

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

```
self.logger = logging.getLogger(__name__)
if path:
    self.file = open(os.path.join(settings['job_dir'], 'fingerprint'), 'w')
    self.file.seek(0)
    Fingerprints.update(fp, settings['job_dir'])
    self.file.write(fp + os.linesep)
    self.file.close()
else:
    settings = Settings(settings)
    if settings.getbool('superuser_fp'):
        job_dir = settings.job_dir(settings)
    else:
        job_dir = settings['job_dir']
    if settings.getbool('seen'):
        f = request_fingerprint(request)
        self.fingerprints.add(f)
        return True
    else:
        self.fingerprints.add(f)
        self.file = open(os.path.join(job_dir, 'fingerprint'), 'w')
        self.file.write(fp + os.linesep)
        self.file.close()
def request_fingerprint(self, request):
    return request_fingerprint(request)
```

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

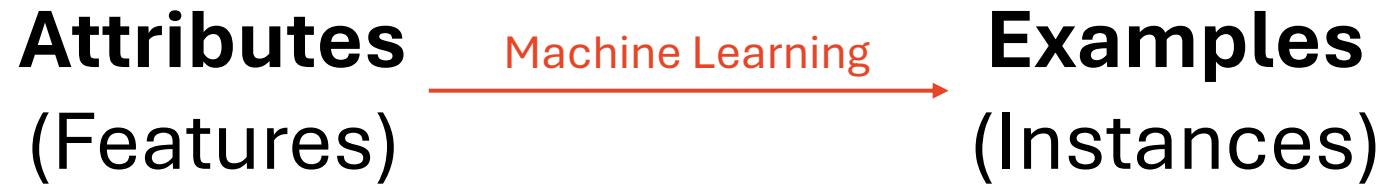
Introduction to Machine Learning

Data Science Process

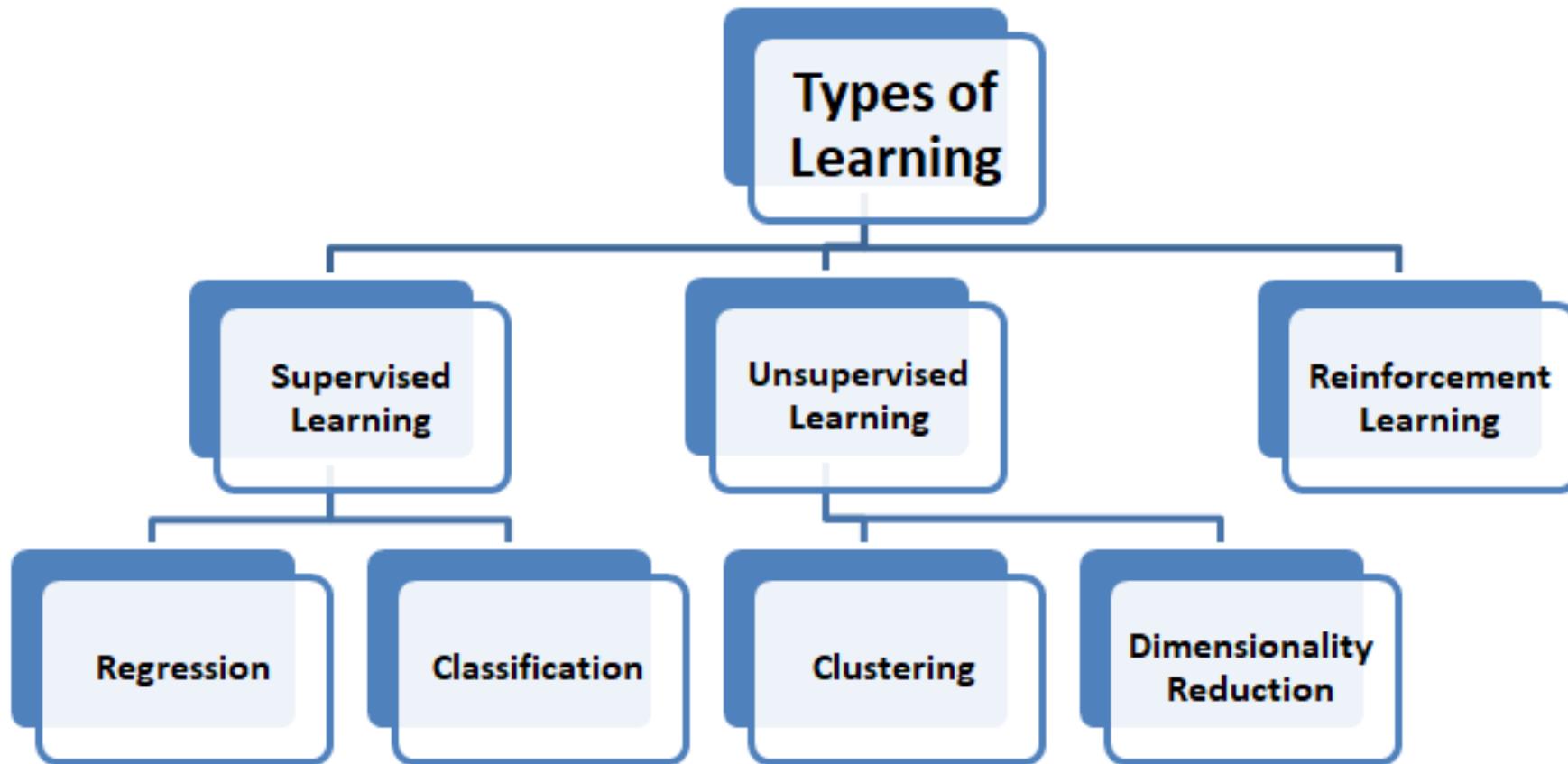
Attributes
(Features)

Examples
(Instances)

