.g., gradient descent)



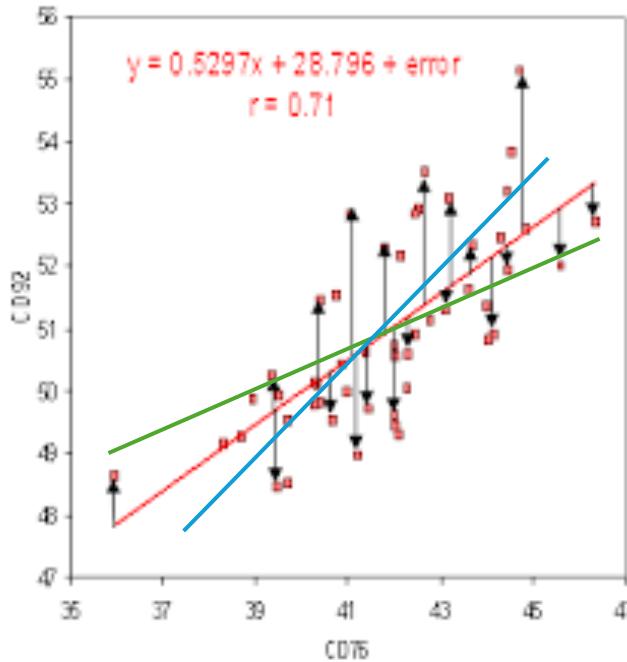
Model Training

Data

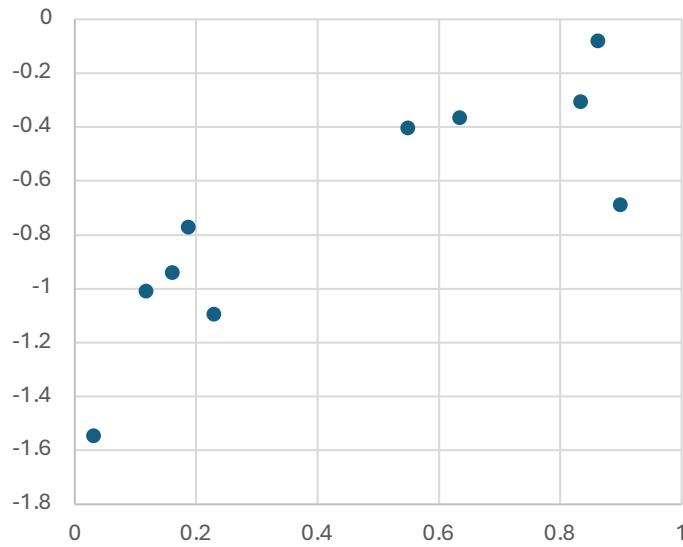
Loss Function

Optimization

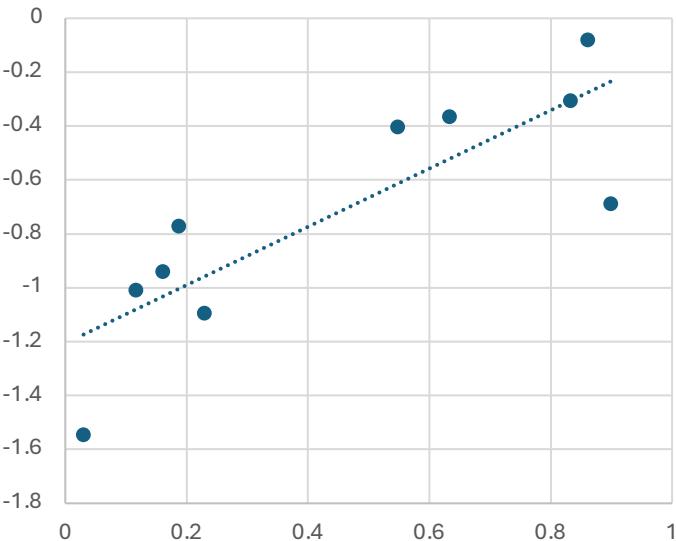
Model training uses data, a loss function, and an optimization method (e.g., gradient descent)



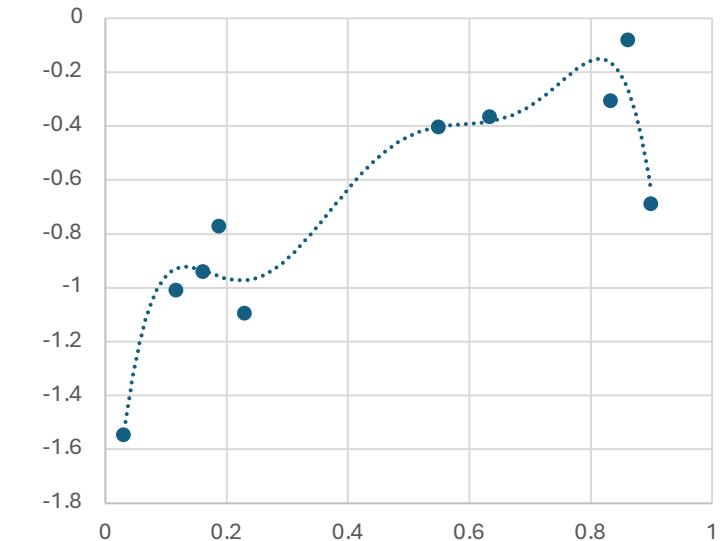
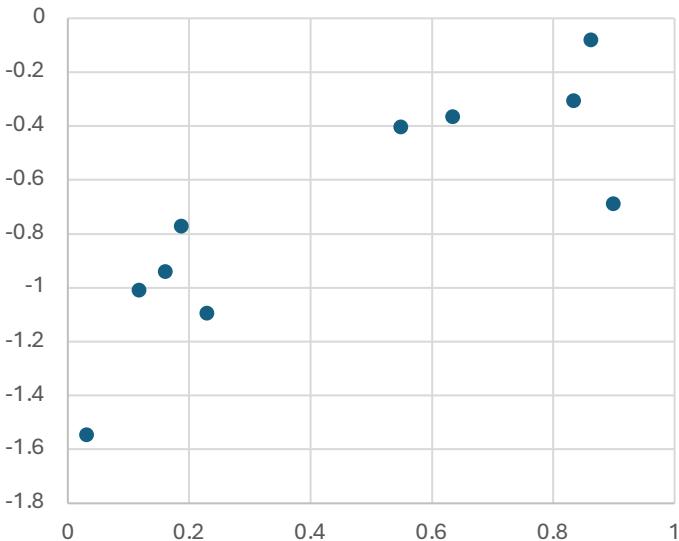
Model Training



