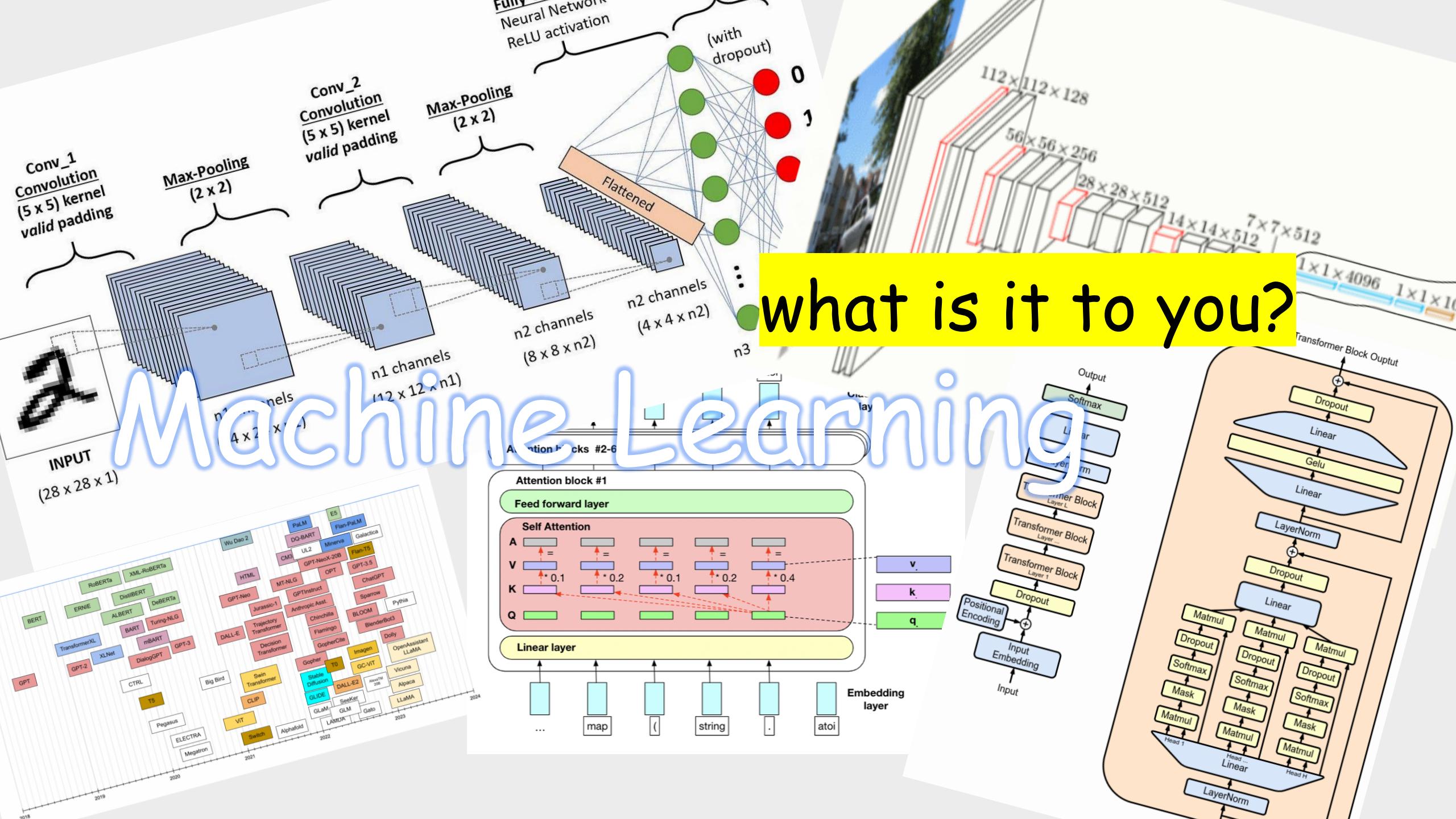


Machine Learning



lesson I

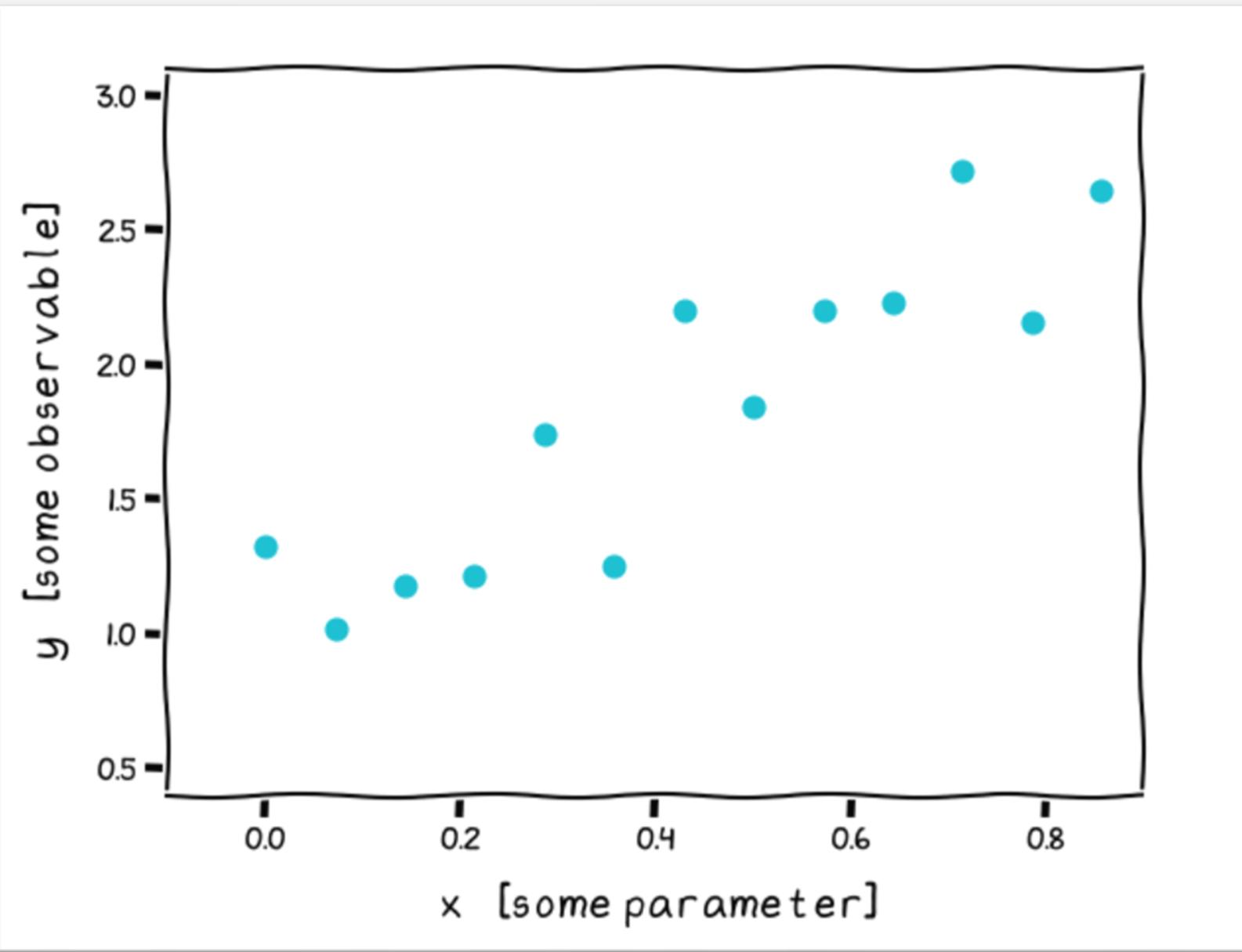
Machine Learning

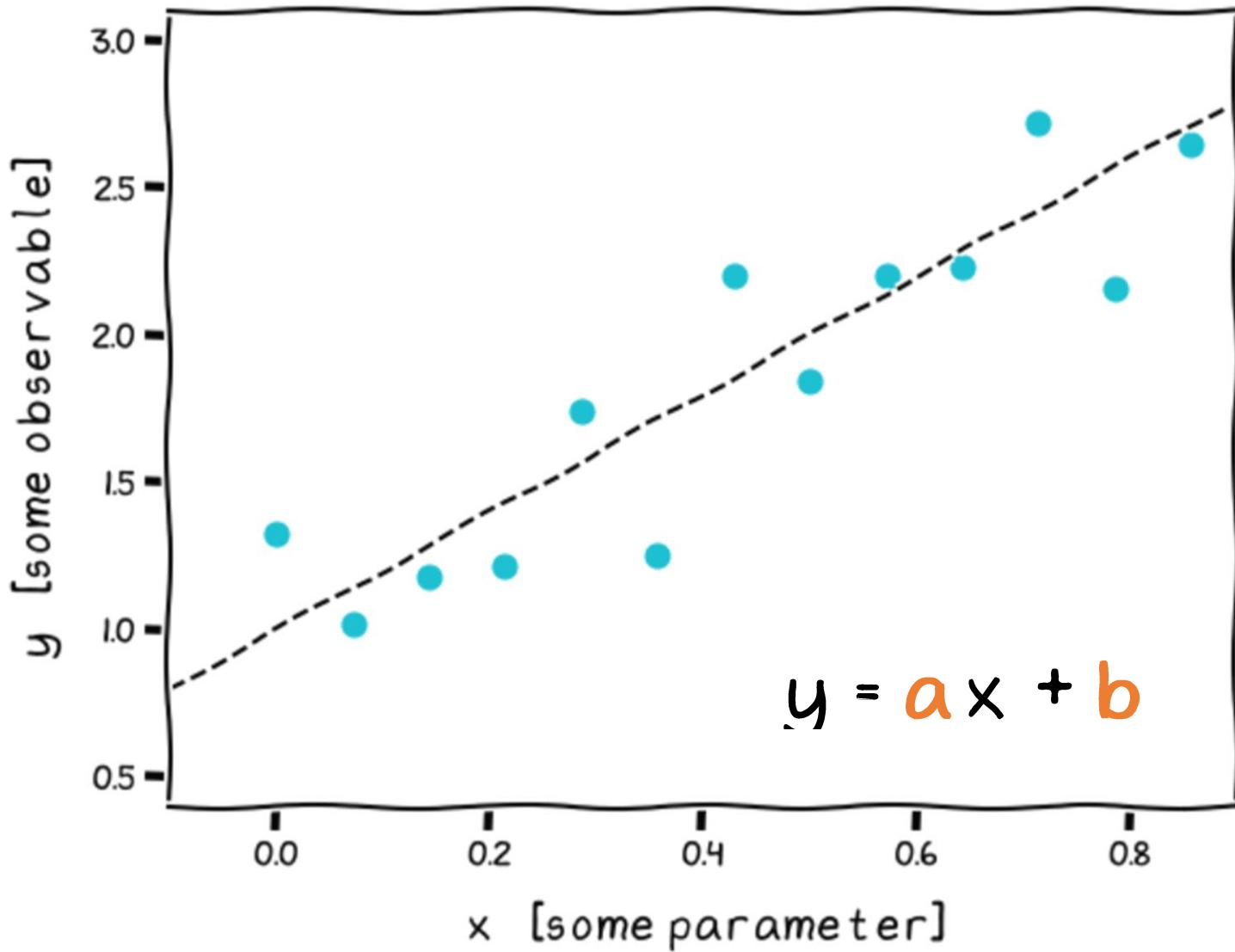
creating a model (with some parameters)

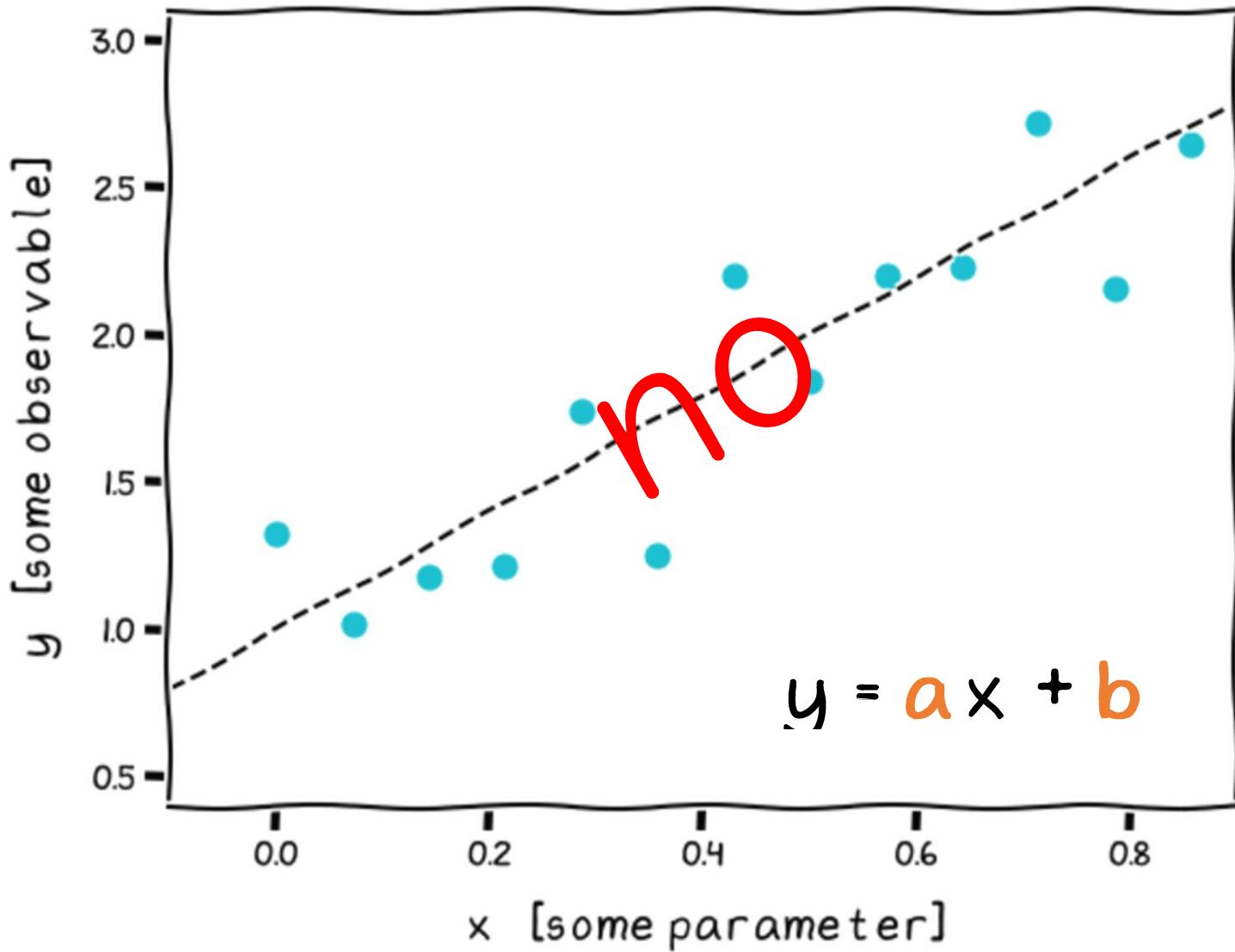
then optimizing it

to predict something
(with a computer algorithm)

©Liza don't quote me



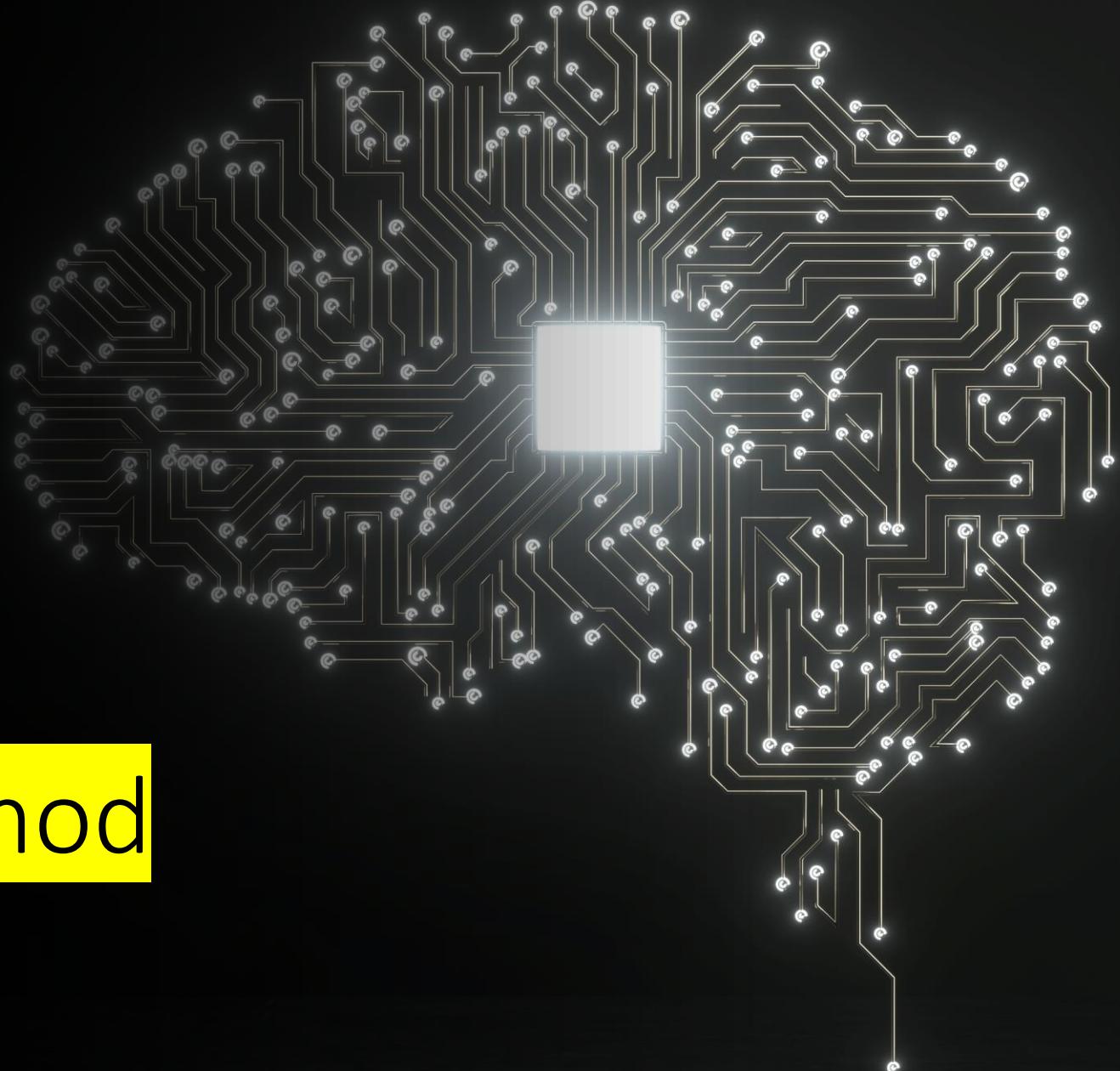






the scientific method

and the philosophy of machine learning



- 1 Form a hypothesis
- 2 Collect data
- 3 Split the data
- 4 Train your model
- 5 Evaluate your model
- 6 Think about your model
- 7 Think of the next project



scientific method
In a nutshell
(ML edition)

1

Form a hypothesis

Before looking at any data, decide for yourself:

- What do you expect to find?
- What is your model?
- Do you plan to compare with other findings? If so, which?
- What do you consider “a good fit” (if modelling)?
- What do you consider “a discrepant result” or “a tension”?
E.g., how many sigmas do you consider significant?

Looking at the data will bias you!

Example *The Black Cloud* (Hoyle, 1958). The Black Cloud appears to be heading for the Earth. The scientific team suggests that this proves the cloud has intelligence. Not so, says the dissenting team member.

Example *The Black Cloud* (Hoyle, 1958). The Black Cloud appears to be heading for the Earth. The scientific team suggests that this proves the cloud has intelligence. Not so, says the dissenting team member. Why? A golf ball lands on a golf course which contains 10^7 blades of grass; it stops on one blade; the chances are 1 in 10^7 of this event occurring by chance. This is not so amazing – the ball had to land somewhere. It would only be amazing if the experiment were re-run to test the newly formulated hypothesis (e.g. the blade being of special attractive character; the golfer of unusual skill) and the event was repeated. However, the importance of deciding if the Black Cloud knew about the Earth cannot await the next event or the sequence of events, and tempts the rush to judgement in which initial data, hypothesis and test data are combined; so in many instances in astronomy and cosmology.

1

Form a hypothesis

Before looking at any data, decide for yourself:

- What do you expect to find?
- What is your model?
- Do you plan to compare with other findings? If so, which?
- What do you consider “a good fit” (if modelling)?
- What do you consider “a discrepant result” or “a tension”?
E.g., how many sigmas do you consider significant?

Looking at the data will bias you!

1

Form a hypothesis

Before looking at any data, decide for yourself:

- What do you expect to find?
- ~~What is your model?~~
- Do you plan to compare with other findings? If so, which?
- What do you consider “a good fit” (if modelling)?
- What do you consider “a discrepant result” or “a tension”?
E.g., how many sigmas do you consider significant?

Looking at the data will bias you!

2

Collect data

No turning on your fancy ML machine yet

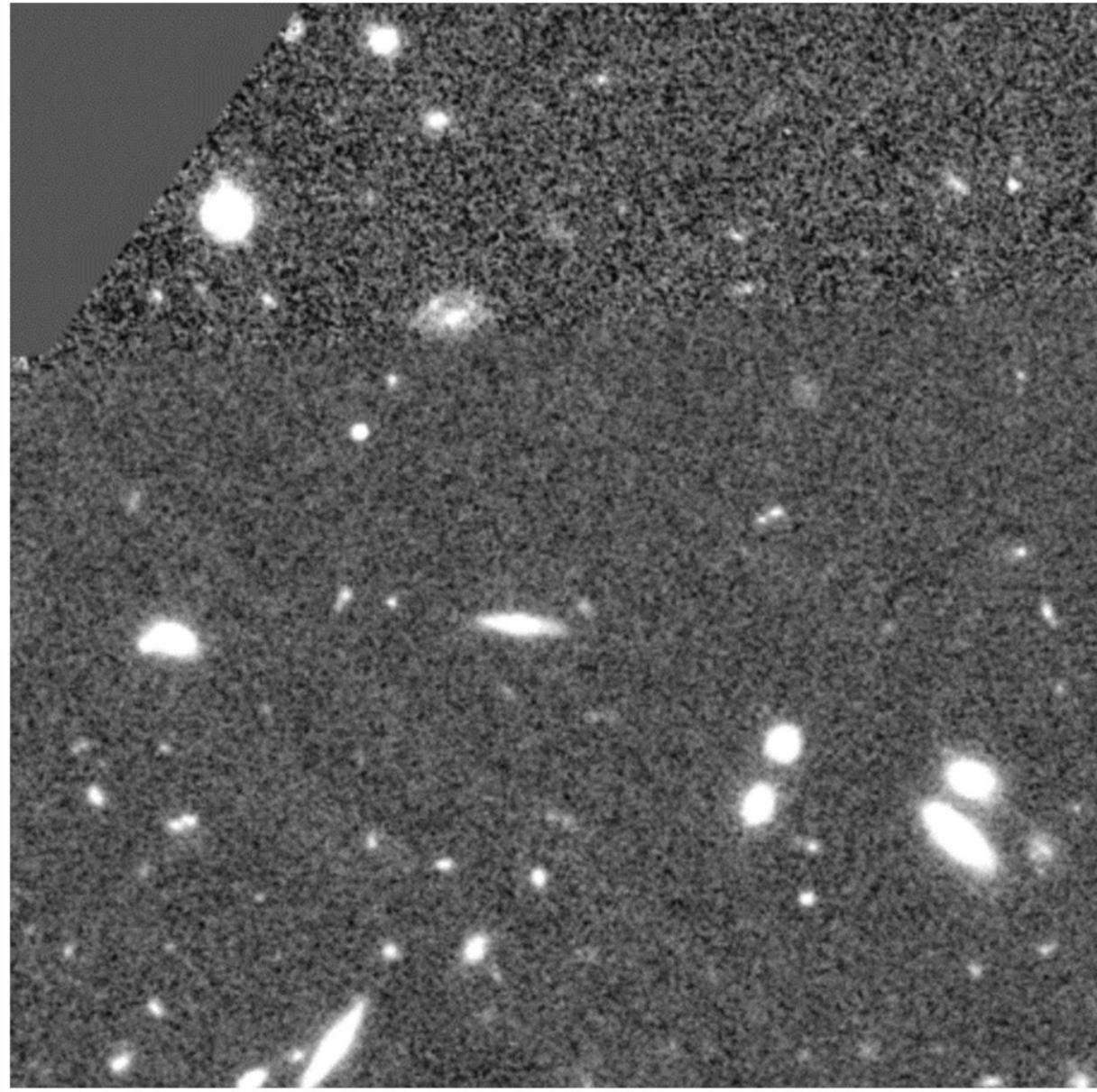
- Run the experiment and preprocess data
- **Preprocess** collected data for analysis

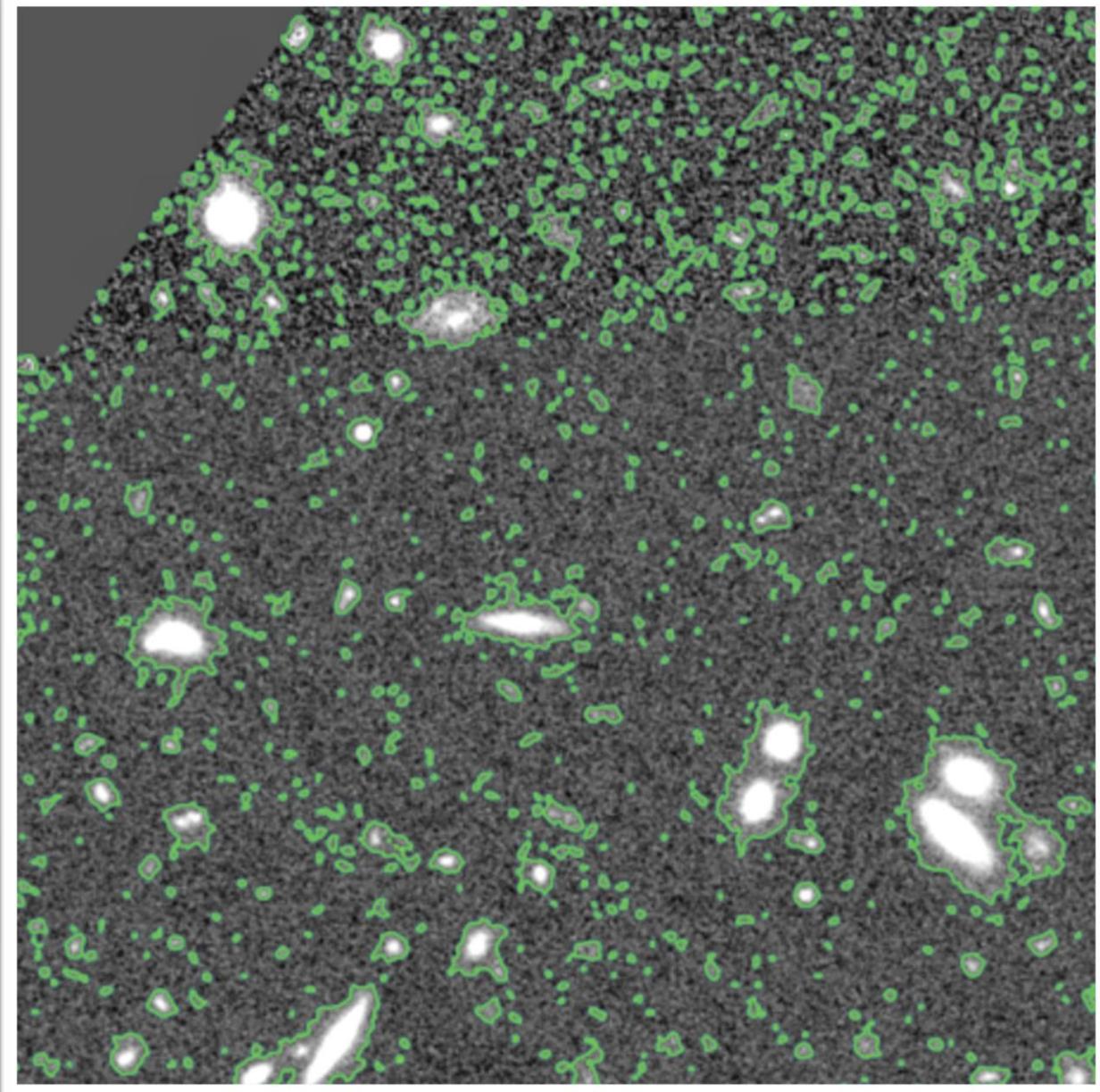
Cleaning?

