

Meta-unsupervised-learning:

A model of unsupervised learning applicable to humans and computers

Vikas Garg – MIT



Adam Tauman Kalai – MSR



Human Computation Theories

H+M

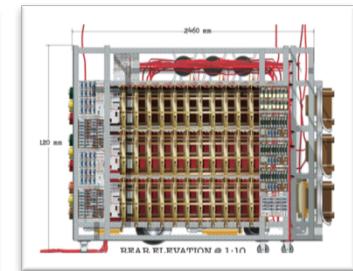
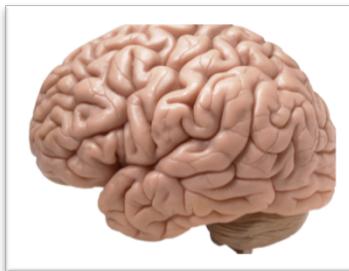
Models of computational systems
that **combine humans + machines**



Models that fit **human** computation
and machine computation

Blum-Vempala HUM (2015)

Turing Machine (1936)



Common characteristics

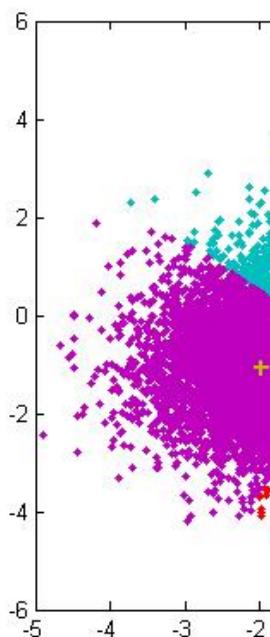
- Time-space-accuracy tradeoffs
- Easier to check than to solve
- Reductions, parallelization, ...

Unsupervised learning

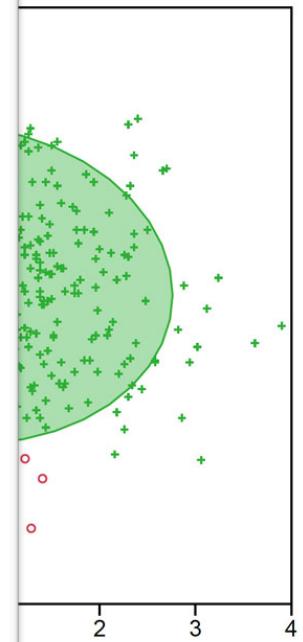
Finding structure in unlabeled data via:

- Clustering
- Dimensionality reduction
- Dictionary learning
- Topic modeling
- Manifold learning

Geometry of clustering



k -means



Gaussian mixture model

Human vs. Machine Clustering

The class was great!

Waste of time ADAM I HOPE YOU DIE.

The class was boring.

Best course I have taken in a long time.

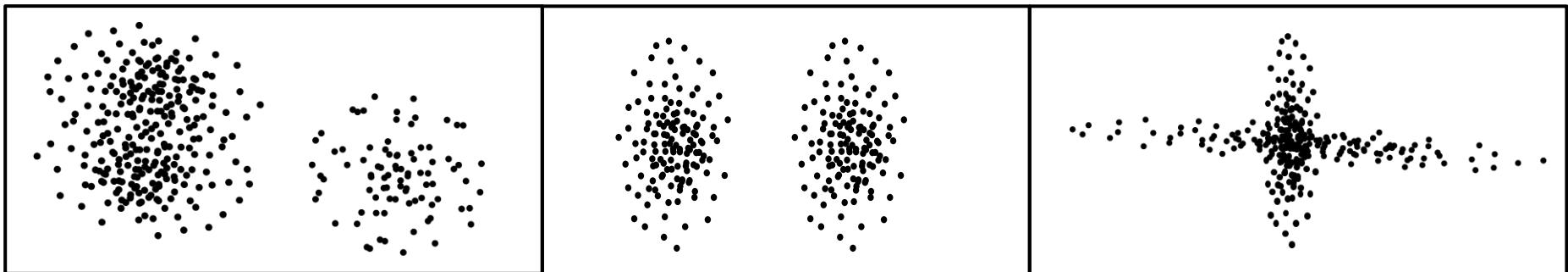
The class was terrible!

I slept the whole time.

The class was great!
The class was boring.
The class was terrible!

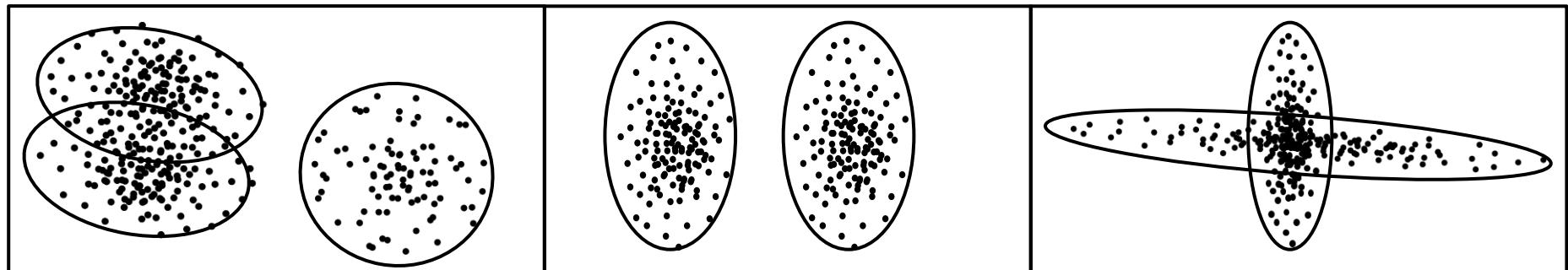
Best course I have taken in a long time.
Waste of time ADAM I HOPE YOU DIE.
I slept the whole time.