Model Training

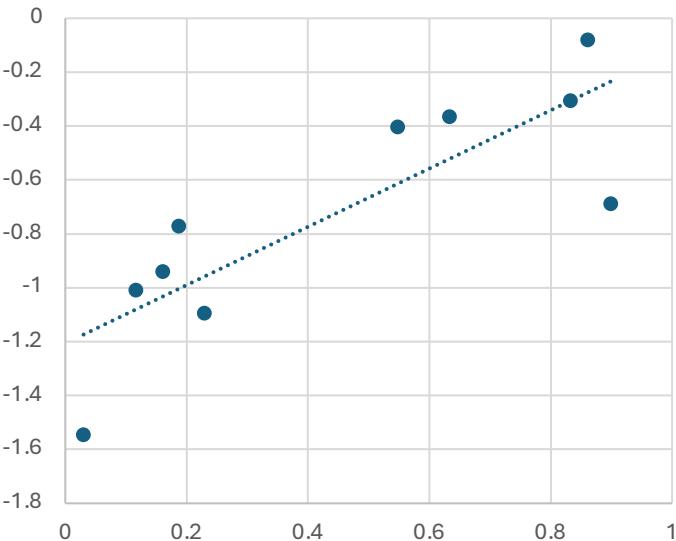


Underfitting

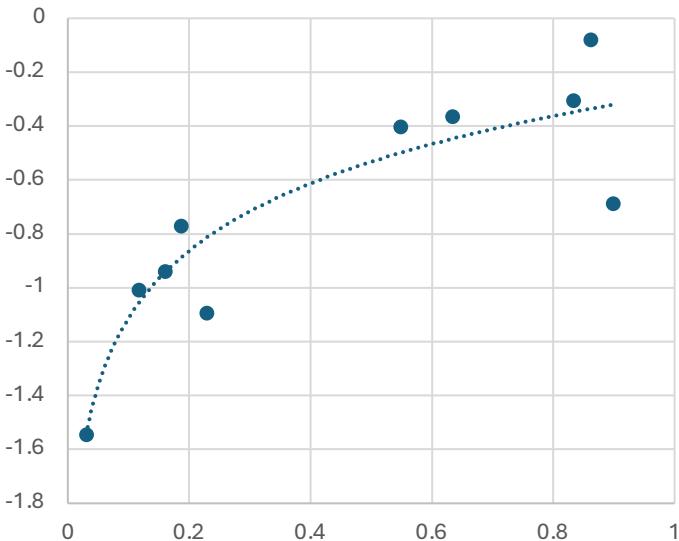


Overfitting

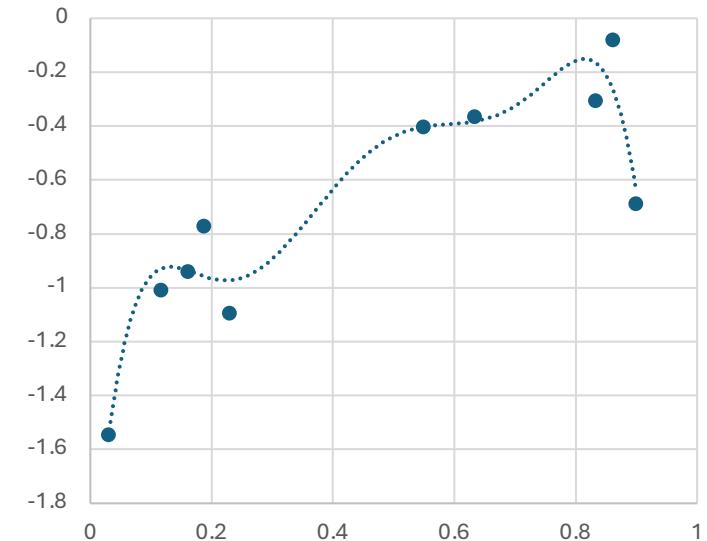
Model Training



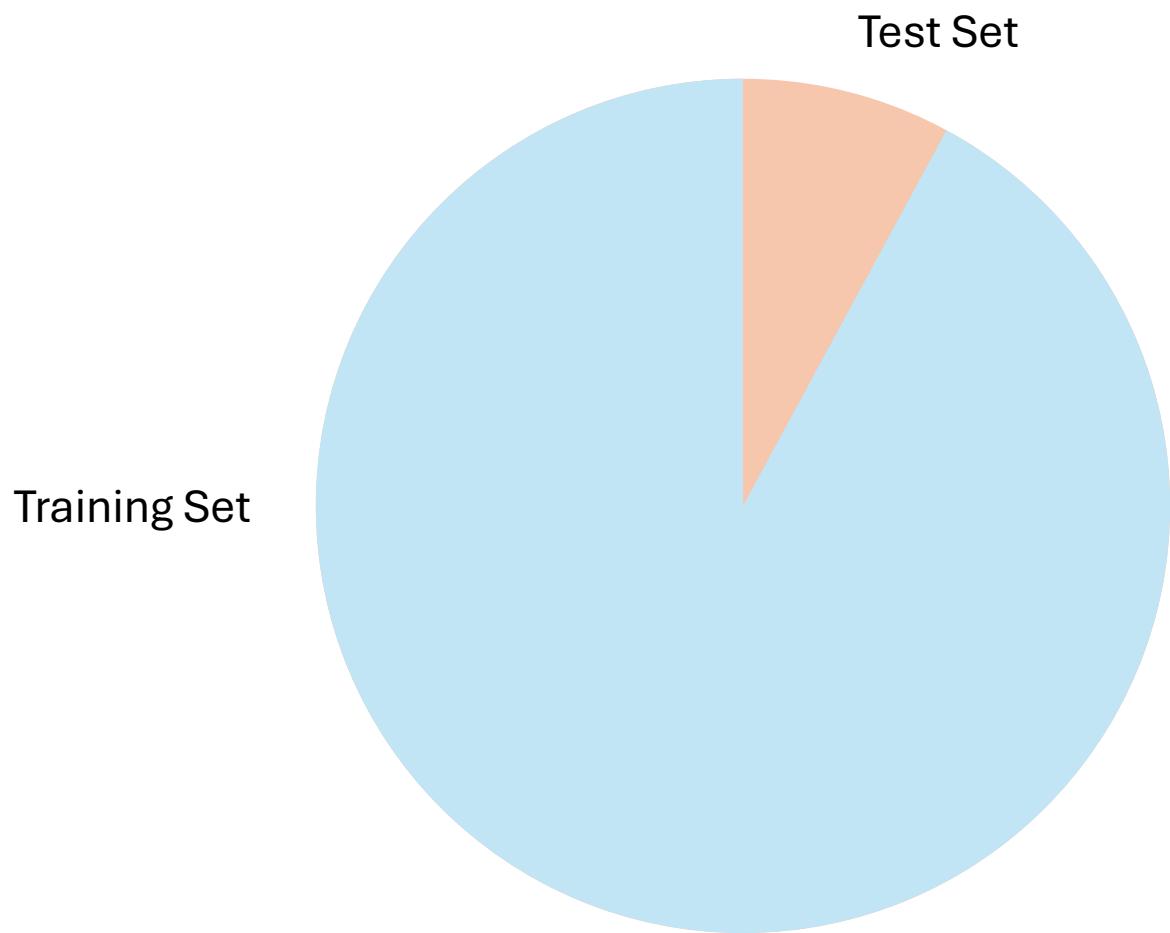
Underfitting



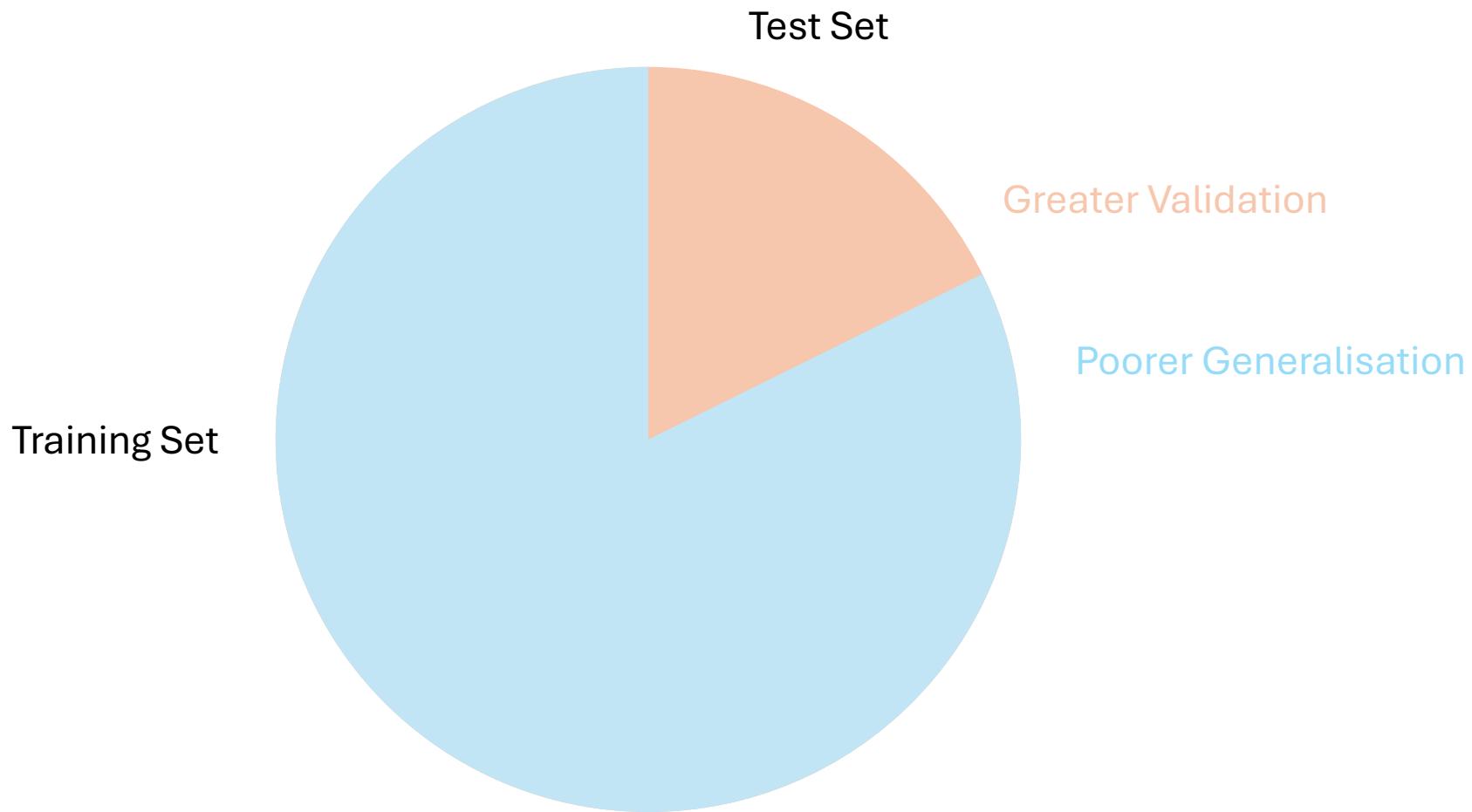
Overfitting



Model Training



Model Training



Content

Clustering

Clustering

THE DEEP SEQUENCING TOOLBOX

Tommy Y. Lu^{1,2}, John Chen³, Colin J. Jackson^{1,3} & Richard J. Payne^{1,2}

¹ Australian Research Council Centre of Excellence for Innovations in Peptide and Protein Science

² School of Chemistry, University of Sydney, NSW, Australia

³ Research School of Chemistry, The Australian National University, Canberra, ACT, Australia

INTRODUCTION

mRNA display is a screening technique that can select peptides for a protein drug target of interest through iterative binding and amplification cycles. [1]

For drug discovery, the sequences that enrich for the target can be optimised for therapeutic use.

The conventional method of choosing peptide hits for optimisation is simply based on their enrichment rank.

Figure 3.

BLAST is an algorithm that finds aligning motifs (k-mers) and extends them.

A sequence similarity network (SSN) is a network built from BLAST pairwise similarities.

Shows an SSN produced from a mRNA display study of CCL22. [3] Edges represent the pairwise similarities. Node colouration is a continuous mapping of enrichment values.

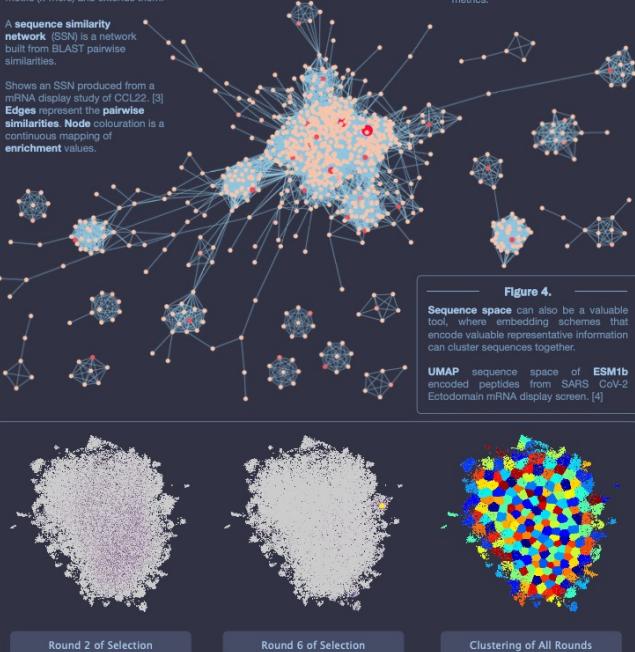


Figure 5.

CD-Hit is a clustering algorithm that uses greedy, length-sorted approach based on global sequence identity. After clustering, sequences can be exported into alignment software that can then backtrack and calculate the key features abstracted in the CD-Hit clustering process.