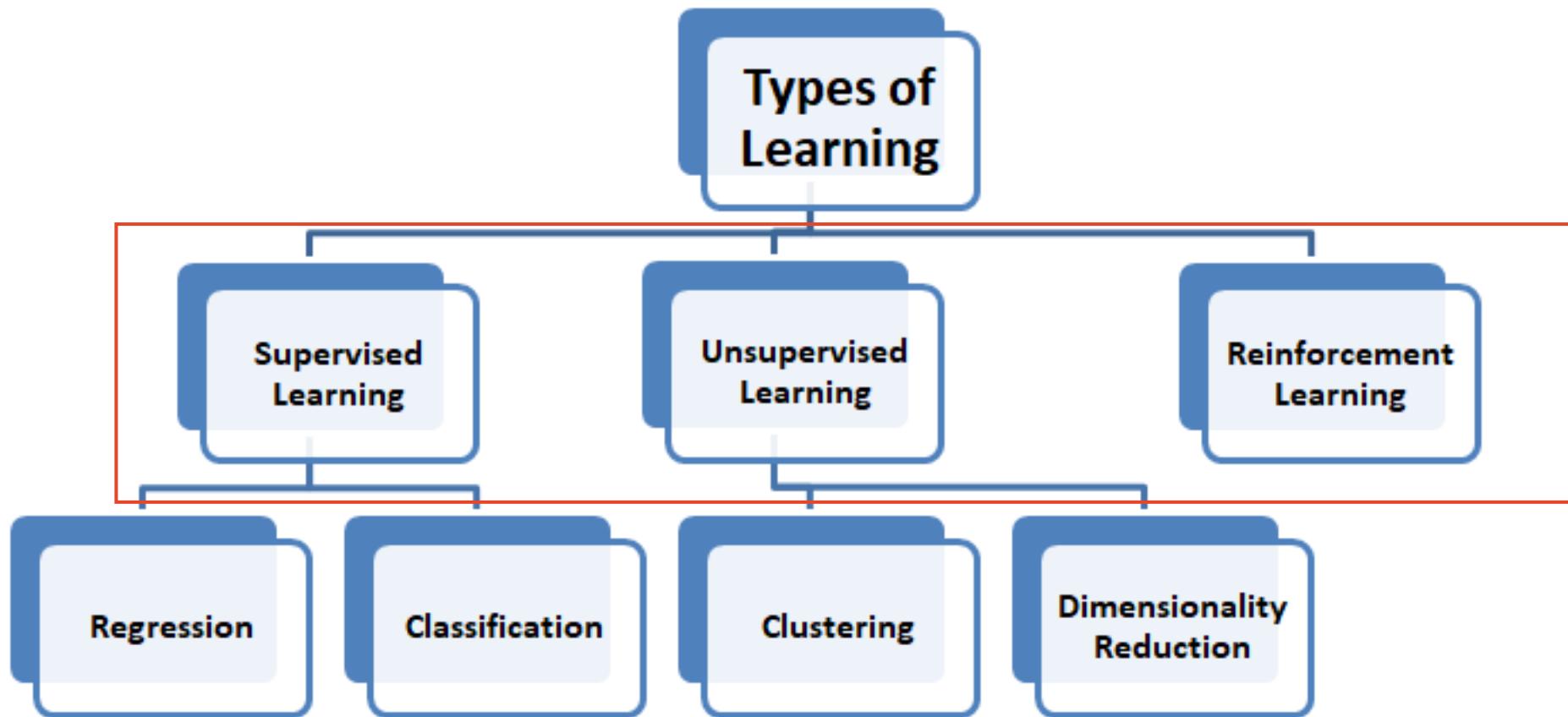
Data Science Process



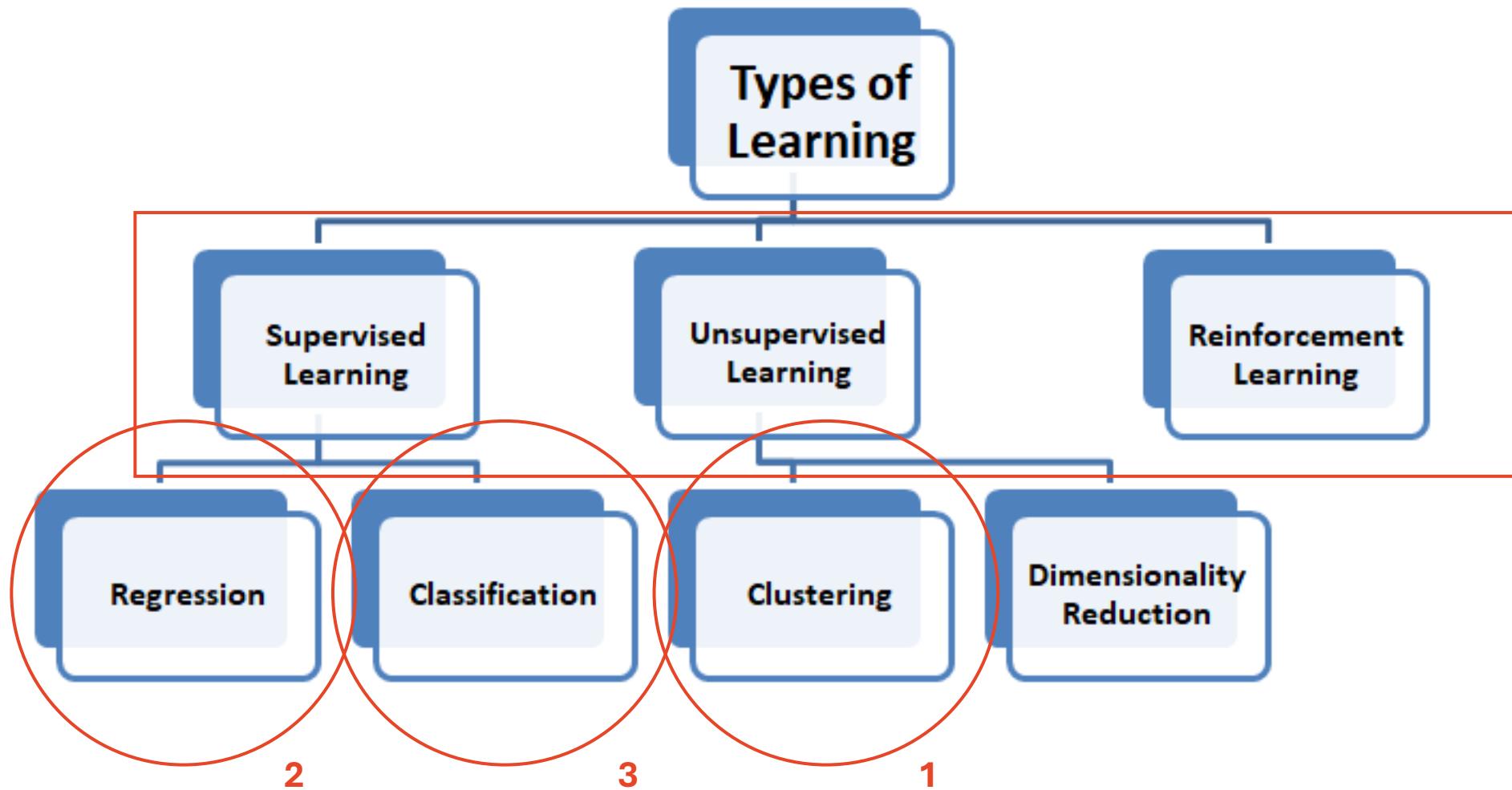
Machine Learning



Machine Learning



Machine Learning



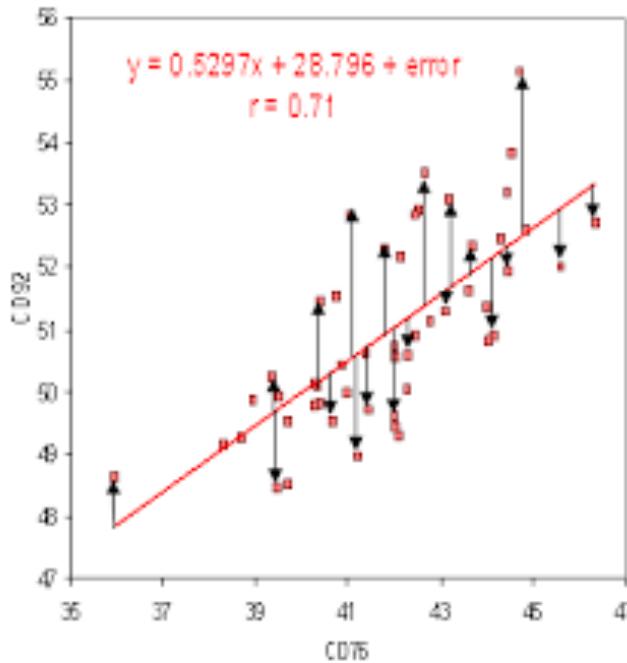
Model Training

Data

Loss Function

Optimization

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



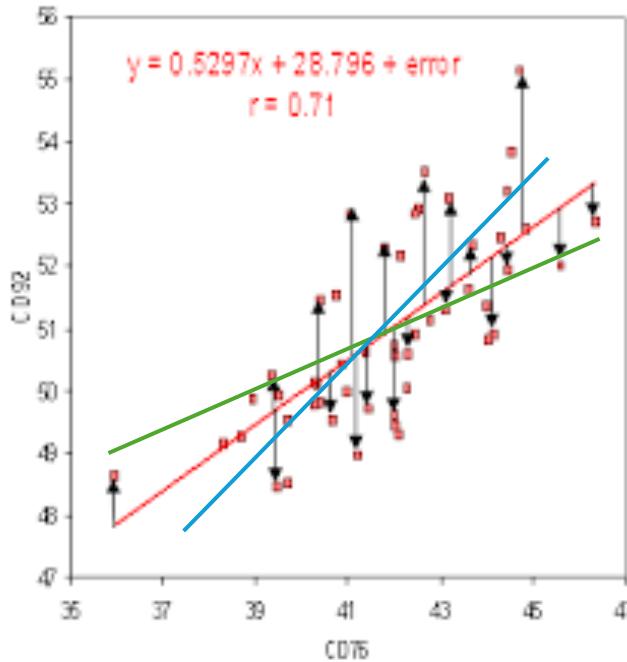
Model Training

Data

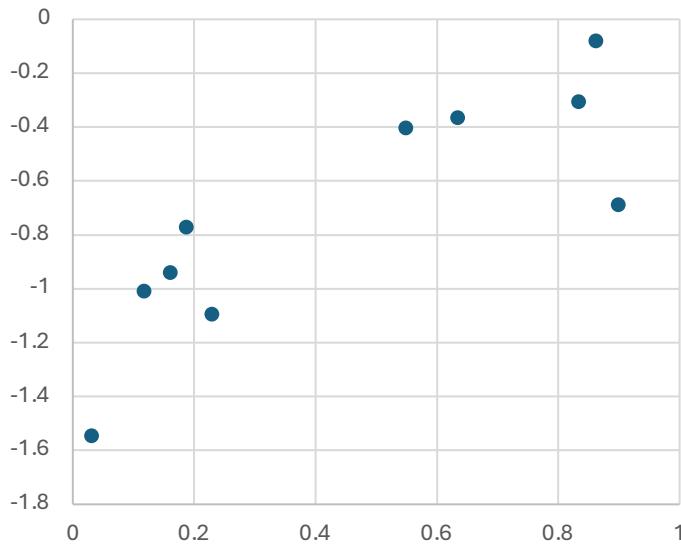
Loss Function

Optimization

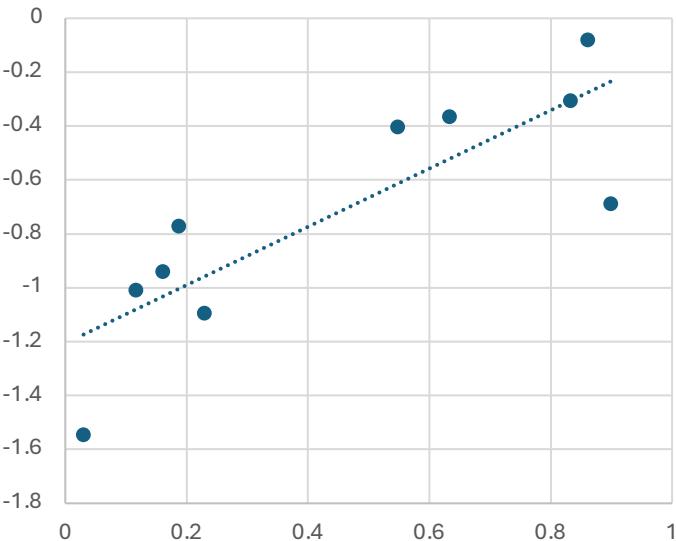
Model training uses data, a loss function, and an optimization method (e.g., gradient descent)



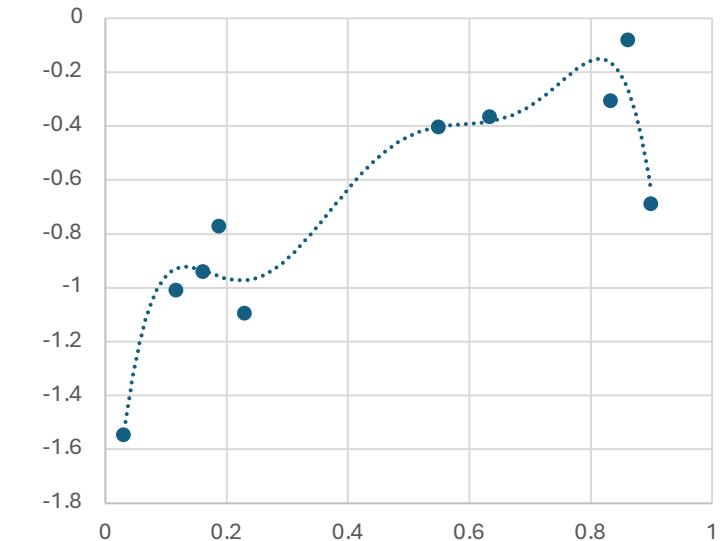
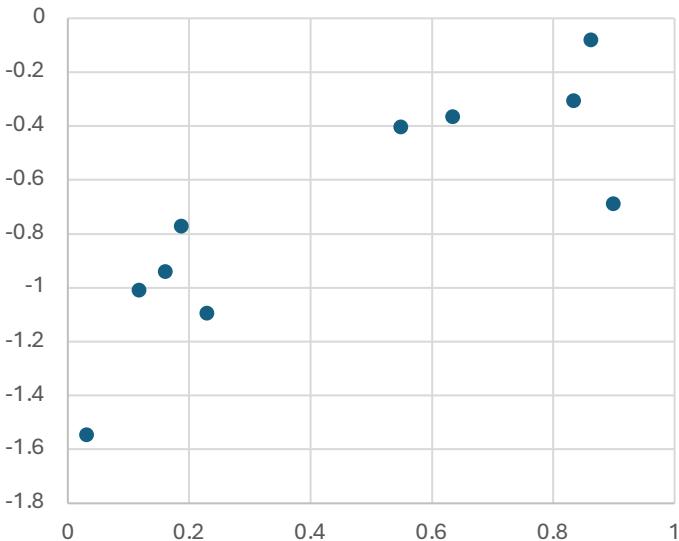
Model Training