Outlier detection?

Normalization?

Look at your data!





3

Split the data

The ultimate ML lifehack

- 70% **training** set

The dataset you will fit everything on

It should not change between re-runs of your code/model fitting/etc
(set `random_state` to a number when sampling your training set)

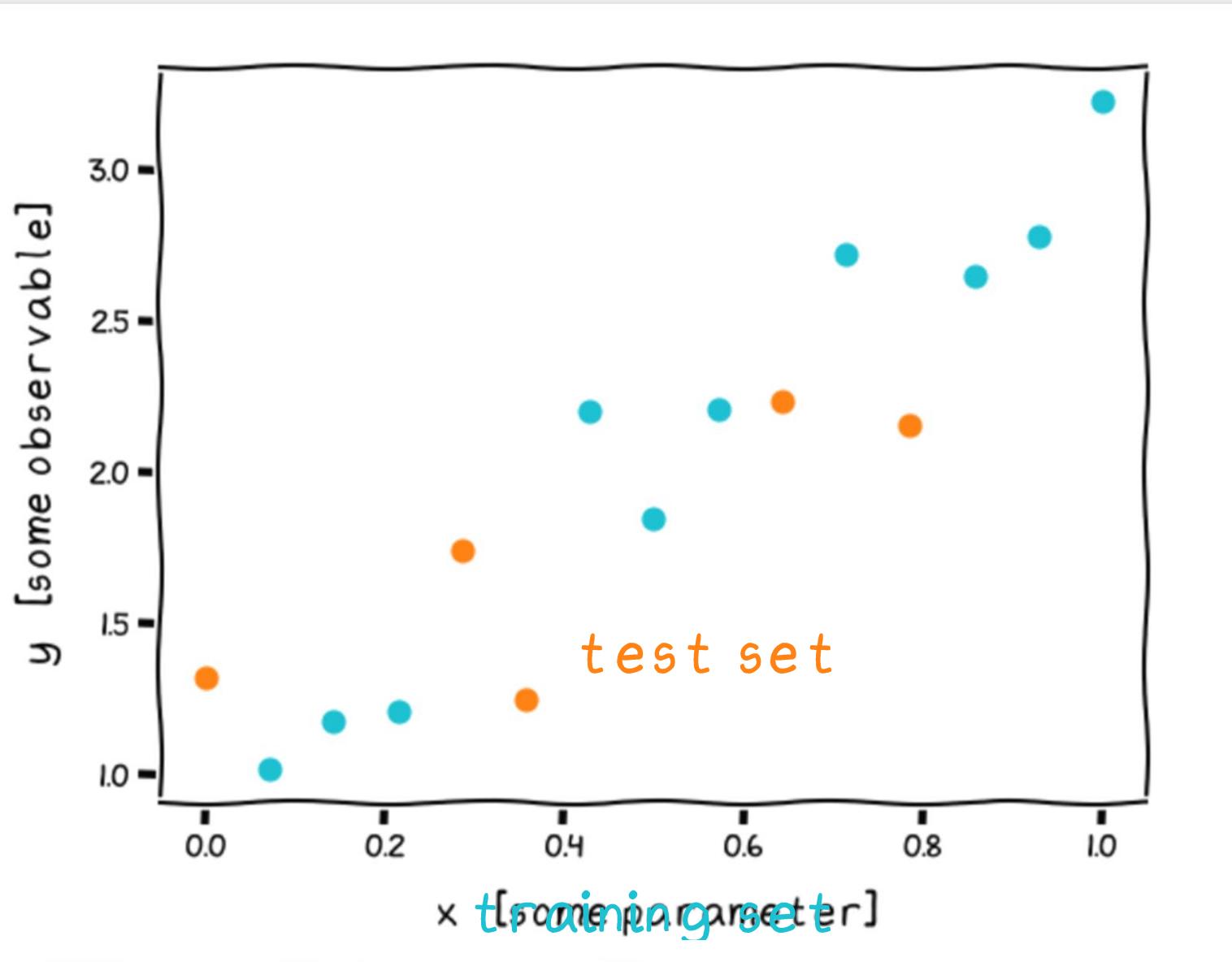
- 30% **validation** set

The dataset you test the model's performance on

- 10% **test** set [optional]

After the model is optimized & you are writing the paper, this is the final test

No need to re-run experiments!



4

“Train” the model

aka optimize your model parameters
using whatever approach you prefer

only look at your training set!

Evaluate the model

Goodness-of-fit, aka “loss”

The **loss function** can be anything you like

- Fitting a curve: χ^2 , reduced χ^2 , RMS loss
- Classification: binary cross-entropy
- More advanced techniques: regularization? Bayesian? Probabilities?

This is also done during training

More on gradient descent + parameter optimization next week



What data do we evaluate the model with?

Evaluate the model

Goodness-of-fit, aka “loss”

The loss function can be anything you like:

- Fitting a curve: χ^2 or R^2
- Classification: Confusion matrix
- More advanced: Bayesian? Probabilities?

Both

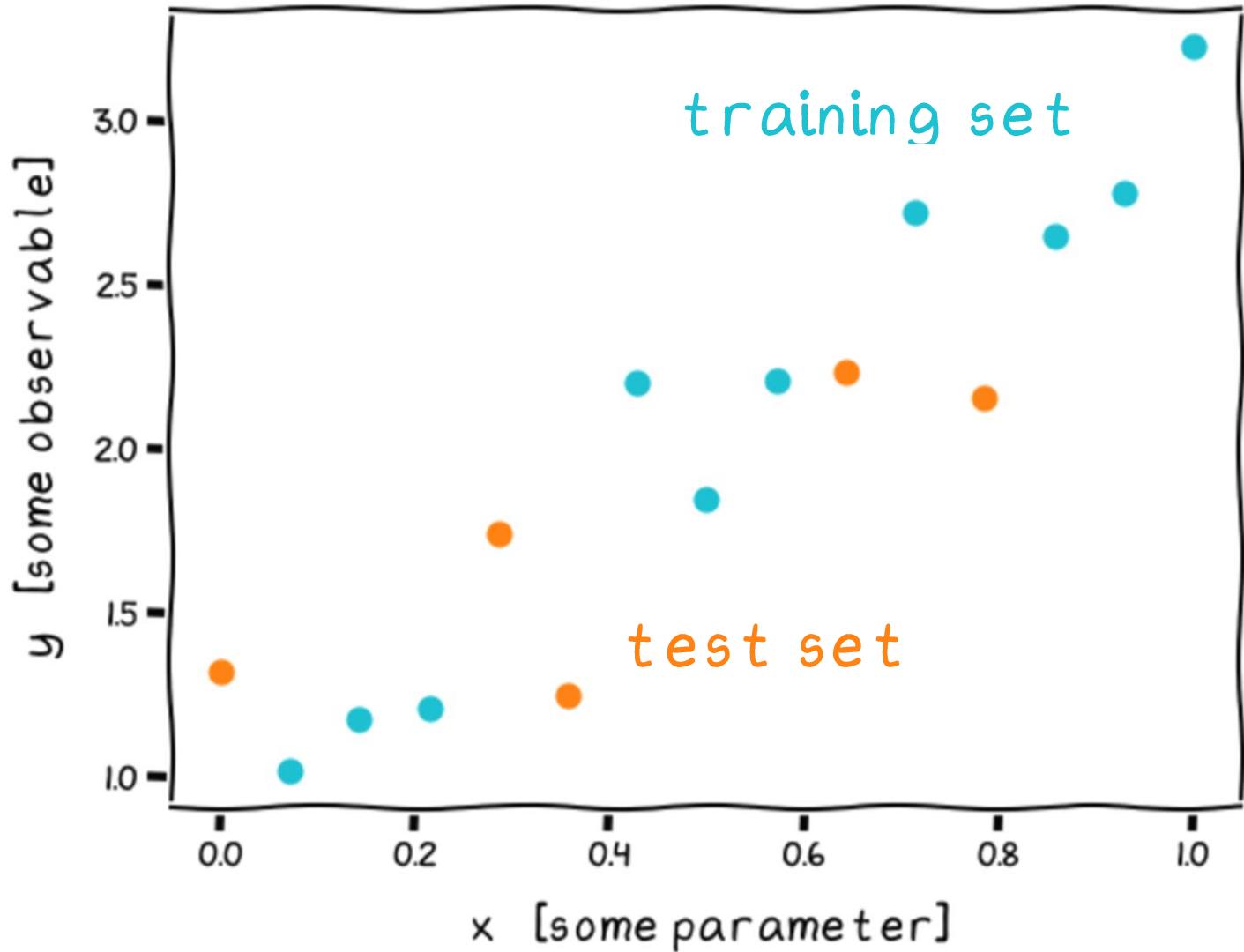
(but mostly test set)

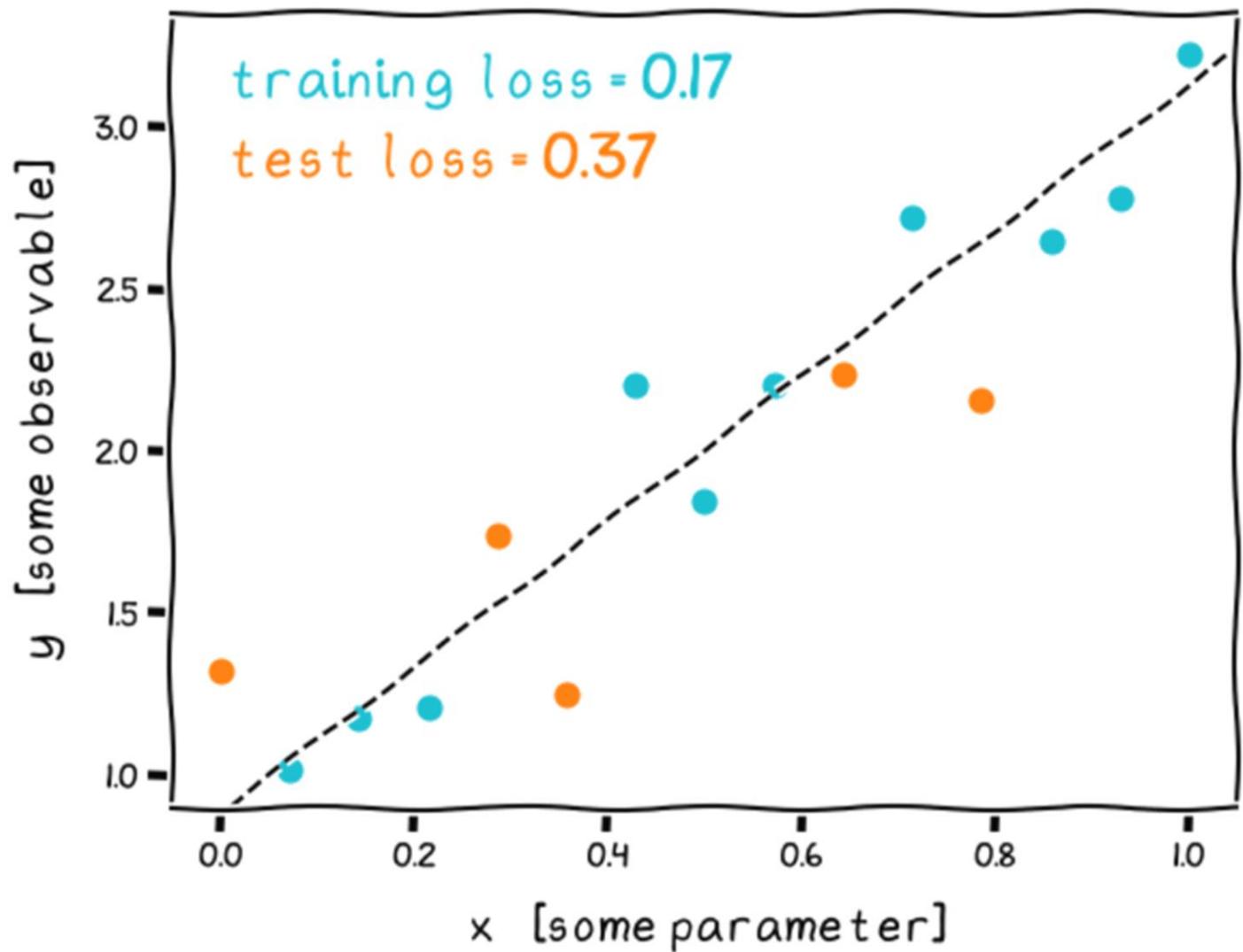
This is also done during training

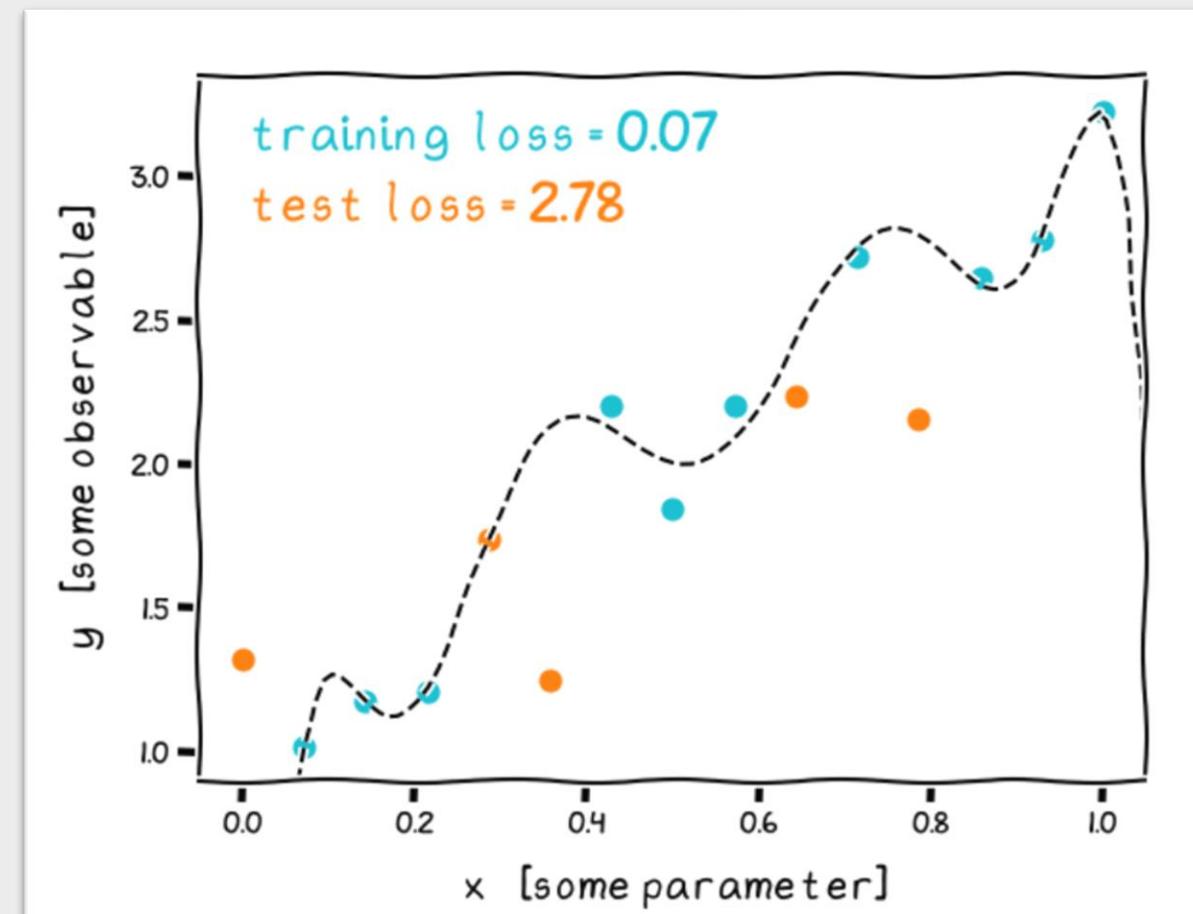
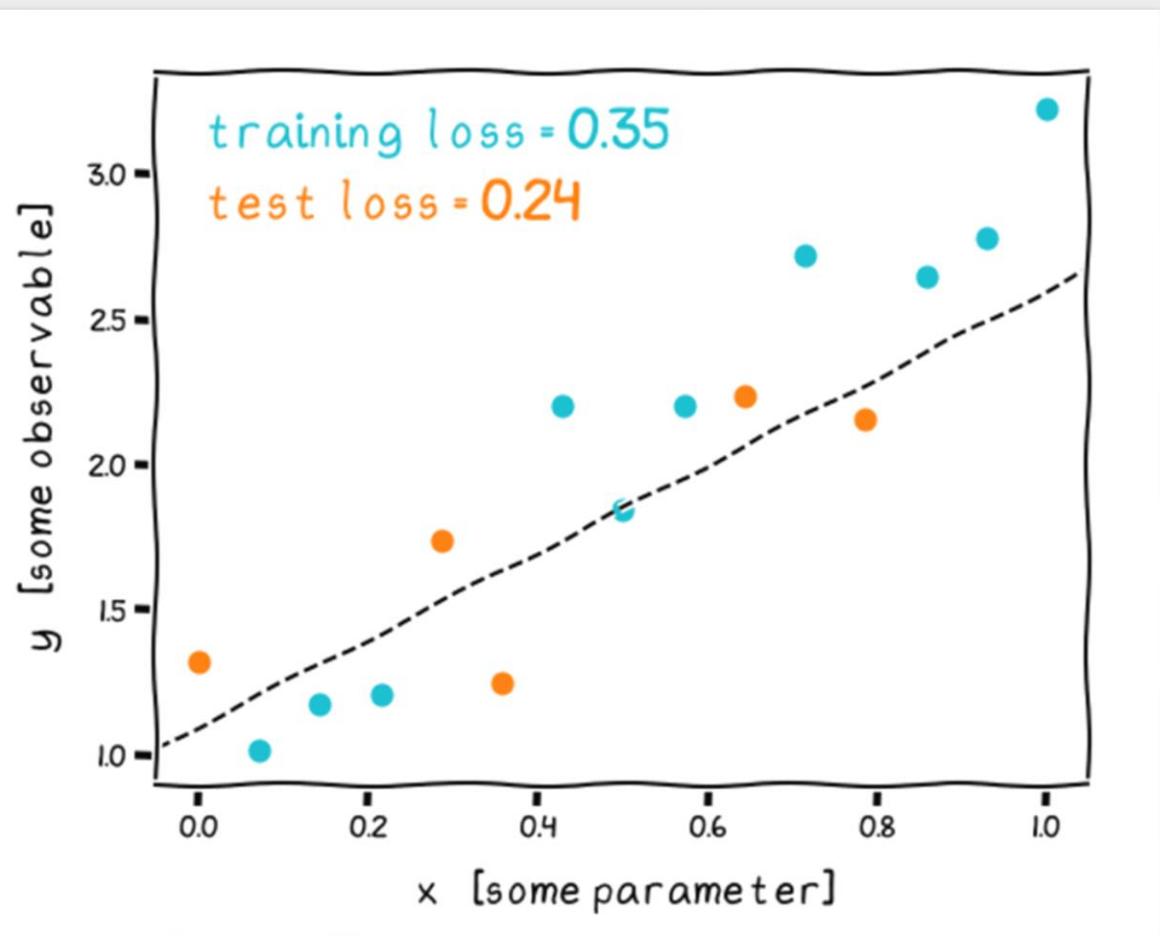
More on gradient descent + parameter optimization next week

What data do we evaluate the model with?

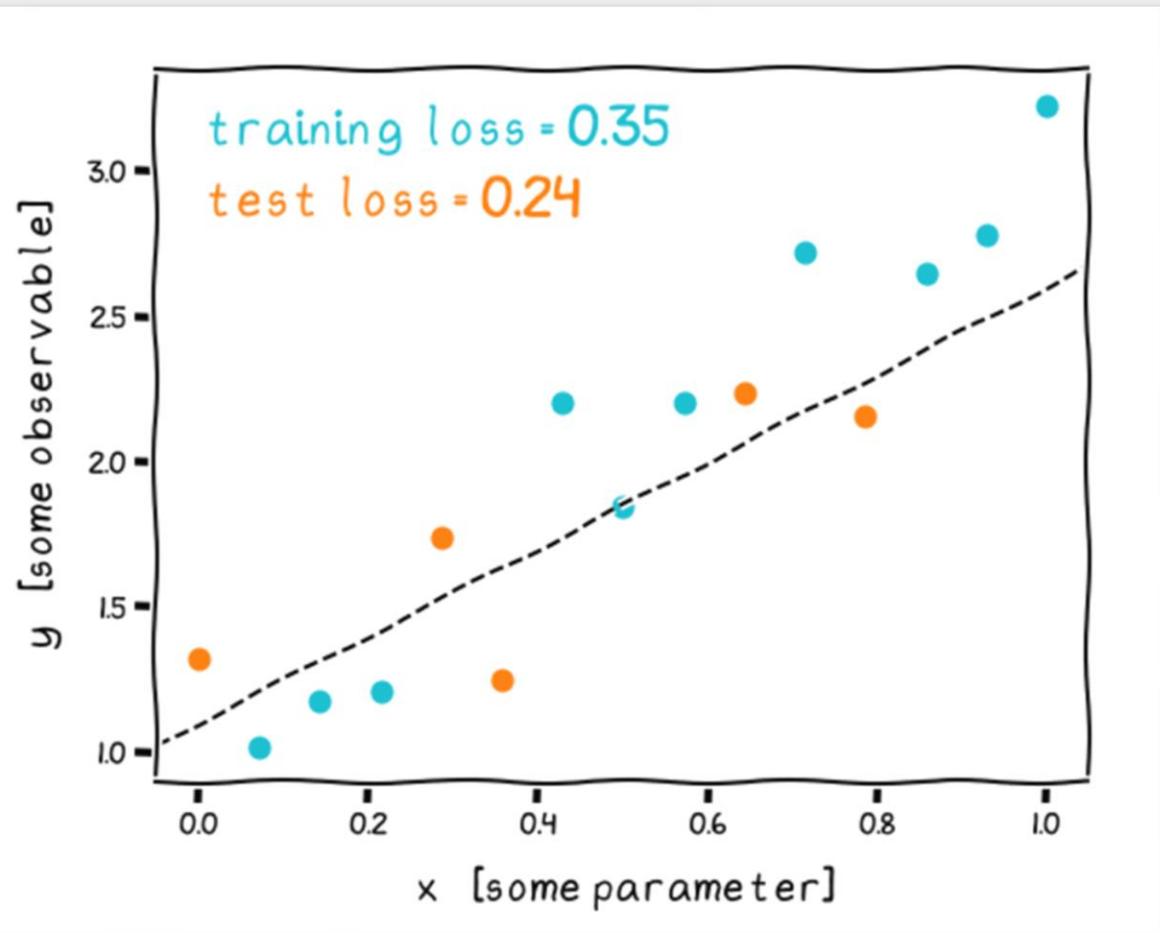




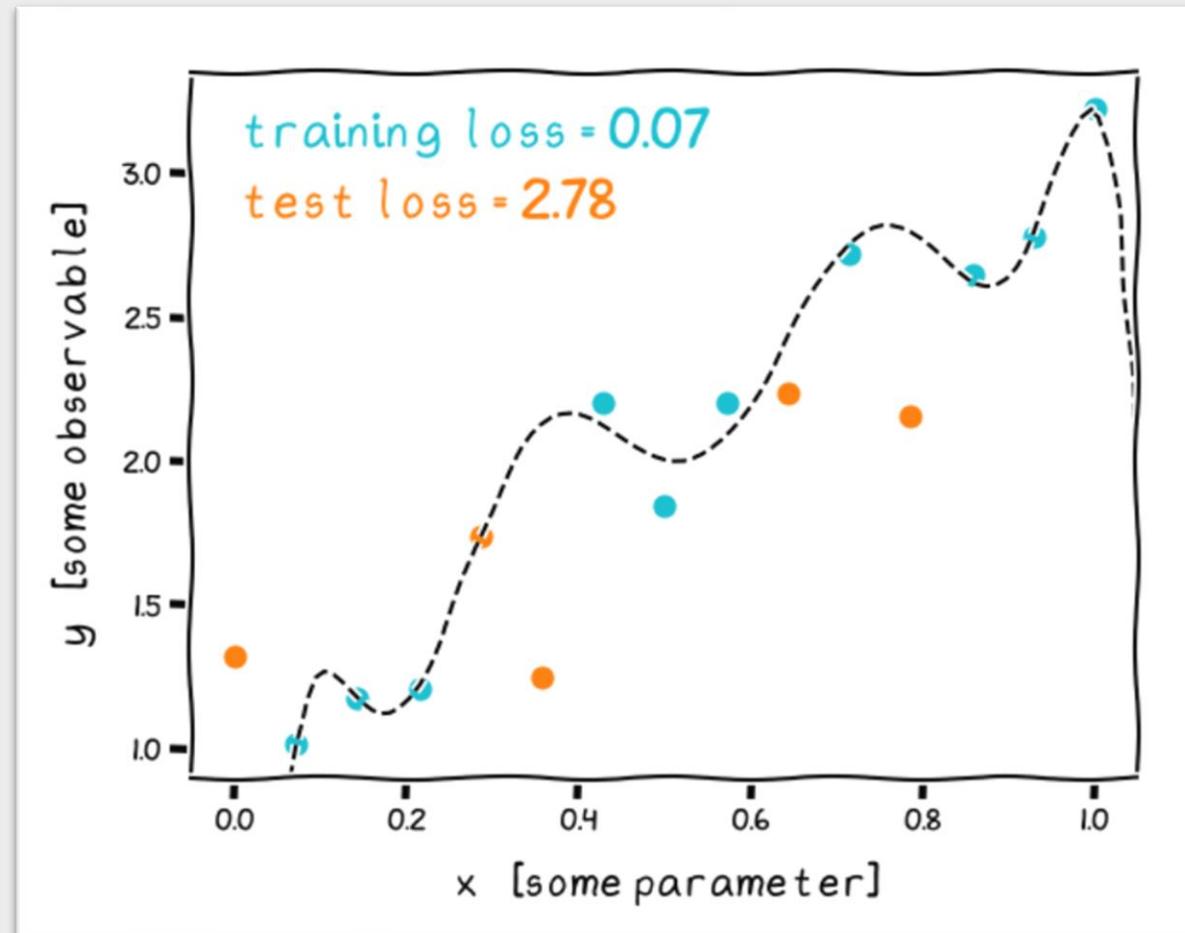




underfitting



overfitting



training loss > test loss
or training loss is high / not convergent

test loss > training loss

Evaluate the model

Goodness-of-fit, aka “loss”

The **loss function** can be anything you like

- Fitting a curve: χ^2 , reduced χ^2 , RMS loss
- Classification: binary cross-entropy
- More advanced techniques: regularization? Bayesian? Probabilities?

This is also done during training

More on gradient descent + parameter optimization next week



What data do we evaluate the model with?

Evaluate the model cheatsheet

`sklearn.metrics`

loss functions
accuracy scores
precision/recall
confusion matrices



Training your model

training loss

Deciding between models

validation loss

Presenting your model

validation loss

or even better, a new test set loss

What data do we evaluate the model with?

Final considerations

What are the limitations of your model?