Figure 1.
Reproducible translation of sequencing data and calculation of enrichment metrics that can be quantified and reported on. Shows four of such key metrics.

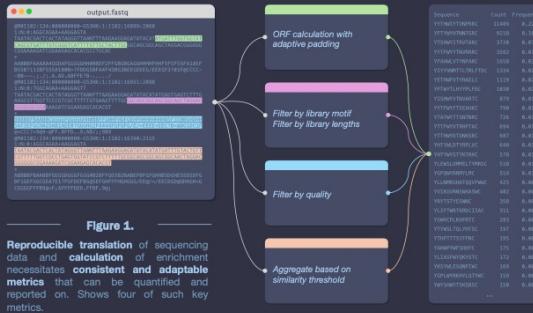
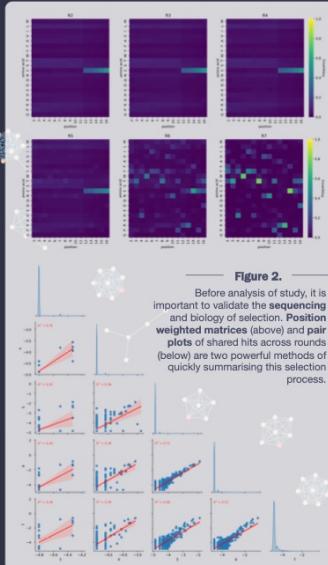


Figure 2.
Before analysis of study, it is important to validate the sequencing and biology of selection. Position weighted matrices (above) and pair plots of shared hits across rounds (below) are two powerful methods of quickly summarising this selection process.



CONCLUSION

Our approach enables the identification of sequence groups that may represent distinct binding modes, increasing the likelihood of uncovering high-quality hits.

Method selection that is deliberate and transparent provides outcome for results that are reproducible and interpretable.

We anticipate that these tools will facilitate more robust hit discovery, improve the comparability of findings between research groups, and ultimately accelerate the translation of mRNA display outputs into therapeutic candidates.



References

- [1] M.S. Newton et al., ACS Synth. Biol., 2020, 9(2), 181-190.
- [2] S.W. Cotton et al., Nat. Protoc., 2011, 6(8), 1163-1185.
- [3] V.T. Zhou et al., ACS Synth. Biol., 2020, 9(2), 3423-3429.
- [4] V. T. Zhou et al., Proc. Natl. Acad. Sci. U.S.A., 2023, 120(6), e230292120.

Clustering

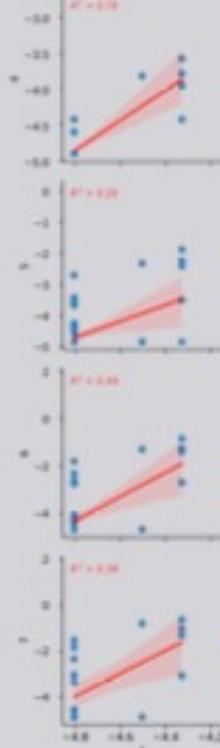


1

output.cdhit.clstr

Y I Y Y V N M T T L T M L F T D C
Y T L T M L F T C - - - - -
Y I Y Y V N M T T L T M L F T N C
Y I Y Y V N M T T L T M L F T D C
Y I Y Y V N M T T L K M L F T D C
Y F Y Y V N M T T L T M L F T D C
Y I Y Y V N M T T L T M L F T A C

JUMAP sequence space of **ESM1b** encoded peptides from SARS CoV-2 Ectodomain mRNA display screen. [4]