Model Training

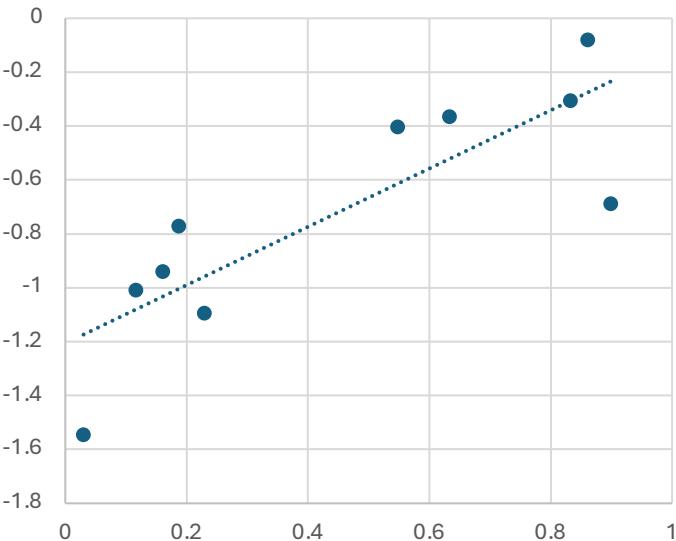


Underfitting

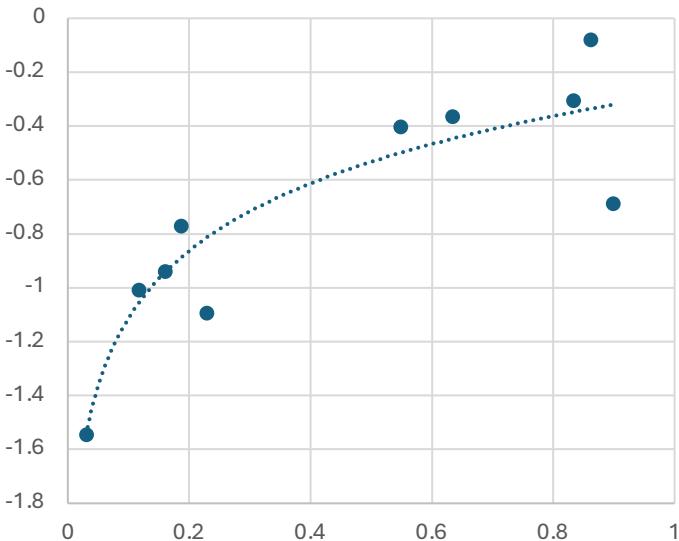


Overfitting

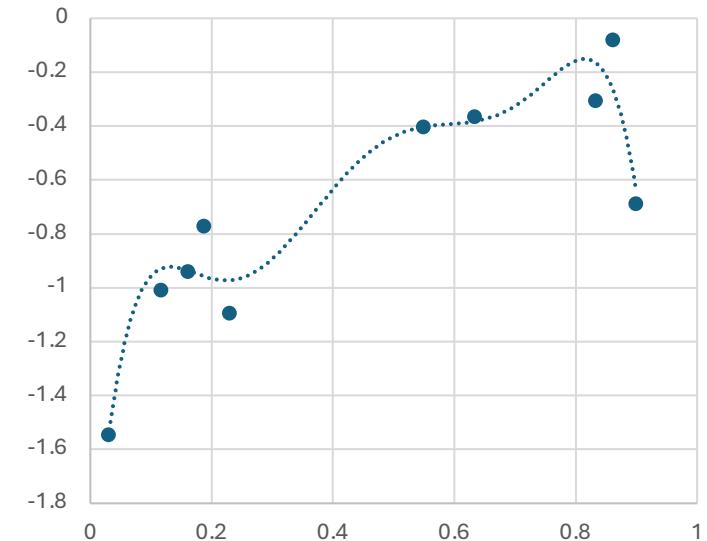
Model Training



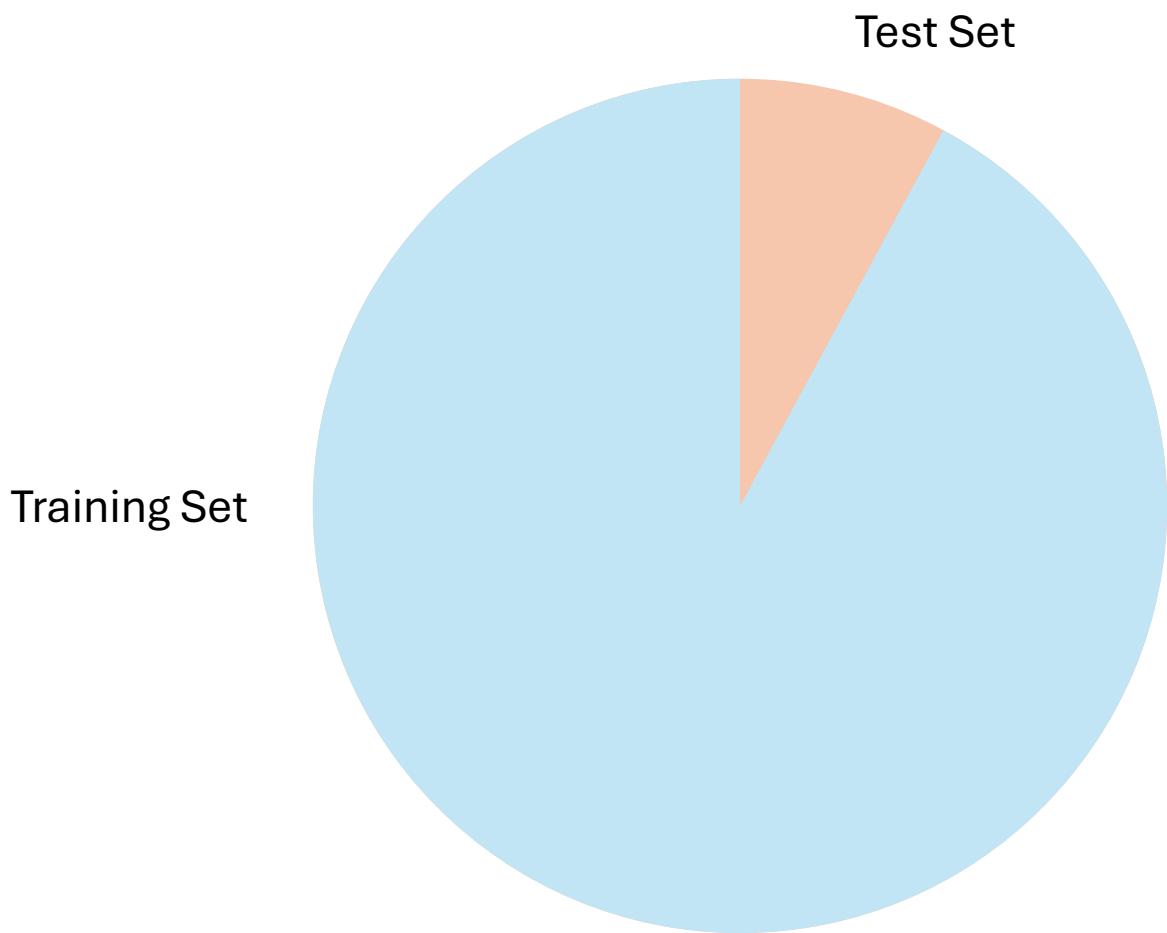
Underfitting



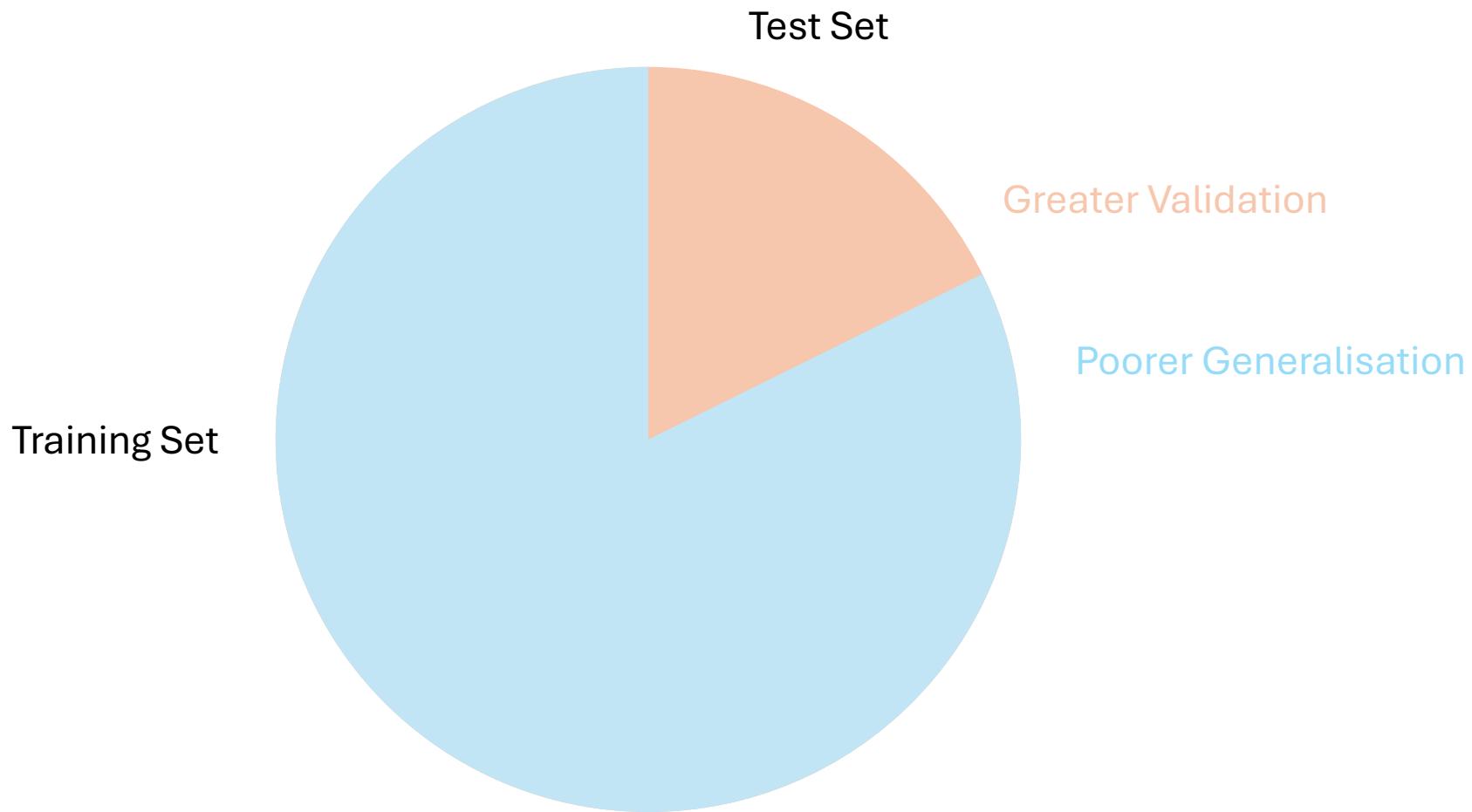
Overfitting



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 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.