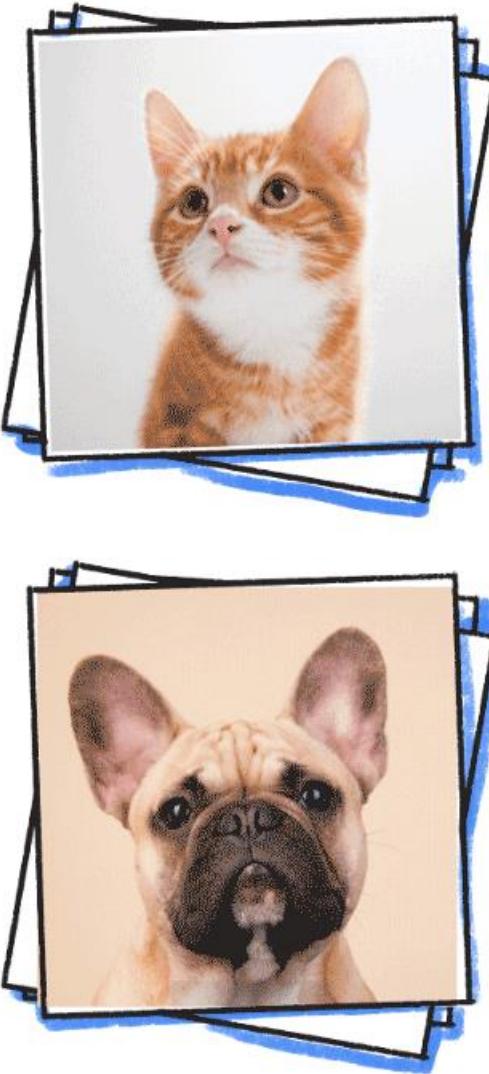
- What was in your training set?
- What was not included?
- What was under / over-represented?
- When would your assumptions break down?

Do not extrapolate!

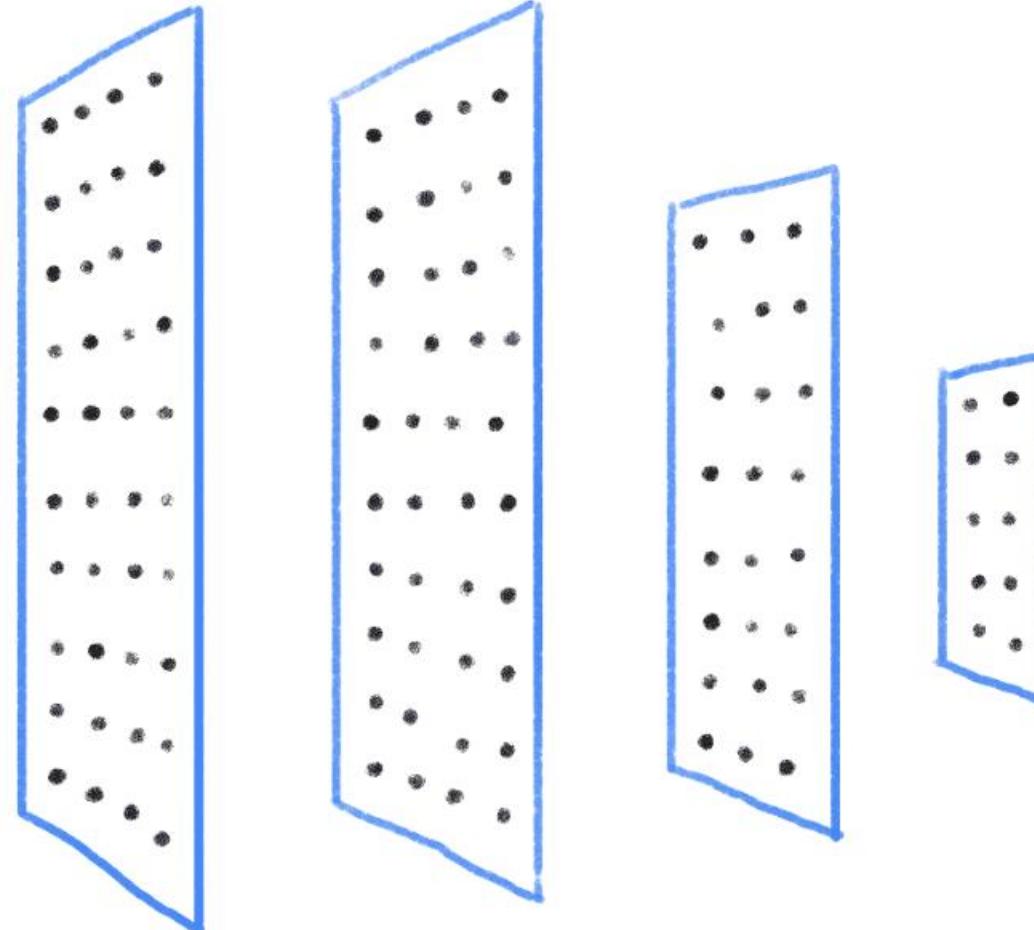
CAT

(LABELED)
PHOTOS

DOG



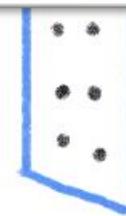
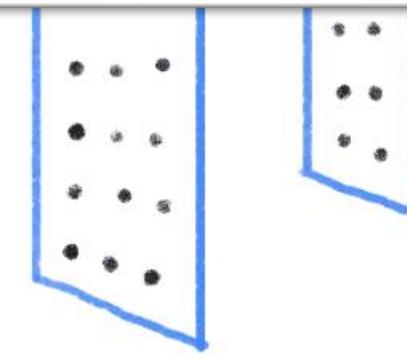
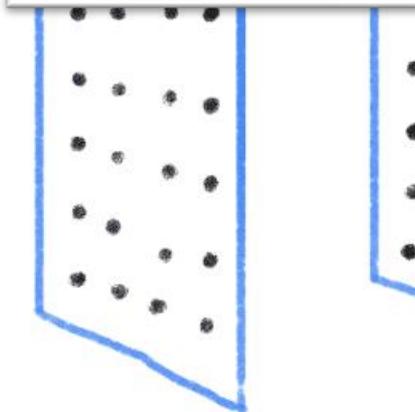
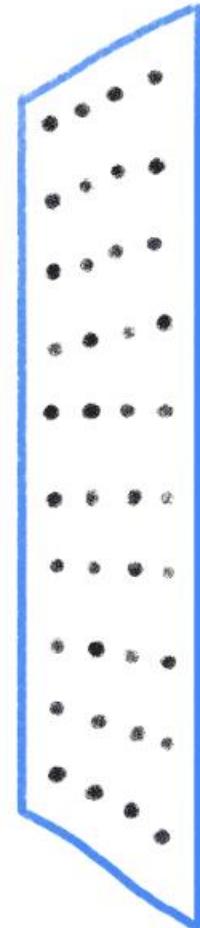
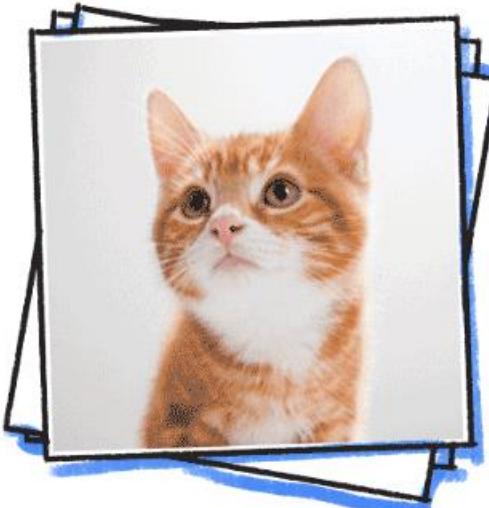
is this a cat?



CAT

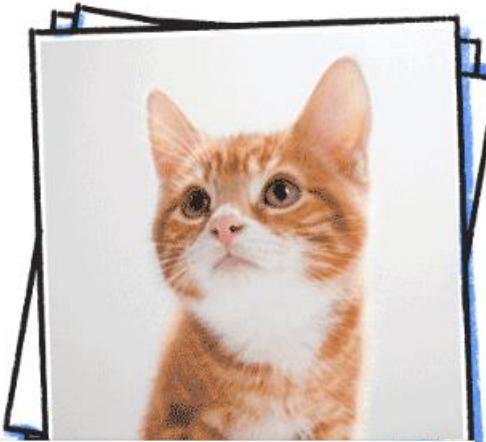
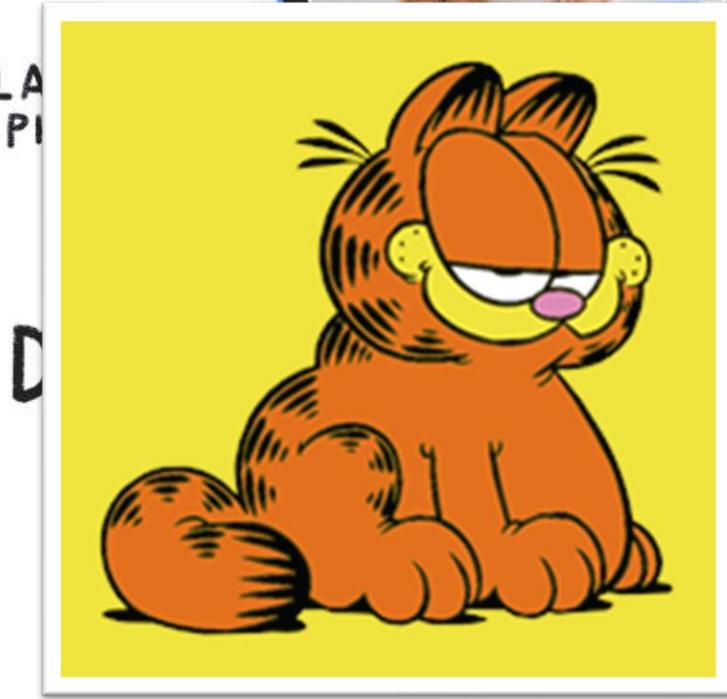
(LABELED)
PHOTOS

DOG

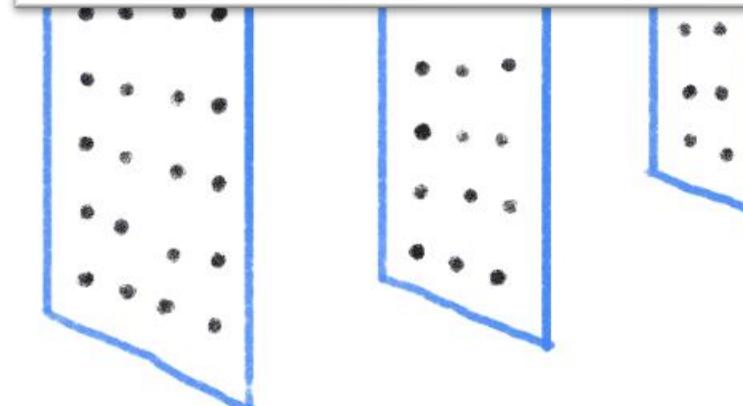
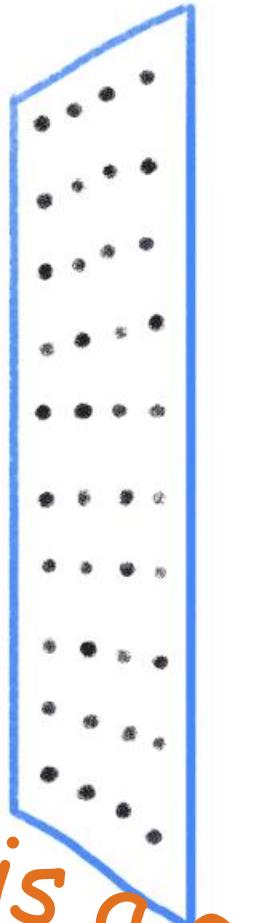


CAT

(LA
PI



is this a cat?



Next steps

Congratulations, you have a model!

- Now you have a new hypothesis – test it
- How does it fit in with other models / what we already know?
- What have you learnt?
- Did you get any cool scientific results? (write a paper)
- Did you provide a cool new dataset for people? (write a paper)

Write a paper, put your
Code on GitHub, make your data public :)

- 1 Form a hypothesis
- 2 Collect data
- 3 Split the data
- 4 Train your model
- 5 Evaluate your model
- 6 Think about your model
- 7 Think of the next project



scientific method
In a nutshell
(ML edition)

Lessons to remember!

1. Always look at your data before doing things
2. Never tell people your model is awesome
using training set metrics
3. Never extrapolate! and don't trust people who do
4. Scikit-learn is your best friend

Types of Machine Learning

Types of Machine Learning

supervised

unsupervised

contrastive

active

reinforcement

semi-supervised

self-supervised

Types of Machine Learning

supervised

unsupervised

contrastive

active

reinforcement

semi-supervised

self-supervised

Types of Machine Learning

supervised

you know the
right answer
(for some of your data)

labelled data

unsupervised

you don't know the
right answer

unlabelled data

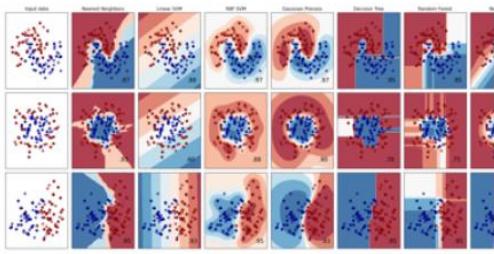
supervised

Classification

Identifying which category an object belongs to.

Applications: Spam detection, image recognition.

Algorithms: Gradient boosting, nearest neighbors, random forest, logistic regression, and more...

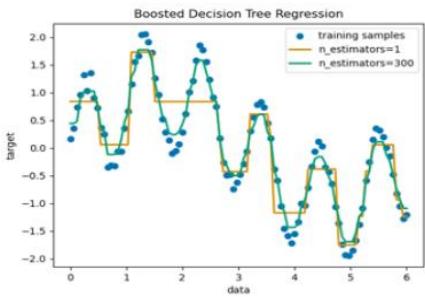


Regression

Predicting a continuous-valued attribute associated with an object.

Applications: Drug response, Stock prices.

Algorithms: Gradient boosting, nearest neighbors, random forest, ridge, and more...

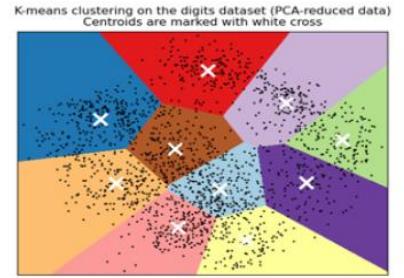


Clustering

Automatic grouping of similar objects into sets.

Applications: Customer segmentation, Grouping experiment outcomes

Algorithms: k-Means, HDBSCAN, hierarchical clustering, and more...



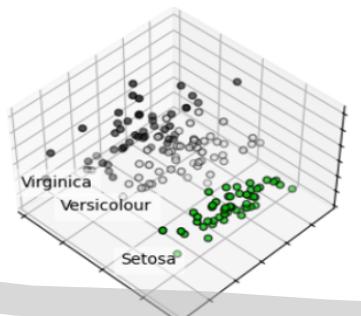
unsupervised

Dimensionality reduction

Reducing the number of random variables to consider.

Applications: Visualization, Increased efficiency

Algorithms: PCA, feature selection, non-negative matrix factorization, and more...

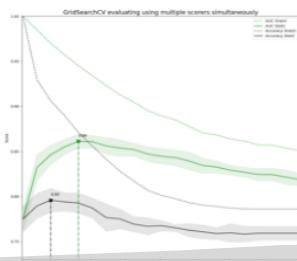


Model selection

Comparing, validating and choosing parameters and models.

Applications: Improved accuracy via parameter tuning

Algorithms: grid search, cross validation, metrics, and more...

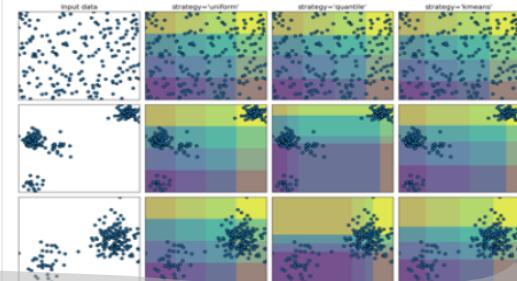


Preprocessing

Feature extraction and normalization.

Applications: Transforming input data such as text for use with machine learning algorithms.

Algorithms: preprocessing, feature extraction, and more...



not ML (this take is a bit hot and probably untrue)

unsupervised learning

A photograph of four cowboys in silhouette against a vibrant orange and yellow sunset sky. They are standing in a row, facing right, holding rifles. The scene is framed by the dark outlines of buildings on the left and right.

no need for
train/test split

try
lots of
models

explore
your data

see what
works

unsupervised learning

data exploration

what type of data am I working with?
how many types of penguins are there?

finding outliers

which points are unusual? why / how?
can I remove artifacts from my galaxy catalog?

unsupervised learning

similarity search

which things are similar, and why?
*which TV shows should I watch
if I like Bojack Horseman?*

finding patterns

what are the patterns or groups in my data?
how do I detect REM sleep from EEG?

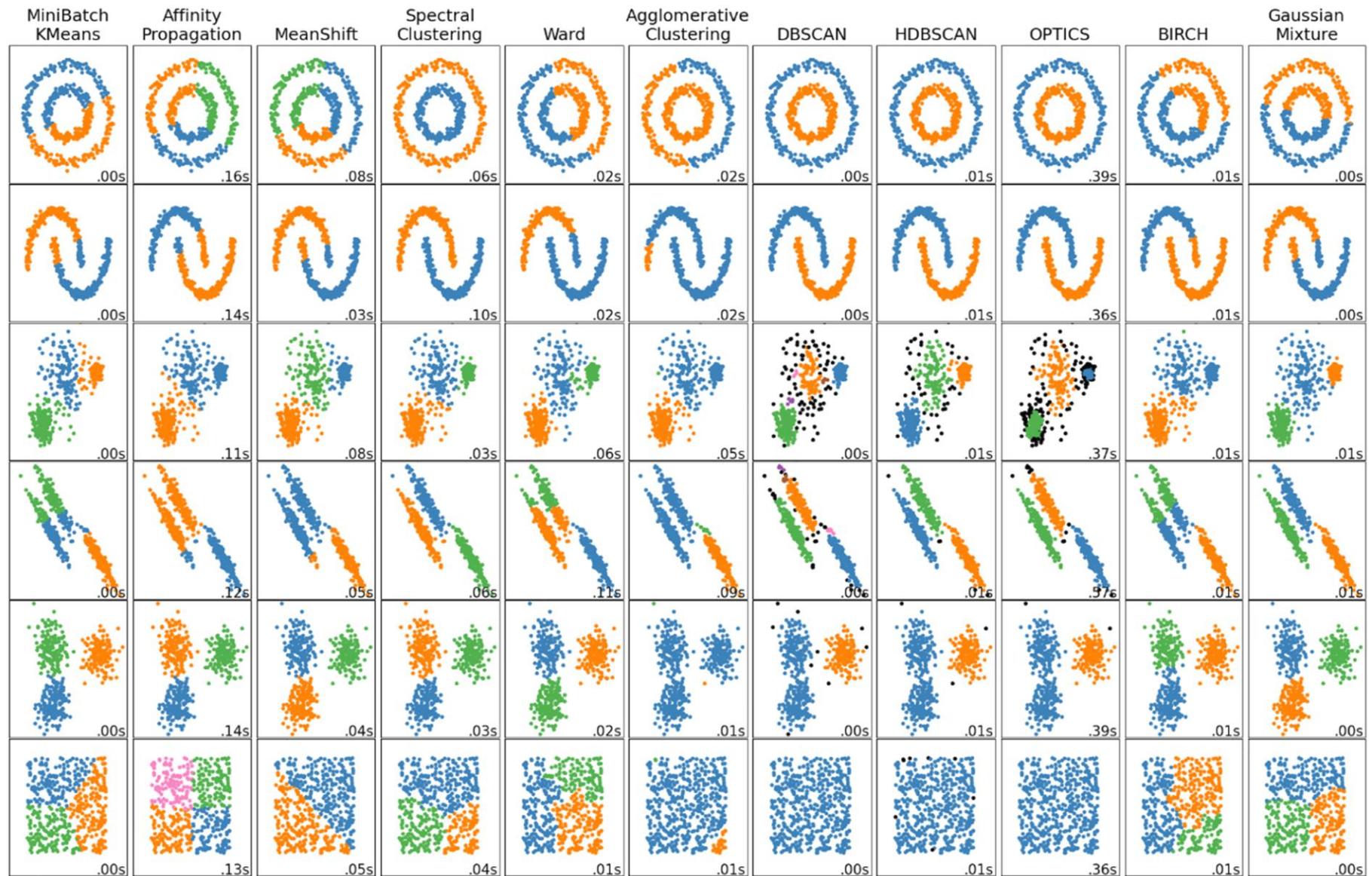
unsupervised learning

look at your data



choose an algorithm

2.3.1. Overview of clustering methods

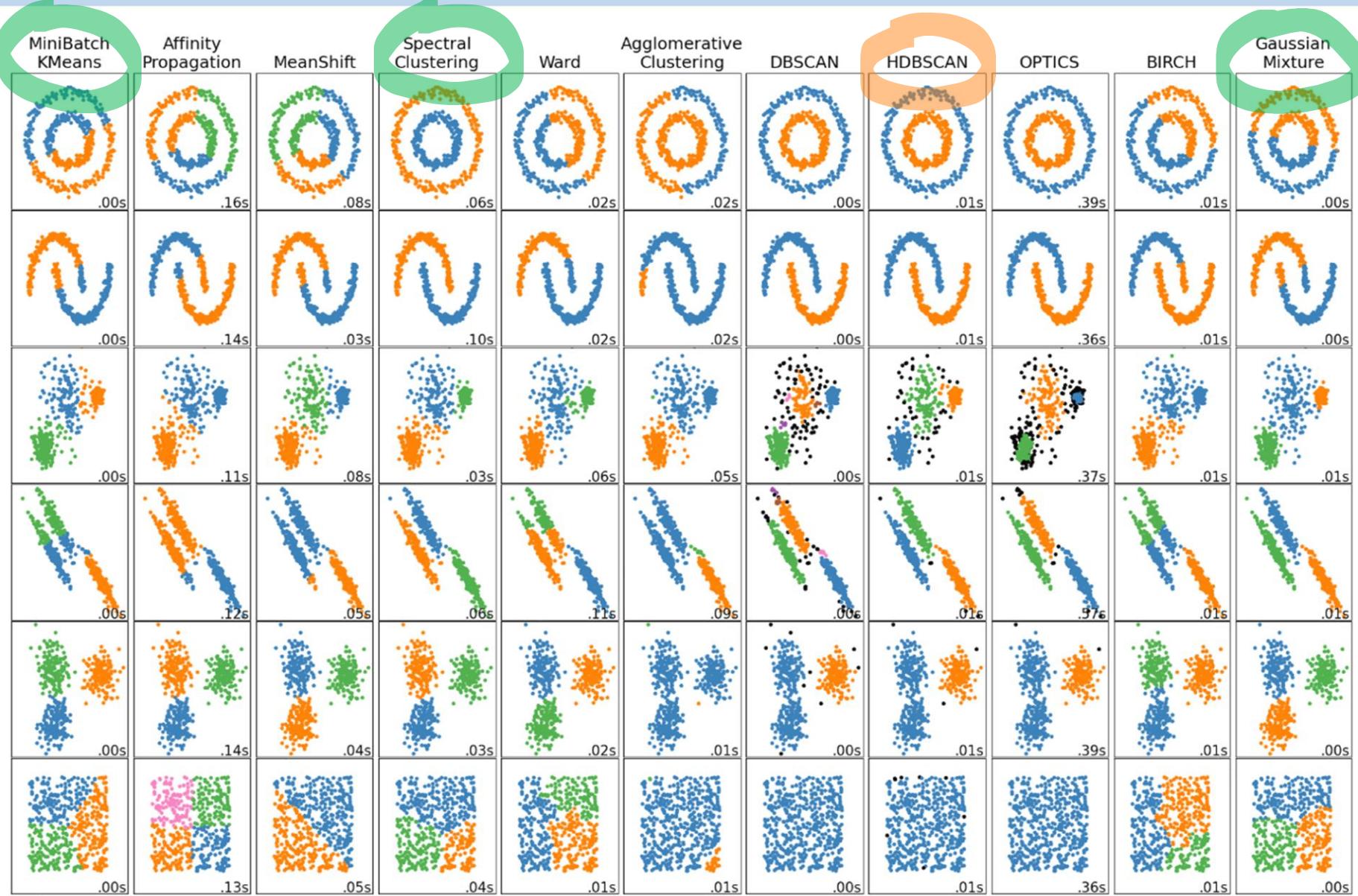


A comparison of the clustering algorithms in scikit-learn

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 <code>MiniBatch</code> code	General-purpose, even cluster size, flat geometry, not too many clusters, inductive	Distances between points
Affinity propagation	damping, sample preference	Not scalable with <code>n_samples</code>	Many clusters, uneven cluster size, non-flat geometry, inductive	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, inductive	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, transductive	Graph distance (e.g. nearest-neighbor graph)
Ward hierarchical clustering	number of clusters or distance threshold	Large <code>n_samples</code> and <code>n_clusters</code>	Many clusters, possibly connectivity constraints, transductive	Distances between points
Agglomerative clustering	number of clusters or distance threshold, linkage type, distance	Large <code>n_samples</code> and <code>n_clusters</code>	Many clusters, possibly connectivity constraints, non Euclidean distances, transductive	Any pairwise distance
DBSCAN	neighborhood size	Very large <code>n_samples</code> , medium <code>n_clusters</code>	Non-flat geometry, uneven cluster sizes, outlier removal, transductive	Distances between nearest points
HDBSCAN	minimum cluster membership, minimum point neighbors	large <code>n_samples</code> , medium <code>n_clusters</code>	Non-flat geometry, uneven cluster sizes, outlier removal, transductive, hierarchical, variable cluster density	Distances between nearest points
OPTICS	minimum cluster membership	Very large <code>n_samples</code> , large <code>n_clusters</code>	Non-flat geometry, uneven cluster sizes, variable cluster density, outlier removal, transductive	Distances between points
Gaussian mixtures	many	Not scalable	Flat geometry, good for density estimation, inductive	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, inductive	Euclidean distance between points
Bisecting K-Means	number of clusters	Very large <code>n_samples</code> , medium <code>n_clusters</code>	General-purpose, even cluster size, flat geometry, no empty clusters, inductive, hierarchical	Distances between points

lesson 2

2.3.1. Overview of clustering methods

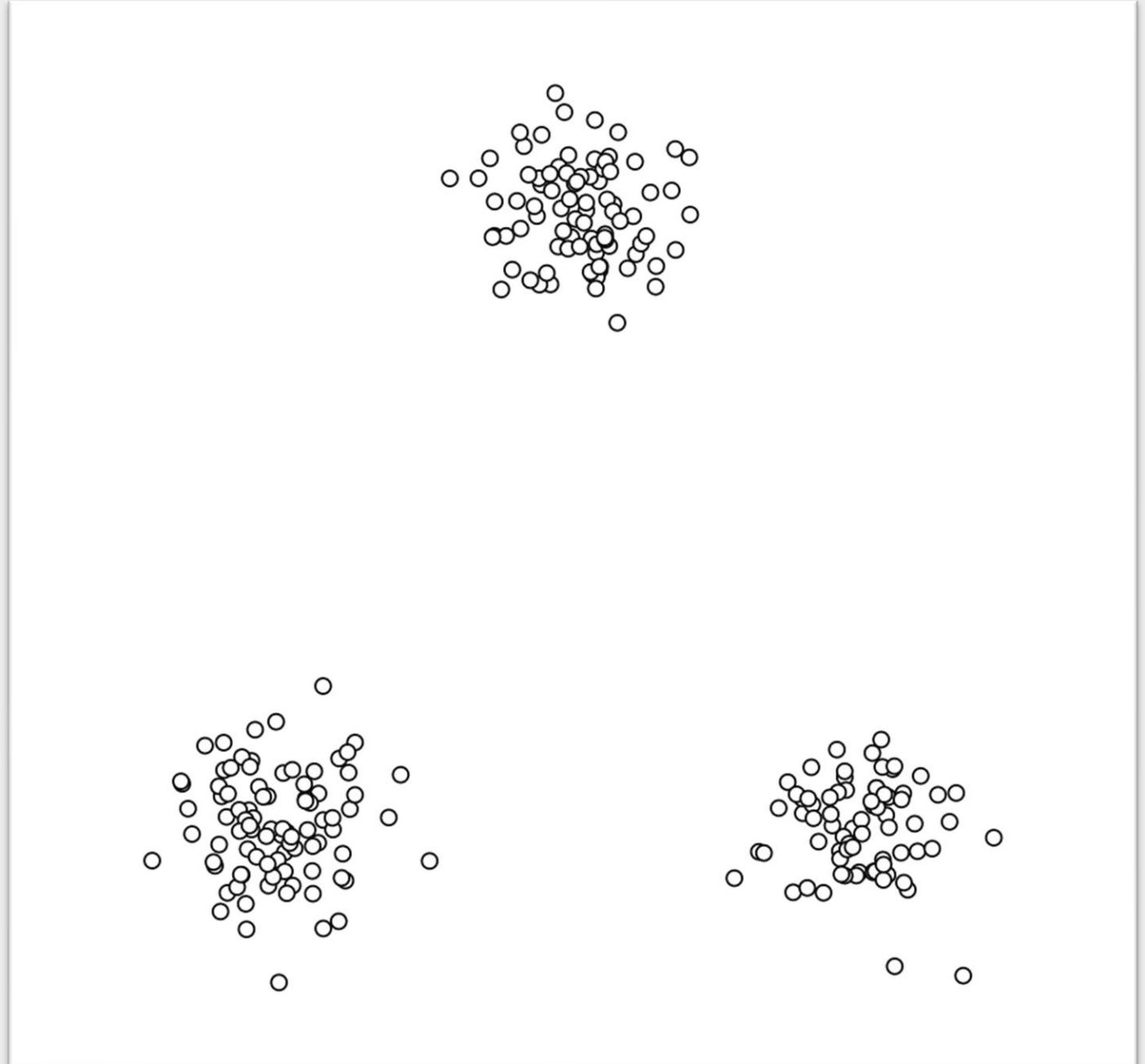


A comparison of the clustering algorithms in scikit-learn

K-means

How many groups
are there?