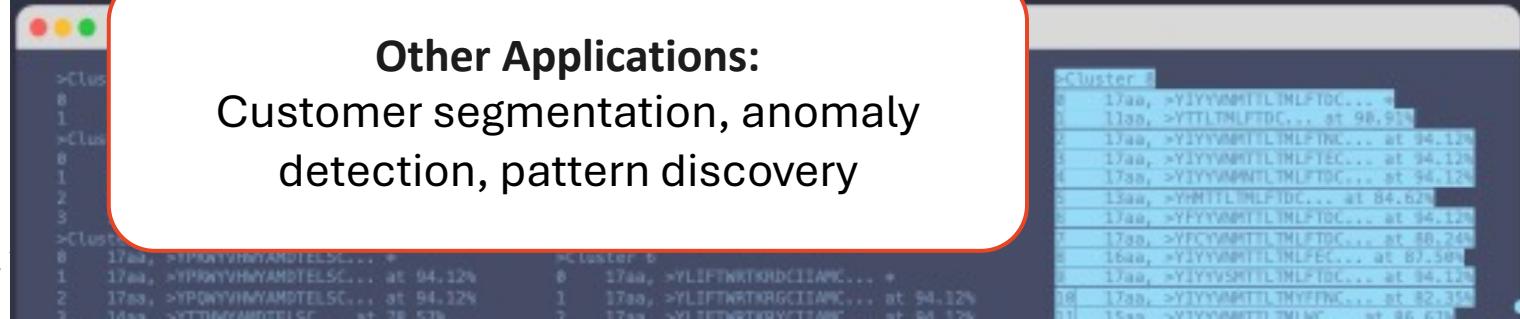


Our approach groups that increasing the

Method selection provides and interprets

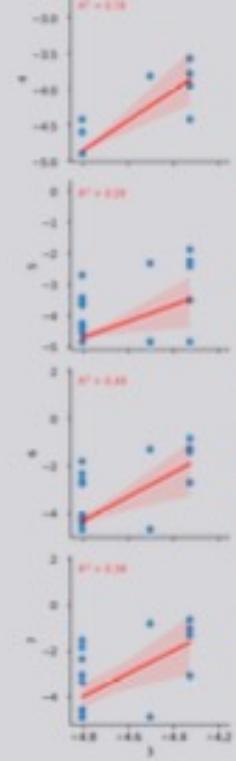
We anticipate the
discovery, implementation

Clustering



encode valuable representative information
can cluster sequences together.

UMAP sequence space of **ESM1b**
encoded peptides from SARS CoV-2
Ectodomain mRNA display screen. [4]



Our approach
groups that
increasing the

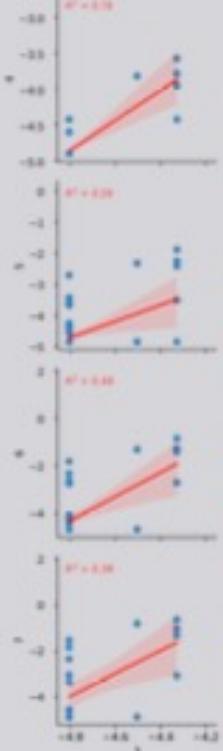
Method selec
provides and
and interpreta

We anticipate t
discovery, imp

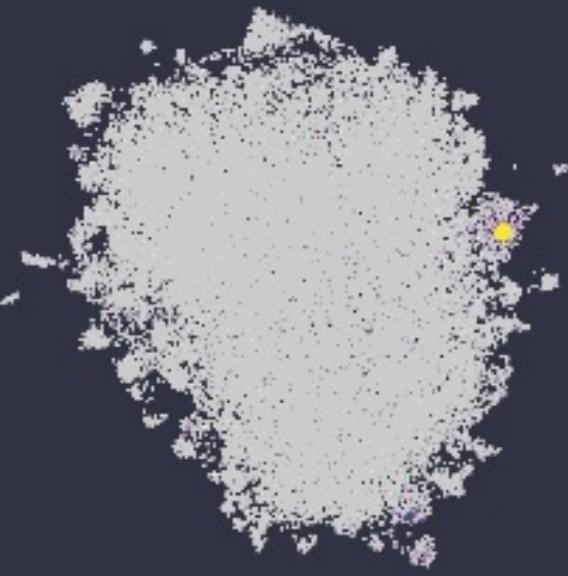
Clustering

encode valuable representative information
can cluster sequences together.

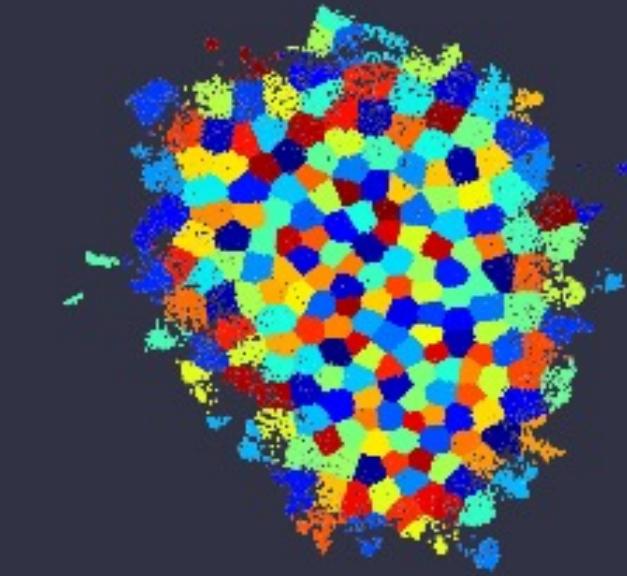
UMAP sequence space of **ESM1b**
encoded peptides from SARS CoV-2
Ectodomain mRNA display screen. [4]



Round 2 of Selection

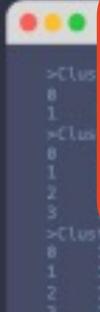


Round 6 of Selection



Clustering of All Rounds

Our approach
groups that
increasing the



Other Applications:
Customer segmentation, anomaly
detection, pattern discovery

Cluster 2		
17aa, >YIYYVNMITLNLFTDC...	at	98.9%
11aa, >YTTLTNLFTDC...	at	98.9%
17aa, >VIYYVNMITLNLFTDC...	at	98.9%
17aa, >VIYYVNMITLNLFTDC...	at	98.9%
13aa, >YHMTTLNLFTDC...	at	84.6%
17aa, >YIYYVNMITLNLFTDC...	at	98.9%
17aa, >YIYYVNMITLNLFTDC...	at	98.9%
16aa, >YIYYVNMITLNLFTDC...	at	87.5%
17aa, >YIYYVNMITLNLFTDC...	at	94.3%
17aa, >YLIFTNTRKRDCLIAMC...	at	94.3%
11aa, >YLIFTNTRKRDCLIAMC...	at	94.3%
17aa, >YIYYVNMITLNLFTDC...	at	98.9%

Similarity metrics:
Euclidean, Manhattan, Cosine

Clustering

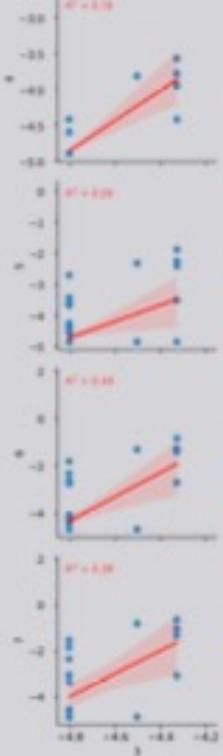


Other Applications:
Customer segmentation, anomaly detection, pattern discovery

Normalisation crucial
Similarity metrics:
Euclidean, Manhattan, Cosine

encode valuable representative information
can cluster sequences together.

UMAP sequence space of **ESM1b**
encoded peptides from SARS CoV-2
Ectodomain mRNA display screen. [4]