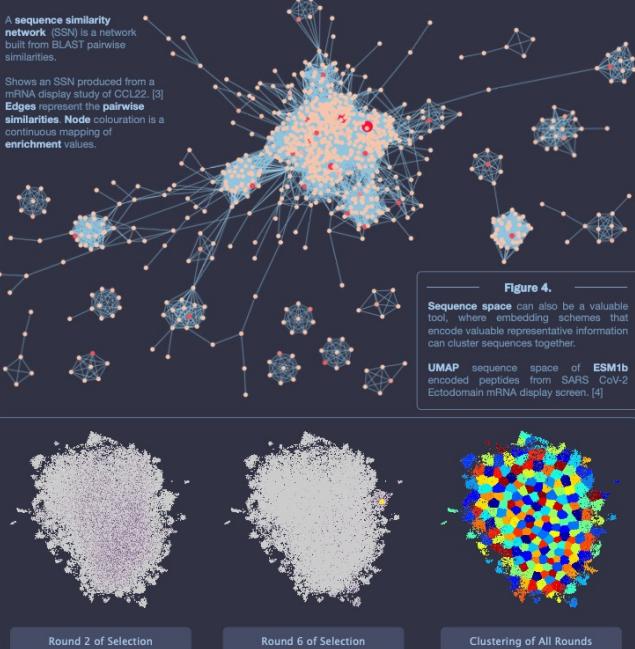


Figure 5.

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.

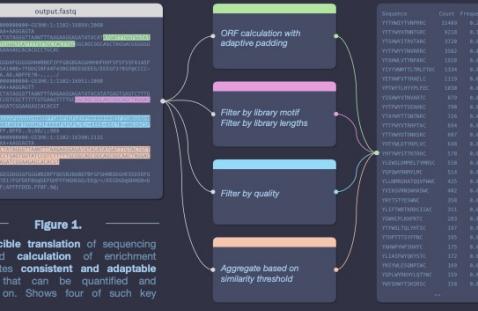


Figure 1.

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

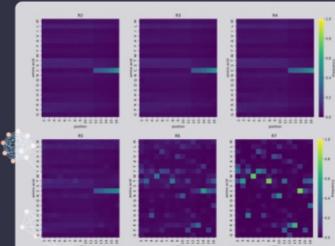


Figure 2.

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



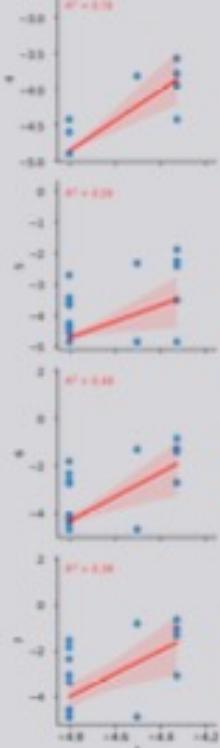
1

output.cdhit.clstr

Cluster	Sequence	Score
1	17aa, >YIYYYMMITLIMLFDC...	4
1	17aa, >YTTLTMLFLDC... at 98.9%	
1	17aa, >YIYYYMMITLIMLFDC...	at 94.1%
1	17aa, >YIYYYMMITLIMLFDC...	at 94.1%
1	17aa, >YIYYYMMITLIMLFDC...	at 94.1%
1	17aa, >YHMTILIMLFDC...	at 84.6%
1	17aa, >YYYYYMMITLIMLFDC...	at 94.1%
1	17aa, >YFCYMMITLIMLFDC...	at 98.2%
1	16aa, >YIYYYMMITLIMLFDC...	at 87.5%
1	17aa, >YIYYYMMITLIMLFDC...	at 94.1%
1	17aa, >YIYYYMMITLIMLFDC...	at 92.3%

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



The screenshot shows a software interface for peptide sequence analysis. It includes a terminal window with command-line input and output, a list of peptides grouped into clusters, and a heatmap visualization. A red box highlights the "Other Applications:" section.

Other Applications:
Customer segmentation, anomaly detection, pattern discovery

Cluster	Peptide Sequence	Similarity (%)
1	>VIVVVNMITLNLFTDC...	98.9%
2	>VIVVVNMITLNLFTDC...	94.12%
3	>VIVVVNMITLNLFTDC...	94.12%
4	>VIVVVNMITLNLFTDC...	94.12%
5	>VIVVVNMITLNLFTDC...	94.12%
6	>VIVVVNMITLNLFTDC...	94.12%
7	>VIVVVNMITLNLFTDC...	94.12%
8	>VIVVVNMITLNLFTDC...	94.12%
9	>VIVVVNMITLNLFTDC...	94.12%
10	>VIVVVNMITLNLFTDC...	94.12%
11	>VIVVVNMITLNLFTDC...	94.12%
12	>VIVVVNMITLNLFTDC...	94.12%
13	>VIVVVNMITLNLFTDC...	94.12%
14	>VIVVVNMITLNLFTDC...	94.12%
15	>VIVVVNMITLNLFTDC...	94.12%
16	>VIVVVNMITLNLFTDC...	94.12%
17	>VIVVVNMITLNLFTDC...	94.12%
18	>VIVVVNMITLNLFTDC...	94.12%
19	>VIVVVNMITLNLFTDC...	94.12%
20	>VIVVVNMITLNLFTDC...	94.12%
21	>VIVVVNMITLNLFTDC...	94.12%
22	>VIVVVNMITLNLFTDC...	94.12%
23	>VIVVVNMITLNLFTDC...	94.12%
24	>VIVVVNMITLNLFTDC...	94.12%
25	>VIVVVNMITLNLFTDC...	94.12%
26	>VIVVVNMITLNLFTDC...	94.12%
27	>VIVVVNMITLNLFTDC...	94.12%
28	>VIVVVNMITLNLFTDC...	94.12%
29	>VIVVVNMITLNLFTDC...	94.12%
30	>VIVVVNMITLNLFTDC...	94.12%
31	>VIVVVNMITLNLFTDC...	94.12%
32	>VIVVVNMITLNLFTDC...	94.12%
33	>VIVVVNMITLNLFTDC...	94.12%
34	>VIVVVNMITLNLFTDC...	94.12%
35	>VIVVVNMITLNLFTDC...	94.12%
36	>VIVVVNMITLNLFTDC...	94.12%
37	>VIVVVNMITLNLFTDC...	94.12%
38	>VIVVVNMITLNLFTDC...	94.12%
39	>VIVVVNMITLNLFTDC...	94.12%
40	>VIVVVNMITLNLFTDC...	94.12%
41	>VIVVVNMITLNLFTDC...	94.12%
42	>VIVVVNMITLNLFTDC...	94.12%
43	>VIVVVNMITLNLFTDC...	94.12%
44	>VIVVVNMITLNLFTDC...	94.12%
45	>VIVVVNMITLNLFTDC...	94.12%
46	>VIVVVNMITLNLFTDC...	94.12%
47	>VIVVVNMITLNLFTDC...	94.12%
48	>VIVVVNMITLNLFTDC...	94.12%
49	>VIVVVNMITLNLFTDC...	94.12%
50	>VIVVVNMITLNLFTDC...	94.12%
51	>VIVVVNMITLNLFTDC...	94.12%
52	>VIVVVNMITLNLFTDC...	94.12%
53	>VIVVVNMITLNLFTDC...	94.12%
54	>VIVVVNMITLNLFTDC...	94.12%
55	>VIVVVNMITLNLFTDC...	94.12%
56	>VIVVVNMITLNLFTDC...	94.12%
57	>VIVVVNMITLNLFTDC...	94.12%
58	>VIVVVNMITLNLFTDC...	94.12%
59	>VIVVVNMITLNLFTDC...	94.12%
60	>VIVVVNMITLNLFTDC...	94.12%
61	>VIVVVNMITLNLFTDC...	94.12%
62	>VIVVVNMITLNLFTDC...	94.12%
63	>VIVVVNMITLNLFTDC...	94.12%
64	>VIVVVNMITLNLFTDC...	94.12%
65	>VIVVVNMITLNLFTDC...	94.12%
66	>VIVVVNMITLNLFTDC...	94.12%
67	>VIVVVNMITLNLFTDC...	94.12%
68	>VIVVVNMITLNLFTDC...	94.12%
69	>VIVVVNMITLNLFTDC...	94.12%
70	>VIVVVNMITLNLFTDC...	94.12%
71	>VIVVVNMITLNLFTDC...	94.12%
72	>VIVVVNMITLNLFTDC...	94.12%
73	>VIVVVNMITLNLFTDC...	94.12%
74	>VIVVVNMITLNLFTDC...	94.12%
75	>VIVVVNMITLNLFTDC...	94.12%
76	>VIVVVNMITLNLFTDC...	94.12%
77	>VIVVVNMITLNLFTDC...	94.12%
78	>VIVVVNMITLNLFTDC...	94.12%
79	>VIVVVNMITLNLFTDC...	94.12%
80	>VIVVVNMITLNLFTDC...	94.12%
81	>VIVVVNMITLNLFTDC...	94.12%
82	>VIVVVNMITLNLFTDC...	94.12%
83	>VIVVVNMITLNLFTDC...	94.12%
84	>VIVVVNMITLNLFTDC...	94.12%
85	>VIVVVNMITLNLFTDC...	94.12%
86	>VIVVVNMITLNLFTDC...	94.12%
87	>VIVVVNMITLNLFTDC...	94.12%
88	>VIVVVNMITLNLFTDC...	94.12%
89	>VIVVVNMITLNLFTDC...	94.12%
90	>VIVVVNMITLNLFTDC...	94.12%
91	>VIVVVNMITLNLFTDC...	94.12%
92	>VIVVVNMITLNLFTDC...	94.12%
93	>VIVVVNMITLNLFTDC...	94.12%
94	>VIVVVNMITLNLFTDC...	94.12%
95	>VIVVVNMITLNLFTDC...	94.12%
96	>VIVVVNMITLNLFTDC...	94.12%
97	>VIVVVNMITLNLFTDC...	94.12%
98	>VIVVVNMITLNLFTDC...	94.12%
99	>VIVVVNMITLNLFTDC...	94.12%
100	>VIVVVNMITLNLFTDC...	94.12%
101	>VIVVVNMITLNLFTDC...	94.12%
102	>VIVVVNMITLNLFTDC...	94.12%
103	>VIVVVNMITLNLFTDC...	94.12%
104	>VIVVVNMITLNLFTDC...	94.12%
105	>VIVVVNMITLNLFTDC...	94.12%
106	>VIVVVNMITLNLFTDC...	94.12%
107	>VIVVVNMITLNLFTDC...	94.12%
108	>VIVVVNMITLNLFTDC...	94.12%
109	>VIVVVNMITLNLFTDC...	94.12%
110	>VIVVVNMITLNLFTDC...	94.12%
111	>VIVVVNMITLNLFTDC...	94.12%
112	>VIVVVNMITLNLFTDC...	94.12%
113	>VIVVVNMITLNLFTDC...	94.12%
114	>VIVVVNMITLNLFTDC...	94.12%
115	>VIVVVNMITLNLFTDC...	94.12%
116	>VIVVVNMITLNLFTDC...	94.12%
117	>VIVVVNMITLNLFTDC...	94.12%
118	>VIVVVNMITLNLFTDC...	94.12%
119	>VIVVVNMITLNLFTDC...	94.12%
120	>VIVVVNMITLNLFTDC...	94.12%
121	>VIVVVNMITLNLFTDC...	94.12%
122	>VIVVVNMITLNLFTDC...	94.12%
123	>VIVVVNMITLNLFTDC...	94.12%
124	>VIVVVNMITLNLFTDC...	94.12%
125	>VIVVVNMITLNLFTDC...	94.12%
126	>VIVVVNMITLNLFTDC...	94.12%
127	>VIVVVNMITLNLFTDC...	94.12%
128	>VIVVVNMITLNLFTDC...	94.12%
129	>VIVVVNMITLNLFTDC...	94.12%
130	>VIVVVNMITLNLFTDC...	94.12%
131	>VIVVVNMITLNLFTDC...	94.12%
132	>VIVVVNMITLNLFTDC...	94.12%
133	>VIVVVNMITLNLFTDC...	94.12%
134	>VIVVVNMITLNLFTDC...	94.12%
135	>VIVVVNMITLNLFTDC...	94.12%
136	>VIVVVNMITLNLFTDC...	94.12%
137	>VIVVVNMITLNLFTDC...	94.12%
138	>VIVVVNMITLNLFTDC...	94.12%
139	>VIVVVNMITLNLFTDC...	94.12%
140	>VIVVVNMITLNLFTDC...	94.12%
141	>VIVVVNMITLNLFTDC...	94.12%
142	>VIVVVNMITLNLFTDC...	94.12%
143	>VIVVVNMITLNLFTDC...	94.12%
144	>VIVVVNMITLNLFTDC...	94.12%
145	>VIVVVNMITLNLFTDC...	94.12%
146	>VIVVVNMITLNLFTDC...	94.12%
147	>VIVVVNMITLNLFTDC...	94.12%
148	>VIVVVNMITLNLFTDC...	94.12%
149	>VIVVVNMITLNLFTDC...	94.12%
150	>VIVVVNMITLNLFTDC...	94.12%
151	>VIVVVNMITLNLFTDC...	94.12%
152	>VIVVVNMITLNLFTDC...	94.12%
153	>VIVVVNMITLNLFTDC...	94.12%
154	>VIVVVNMITLNLFTDC...	94.12%
155	>VIVVVNMITLNLFTDC...	94.12%
156	>VIVVVNMITLNLFTDC...	94.12%
157	>VIVVVNMITLNLFTDC...	94.12%
158	>VIVVVNMITLNLFTDC...	94.12%
159	>VIVVVNMITLNLFTDC...	94.12%
160	>VIVVVNMITLNLFTDC...	94.12%
161	>VIVVVNMITLNLFTDC...	94.12%
162	>VIVVVNMITLNLFTDC...	94.12%
163	>VIVVVNMITLNLFTDC...	94.12%
164	>VIVVVNMITLNLFTDC...	94.12%
165	>VIVVVNMITLNLFTDC...	94.12%
166	>VIVVVNMITLNLFTDC...	94.12%
167	>VIVVVNMITLNLFTDC...	94.12%
168	>VIVVVNMITLNLFTDC...	94.12%
169	>VIVVVNMITLNLFTDC...	94.12%
170	>VIVVVNMITLNLFTDC...	94.12%
171	>VIVVVNMITLNLFTDC...	94.12%
172	>VIVVVNMITLNLFTDC...	94.12%
173	>VIVVVNMITLNLFTDC...	94.12%
174	>VIVVVNMITLNLFTDC...	94.12%
175	>VIVVVNMITLNLFTDC...	94.12%
176	>VIVVVNMITLNLFTDC...	94.12%
177	>VIVVVNMITLNLFTDC...	94.12%
178	>VIVVVNMITLNLFTDC...	94.12%
179	>VIVVVNMITLNLFTDC...	94.12%
180	>VIVVVNMITLNLFTDC...	94.12%
181	>VIVVVNMITLNLFTDC...	94.12%
182	>VIVVVNMITLNLFTDC...	94.12%
183	>VIVVVNMITLNLFTDC...	94.12%
184	>VIVVVNMITLNLFTDC...	94.12%
185	>VIVVVNMITLNLFTDC...	94.12%
186	>VIVVVNMITLNLFTDC...	94.12%
187	>VIVVVNMITLNLFTDC...	94.12%
188	>VIVVVNMITLNLFTDC...	94.12%
189	>VIVVVNMITLNLFTDC...	94.12%
190	>VIVVVNMITLNLFTDC...	94.12%
191	>VIVVVNMITLNLFTDC...	94.12%
192	>VIVVVNMITLNLFTDC...	94.12%
193	>VIVVVNMITLNLFTDC...	94.12%
194	>VIVVVNMITLNLFTDC...	94.12%
195	>VIVVVNMITLNLFTDC...	94.12%
196	>VIVVVNMITLNLFTDC...	94.12%
197	>VIVVVNMITLNLFTDC...	94.12%
198	>VIVVVNMITLNLFTDC...	94.12%
199	>VIVVVNMITLNLFTDC...	94.12%
200	>VIVVVNMITLNLFTDC...	94.12%