`sklearn.cluster.KMeans`



<https://www.naftaliharris.com/blog/visualizing-k-means-clustering/>

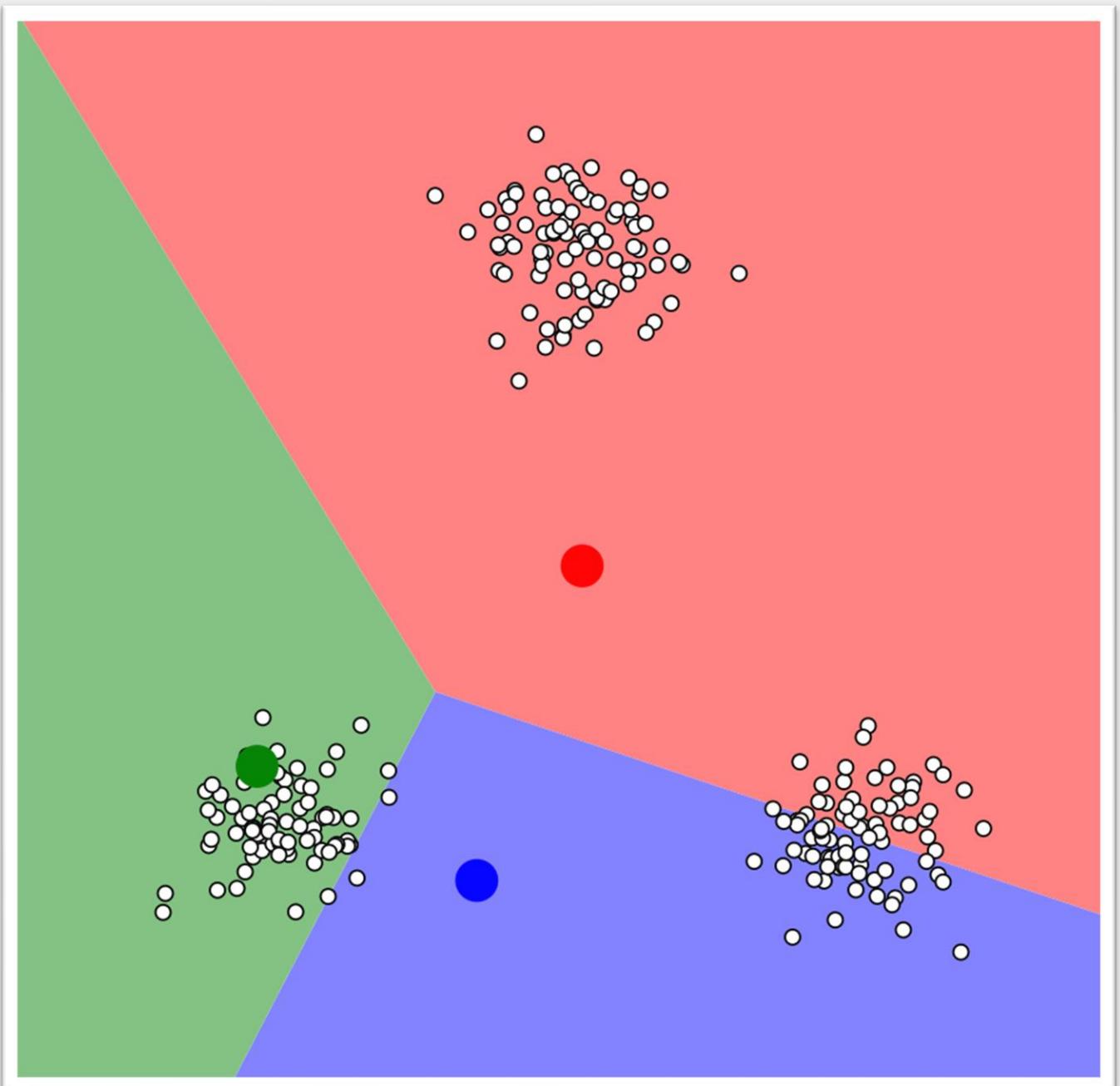
K-means

Take a guess
on where the cluster
centers are

Options

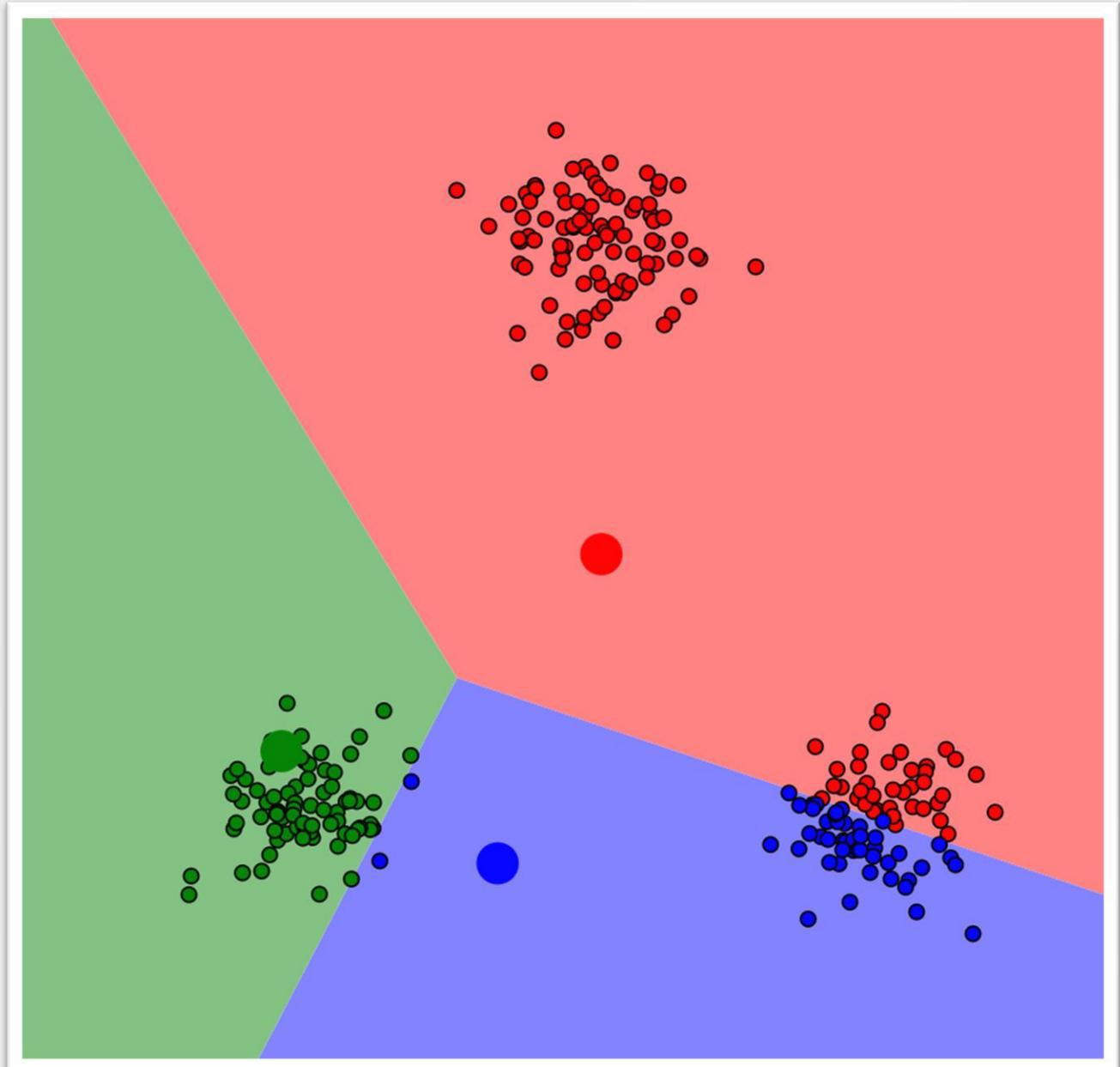
- o Random
- o User-defined
- o Optimal ([kmeans++](#))

[`sklearn.cluster.KMeans`](#)



K-means

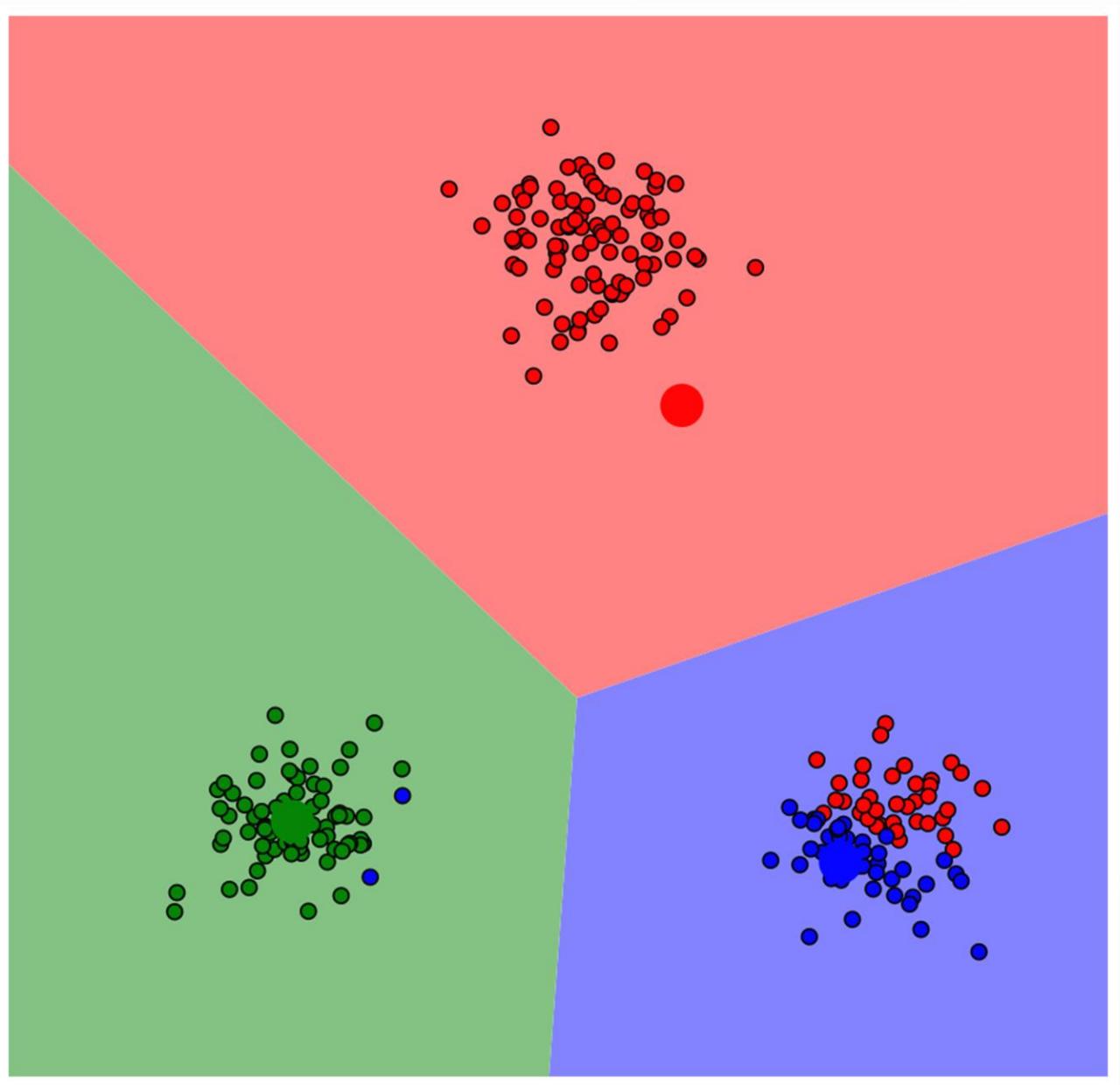
Assign points
based on their
distance to the cluster
center



`sklearn.cluster.KMeans`

K-means

Update centroids
based on the average
position of all
the points in the cluster

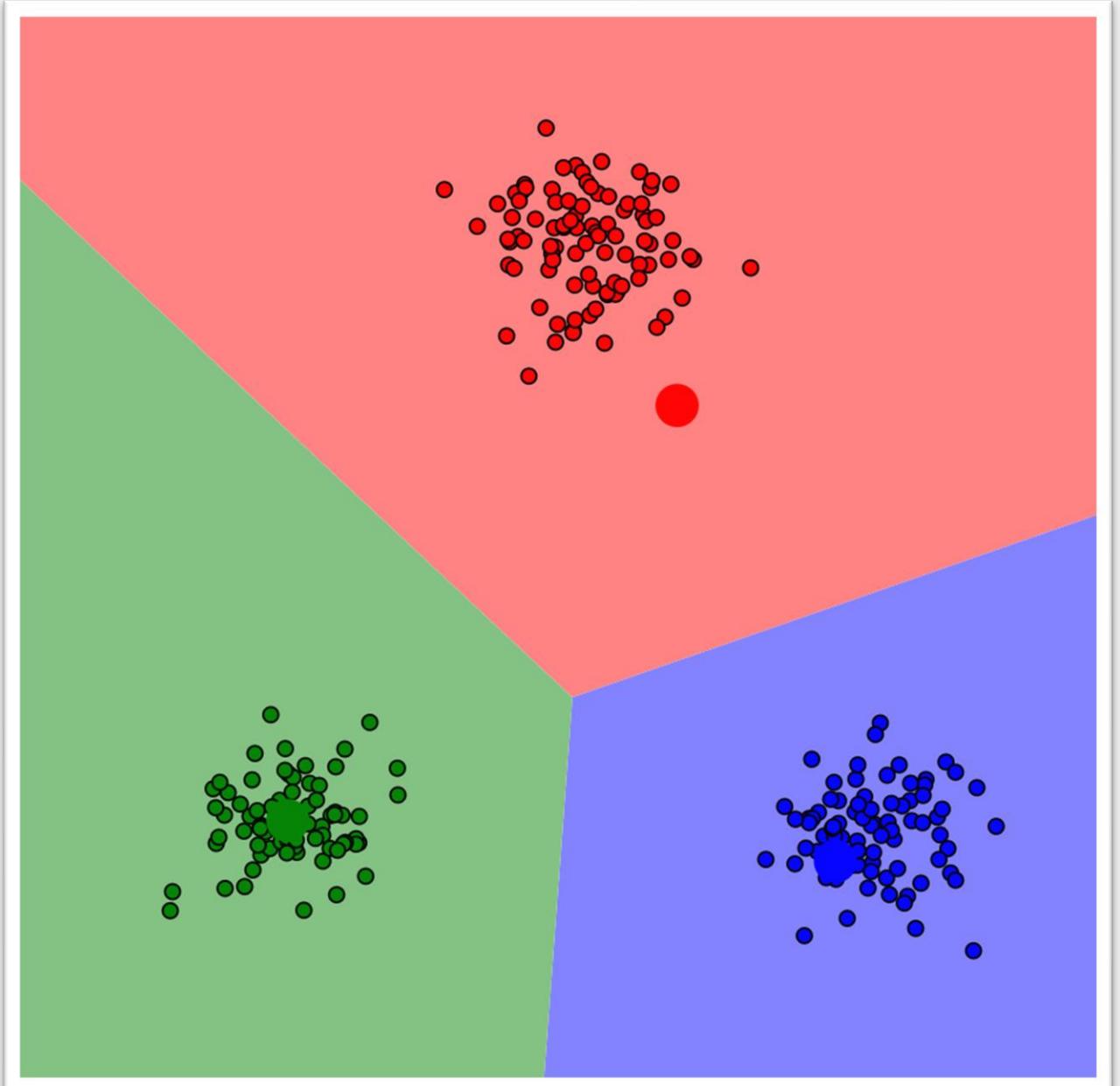


`sklearn.cluster.KMeans`

K-means

Reassign points...

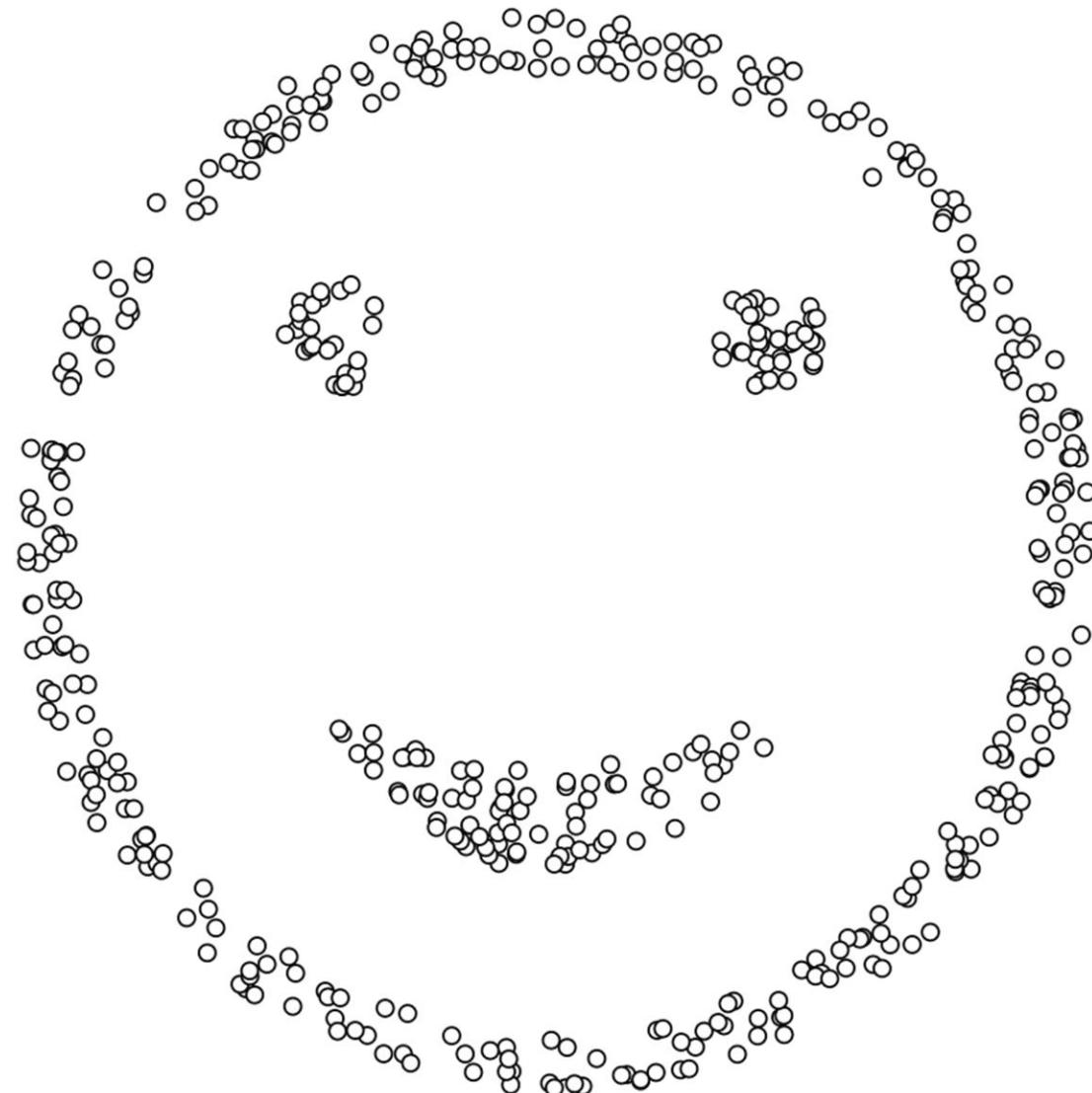
Rinse and repeat N times...



`sklearn.cluster.KMeans`

K-means

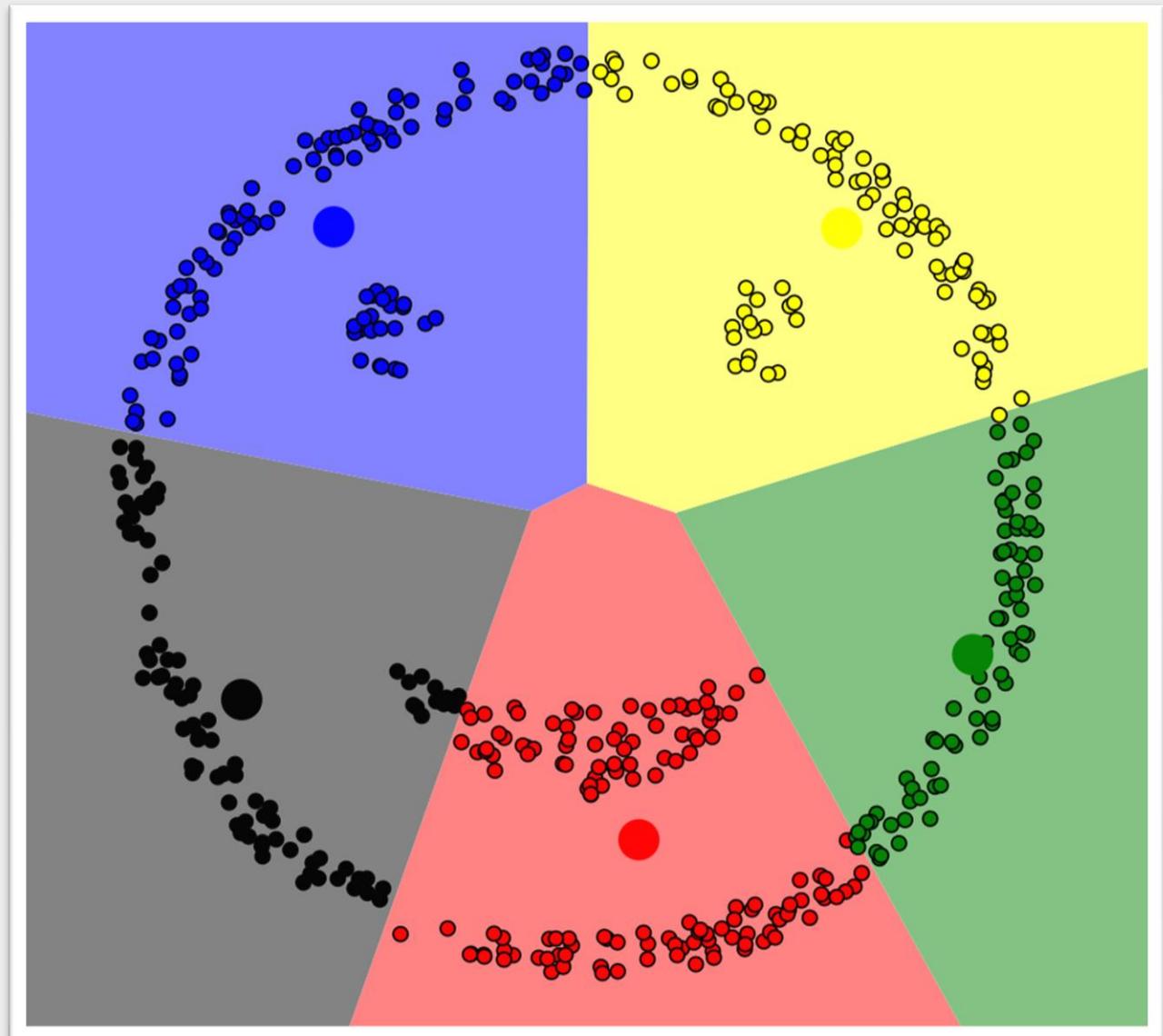
How many groups
are there?



`sklearn.cluster.KMeans`

K-means

What about this data?