H&M idea: use prior data

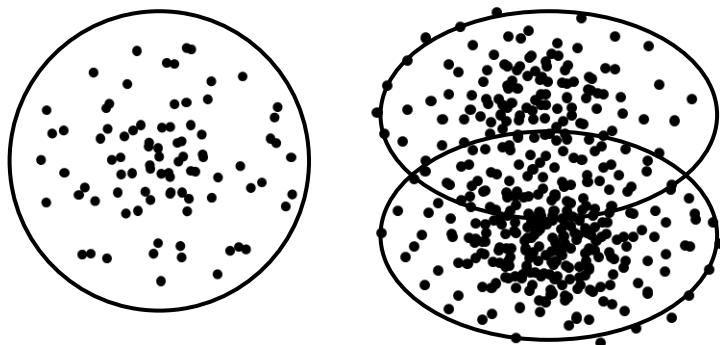


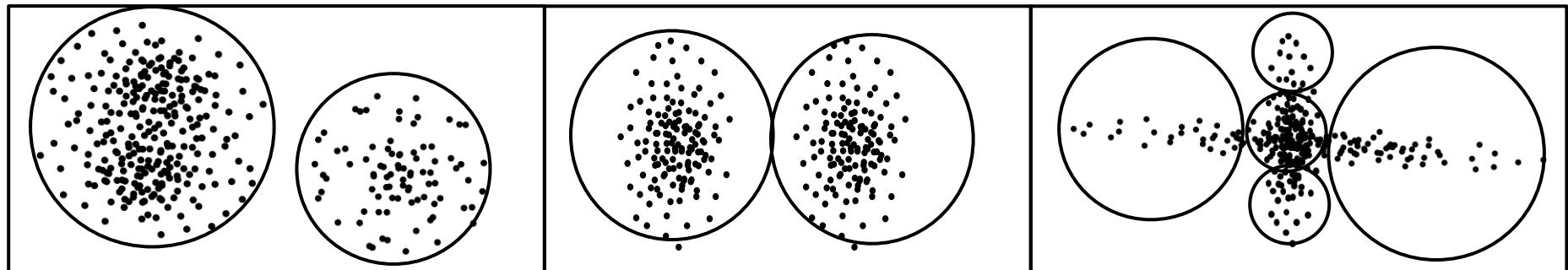
Clustering with prior labeled data



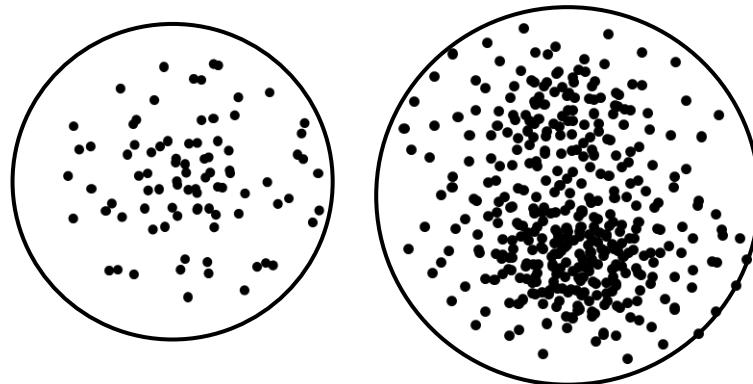


Clustering with prior labeled data





Clustering with prior labeled data



META-MODEL

Meta-clustering

#dim #data

- Unlabeled data $X \in U_{d \geq 1, m \geq 1} \uparrow \mathbb{R}^{d \times m}$
- Ground truth clustering $Y \in \Pi(X)$
- Predicted clustering $Z \in \Pi(X)$
- Accuracy $\text{acc}(Y, Z) \geq 0$
- Meta-dist. μ over labeled problems (X, Y)
- Training problems $T = (X \downarrow 1, Y \downarrow 1), \dots, (X \downarrow n, Y \downarrow n) \sim \mu^{\uparrow n}$
- Meta-algorithm $M(T)$ outputs clustering alg. C

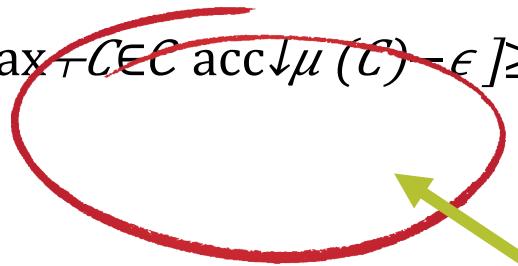
Def: Efficient meta-clustering

Fix a class of clustering algorithms c .