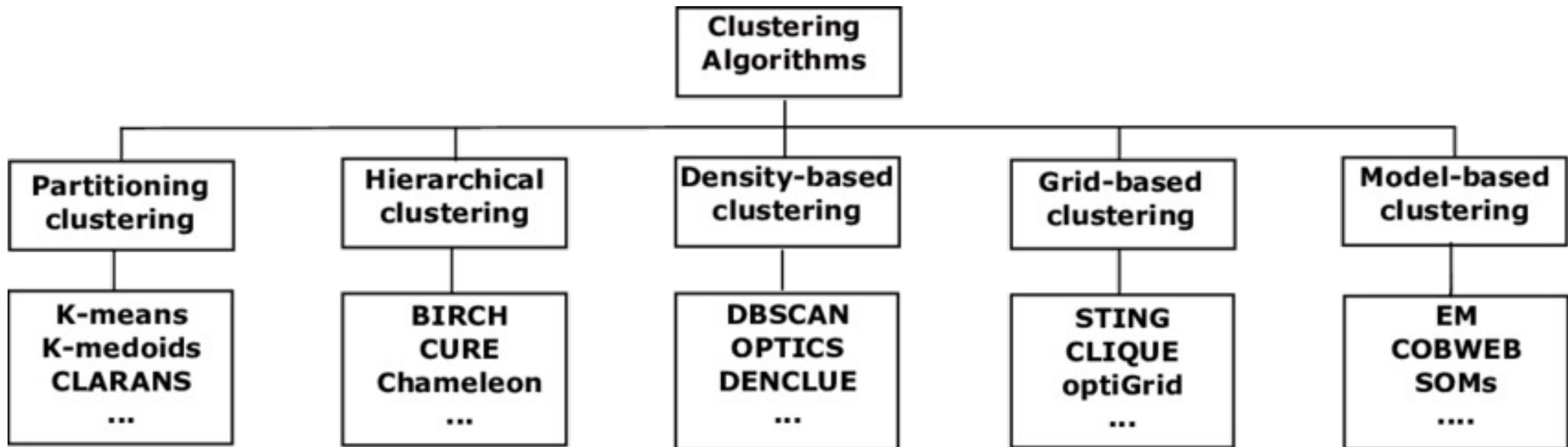


Our approach groups that increasing the

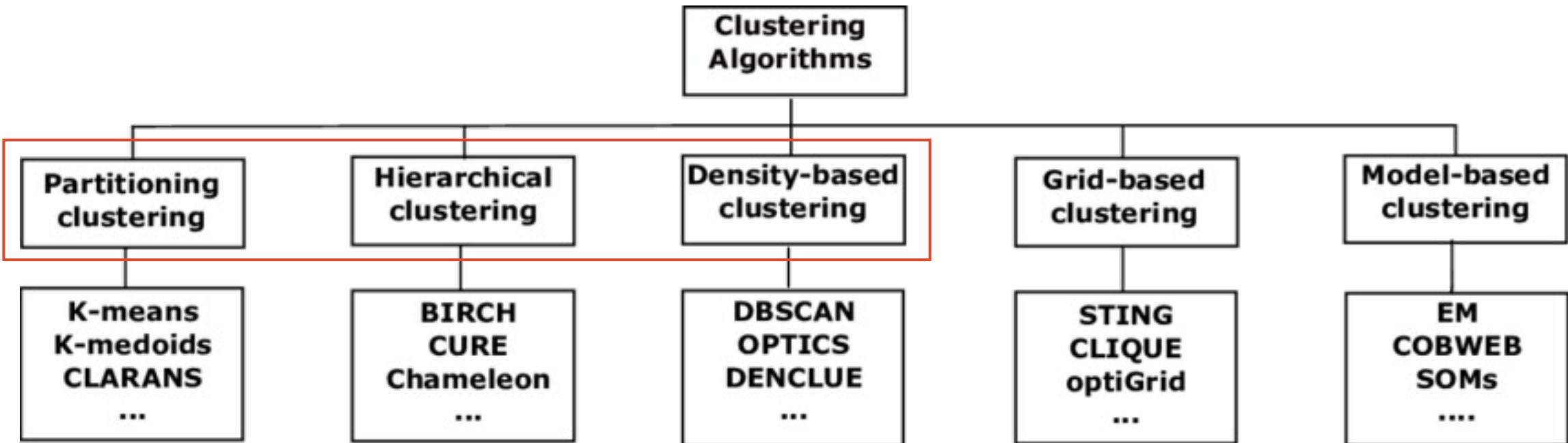
select
and
retai

discovery, imp

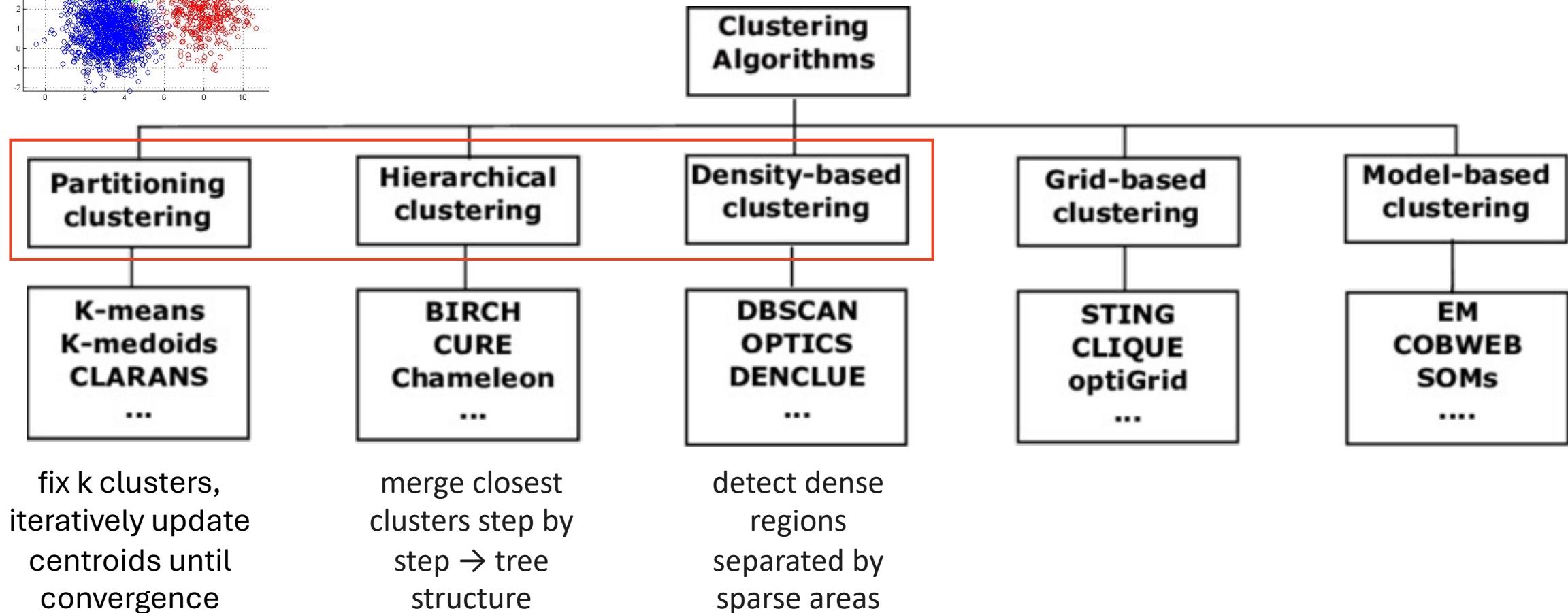
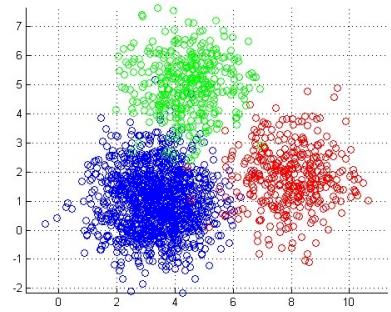
Types of Clustering