`sklearn.cluster.KMeans`

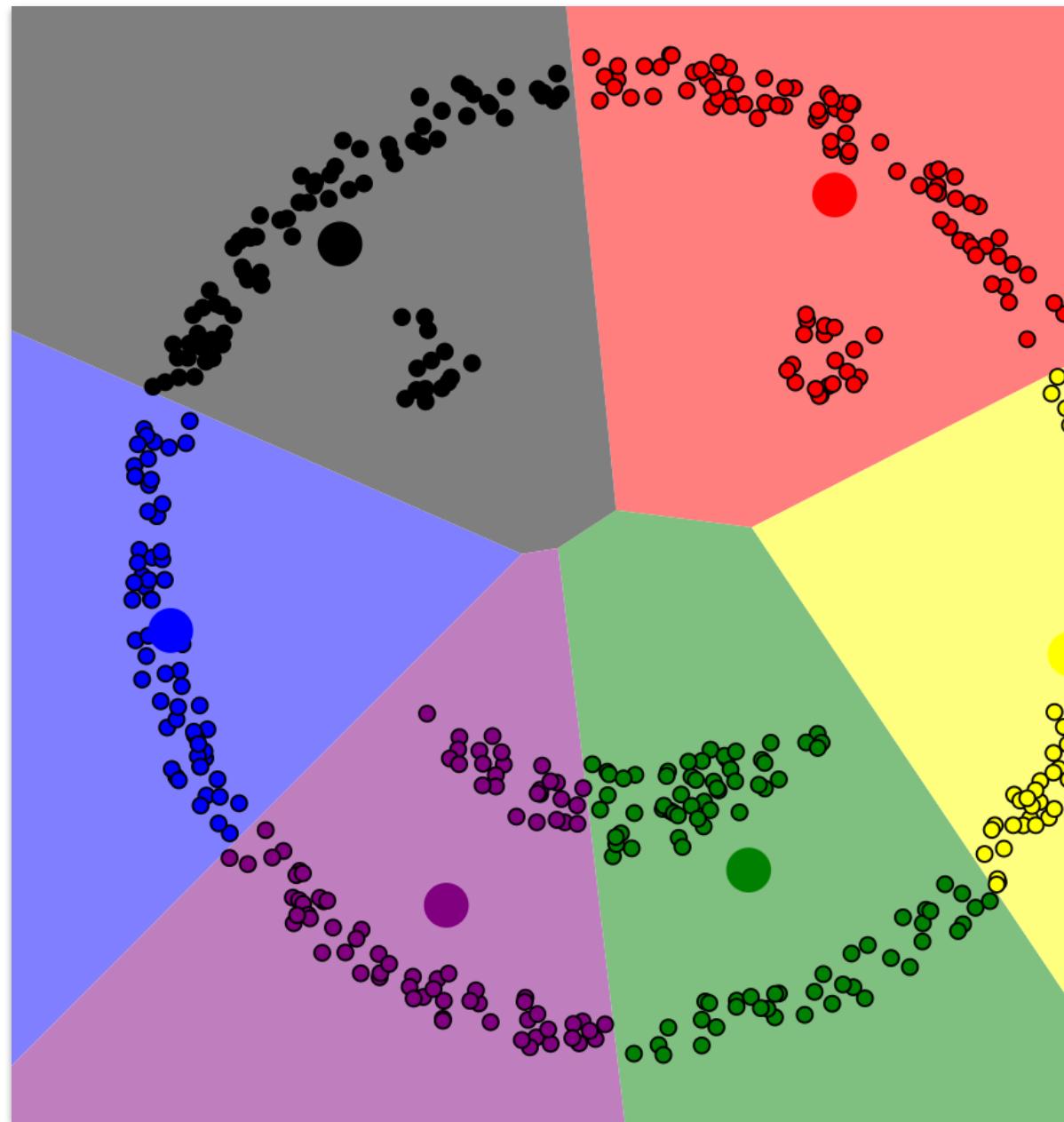
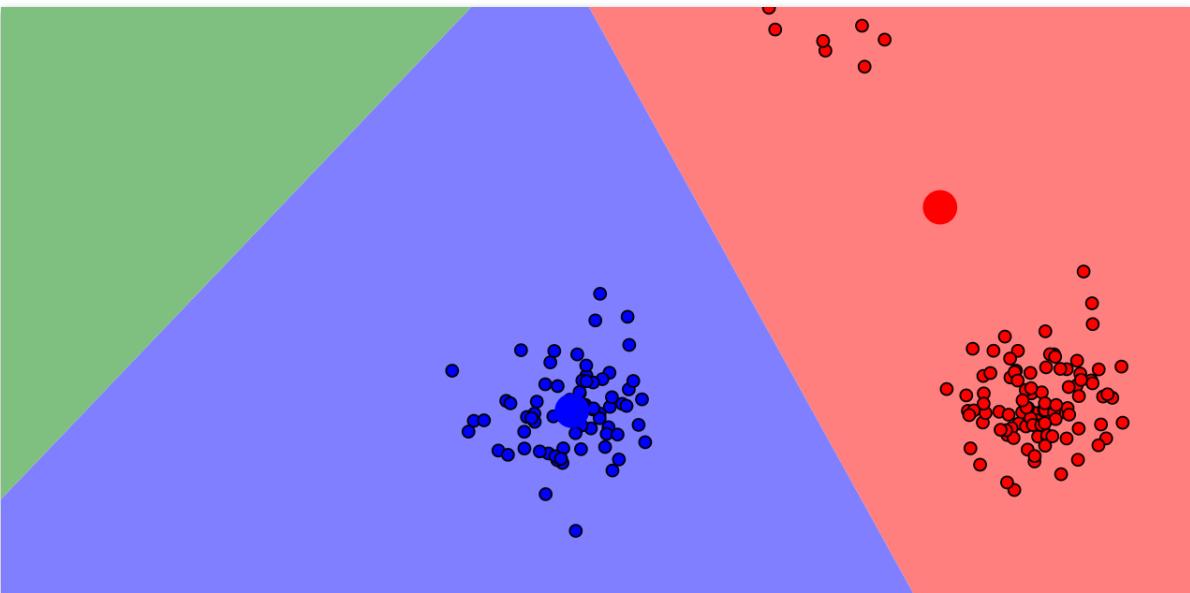
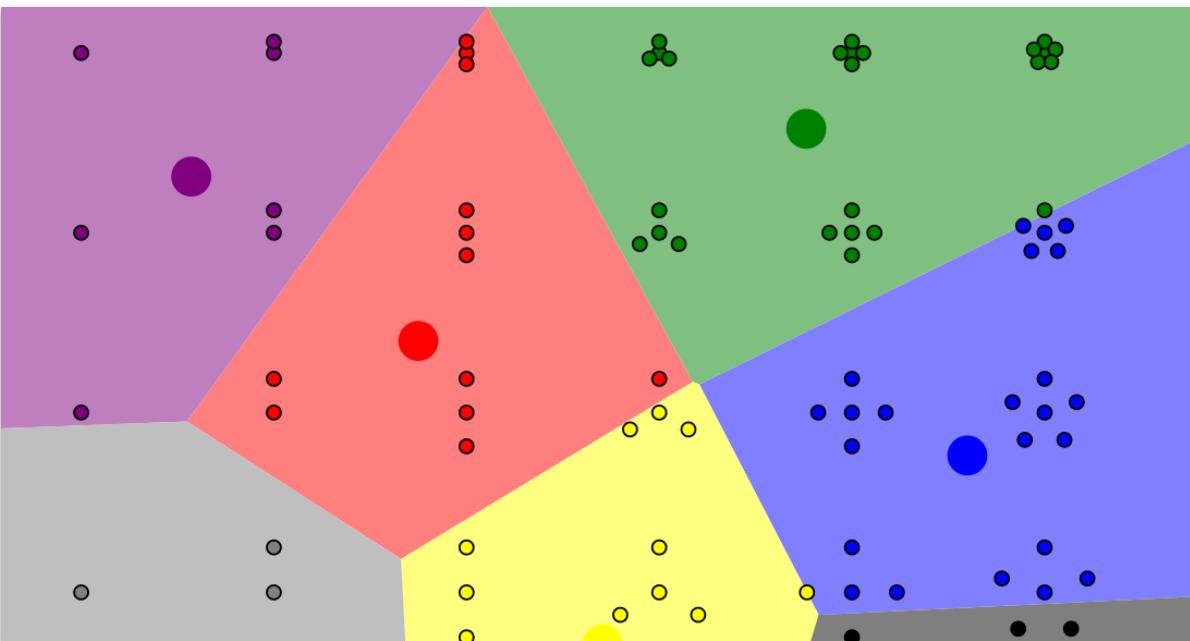
K-means

pros

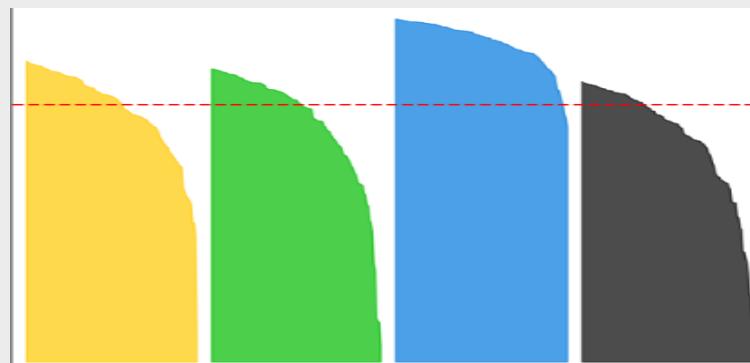
- easy
- really fast!
- easy to interpret
- great in many dims
- no assumption
of cluster shape**

cons

- need to decide on k
- depends on initial values
- not great for uneven clusters
- outliers
- complex patterns

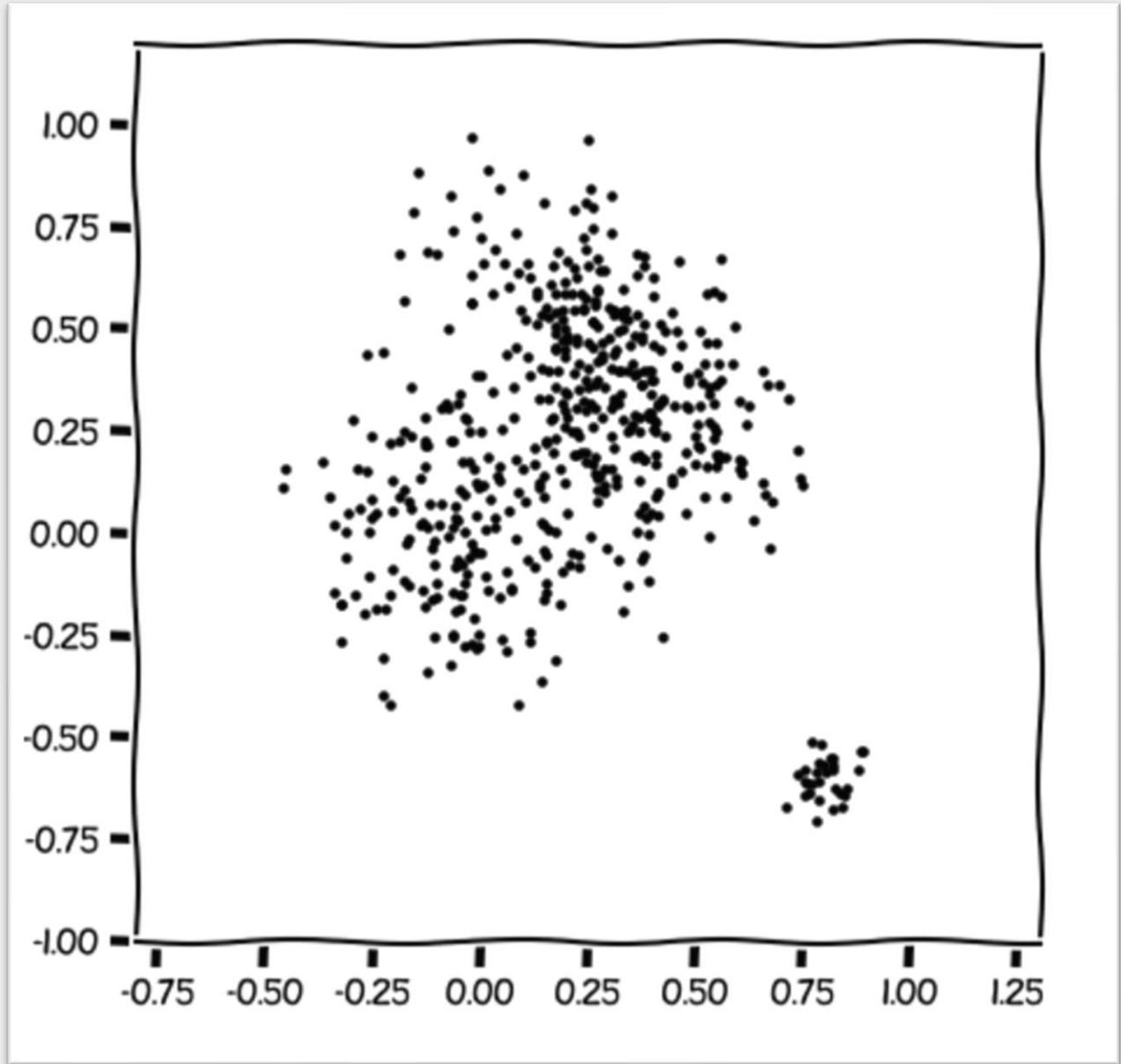


k-means silhouette score



K-means

How many groups
are there?



`numpy.random.multivariate_normal`

K-means

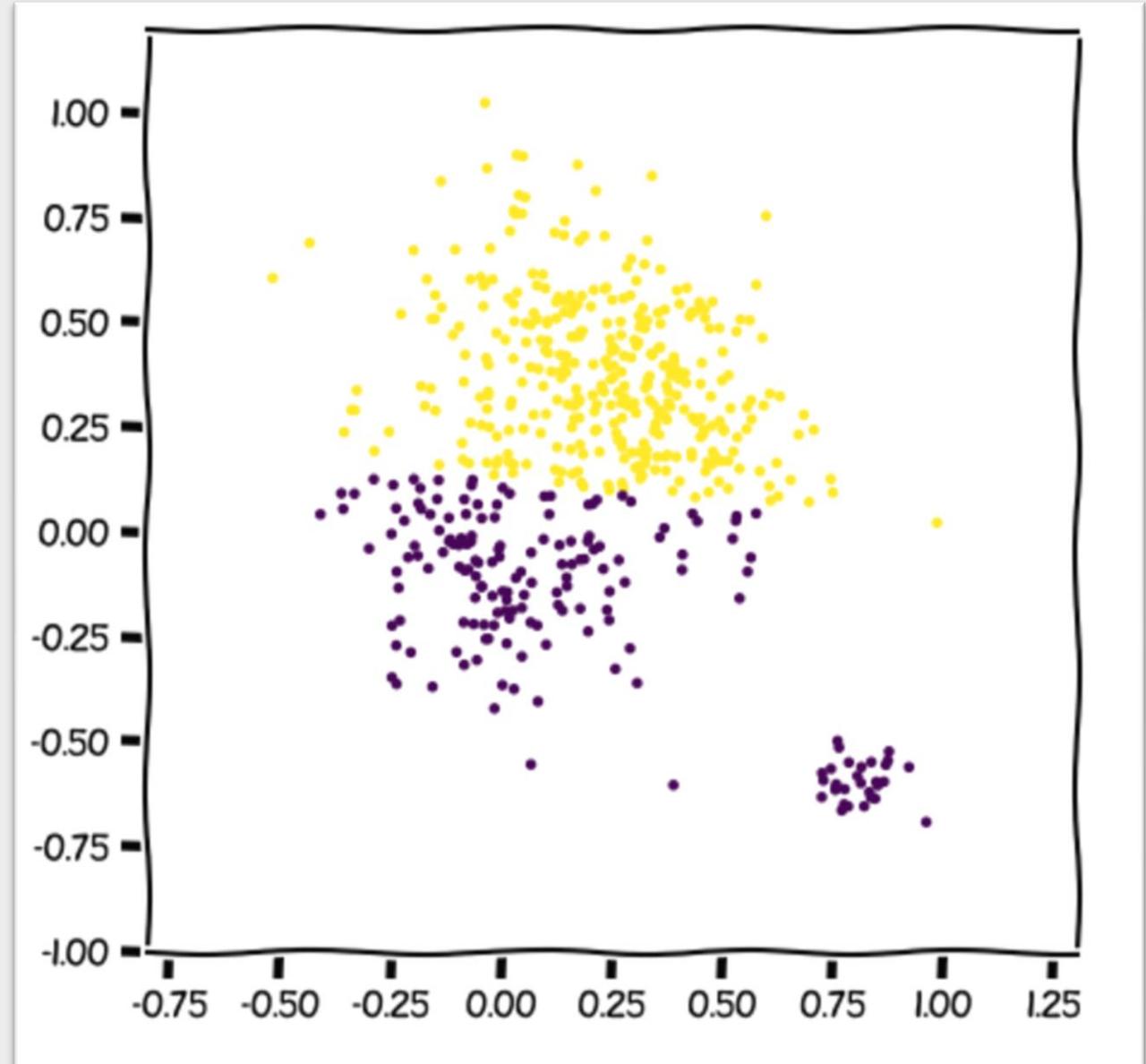
Fit results

Parameters

- Two clusters
- Automatic centroids (kmeans++)

Lesson learnt

Do not blindly trust your kmeans results when making a lesson plan...



`numpy.random.multivariate_normal`

K-means

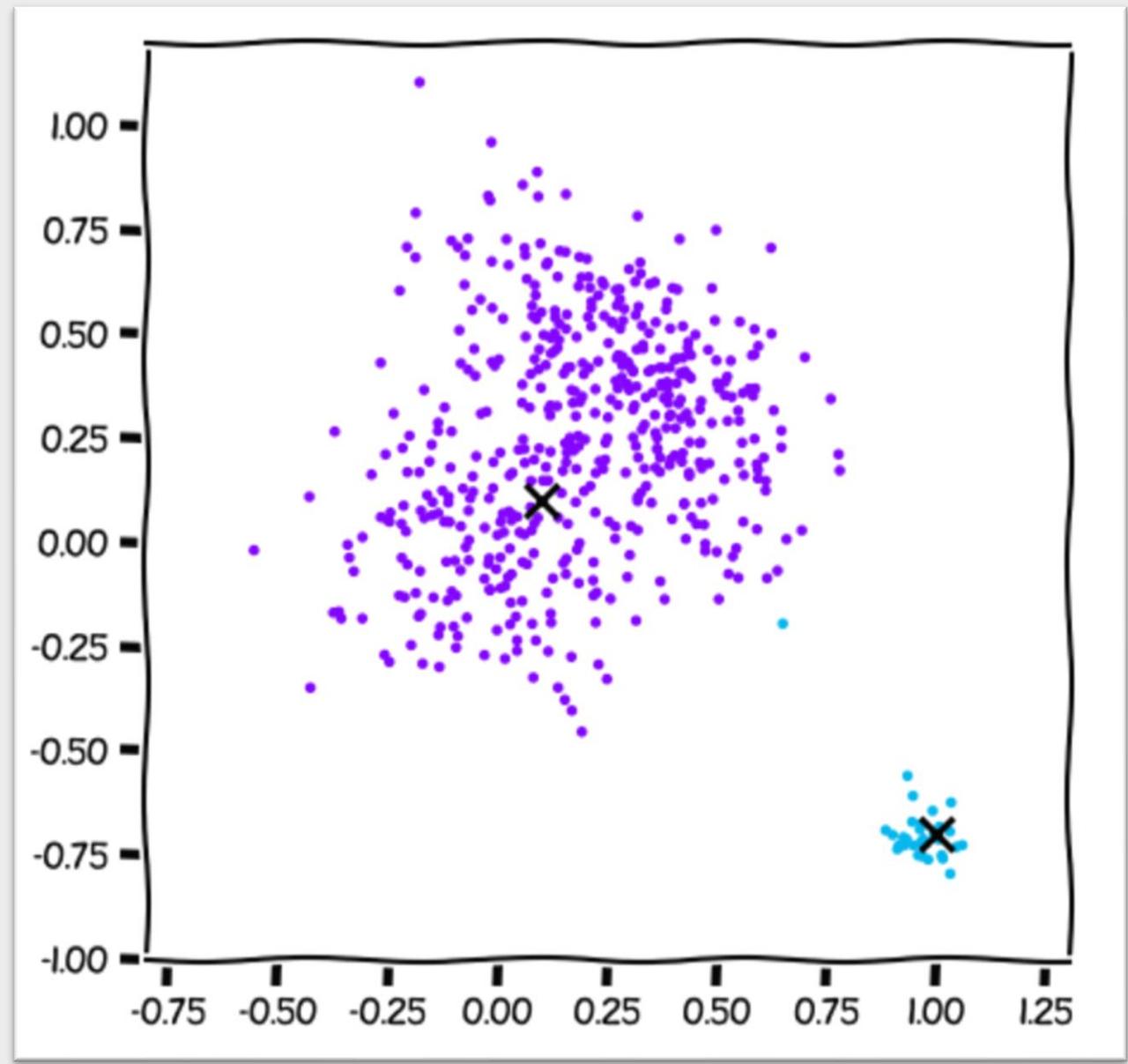
Fit results

Parameters

- Two clusters
- Custom centroids
- Moved the small cluster a little further away to help K-means do its thing

Lesson learnt

It really doesn't like having uneven clusters...

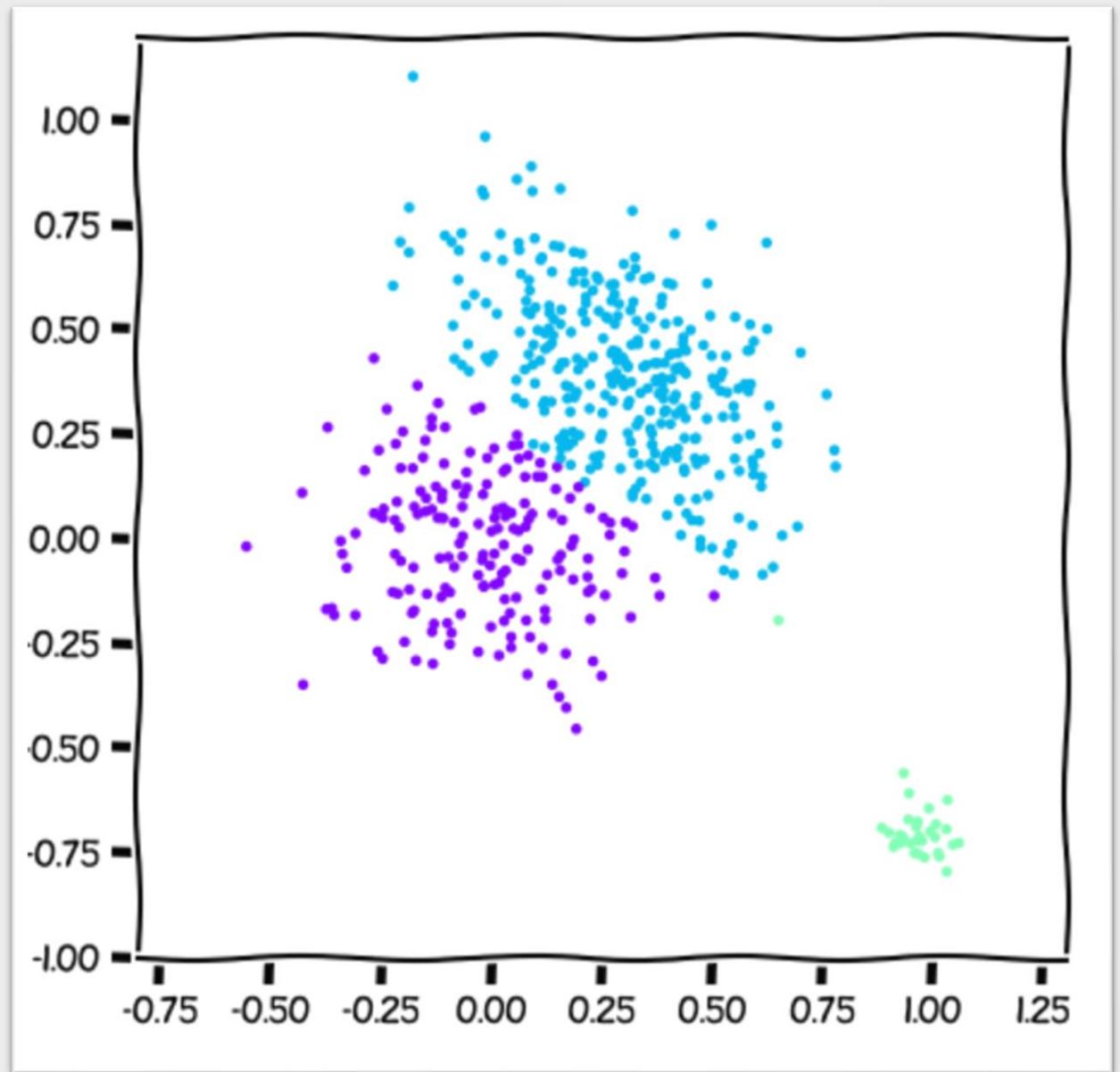


K-means

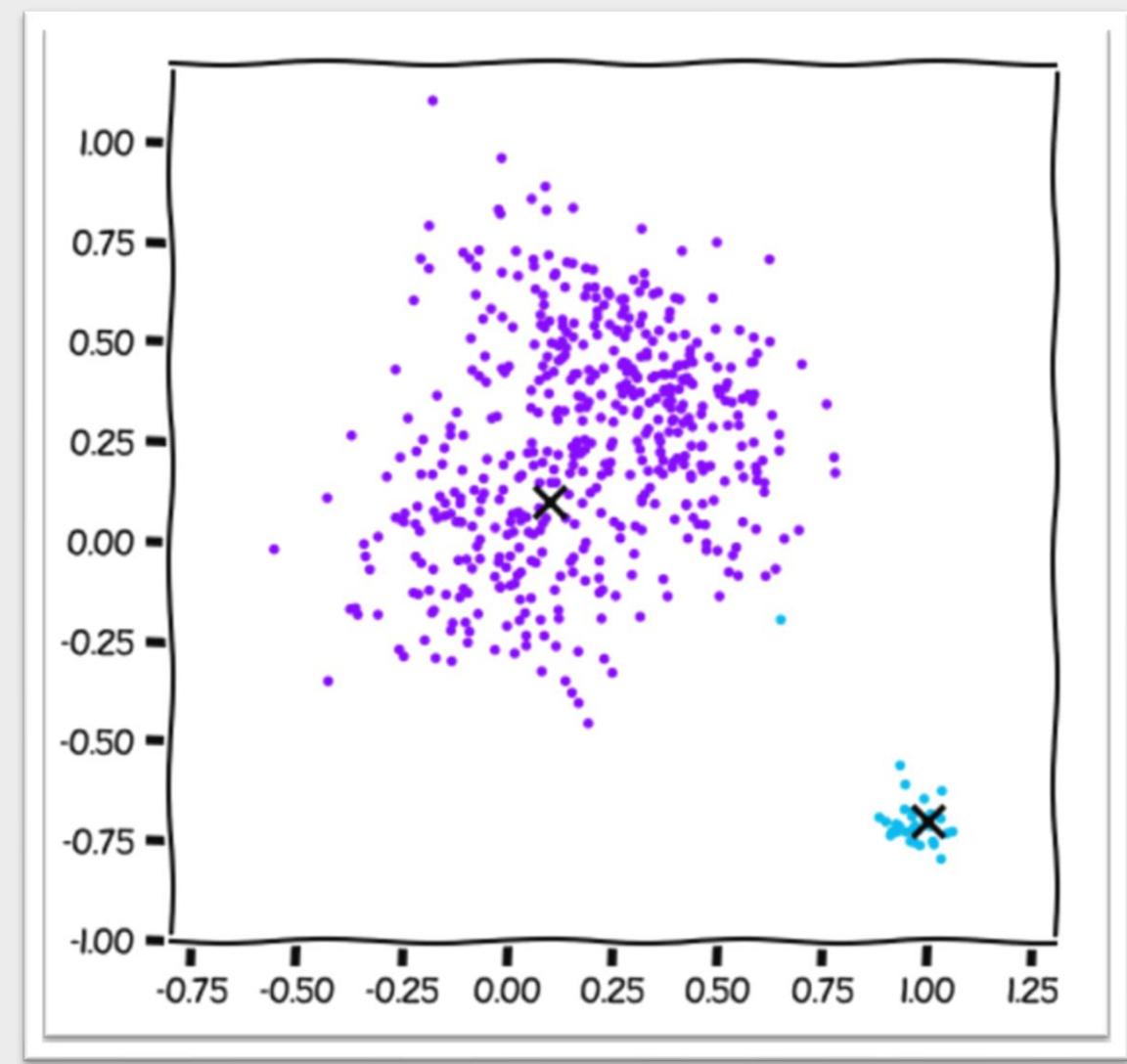
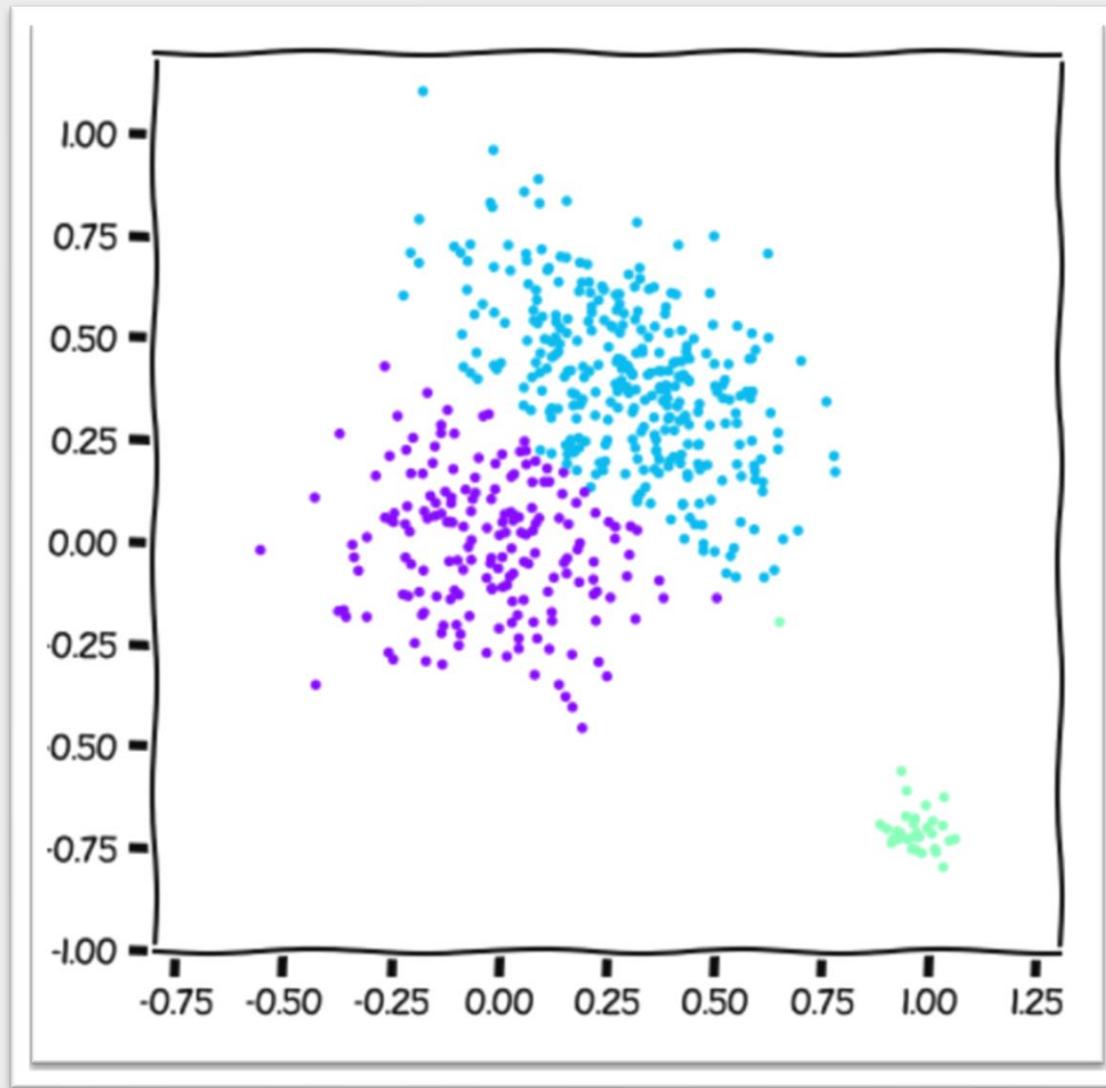
Fit results