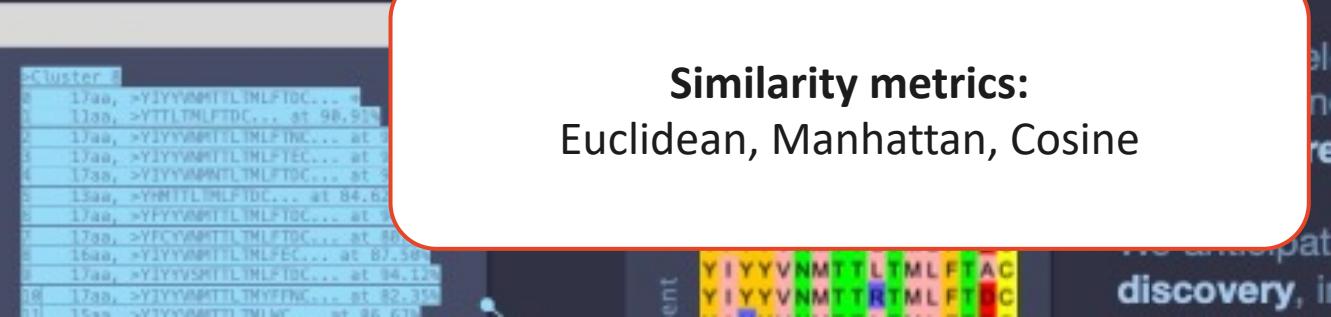
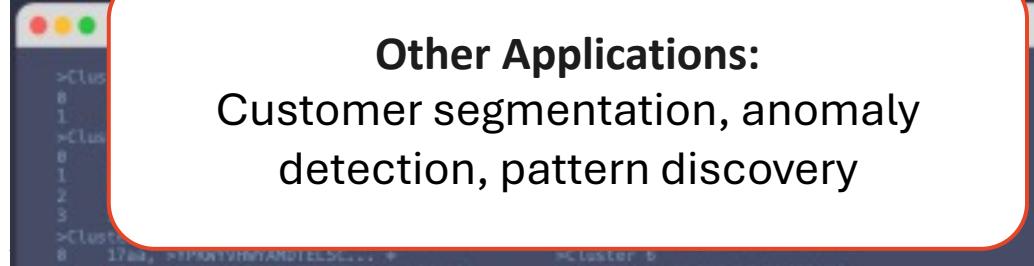
encode valuable representative information
can cluster sequences together.

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

We anticipate the method provides and
and interpretation for discovery, imp...

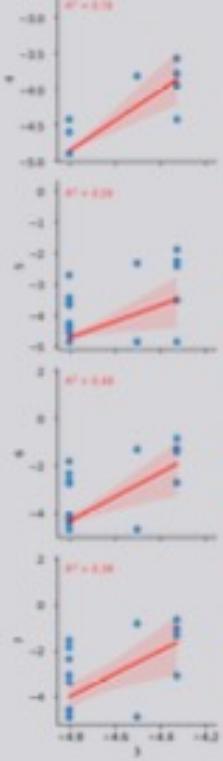
Our approach groups that increasing the

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

Clustering

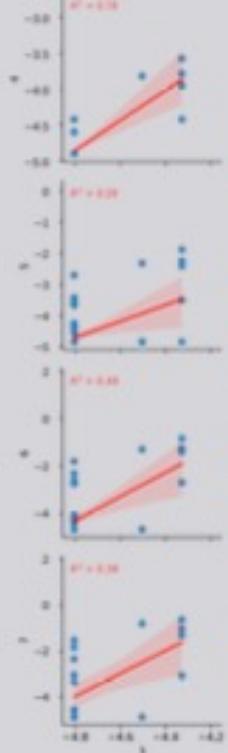


Other Applications:
Customer segmentation, anomaly detection, pattern discovery

```
>Clus  
8  
1  
>Clus  
8  
1  
>Clus  
8  
1  
Cluster 2  
17aa, >YIYVNMITLNLFTDC... at 98.9%  
11aa, >YTTLTNLFTDC... at 98.9%  
17aa, >VIYVNMITLNLFTDC... at 98.9%  
17aa, >VIYVNMITLNLFTDC... at 98.9%  
13aa, >YHITLNLFTDC... at 84.6%  
17aa, >YIYVNMITLNLFTDC... at 98.9%  
17aa, >YIYVNMITLNLFTDC... at 98.9%  
16aa, >YIYVNMITLNLFTDC... at 87.5%  
17aa, >YIYVNMITLNLFTDC... at 94.32%  
17aa, >YIYVNMITLNLFTDC... at 94.32%  
17aa, >YIYVNMITLNLFTDC... at 94.32%  
17aa, >YIYVNMITLNLFTDC... at 94.32%
```

encode valuable representative information
can cluster sequences together.