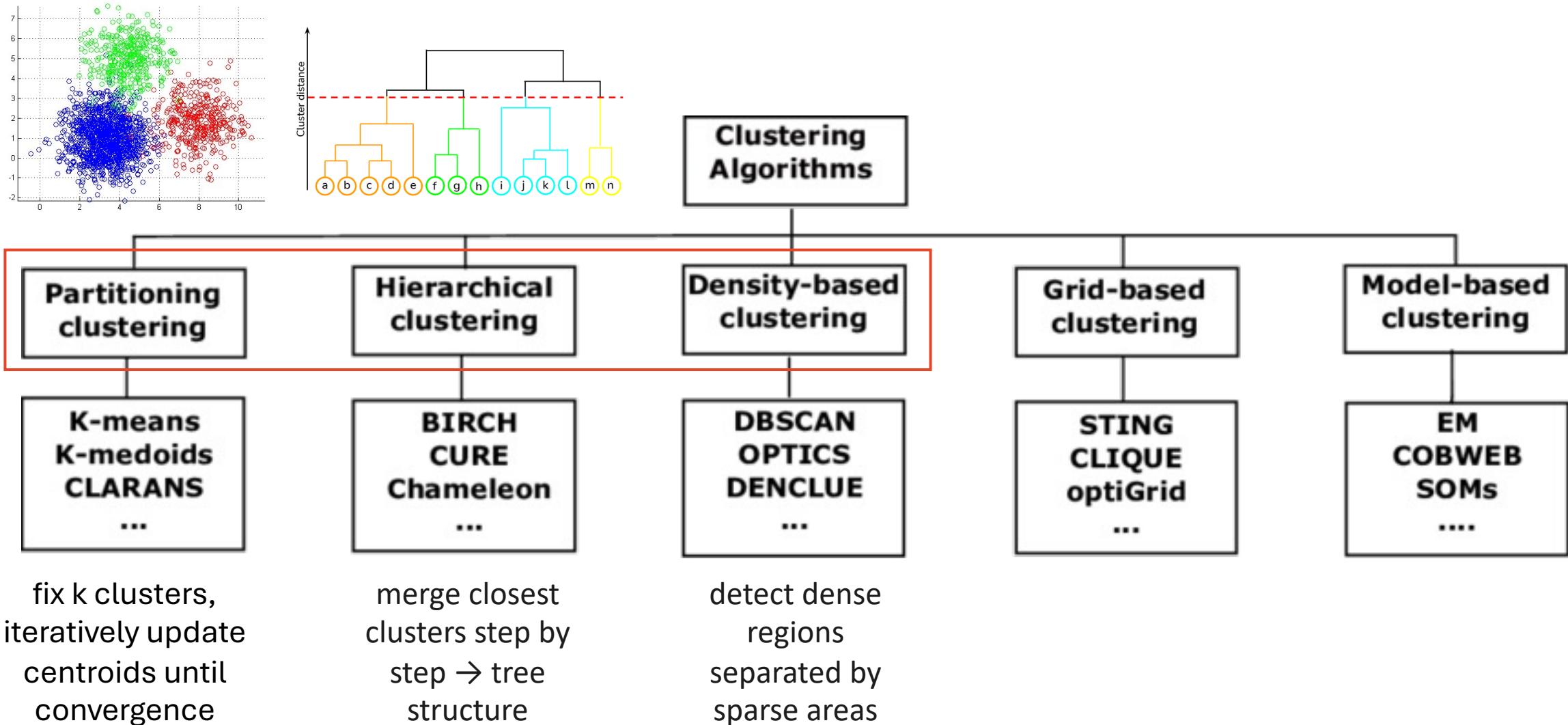
Types of Clustering



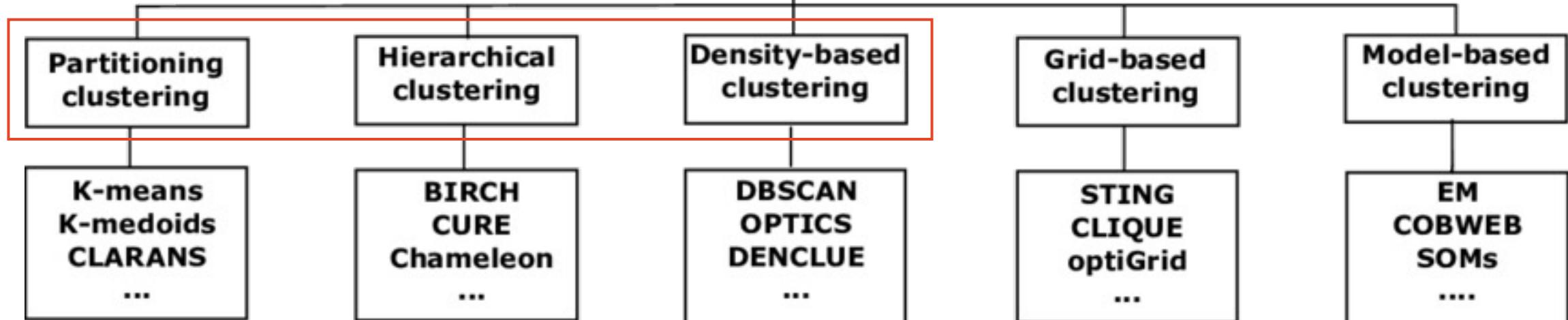
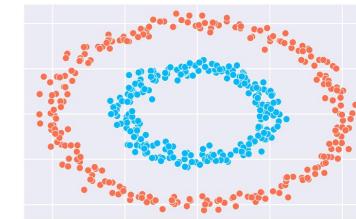
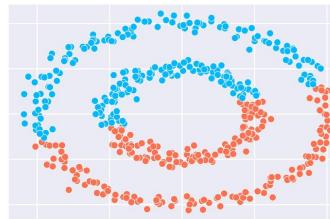
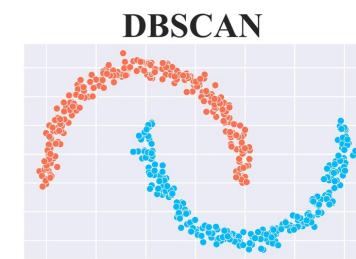
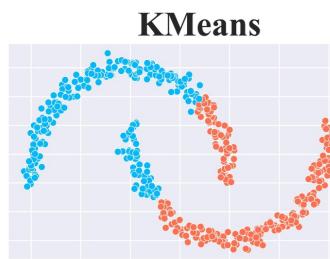
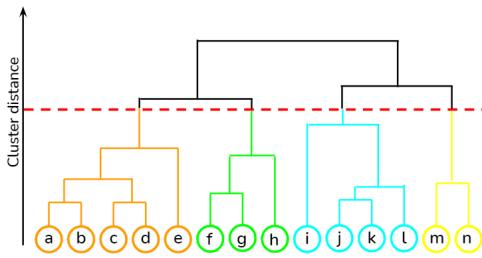
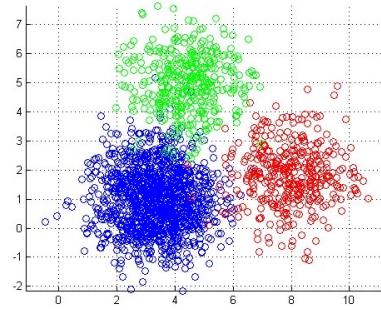
Types of Clustering