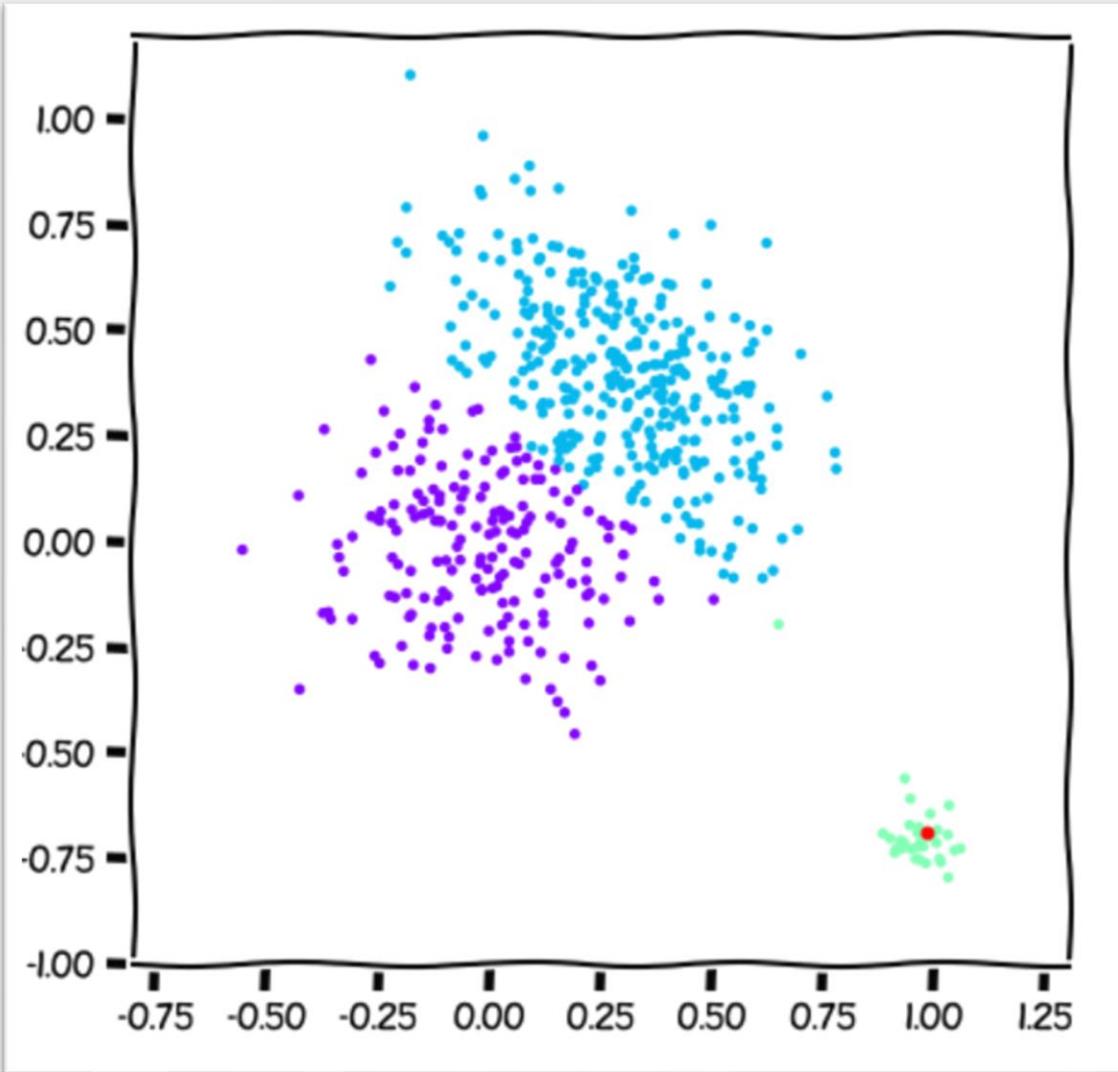
Parameters

- Three clusters
- kmeans++ centroids



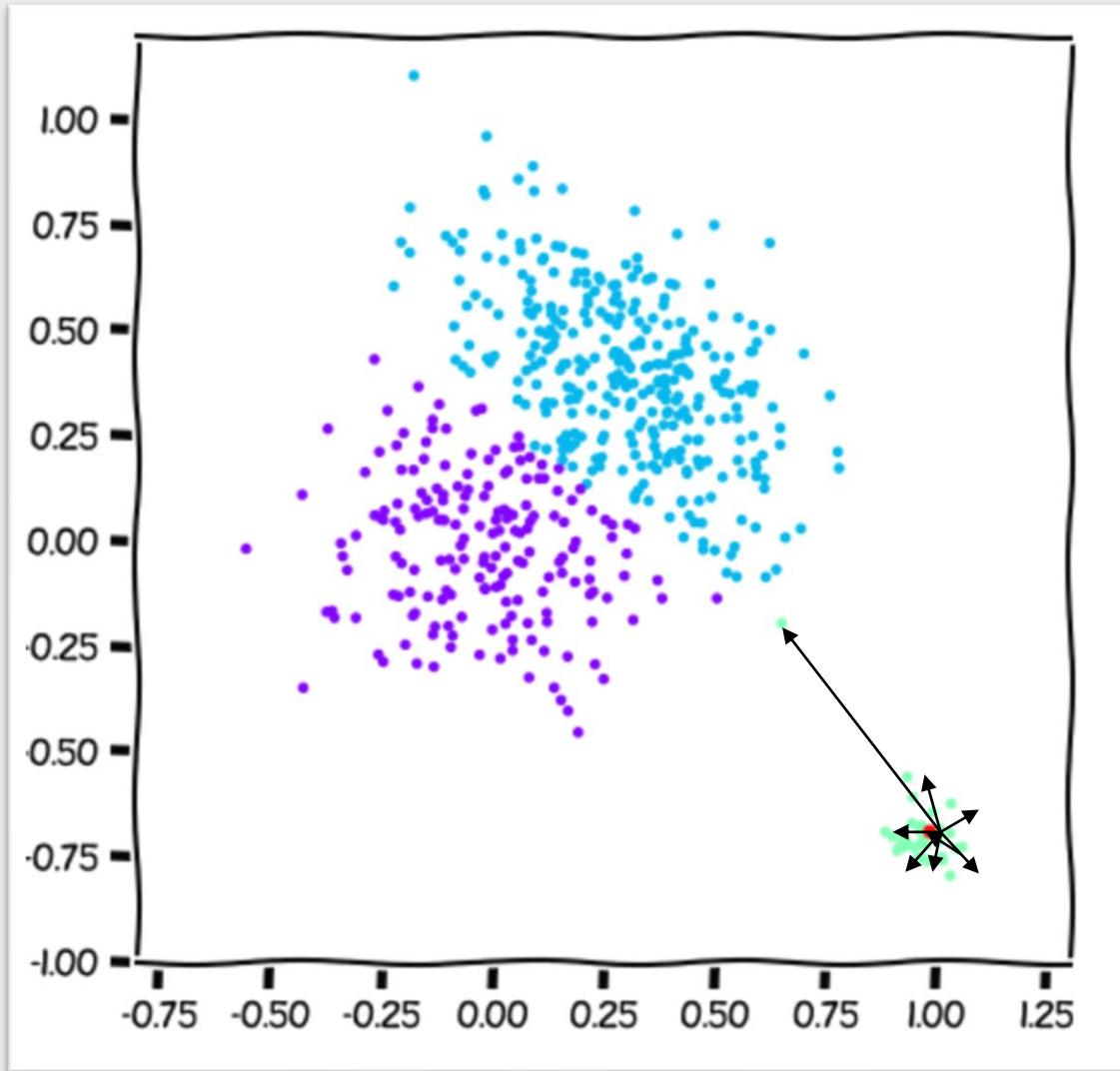
which fit is better?





Silhouette score

Calculated for each point in each cluster

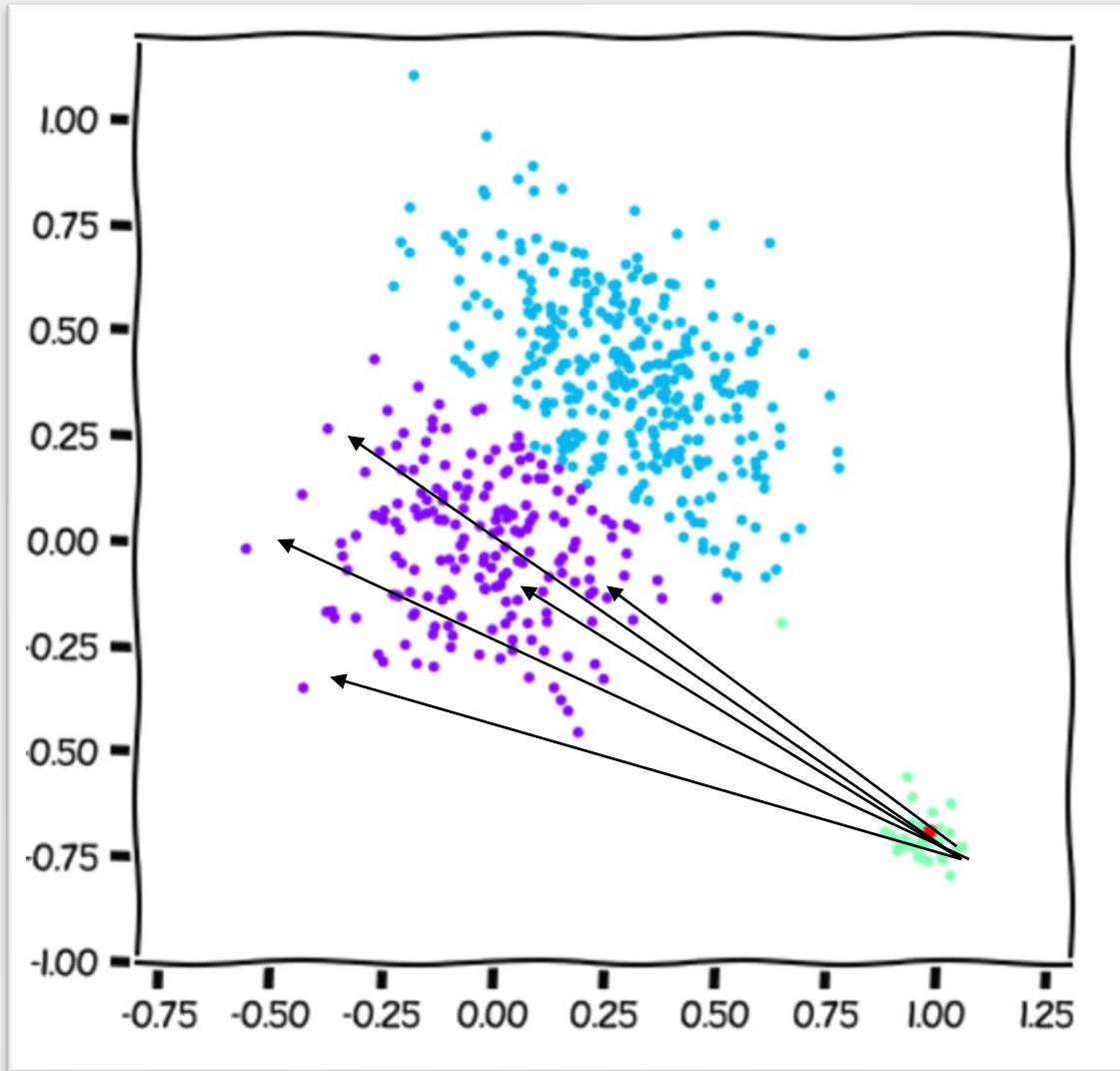


Silhouette score

Calculated for each point in each cluster

A

average distance from the point
to each other point in its cluster

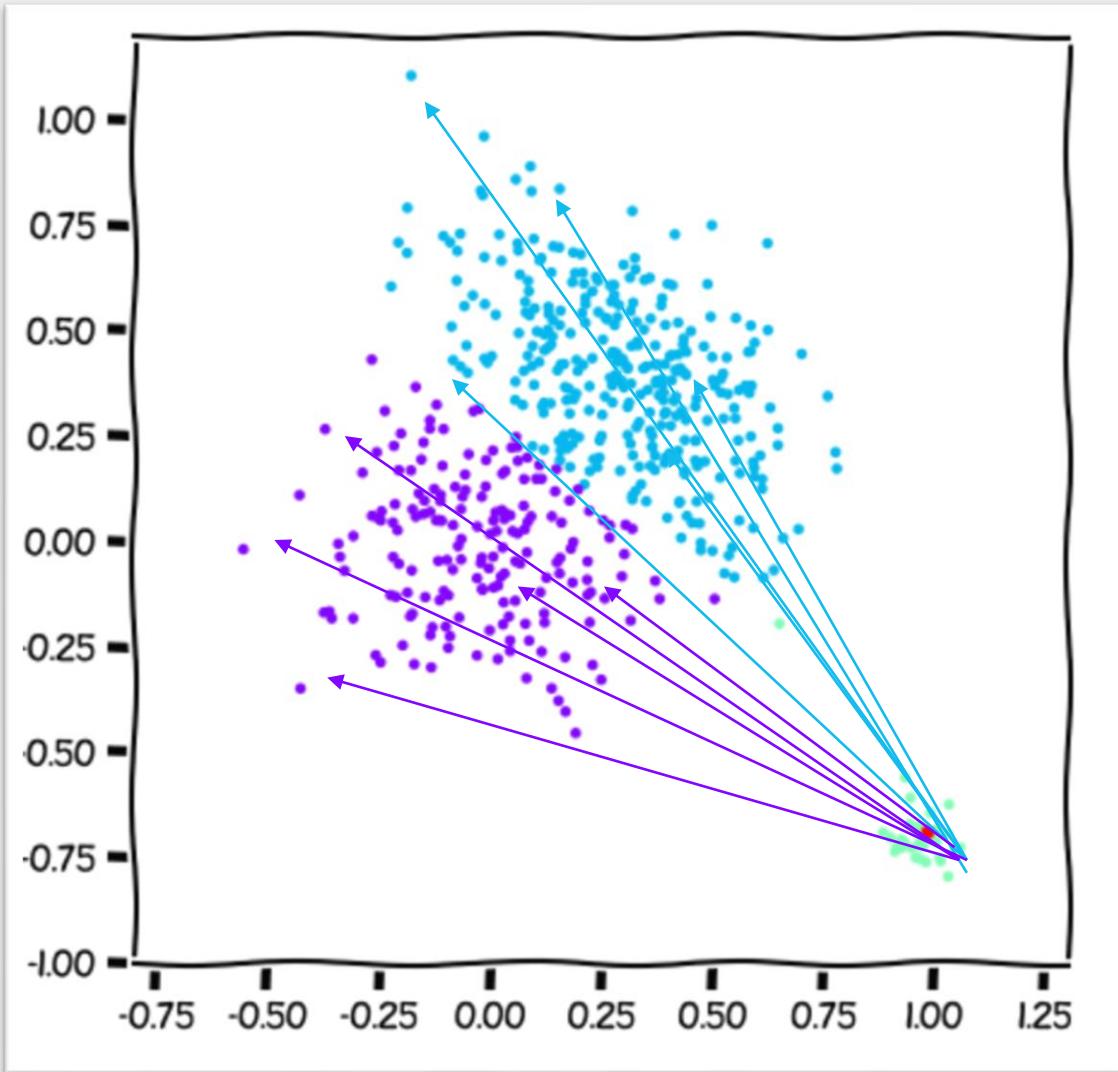


Silhouette score

Calculated for each point in each cluster

A average distance from the point
to each other point in its cluster

B average distance from the point
to each point of another cluster



Silhouette score

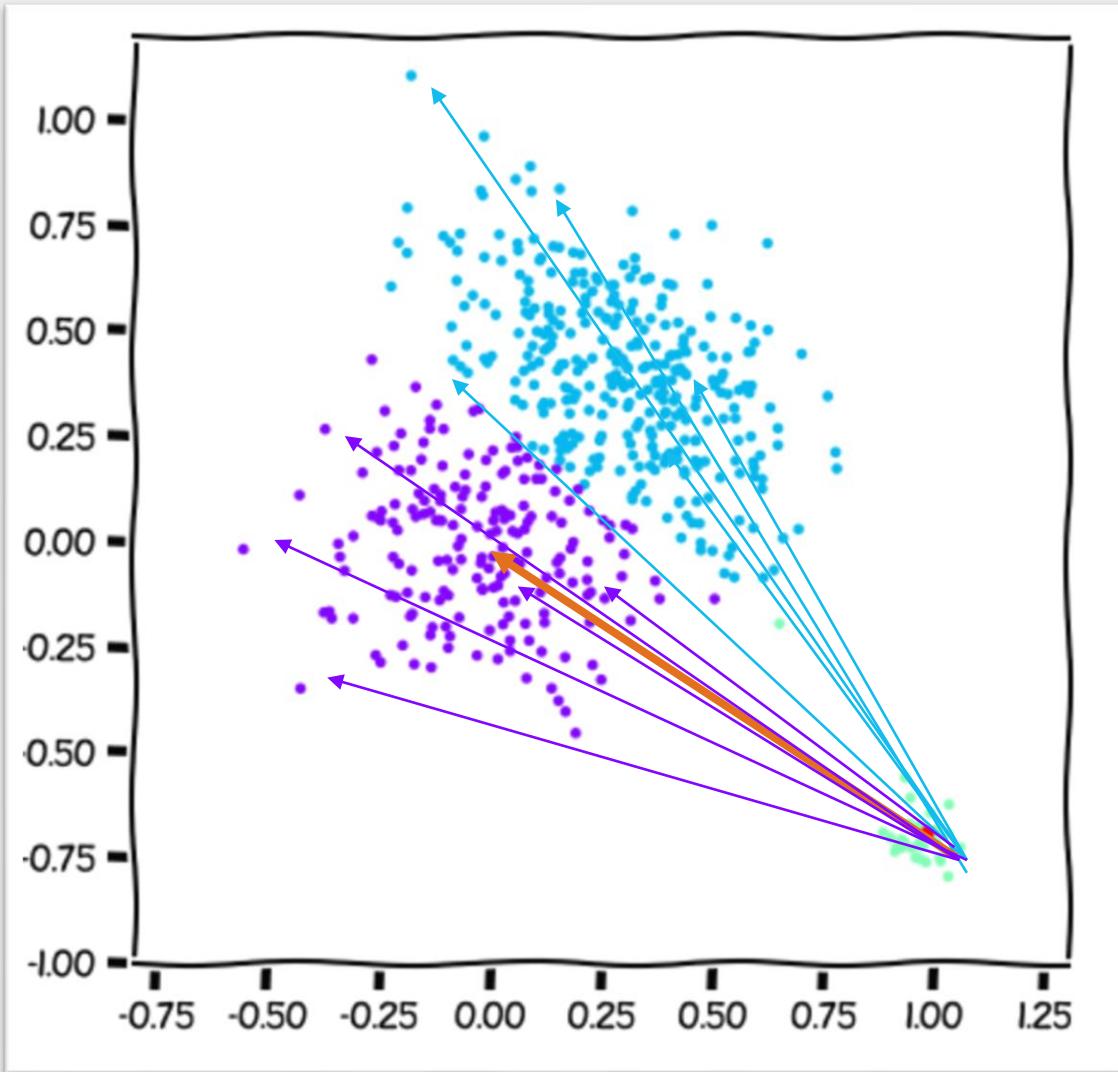
Calculated for each point in each cluster

A

average distance from the point
to each other point in **its cluster**

B₁, B₂

average distance from the point
to each point of **another cluster**



Silhouette score

Calculated for each point in each cluster

A

average distance from the point
to each other point in its cluster

B₁, B₂

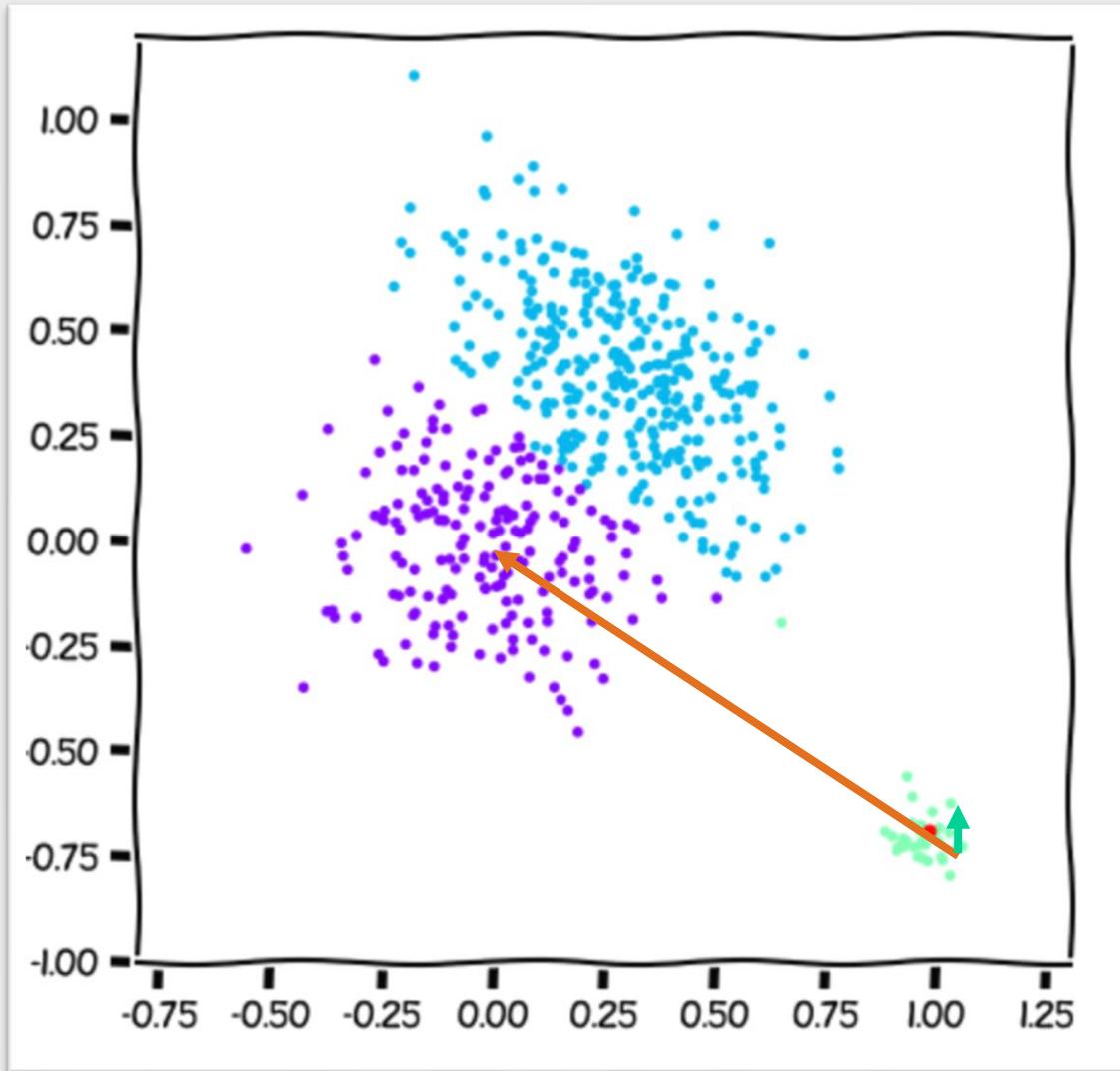
average distance from the point
to each point of another cluster