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



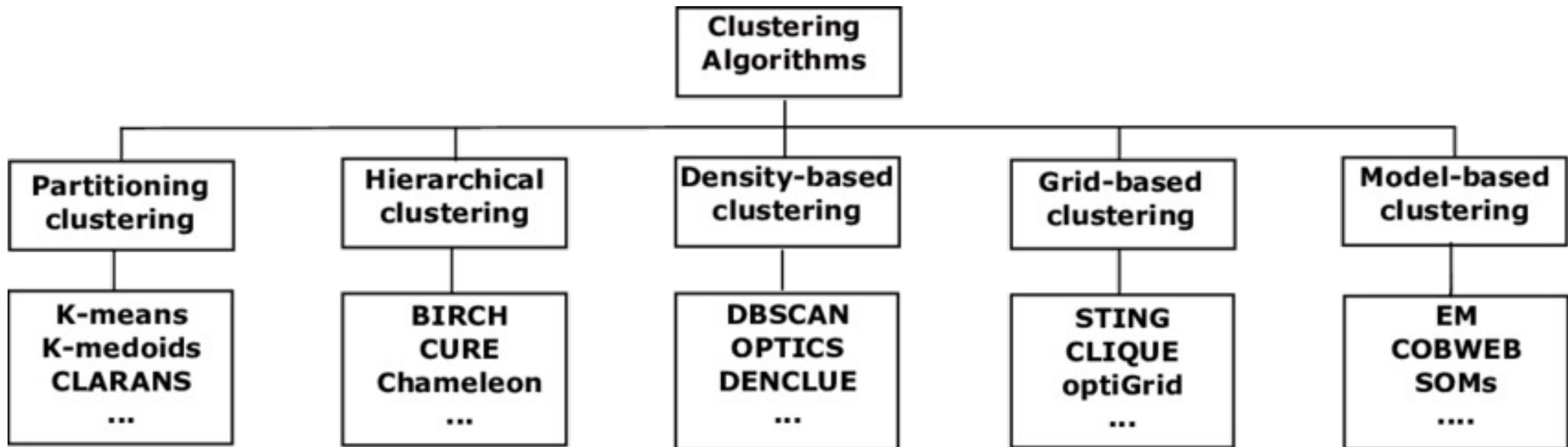
Our approach groups that increasing the

Normalisation crucial

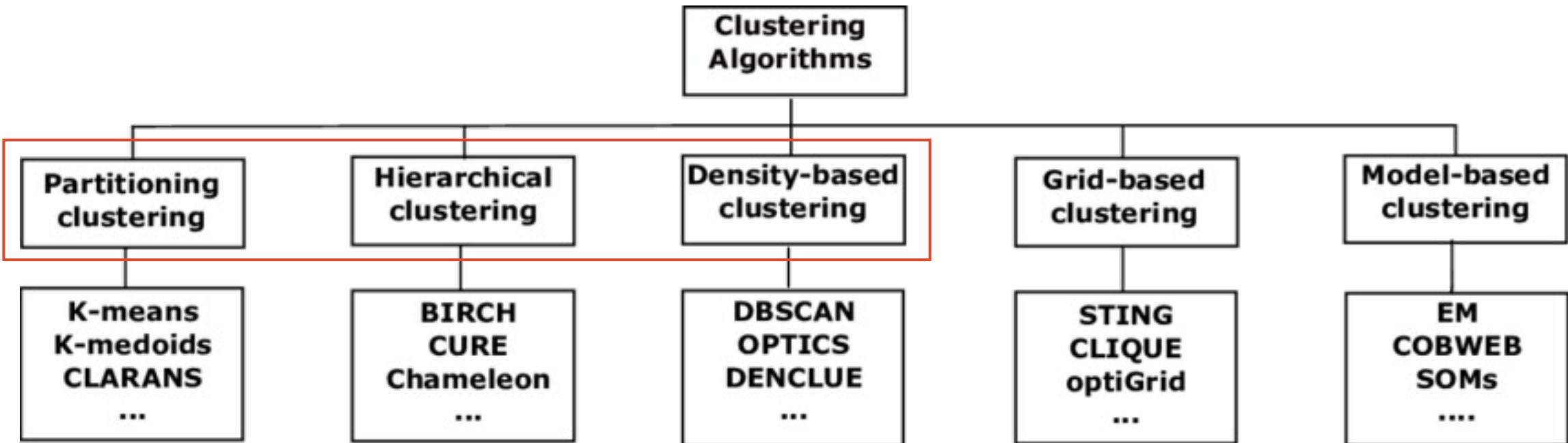
Similarity metrics:
Euclidean, Manhattan, Cosine

select
and
retai
discovery, imp

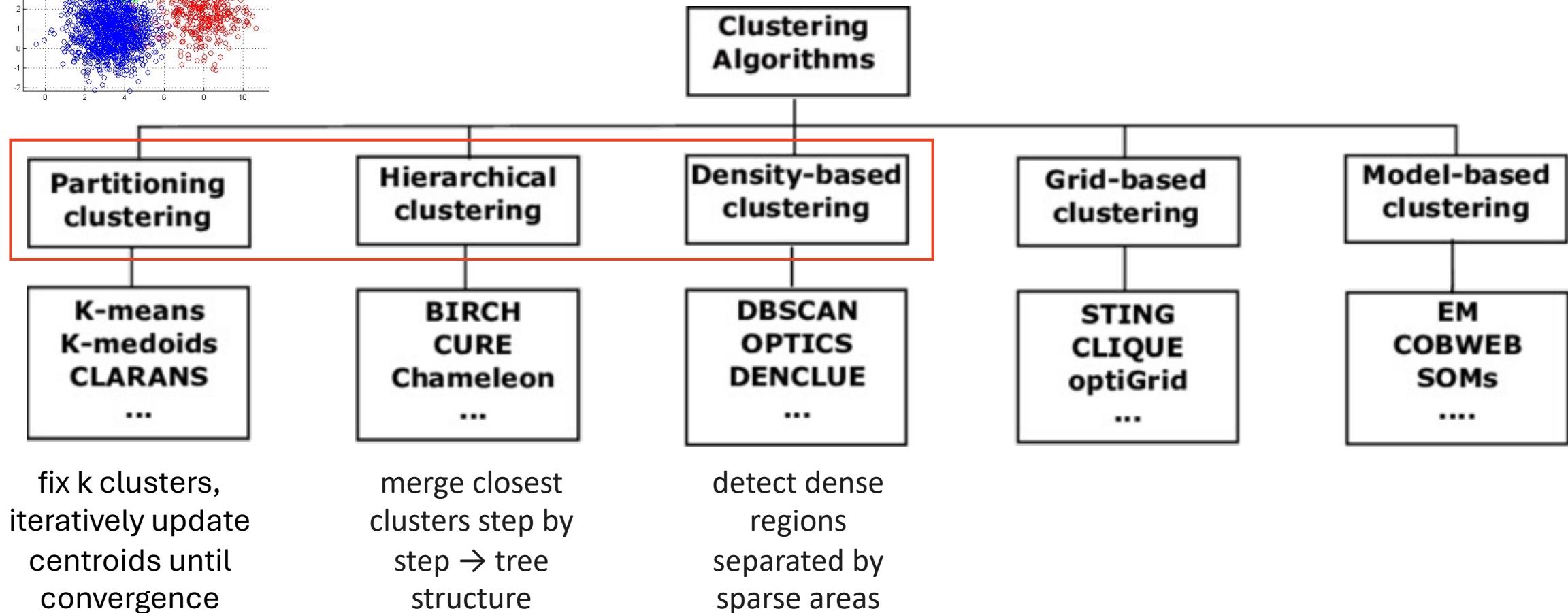
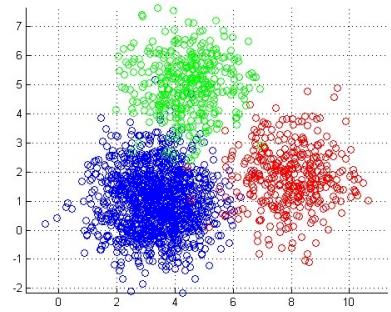
Types of Clustering