M efficiently meta-clusters as well as c if:

for any meta-dist. μ over clustering problems,
for any $\epsilon, \delta > 0$, $n \geq \text{poly}(1/\epsilon, 1/\delta)$,

$$\Pr_{T \sim \mu \uparrow n} [\text{acc}_{\downarrow \mu}(M(T)) \geq \max_{C \in \mathcal{C}} \text{acc}_{\downarrow \mu}(C) - \epsilon] \geq 1 - \delta$$



Also, $M, M(T)$ must be poly-time.

$$\mathbf{E}_{\tau}[\text{acc}(Y, C(X))]$$

$$(X, Y) \sim \mu$$

Meta-unsupervised learning

Similar to theory of supervised learning:

PAC (Valiant 84) Agnostic (Kearns, Schapire, Sellie 94)



Meta-clustering Occam bound

$C \downarrow \text{EMP} \in \mathcal{C}$ maximizes empirical accuracy on training data τ

$$\Pr_{\tau \sim \mu \uparrow n} [\text{acc} \downarrow \mu (C \downarrow \text{EMP}) \geq \max_{\tau \in \mathcal{C}} \text{acc} \downarrow \mu (C) + \sqrt{2/n \log |\mathcal{C}| / \delta}] \geq 1 - \delta$$



#bits

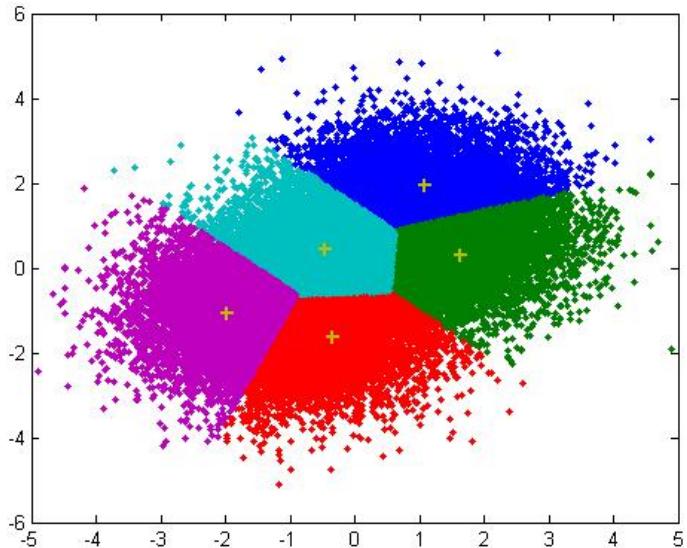
If M efficiently finds $C \downarrow \text{EMP}$,
then M efficiently meta-clusters as well as c
using $\#\text{problems} \approx \#\text{bits}$ to represent c

EXPERIMENTAL EVALUATION

Silhouette heuristic (intrinsic clustering quality)

Rand Index (accuracy metric w.r.t. ground-truth/gold-standard)

Silhouette clustering quality heuristic



How much closer, on average, are points to other points in their same cluster than to points in the nearest other cluster?
(normalized to be in $[-1, 1]$)

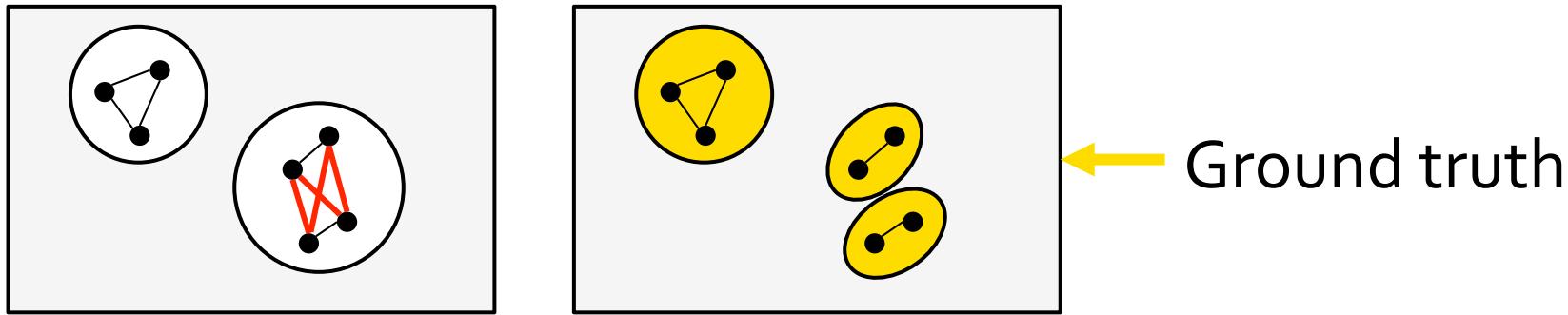
$$a = \text{Avg}_{x' \in C(x)}$$

$$b = \min_{x' \neq x} \text{Avg}_{x' \in C(x')} / |x - x'|$$

$$\text{Silh} = \frac{b(x) - a(x)}{\max\{a(x), b(x)\}}$$

Often used to pick #clusters k

Rand Index: ground truth *accuracy*



- Fraction of edge/non-edge agreement pairs
- Adjusted Rand Index (ARI)
 - Corrects for random chance
 - Common empirical clustering accuracy metric



International Federation of Classification Societies

Cluster Benchmark Data Repository

Search ...