B

$\min(B_1, B_2, \dots)$

average distance from the point
to the closest cluster

here B = B₁



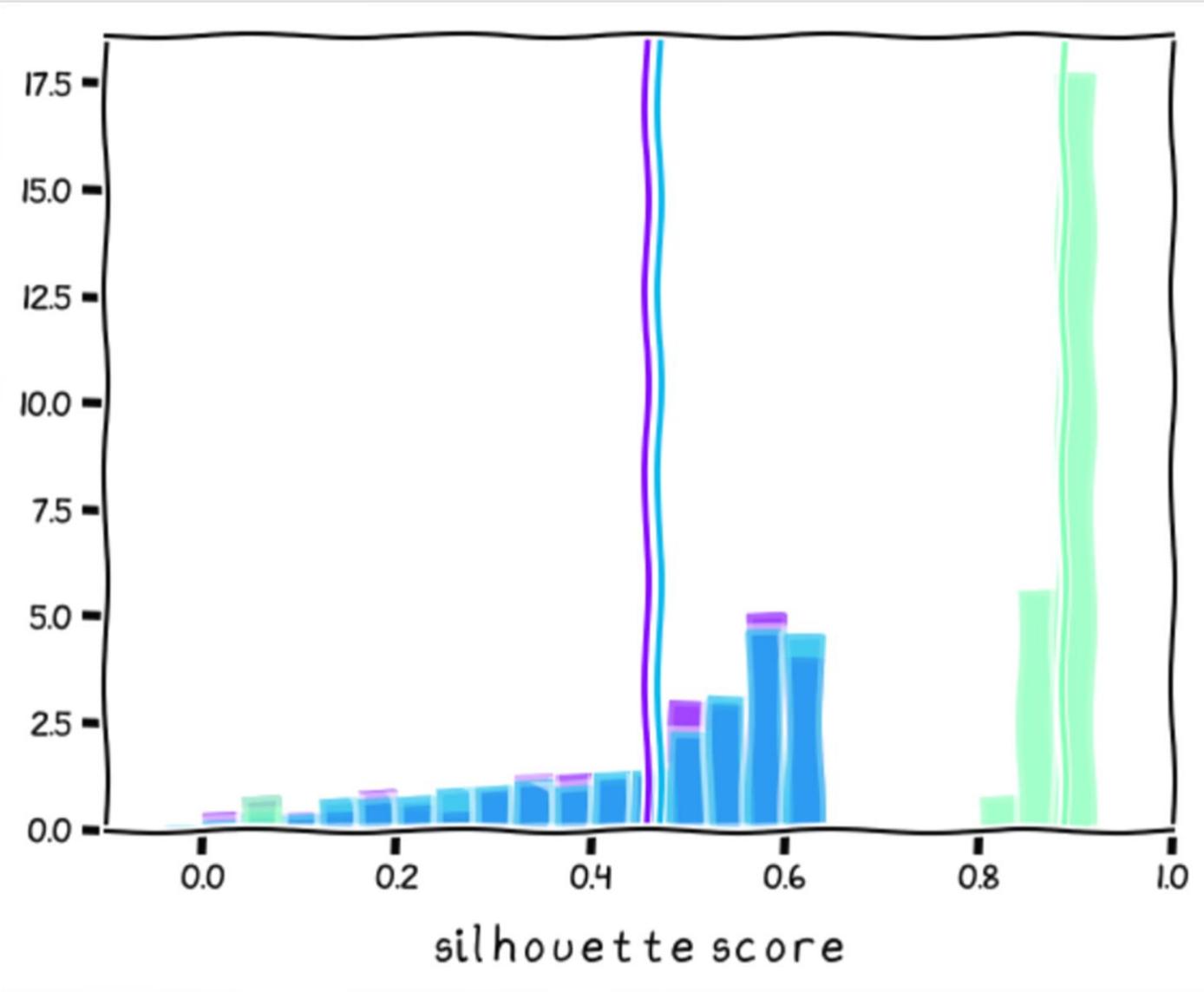
Silhouette score

Calculated for each point in each cluster

A average distance from the point to each other point in **its cluster**

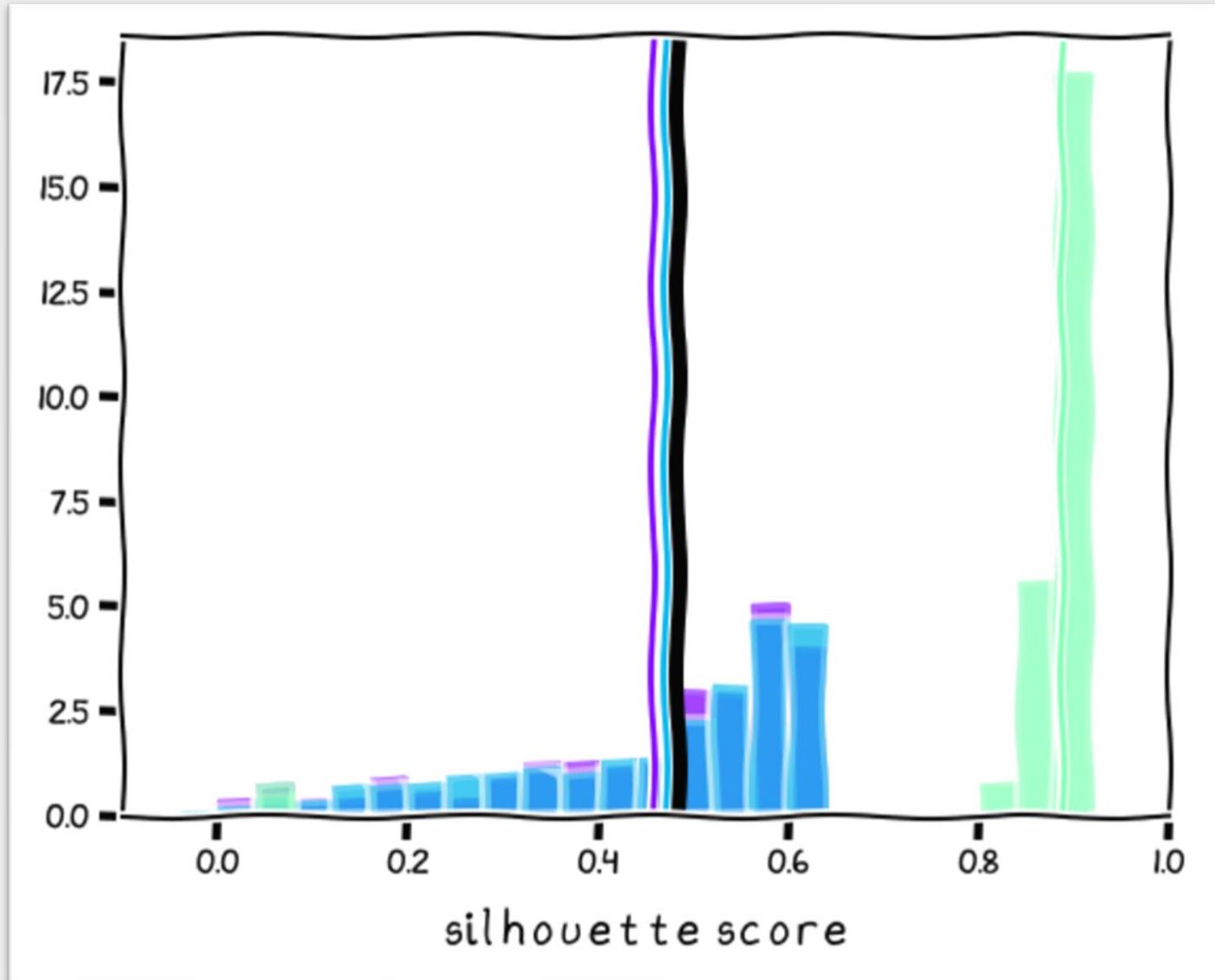
B $\min(B_1, B_2, \dots)$
average distance from the point to the **closest cluster**
here $B = B_1$

$$\text{score} = \frac{B - A}{\max[B, A]} = 0.9$$



Do this for every point in
one cluster

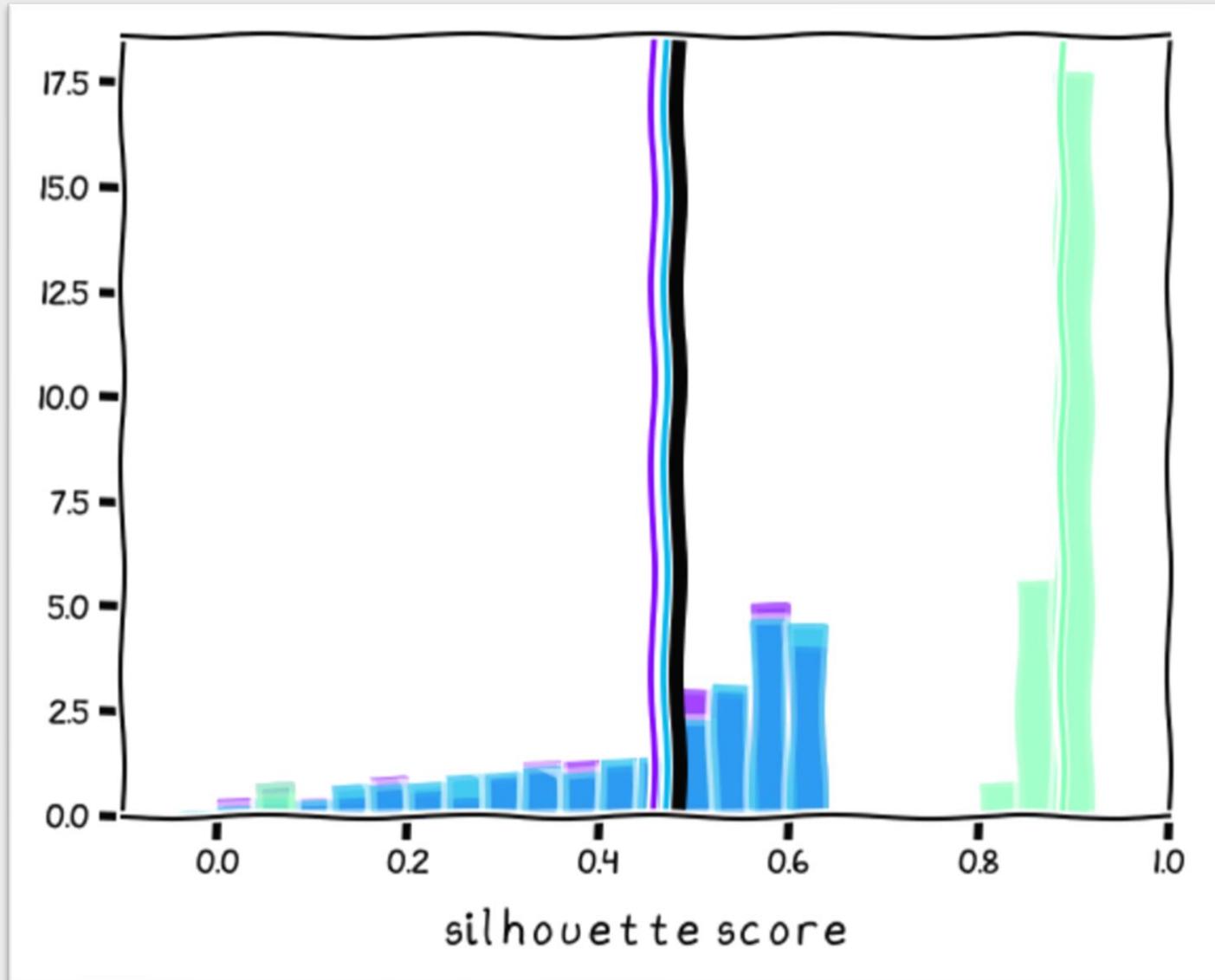
...and then every point in
the other clusters...



Do this for every point in one cluster

...and then every point in the other clusters...

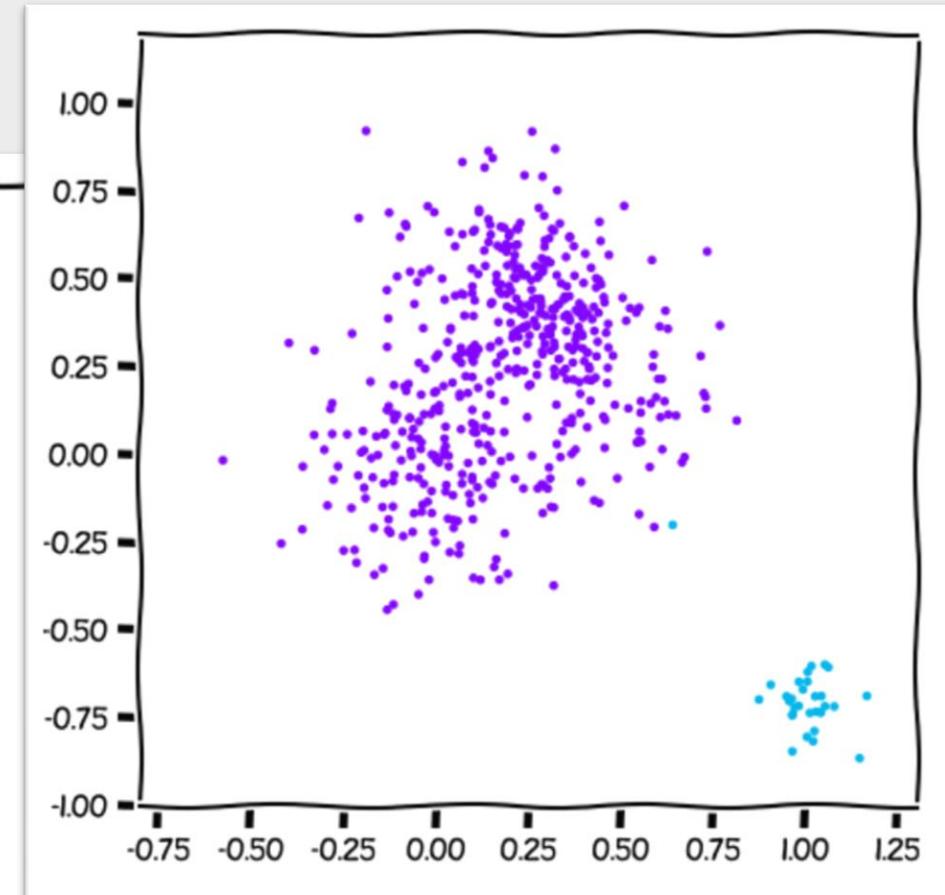
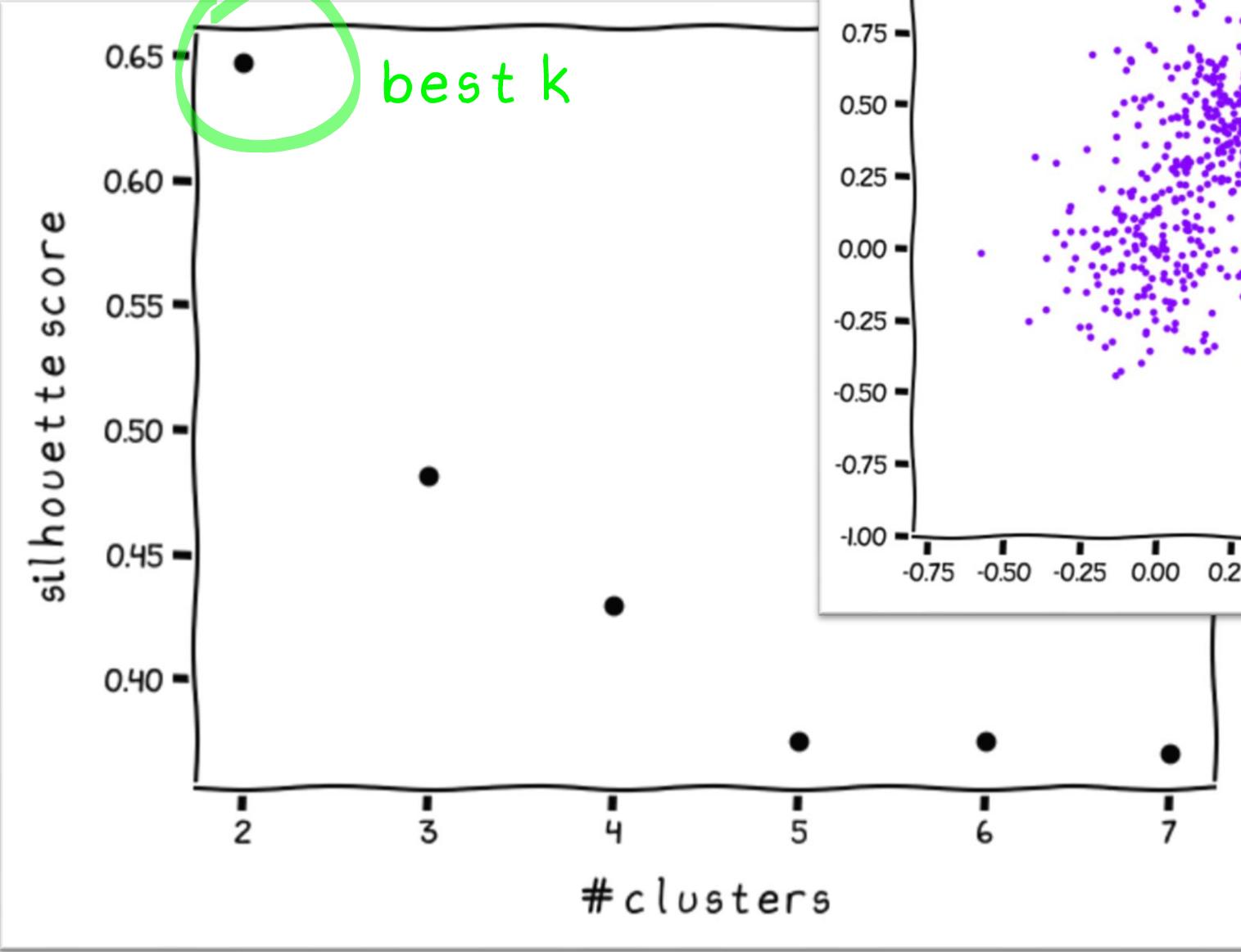
and find the **mean score** for the entire dataset



Do this for every point in one cluster

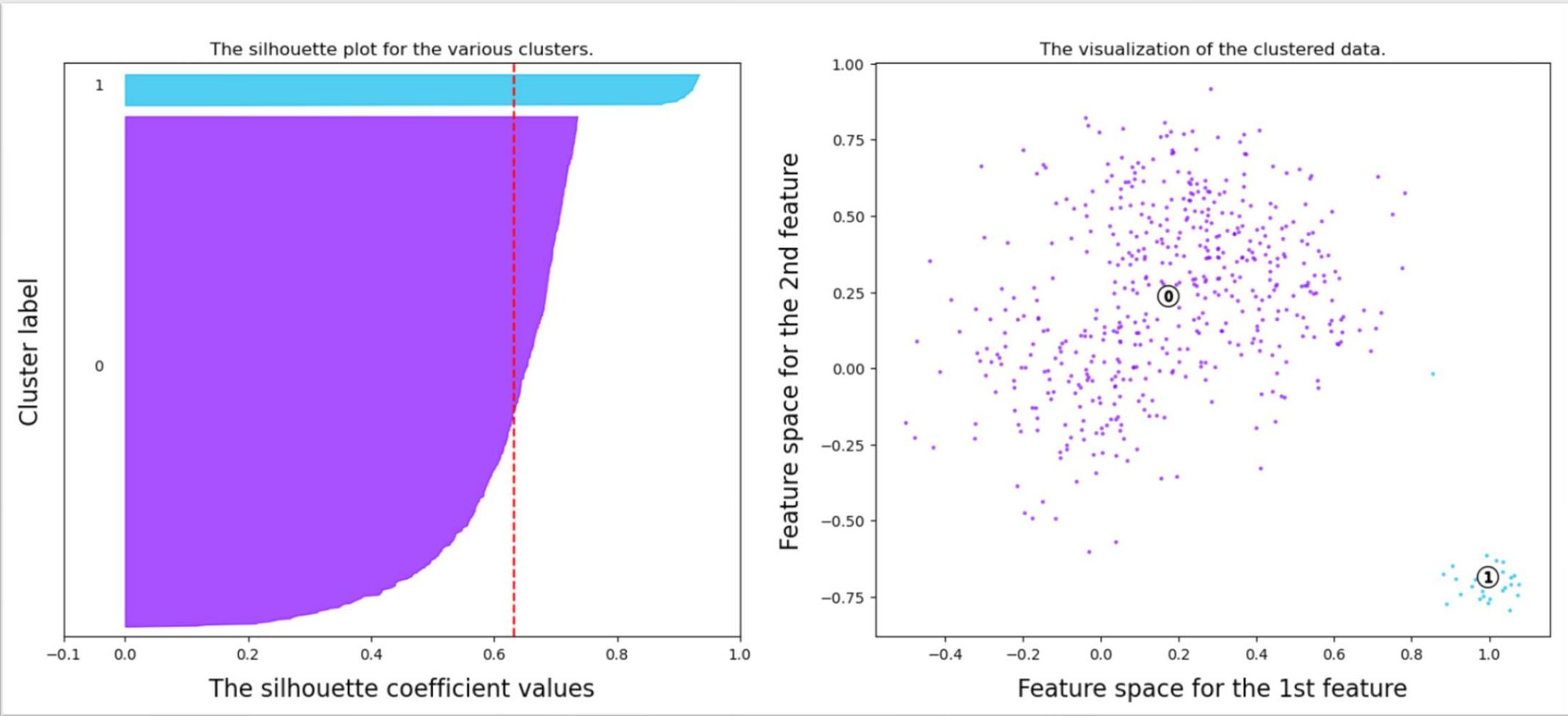
...and then every point in the other clusters...

and find the **mean score** for the entire dataset



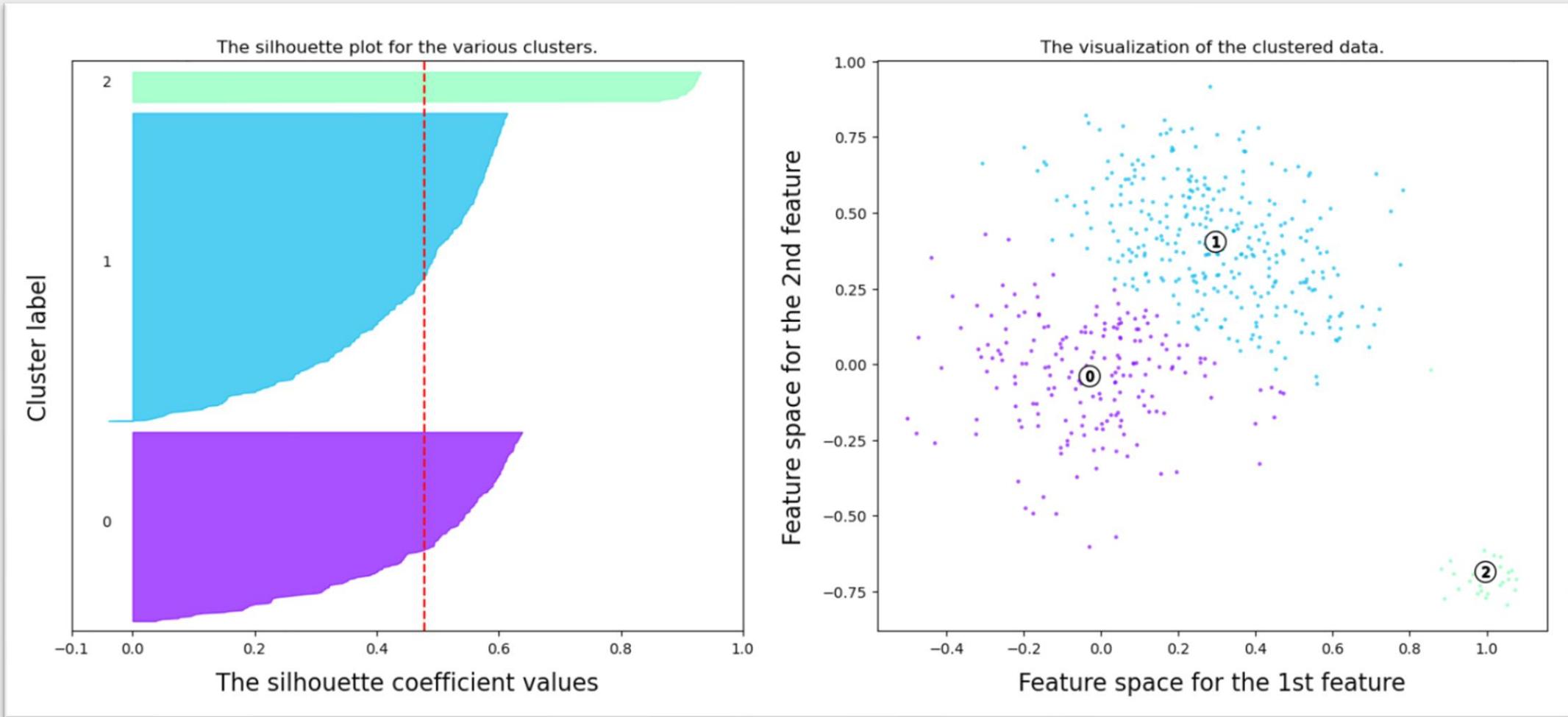
K-means silhouette score

`sklearn.metrics.silhouette_score`



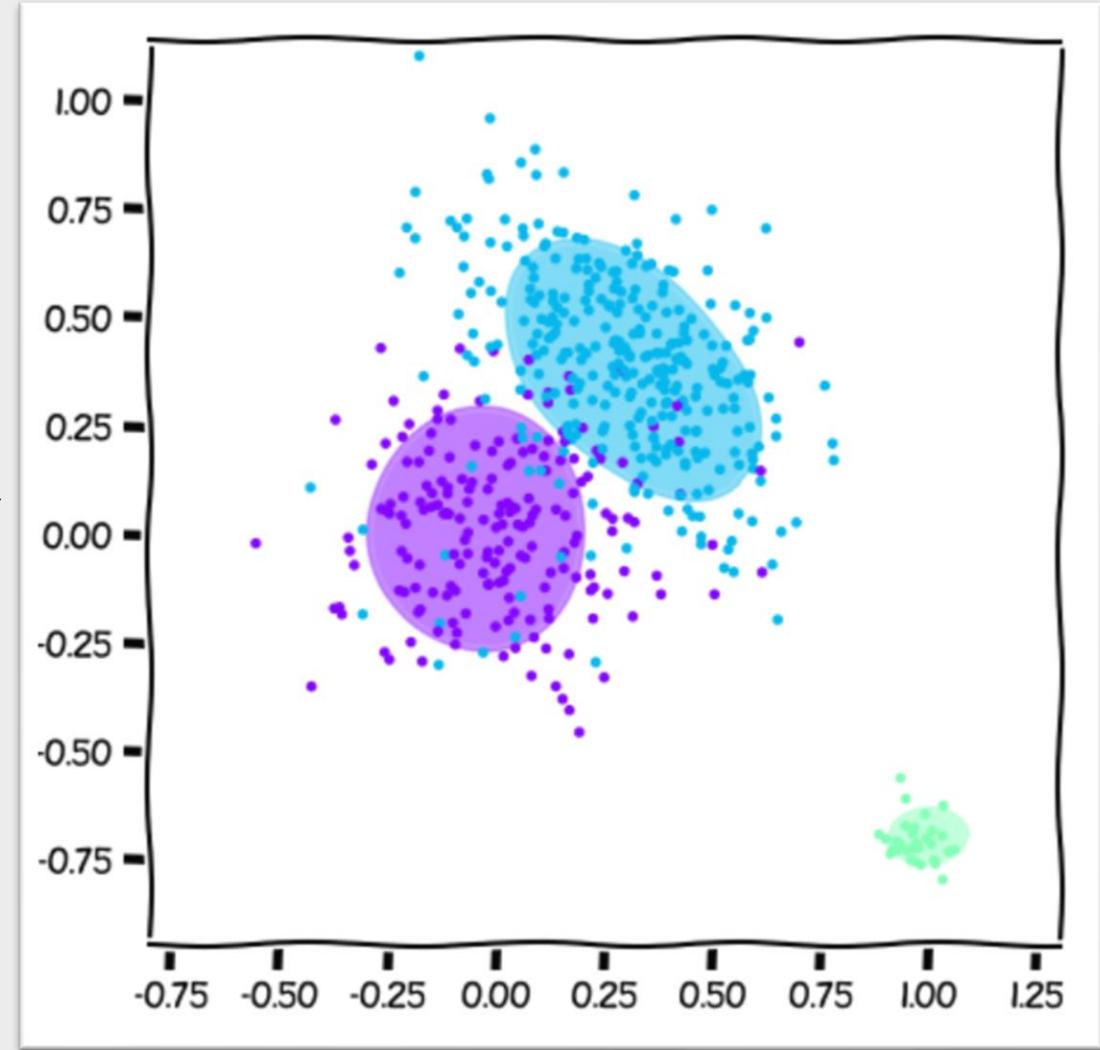
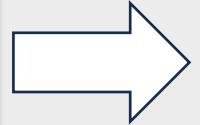
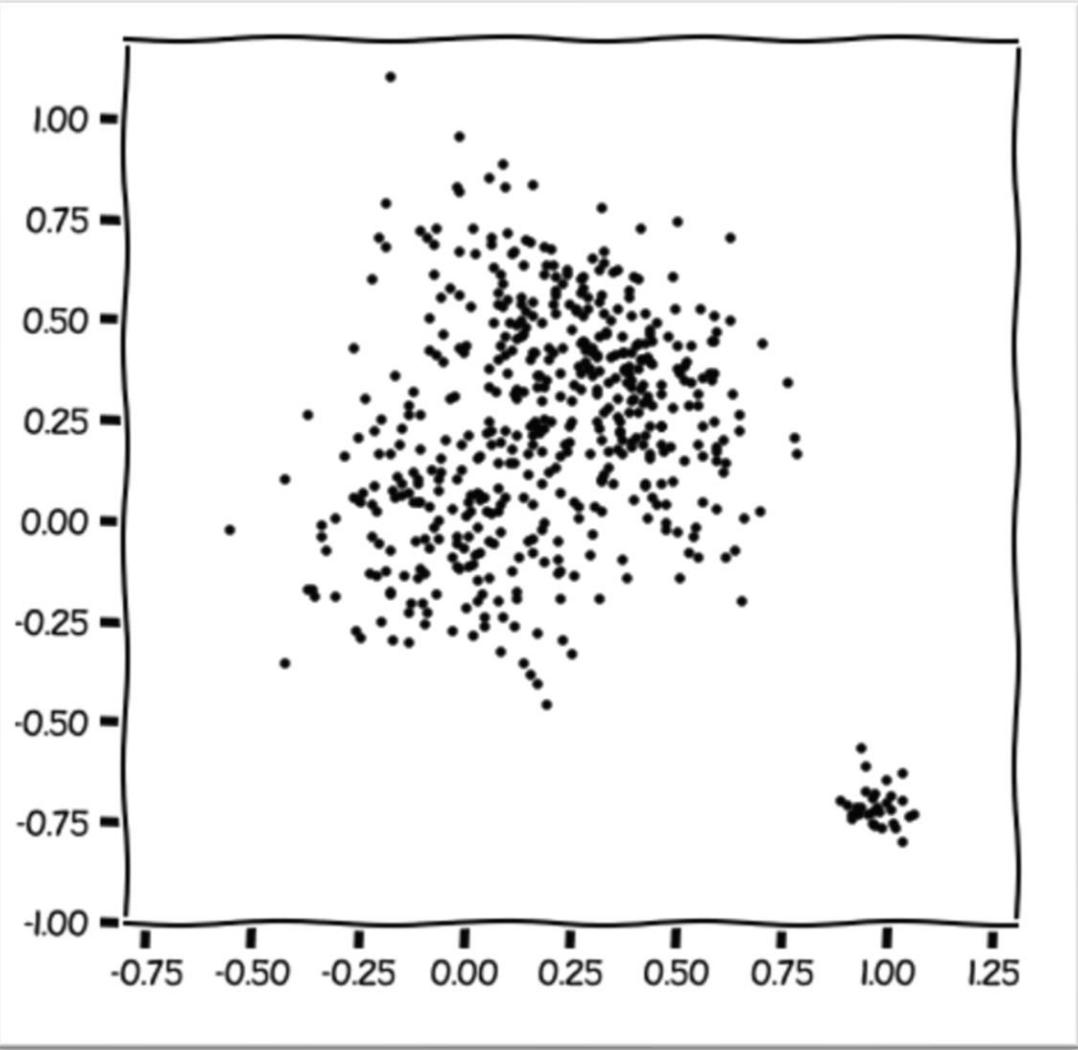
K-means silhouette score