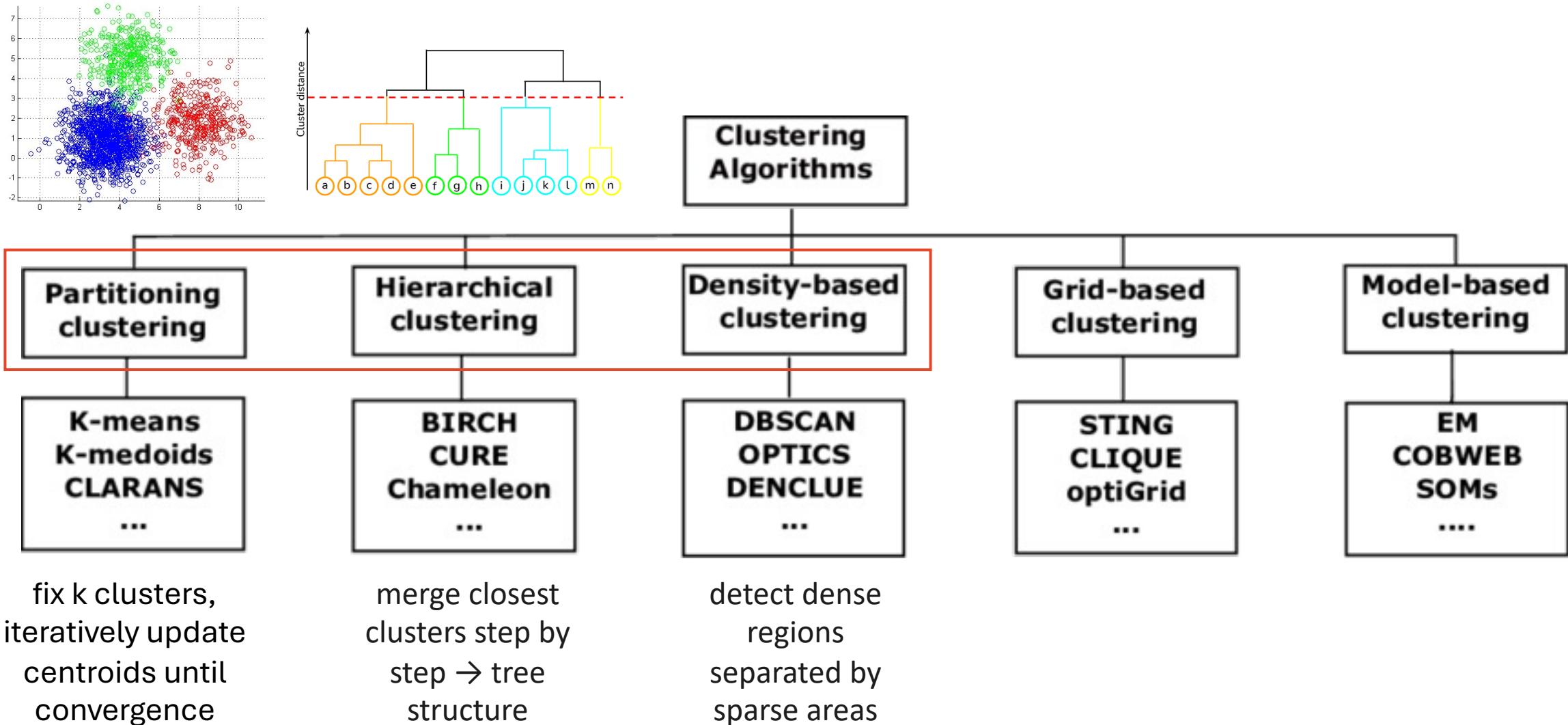
Types of Clustering



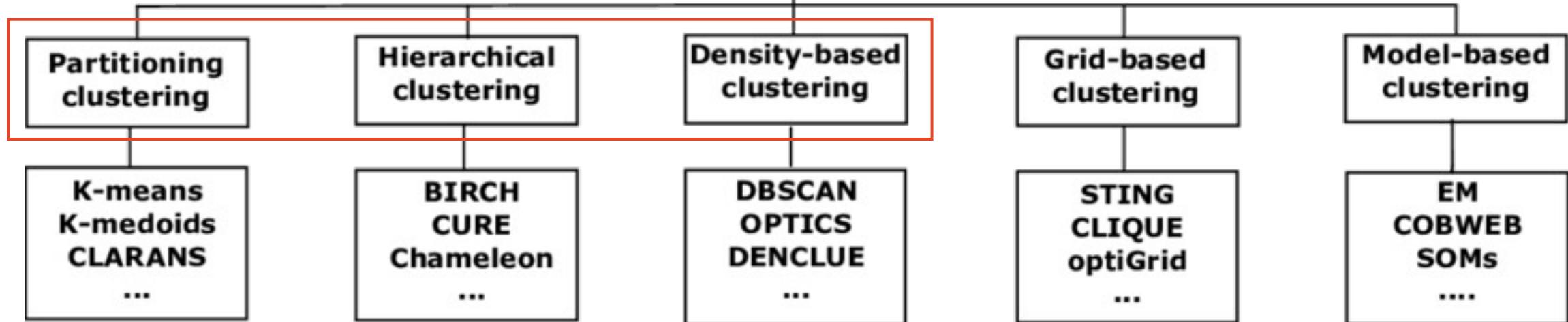
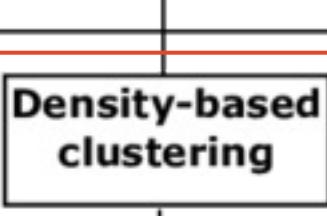
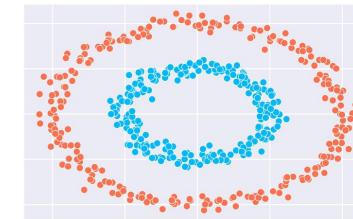
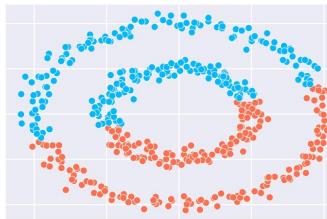
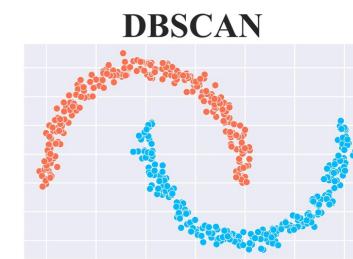
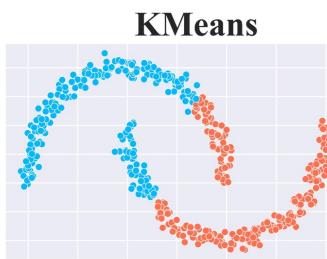
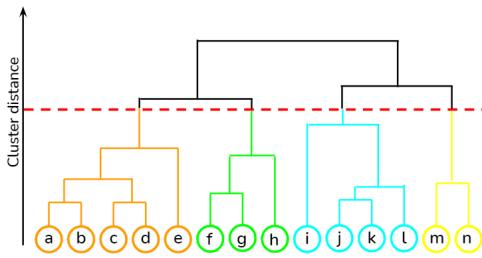
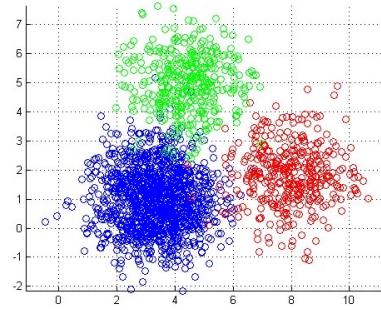
Types of Clustering



Types of Clustering



Types of Clustering

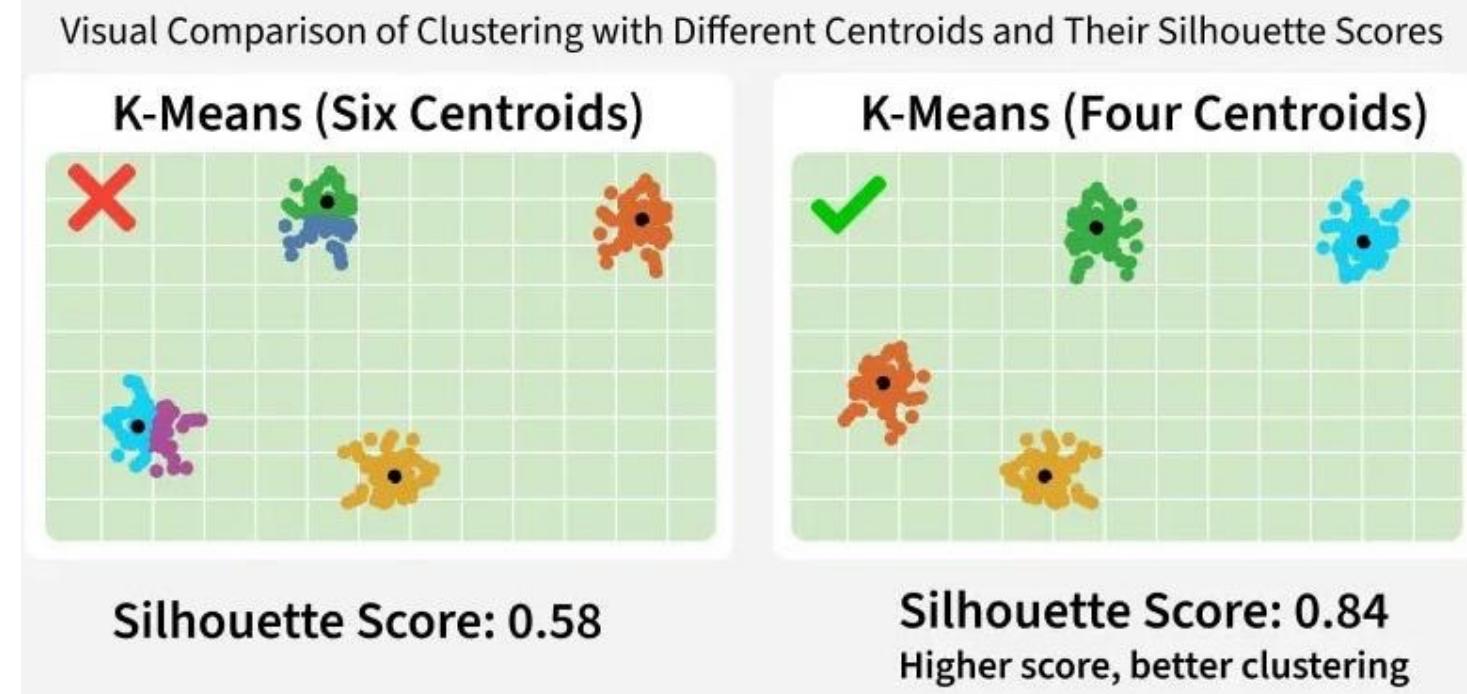


fix k clusters,
iteratively update
centroids until
convergence

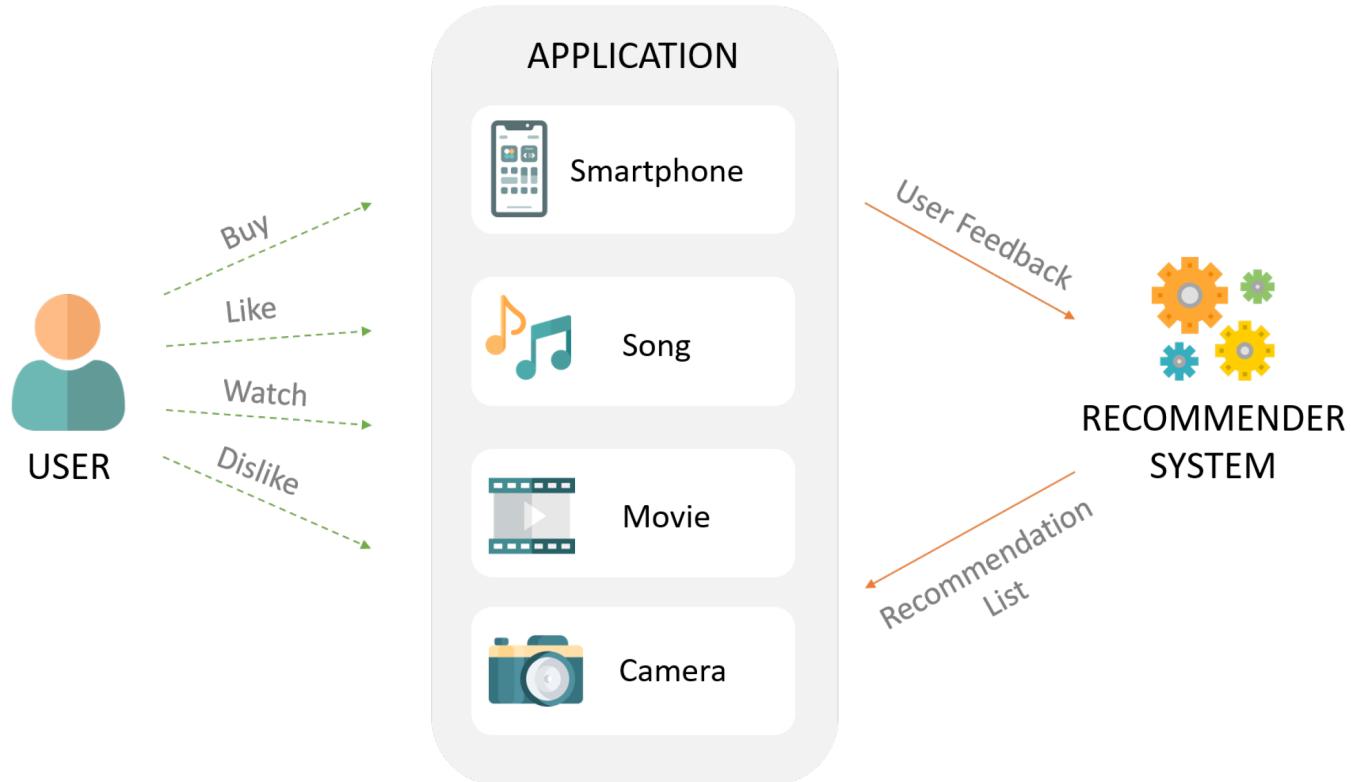
merge closest
clusters step by
step → tree
structure