- [Welcome to Repository](#)
- [Data Sets](#)
- [Philosophy](#)
- [Citation Policy](#)
- [Disclaimer](#)
- [Contribute a Data Set](#)
- [Licenses](#)
- [Contact](#)
- [IFCS Homepage](#)

Welcome to the IFCS Cluster Benchmark Data Repository

The aim of this Repository is to stimulate better practice in benchmarking (performance comparison of methods) for cluster analysis by providing a variety of well documented high quality datasets and simulation routines for use in practical benchmarking.

The repository collects datasets with and without given "true" clusterings. A particular feature of the repository is that every dataset comes with a comprehensive documentation, including information on the specific nature of the clustering problem in this dataset and the characteristics that useful clusters should fulfill, with scientific justification.

Currently < 10 datasets



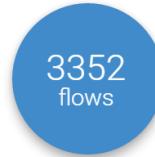
Exploring machine learning better, together



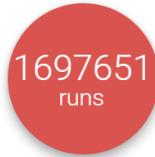
Find or add **data** to analyse



Download or create scientific
tasks



Find or add data analysis
flows



Upload and explore all **results**
online.

Classification data → Clustering data

x								y
1	2.2	2	1.3	1.7	4	-1		Cat
1	1.7	1	2.4	3.9	5	-1		Dog
0	3.9	4	17	2.8	-1	1		Cat
1	2.8	5	17	1.1	2	-1		Bird
0	1.1	-1	17	3.8	0	-1		Bird
1	4.5	2	17	1.1	8	1		Dog

Classification data → Clustering data

x								y
1	2.2	2	1.3	1.7	4	-1	Cat	
1	1.7	1	2.4	3.9	5	-1	Dog	
0	3.9	4	17	2.8	-1	1	Cat	
1	2.8	5	17	1.1	2	-1	Bird	
0	1.1	-1	17	3.8	0	-1	Bird	
1	4.5	2	17	1.1	8	1	Dog	