Types of Clustering



Types of Clustering

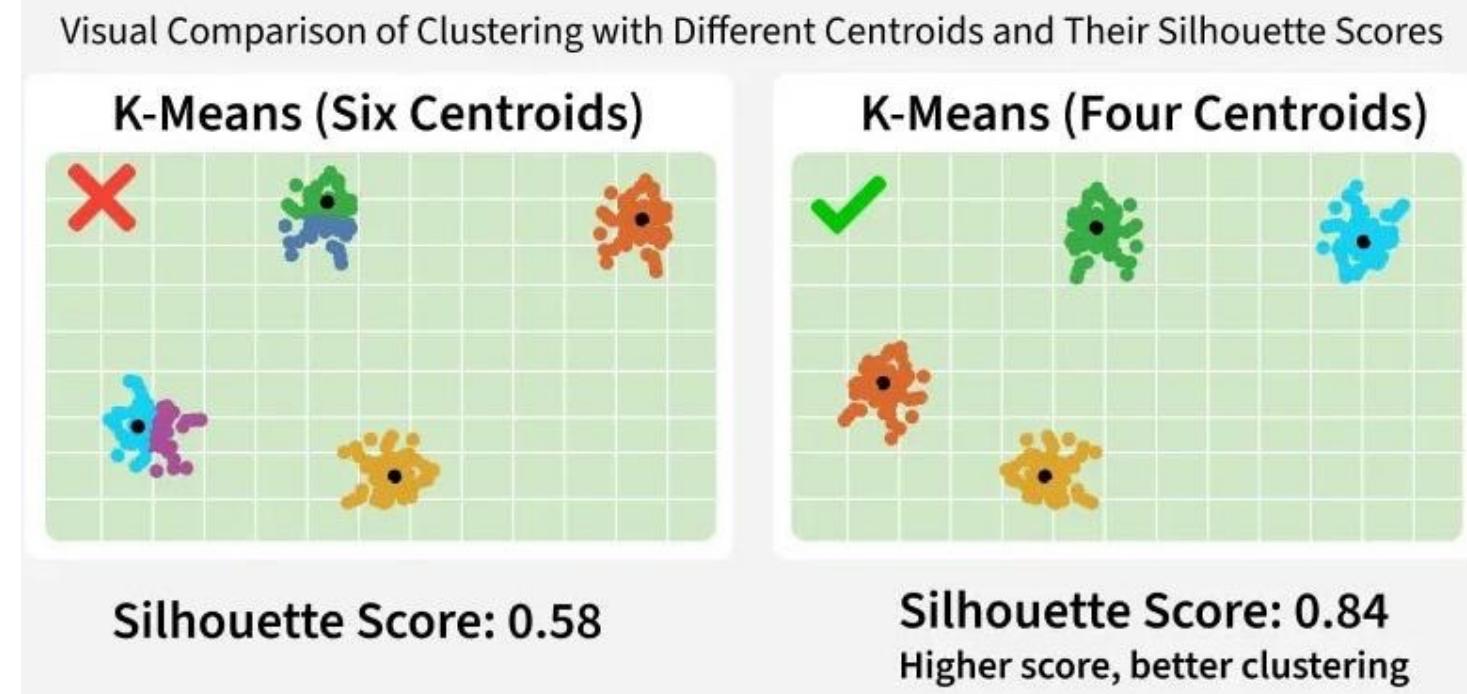


fix k clusters,
iteratively update
centroids until
convergence

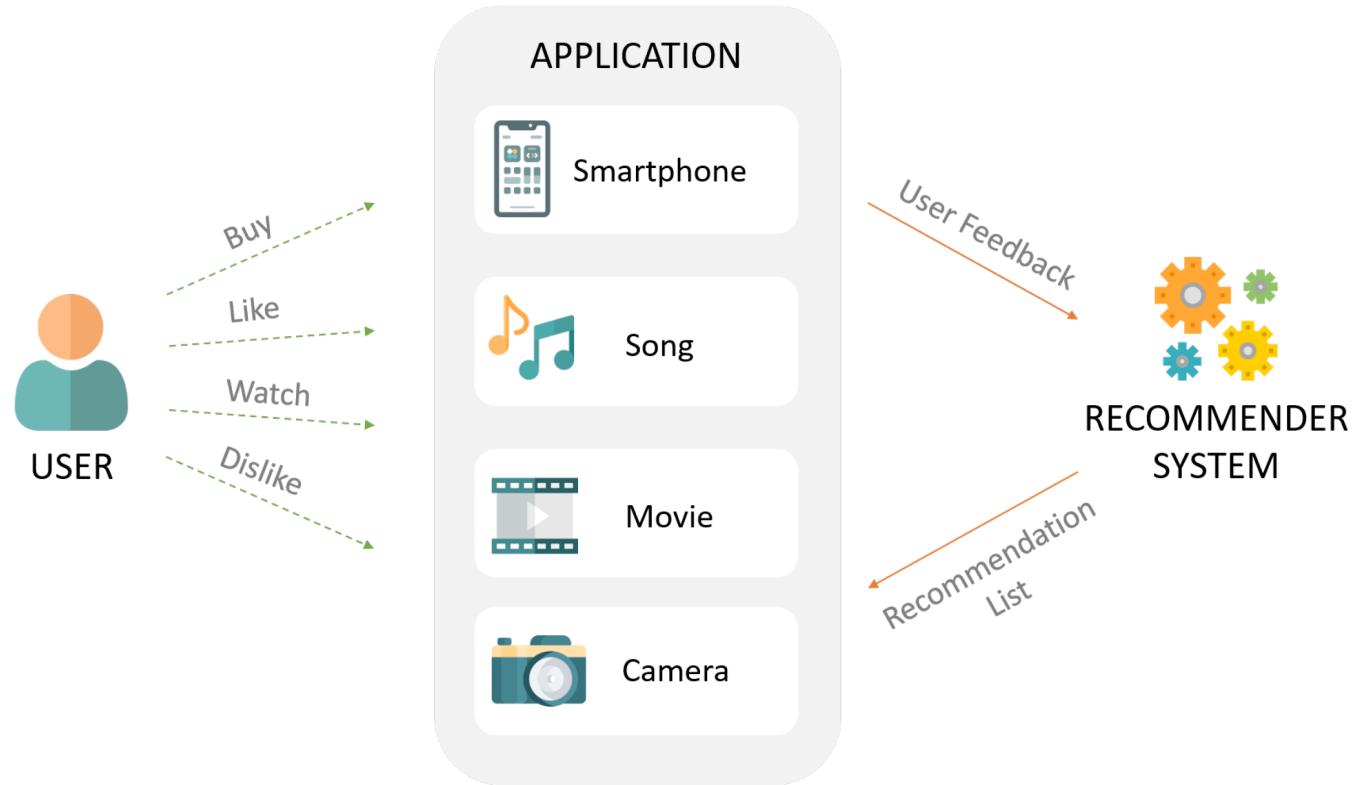
merge closest
clusters step by
step → tree
structure

detect dense
regions
separated by
sparse areas

Evaluation



Recommender



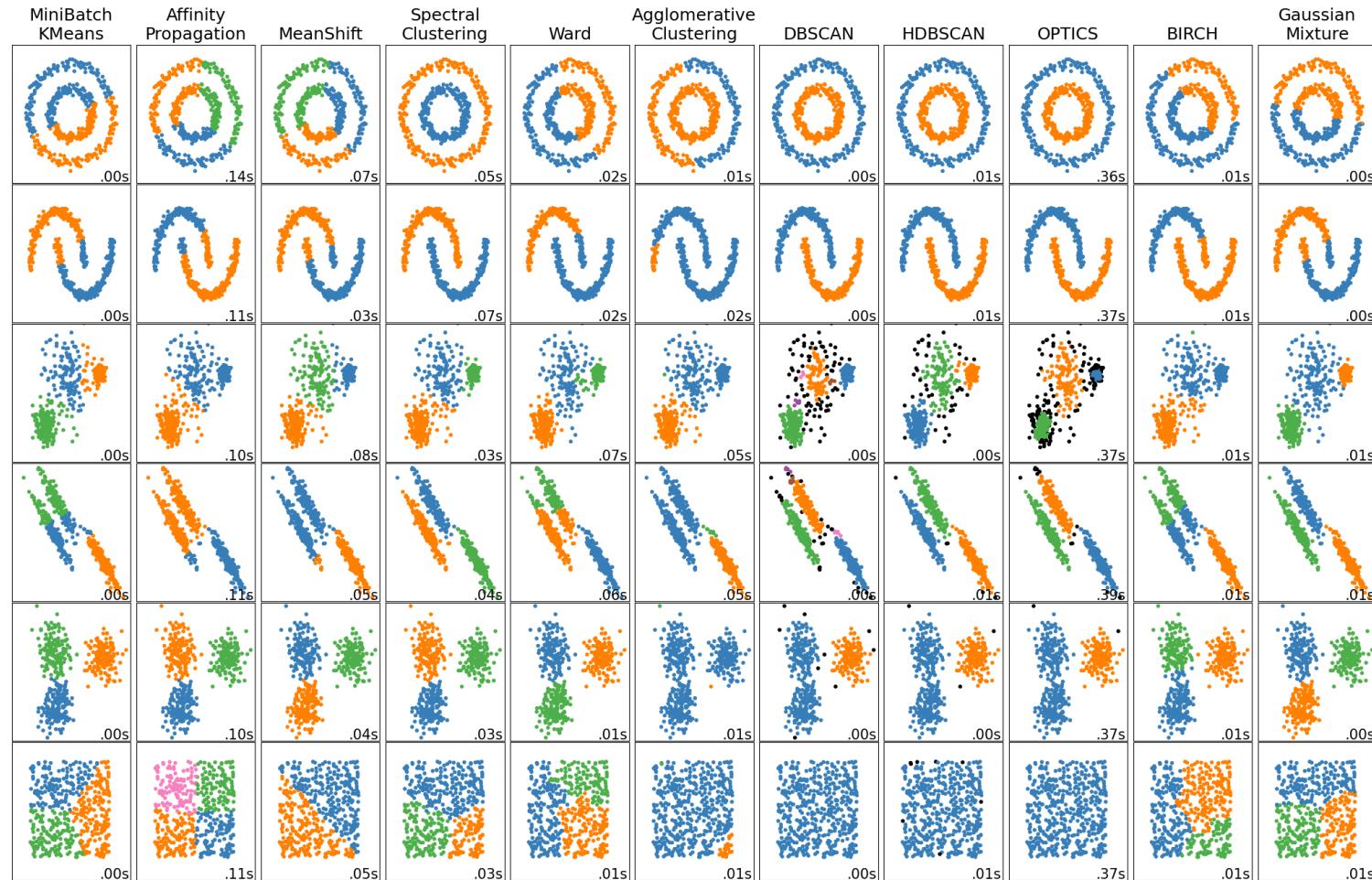
Python Demonstration

Clustering

Try running clustering algorithm

<https://scikit-learn.org/stable/modules/clustering.html>

See `cluster.py` on the Ed Workspace



Content

Regression

Regression Model

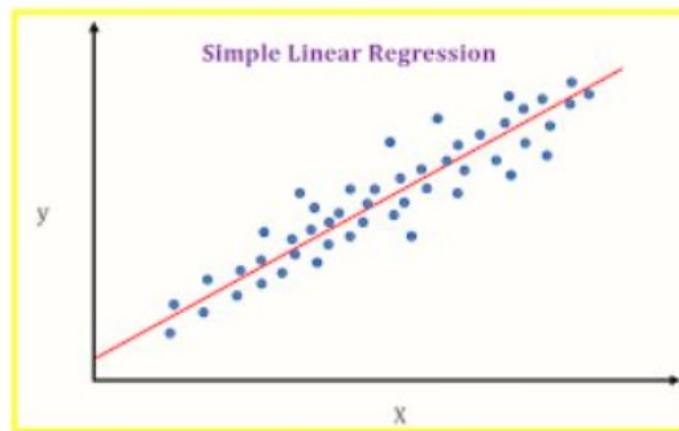
- It is a statistical technique used in machine learning, statistics, and data analysis to quantify the relationship between a dependent variable (also called the target) and one or more independent variables (predictors or features)
- Its primary purpose is to make predictions or estimate the value of the dependent variable based on the values of the independent variables