detect dense
regions
separated by
sparse areas

Evaluation



Recommender



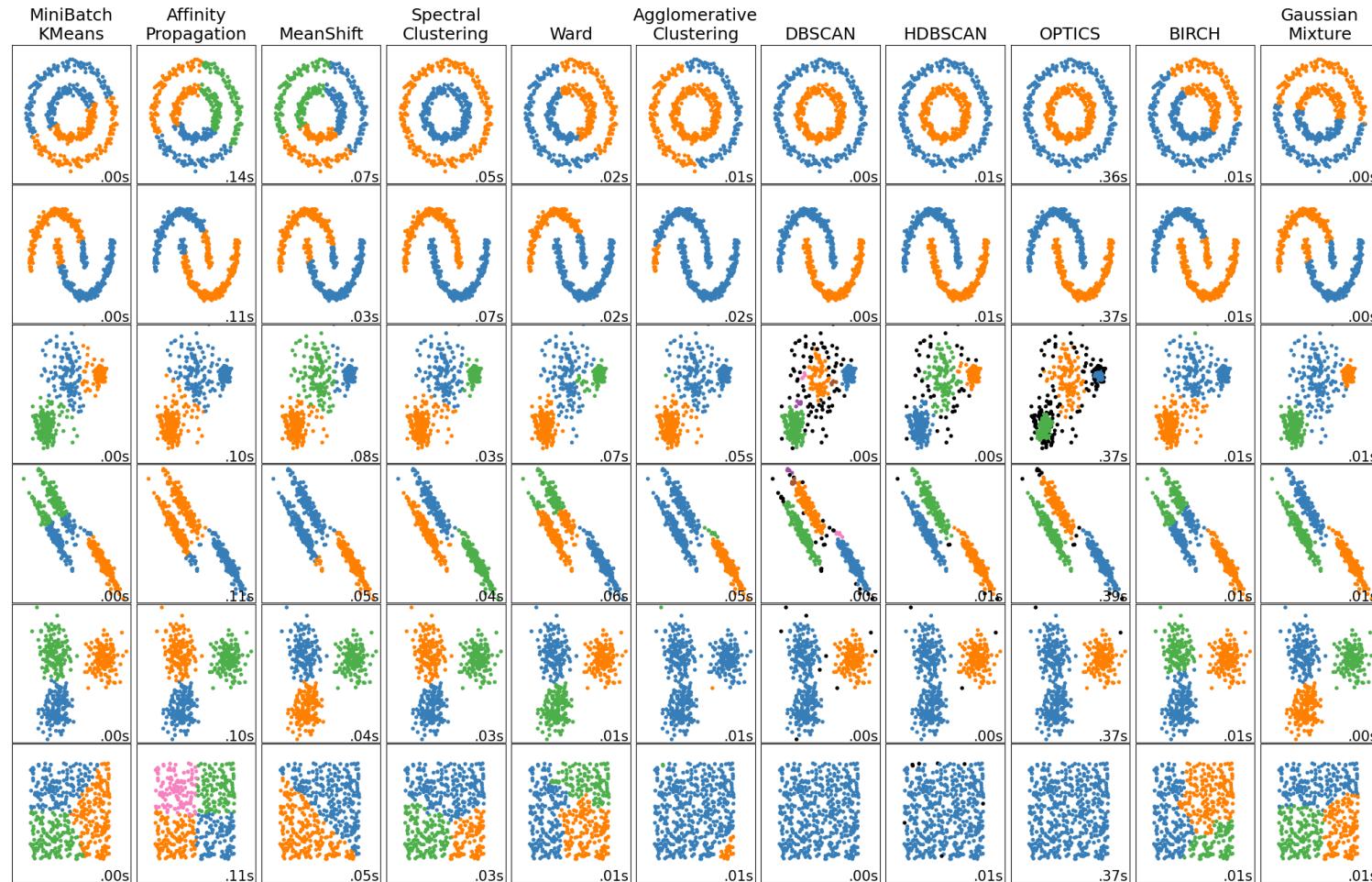
Python Demonstration

Clustering

Try running clustering algorithm

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

See `cluster.py` on the Ed Workspace



Content

Regression

Regression Model

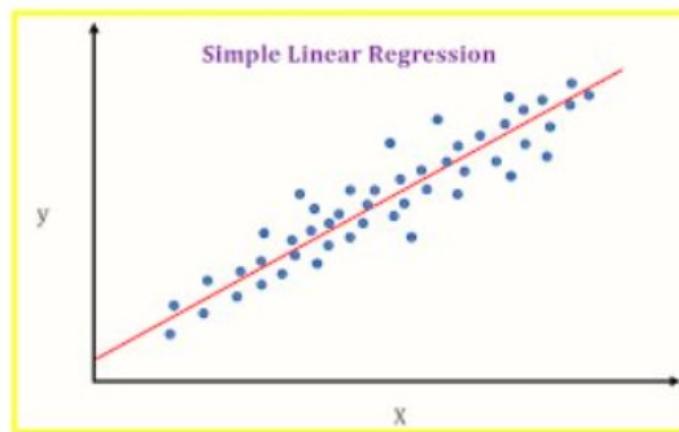
- It is a statistical technique used in machine learning, statistics, and data analysis to quantify the relationship between a dependent variable (also called the target) and one or more independent variables (predictors or features)
- Its primary purpose is to make predictions or estimate the value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

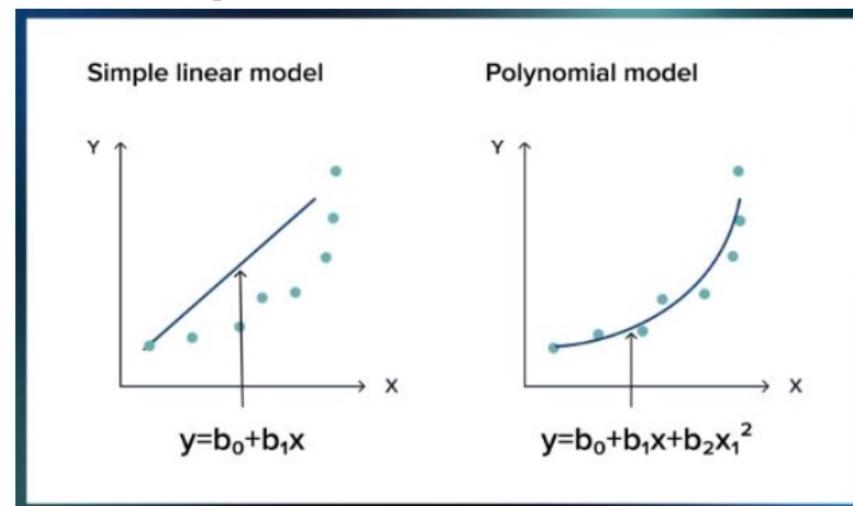


Linear Regression Model

- **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 2

Activity

The first main goal is to agree on the dataset and which attribute you plan to predict (each person using a different approach for producing a predictive model). First though, make sure that you have several ways to communicate with one another between labs.

In many cases, the dataset will be carried over from previous Stage, but even so, you need to consider which attribute you will be predicting. Notice that if you want to predict a nominal attribute, then you will all be using classification techniques; on the other hand, if you predict a quantitative attribute then you all must use various regression approaches.

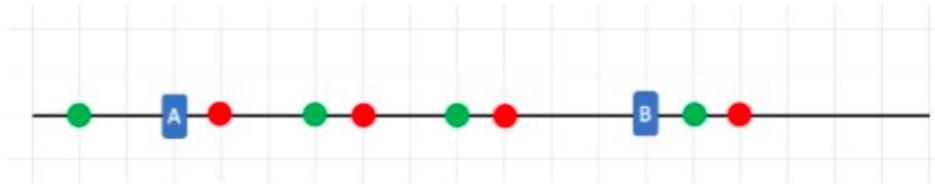
You should at this time divide the dataset into a training set, a validation set and a test set.

Now start choosing the ways you will each produce a predictive model. In the slides, we presented a few approaches for regression and classification (namely linear regression and kNN regression, and decision tree classifier and “logistic regression” classifier) but there are plenty of variants possible within these (for example, which distance function to use in kNN, how to encode any categorical attributes for linear regression, whether to transform any quantitative attributes by re-scaling them, which attributes to include in making the prediction, whether to apply regularization, how to tune any hyper-parameters, etc). You may also want to do some reading to find more ideas to try. But at least pick a few things to try initially.

Please make sure that you document your choice of dataset (including indication of where the original dataset, and the training and test subsets, are available), and which attribute to predict, by a post in the Canvas page of your group. Also post about the tasks people are given, and a due date for each task, to keep the group’s progress going.

Exam-Style Questions

- What would A and B be classified as (green or red) under different values of k in a k-NN algorithm?
- What will be your default colour if there is a tie?



k	A	B
1		
2		
3		
4		

Exam-Style Questions

- Build a decision tree from the following data using Gini Impurity Index.

outlook	temp.	humidity	windy	play
sunny	hot	high	false	no
sunny	hot	high	true	no
overcast	hot	high	false	yes
rainy	mild	high	false	yes
rainy	cool	normal	false	yes
rainy	cool	normal	true	no
overcast	cool	normal	true	yes
sunny	mild	high	false	no
sunny	cool	normal	false	yes
rainy	mild	normal	false	yes
sunny	mild	normal	true	yes
overcast	mild	high	true	yes
overcast	hot	normal	false	yes
rainy	mild	high	true	no

That's it folks!

Remaining Ed Lessons, Questions, Assignment etc.