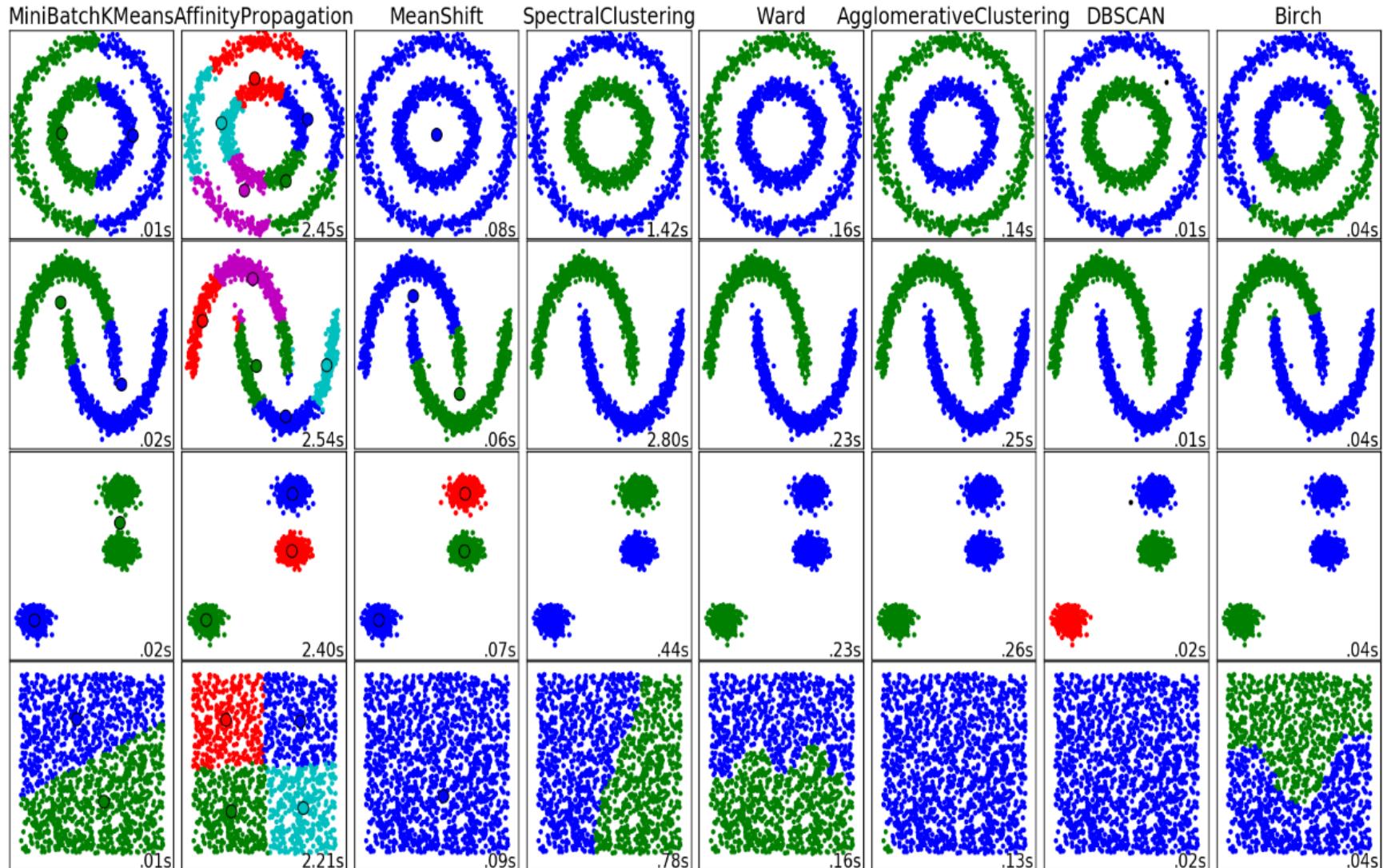
Cluster

Cluster

Cluster

WHICH ALGORITHM?

2.3.1. Overview of clustering methods



A comparison of the clustering algorithms in scikit-learn

Method name	Parameters	Scalability	Usecase	Geometry (metric used)
K-Means	number of clusters	Very large <code>n_samples</code> , medium <code>n_clusters</code> with MiniBatch code	General-purpose, even cluster size, flat geometry, not too many clusters	Distances between points
Affinity propagation	damping, sample preference	Not scalable with <code>n_samples</code>	Many clusters, uneven cluster size, non-flat geometry	Graph distance (e.g. nearest-neighbor graph)
Mean-shift	bandwidth	Not scalable with <code>n_samples</code>	Many clusters, uneven cluster size, non-flat geometry	Distances between points
Spectral clustering	number of clusters	Medium <code>n_samples</code> , small <code>n_clusters</code>	Few clusters, even cluster size, non-flat geometry	Graph distance (e.g. nearest-neighbor graph)
Ward hierarchical clustering	number of clusters	Large <code>n_samples</code> and <code>n_clusters</code>	Many clusters, possibly connectivity constraints	Distances between points
Agglomerative clustering	number of clusters, linkage type, distance	Large <code>n_samples</code> and <code>n_clusters</code>	Many clusters, possibly connectivity constraints, non Euclidean distances	Any pairwise distance
DBSCAN	neighborhood size	Very large <code>n_samples</code> , medium <code>n_clusters</code>	Non-flat geometry, uneven cluster sizes	Distances between nearest points
Gaussian mixtures	many	Not scalable	Flat geometry, good for density estimation	Mahalanobis distances to centers
Birch	branching factor, threshold, optional global clusterer.	Large <code>n_clusters</code> and <code>n_samples</code>	Large dataset, outlier removal, data reduction.	Euclidean distance between points

Which clustering algorithm?