Types of Regression Model

- **Linear Regression Model**
- **Polynomial Regression Model**
- **Multiple Regression Model**
- **Logistic Regression Model**

Linear Regression Model

- Simple and effective way to model linear relationships between two variables
- Increase or decrease of dependent variable is proportional to the increase of the predictor – i.e., independent variable

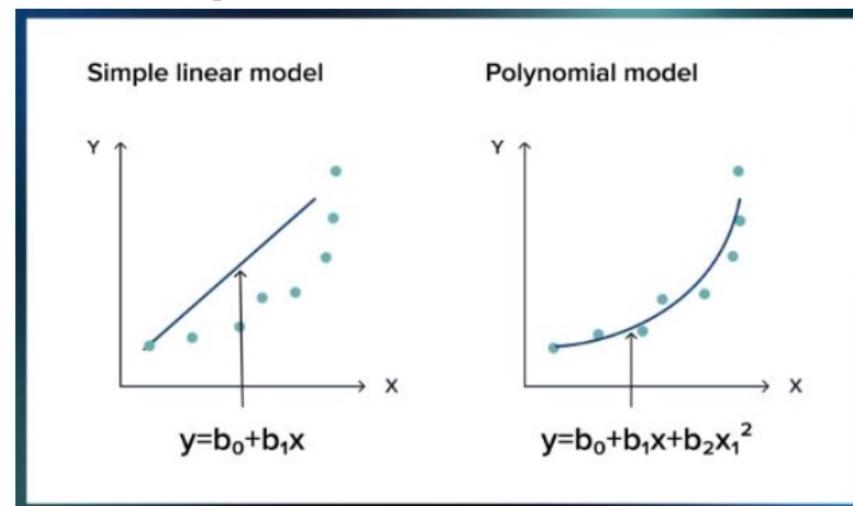


Linear Regression Model

- **Linearity:** It assumes that a **linear relationship** between predictors and the target variable
- **Least Squares Method:** The coefficient is estimated by **minimising the sum of squared differences** between predicted and actual values
- **Assumptions:** Linear regression assumes that **residuals** (differences between predicted and actual values) are **normally distributed**. The residuals have constant variance, and there is no correlation between predictors (**multicollinearity**)
- **Inference:** This regression model allows **hypothesis testing** on **coefficients** to determine their significance and assess the model's explanatory power

Polynomial Regression Model

- Models the relationship between a dependent variable and one or more independent variables using a **polynomial function**
 - Polynomial functions: functions in which the **independent variable is raised to a power**



Polynomial Regression Model

- It is an extension of the linear regression model by allowing for **nonlinear relationships** between **predictors** and the **target variable**
- It uses **polynomial terms** (e.g., x^2, x^3) to capture curved patterns in the data.
- It aims to fit a polynomial equation to the data using the **method of least squares**, **minimising the sum of squared differences** between predicted and actual values.
- **Warning:** only do regression with powers of x, when you have good reason to expect a polynomial impact of the degree you are including
 - It is very easy to overfit training data with polynomial model

Which Model To Use?

- It depends on the specific problem you are trying to solve
 - If you are trying to model a linear relationship between variables, then linear regression is a good choice
 - If you are trying to model a non-linear relationship between variables, then polynomial regression or multiple regression may be a better choice
 - If you are trying to classify data, then logistic regression is a good choice
- All regression models are statistical models, and they have limitations
 - It is important to consider the assumptions of the model and the quality of the data before using the model to make predictions

Research Task

Regression

Regression Models

Pretend you're doing your assignment and look up all the different types of regression and try and get at least one running in Python

https://scikit-learn.org/stable/supervised_learning.html

Content

Regression vs Classification

What is classification?

Classification is a supervised learning task where the goal is **to predict a discrete class label for a given input**. The model learns from labeled training data to make predictions on new, unseen data.

Key aspects of classification:

Types of classification:

- Binary classification: Two possible classes (e.g., spam or not spam)
- Multi-class classification: More than two classes (e.g., classifying animals into species)
- Multi-label classification: Each instance can belong to multiple classes simultaneously

Examples of common algorithms:

- Logistic Regression: Despite its name, it's used for binary classification
- Decision Trees: Tree-like model of decisions
- Random Forests: Ensemble of decision trees

Output of these algorithms are **categorical**

Evaluation Metrics

Accuracy: Proportion of correct predictions

Precision: Proportion of true positive predictions among all positive predictions

- True Positives / (True Positives + False Positives)

Recall: Proportion of true positive predictions among all actual positive instances

- Recall = True Positives / (True Positives + False Negatives)

F1-score: Harmonic mean of precision and recall

- $F1 = 2 * (\text{Precision} * \text{Recall}) / (\text{Precision} + \text{Recall})$

Regression vs. Classification

Output type:

- Classification produces categorical outputs
- Regression produces numerical outputs.

Evaluation metrics:

- Classification uses metrics like accuracy, precision, recall, and F1-score.
- Regression uses metrics like mean squared error (MSE), root mean squared error (RMSE), and R-squared.

Decision boundaries:

- Classification algorithms try to find decision boundaries between classes
- Regression algorithms aim to fit a line or curve to the data.

Lab Activities

Working on Assignment 1

Activity

For Data1002, the main goal of this meeting is to finalize the report. Specifically, check each item in the marking table at <https://canvas.sydney.edu.au/courses/65392/assignments/621825>. Create a to-do list and make sure all items have been covered correctly. If you are unsure, ask your tutor.

Exam-Style Questions

Question 1:

How can visualising data using different types of charts improve the understanding of complex datasets in data science projects?

Provide examples of suitable chart types for various data analysis tasks.

Exam-Style Questions

Visualising data using different types of charts can significantly enhance the comprehension of complex datasets by presenting data in an accessible and interpretable manner.

For example, line charts are ideal for **showing trends over time**, such as tracking sales figures over months. Scatter plots effectively **display relationships between two variables**, useful for identifying correlations in datasets, like age and income levels. Bar charts are excellent for comparing categorical data, such as sales across different regions. Pie charts, though less effective for precise comparisons, **can illustrate proportions within a whole**, like **market share distribution**.

By selecting the appropriate chart type, data scientists can **highlight key insights, making complex data more understandable and actionable**.

Exam-Style Questions

Question 2 [DATA1002]:

Discuss the role of ethical considerations in data visualisation within data science projects.

How can misleading charts impact decision-making, and what steps can be taken to ensure ethical visualisation practices?

Exam-Style Questions

Ethical considerations in data visualisation are crucial to maintain trust and integrity in data science projects. Misleading charts can distort data interpretation, leading to poor decision-making.

For instance, truncated y-axes can exaggerate differences between data points, while improper scaling can misrepresent trends. To ensure ethical visualisation practices, data scientists should adhere to principles of clarity and accuracy, such as using appropriate scales, avoiding deceptive design choices, and providing necessary context through labels and legends.

Additionally, transparency about data sources and methods used to create visualisations helps stakeholders understand the limitations and reliability of the presented data. Ethical visualisation fosters informed decision-making and upholds the credibility of the data science profession.

That's it folks!

Remaining Ed Lessons, Questions, Assignment etc.