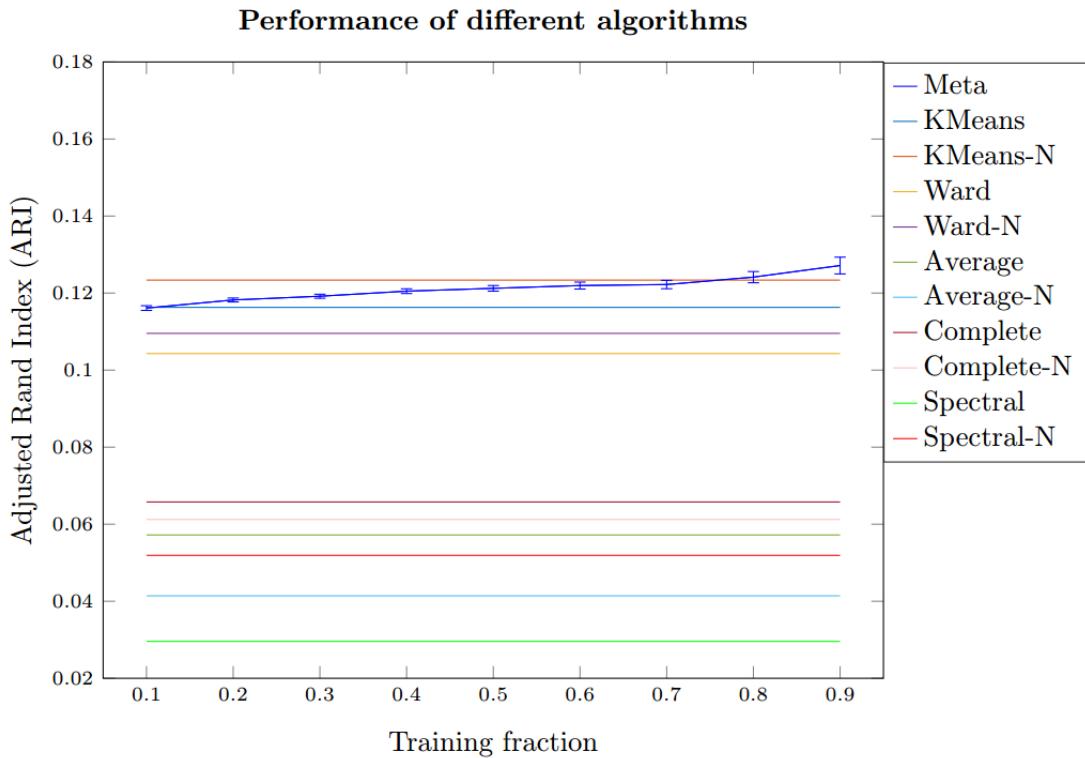
Fix $k=2$ clusters

250 binary classification problems from openml.org

5 algorithms * 2 normalizations

Meta takes max ARI estimate from:

- # dimensions
- # examples
- **Silhouette score**
- max/min singular value



HOW MANY CLUSTERS?

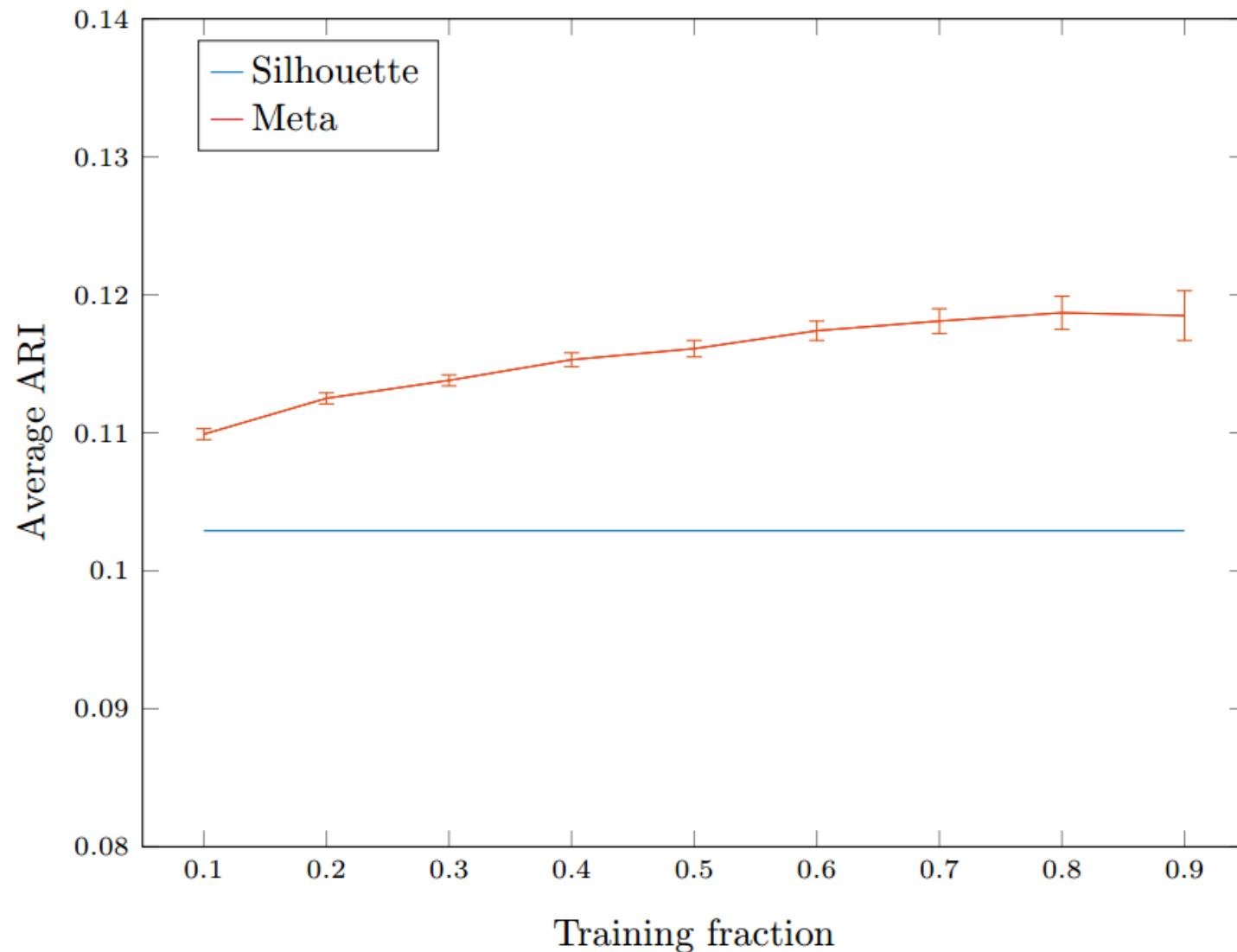
Std. Silh approach to pick k

- Given new clustering problem, run k -means for $k=2,\dots,10$ and compute $\text{Silh} \downarrow 2, \dots, \text{Silh} \downarrow 10$
- Take k to maximize $\text{Silh} \downarrow k$

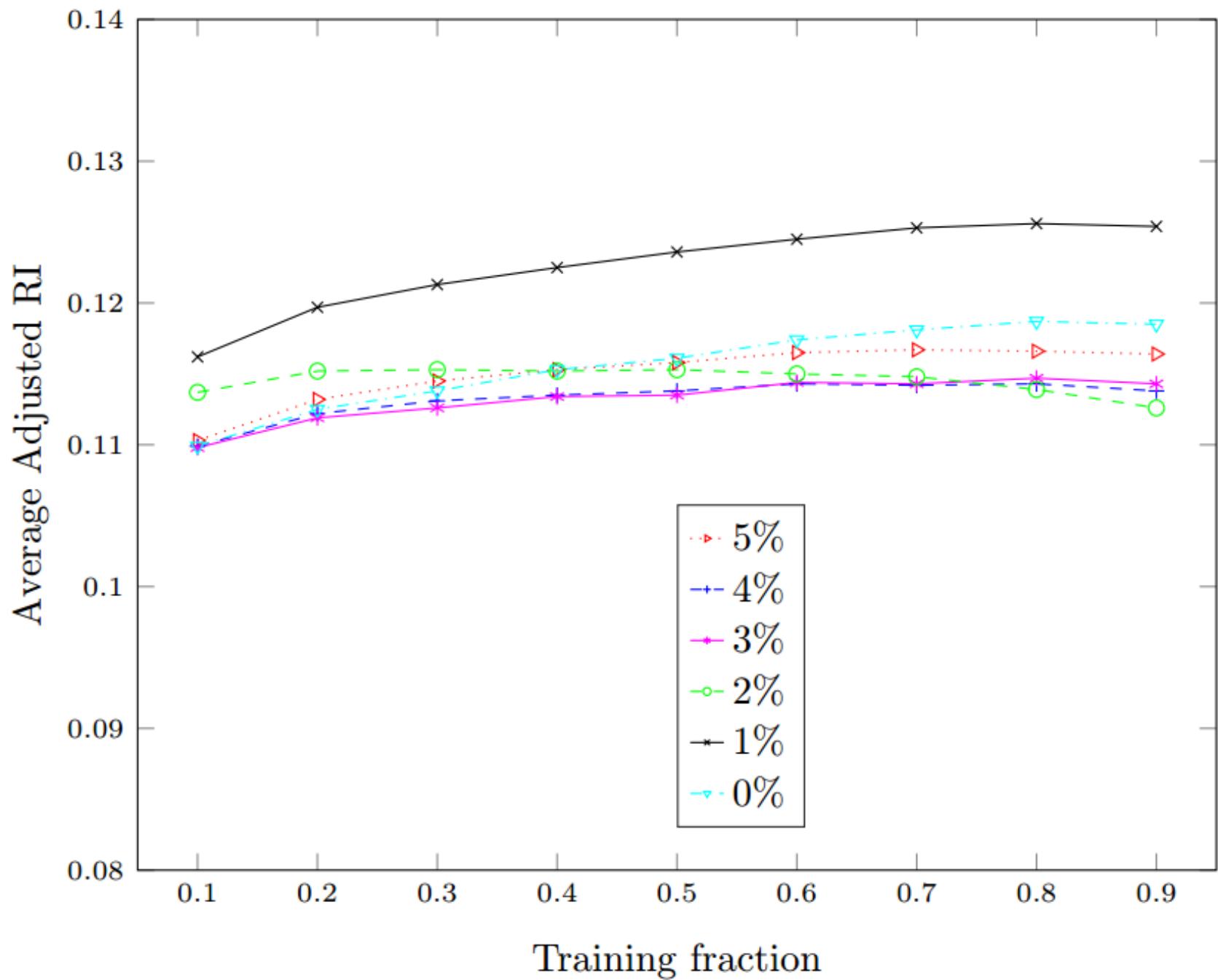
Meta- k

- For each $k=2, \dots, 10$, fit $\text{ARI} \downarrow k = \alpha \downarrow k + \beta \downarrow k \text{Silh} \downarrow k$
 - Use least-squares linear regression
 - Data is ground-truth problems from openml.org
- Given new clustering problem, compute $\text{Silh} \downarrow 2, \dots, \text{Silh} \downarrow 10$ and estimate $\text{ARI} \downarrow 2, \dots, \text{ARI} \downarrow 10$
- Take k to maximize $\text{ARI} \downarrow k$

Choosing # clusters k

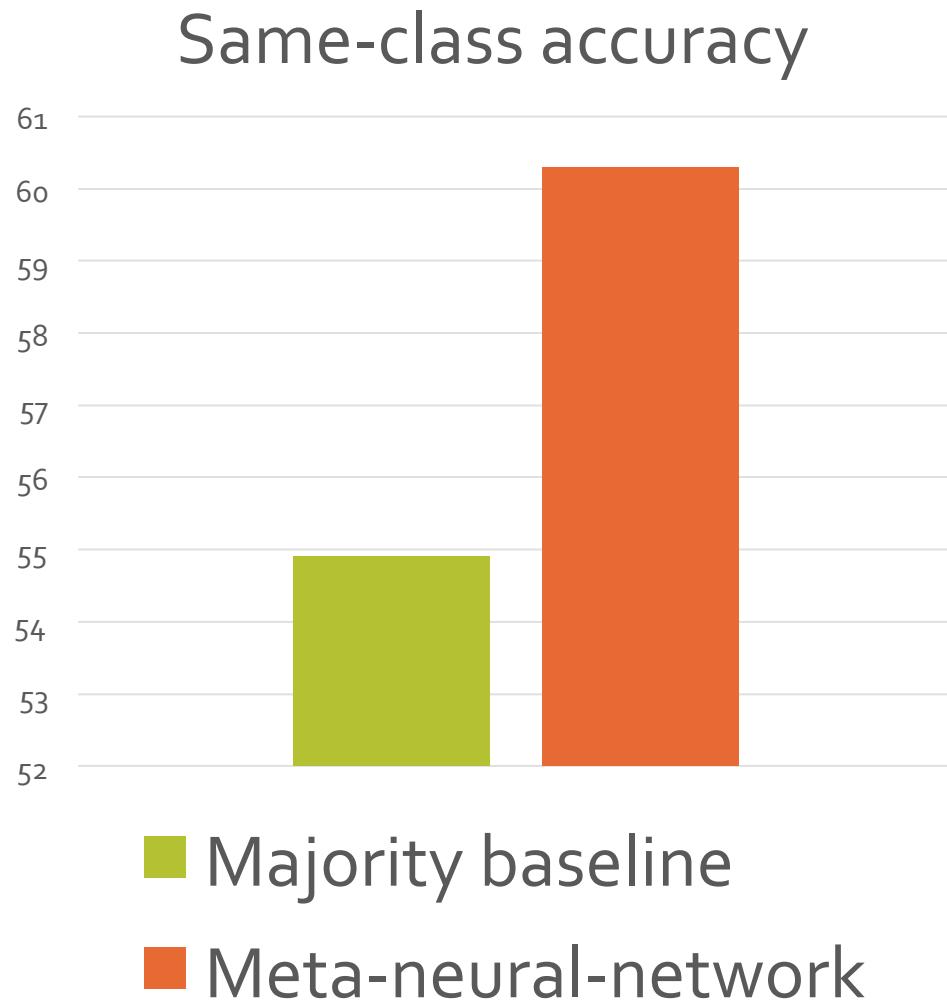


META-OUTLIER- REMOVAL



Meta-neural-network-unsupervised-learning

Lots of little
data makes
big data



HUMAN EXPERIMENTS?

Cluster this

✿❀Ⓜ️ ❁●❁◆ ♦❁ ❁ ⓘ□Ⓜ️❁♦❁

❖❁◆Ⓜ️ ❁□❖ ❁♦❁○Ⓜ️ ❁❖❖❖❖❖ ❁❖ ❁❖❖❖❖ ❁❖ ❁❖❖❖❖ ❁❖ ❁❖❖❖❖

✿❀Ⓜ️ ❁●❁◆ ♦❁ ❁ ⓘ□□❖■■❖

❖Ⓜ️♦❁ ❁□♦□♦Ⓜ️ ❁❖ ❁❖❖Ⓜ️ ❁♦❁○Ⓜ️ ■■ ❁❖ ❁❖ ❁❖ ❁❖ ❁❖ ❁❖

✿❀Ⓜ️ ❁●❁◆ ♦❁ ❁♦❁□❖❖●Ⓜ️ ❁

❖ ❁♦●Ⓜ️□♦❁ ❁♦❀Ⓜ️ ❁♦❀□●Ⓜ️ ❁♦❁○Ⓜ️

Cluster this

＊＊＊＊＊ ＊＊＊＊＊ ♦＊♦ ＊＊＊＊＊♦＊＊

＊＊♦＊＊ ＊＊＊ ♦＊＊＊＊＊ ♦＊＊＊＊＊＊＊＊

＊＊＊＊＊ ＊＊＊＊＊ ♦＊♦ ＊＊＊＊＊＊＊＊＊

＊＊＊＊＊ ♦＊＊＊＊＊＊＊＊＊＊＊＊＊＊＊＊＊＊＊＊

＊＊＊＊＊ ＊＊＊＊＊ ♦＊＊＊＊＊＊＊＊＊＊＊＊＊＊

＊＊＊＊＊ ♦＊＊＊＊＊＊＊＊＊＊＊＊＊＊＊＊＊＊＊＊

Cluster this

The class was great!

Waste of time ADAM I HOPE YOU DIE.

The class was boring.

Best course I have taken in a long time.

The class was terrible!

I slept the whole time.

Summary

Don't miss the meta-boat

- Meta-clustering
- Meta-topic-modeling
- Meta-feature-selection
- **Meta-classification**
- **Meta-algorithms**



H&M theories that apply to Humans & Machines

An open meta-problem

Humans solve problems, write programs, come up with and improve strategies

8			9	6
7		1	8	
6	3			1
5	2			
			2	

How can computers do the same?

Meta-thank-you