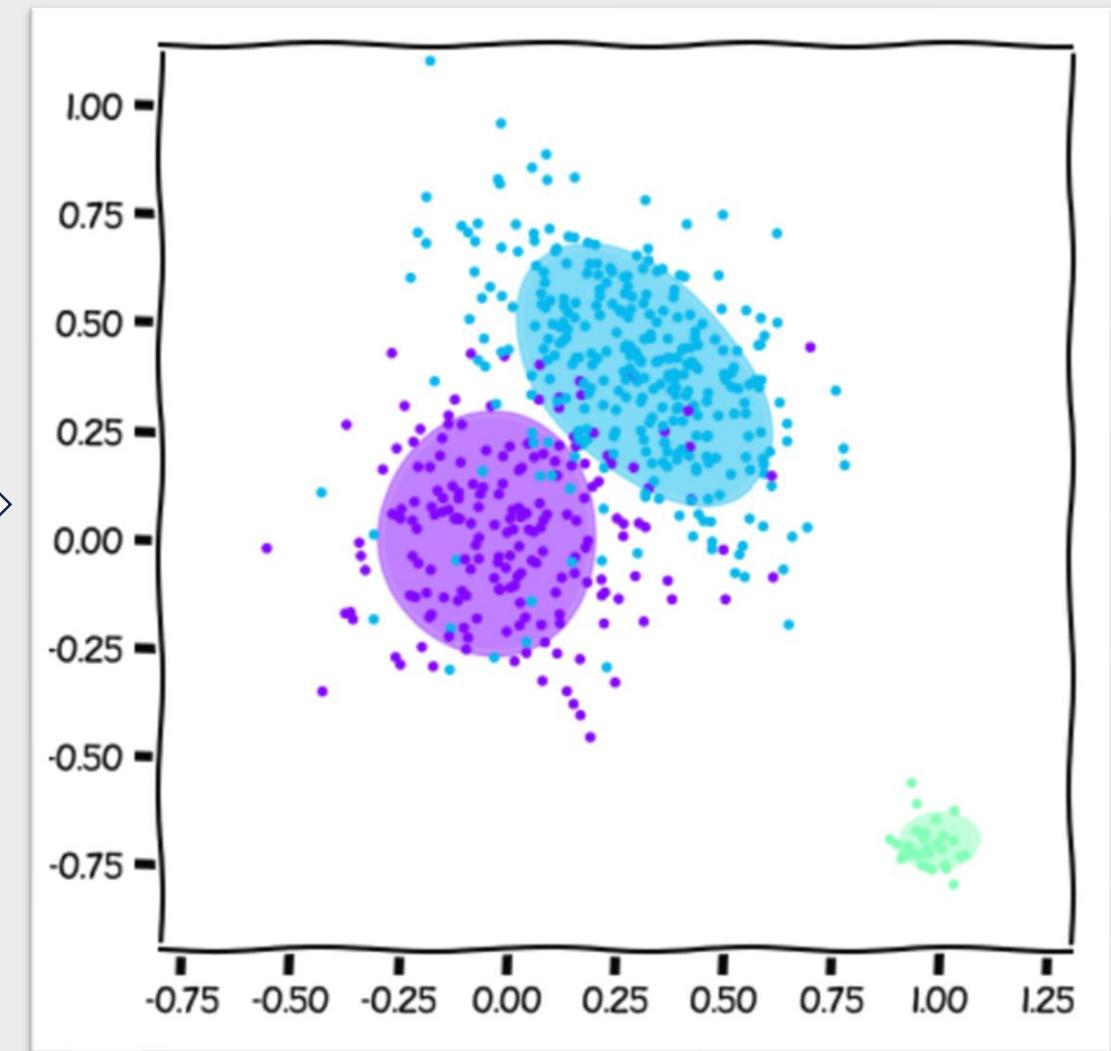
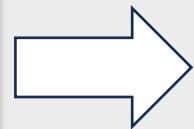
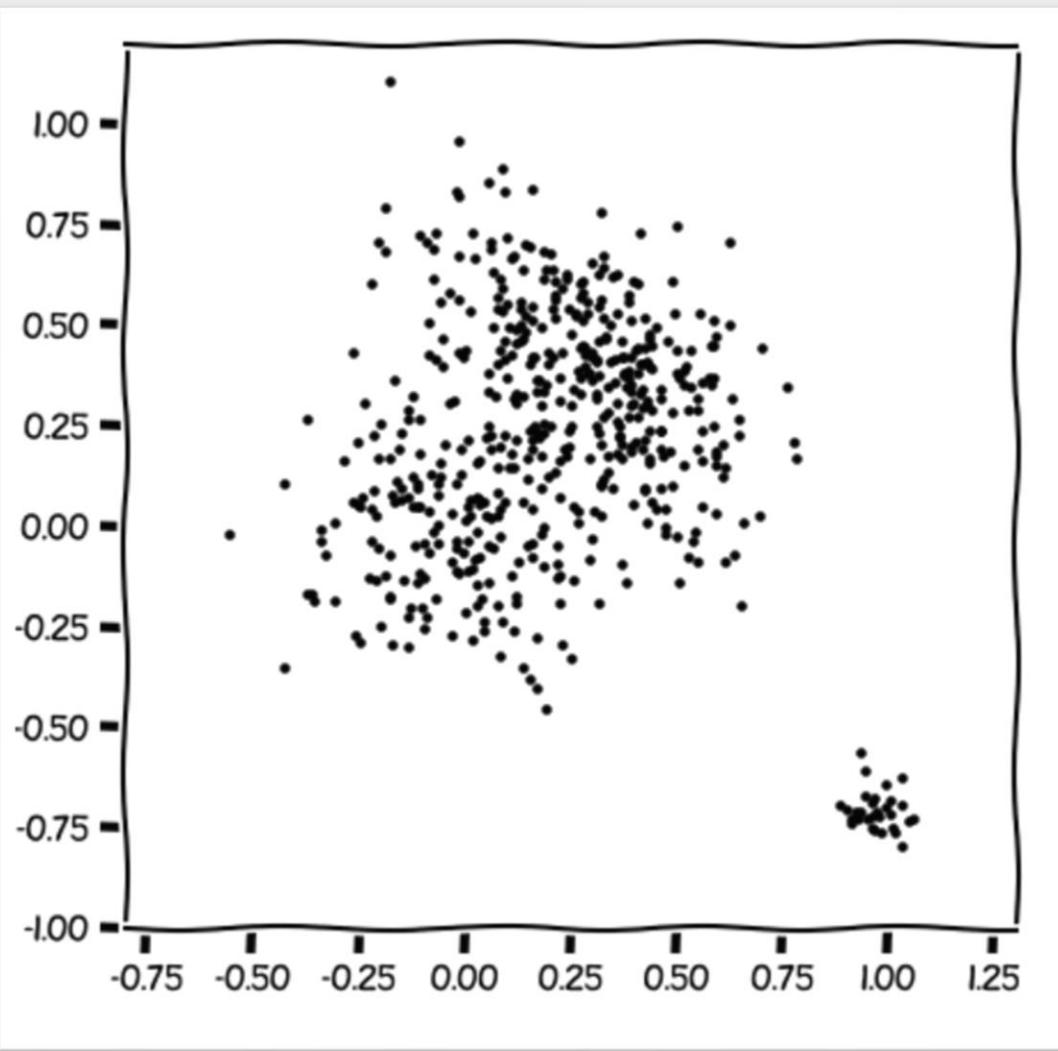
`sklearn.metrics.silhouette_score`



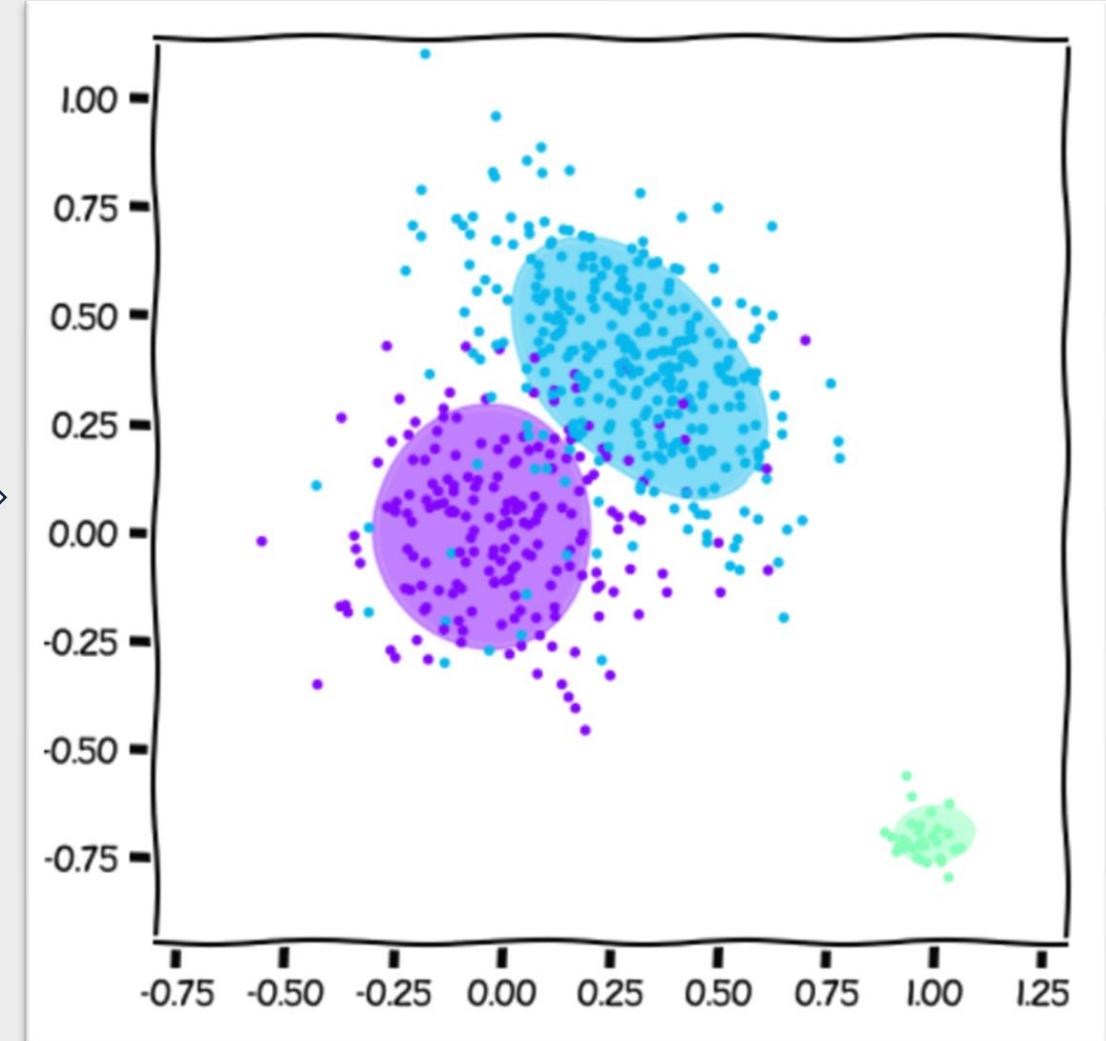
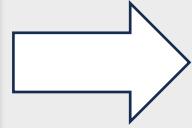
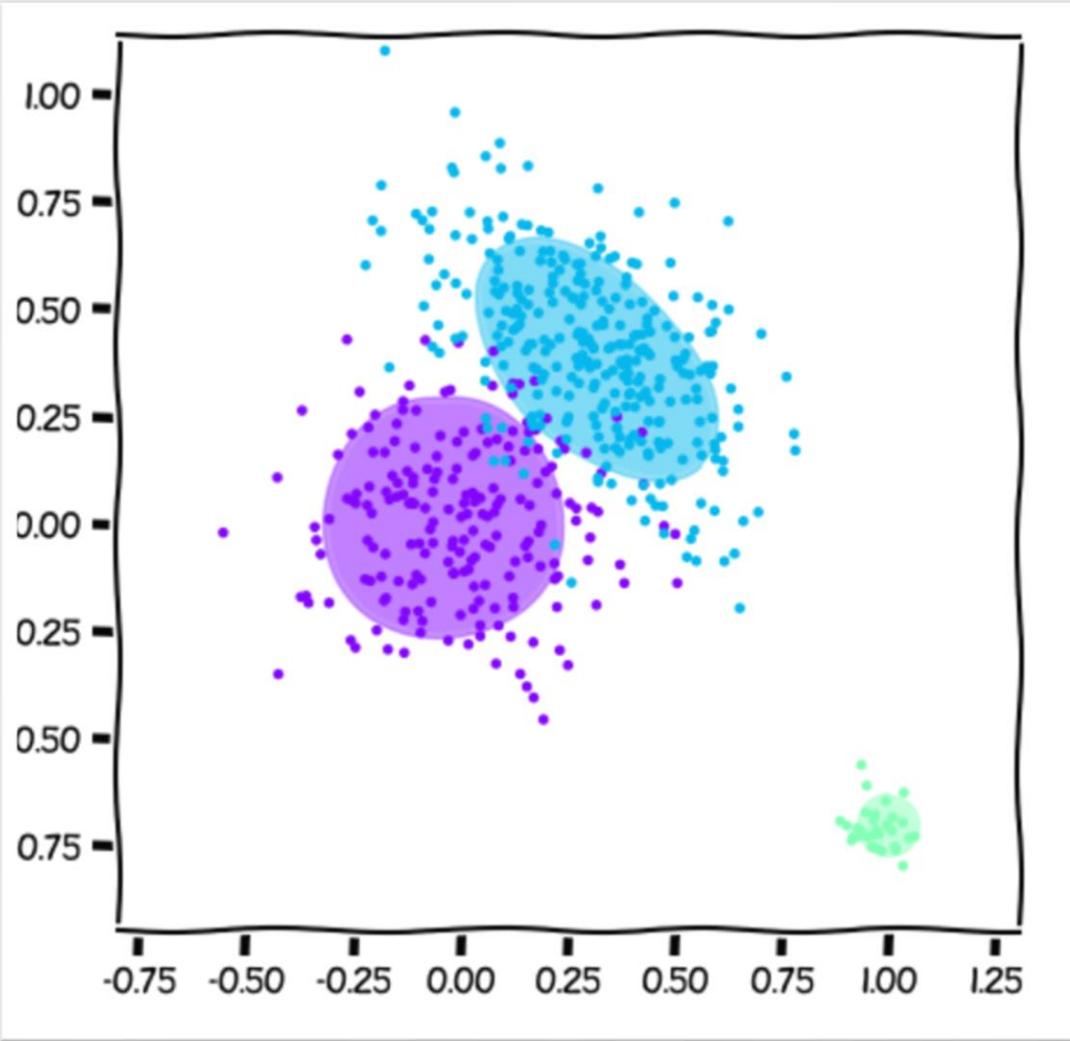
Gaussian Mixture Models

`sklearn.mixture.GaussianMixture`

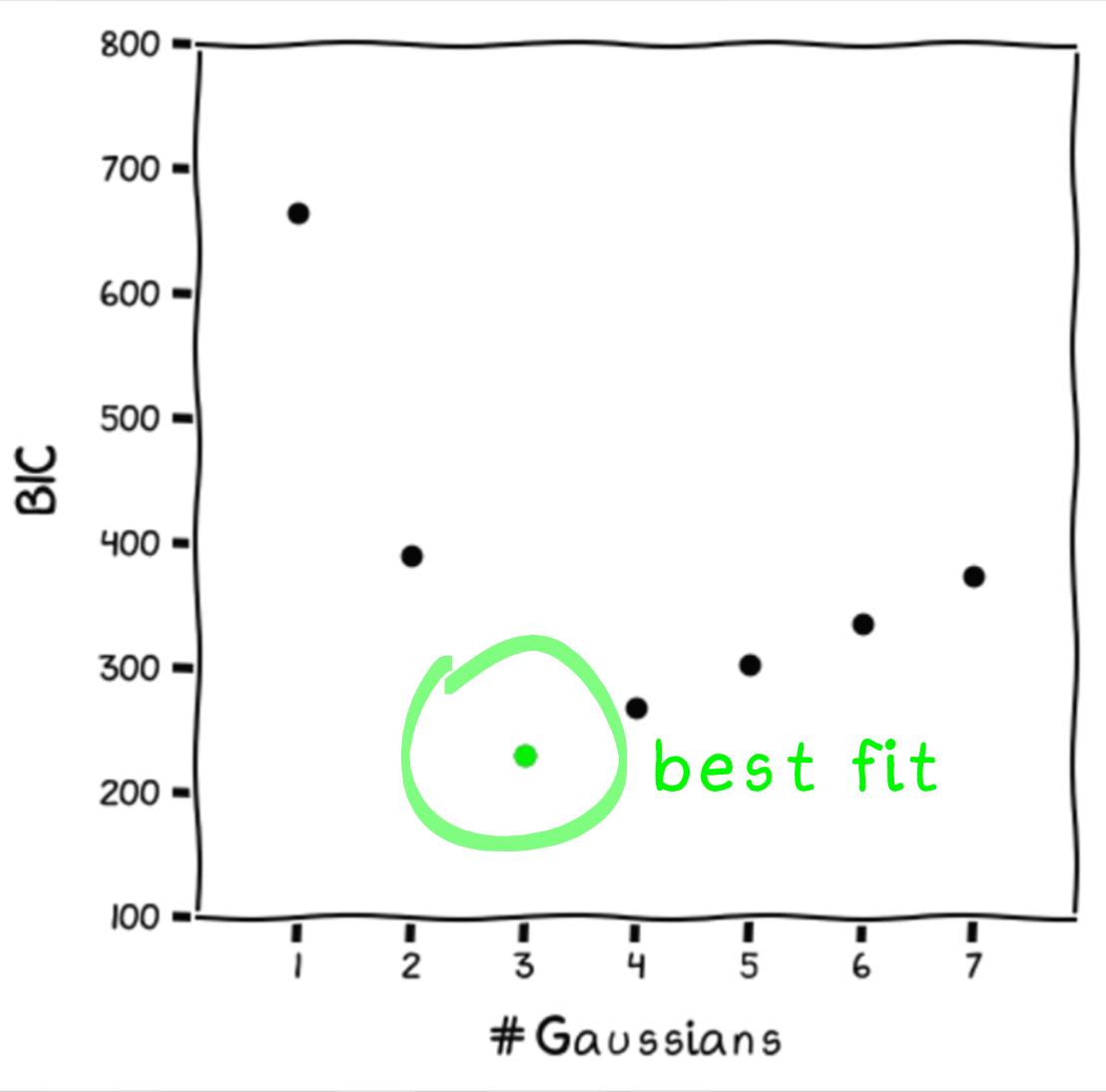




works amazingly if your data is Gaussian!



works amazingly if your data is Gaussian!



GMM

Bayesian Information Criterion

$$\text{BIC} = k \ln n - 2 \ln L$$

of data points n
of params k maximum likelihood

GMMs

pros

- easy
- really fast!
- amazing
- if your data is Gaussian
- good for uneven clusters
- provides statistics
- provides probabilities

cons

- need to decide on N
- outliers
- horrible
- if your data isn't Gaussian
- scales ok but not too well

means_ : array-like of shape (n_components, n_features)

The mean of each mixture component.

covariances_ : array-like

The covariance of each mixture component. The shape depends on `covariance_type`:

<code>aic(X)</code>	Akaike information criterion for the current model on the input X.
<code>bic(X)</code>	Bayesian information criterion for the current model on the input X.
<code>predict_proba(X)</code>	Evaluate the components' density for each sample.

Spectral Clustering

`sklearn.cluster.SpectralClustering`

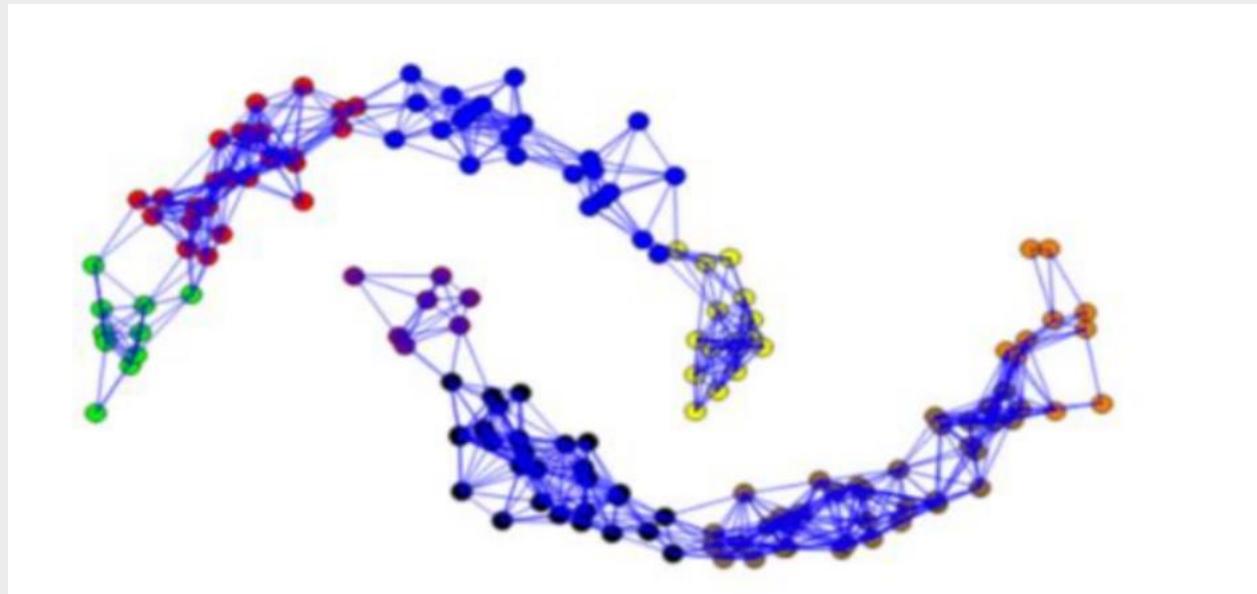
Apply clustering to a projection of the normalized Laplacian.

In practice Spectral Clustering is very useful when the structure of the individual clusters is highly non-convex, or more generally when a measure of the center and spread of the cluster is not a suitable description of the complete cluster, such as when clusters are nested circles on the 2D plane.

Spectral Clustering

`sklearn.cluster.SpectralClustering`

1. Calculate the distances between points
2. Construct a similarity graph



Affinity Matrix

How *close* are the data points i and j with a hard cut-off
Typically doesn't use a Gaussian distance metric

$$A_{ij} = \exp \left[-\epsilon |x_i - x_j|^2 \right]$$

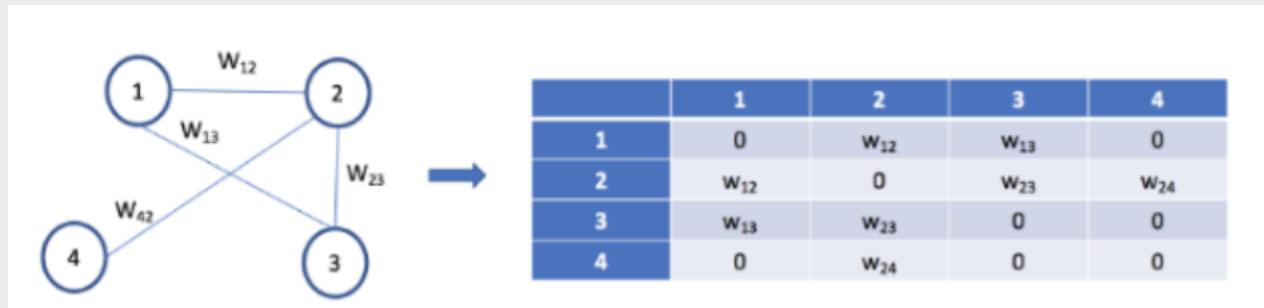
Degree Matrix

How many neighbors each point has

$$D_{ij} = \text{diag}[N]$$

Laplacian Matrix

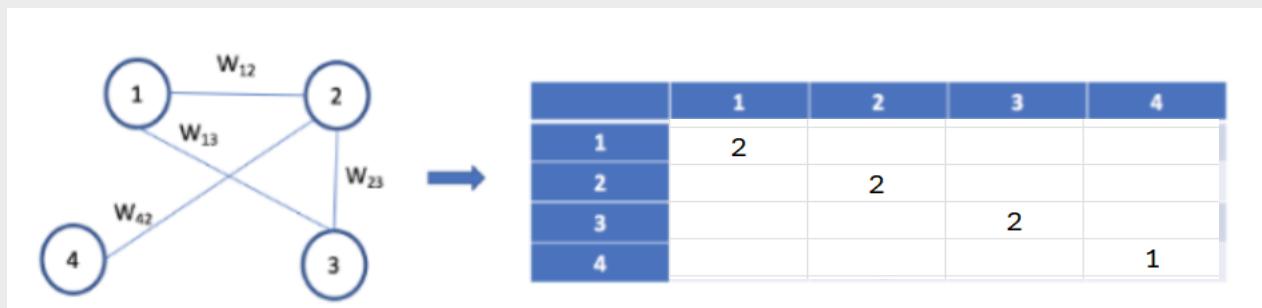
$$L = D - A$$



Affinity Matrix

How *close* are the data points i and j with a hard cut-off
Typically doesn't use a Gaussian distance metric

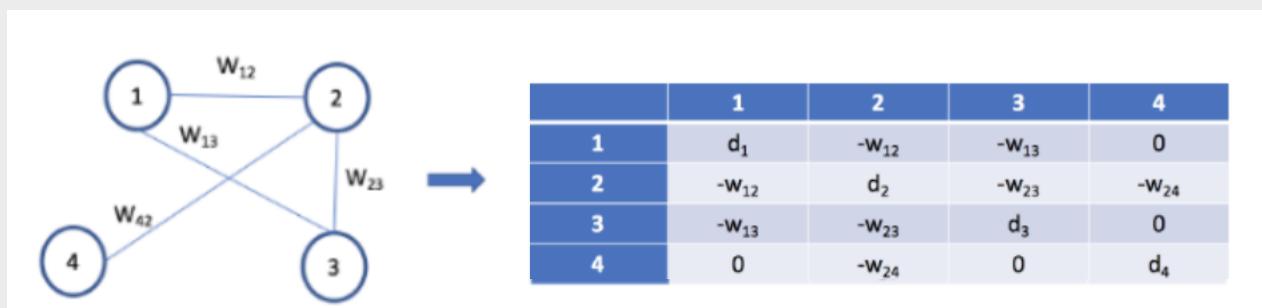
$$A_{ij} = \exp \left[-\epsilon |x_i - x_j|^2 \right]$$



Degree Matrix

How many neighbors each point has

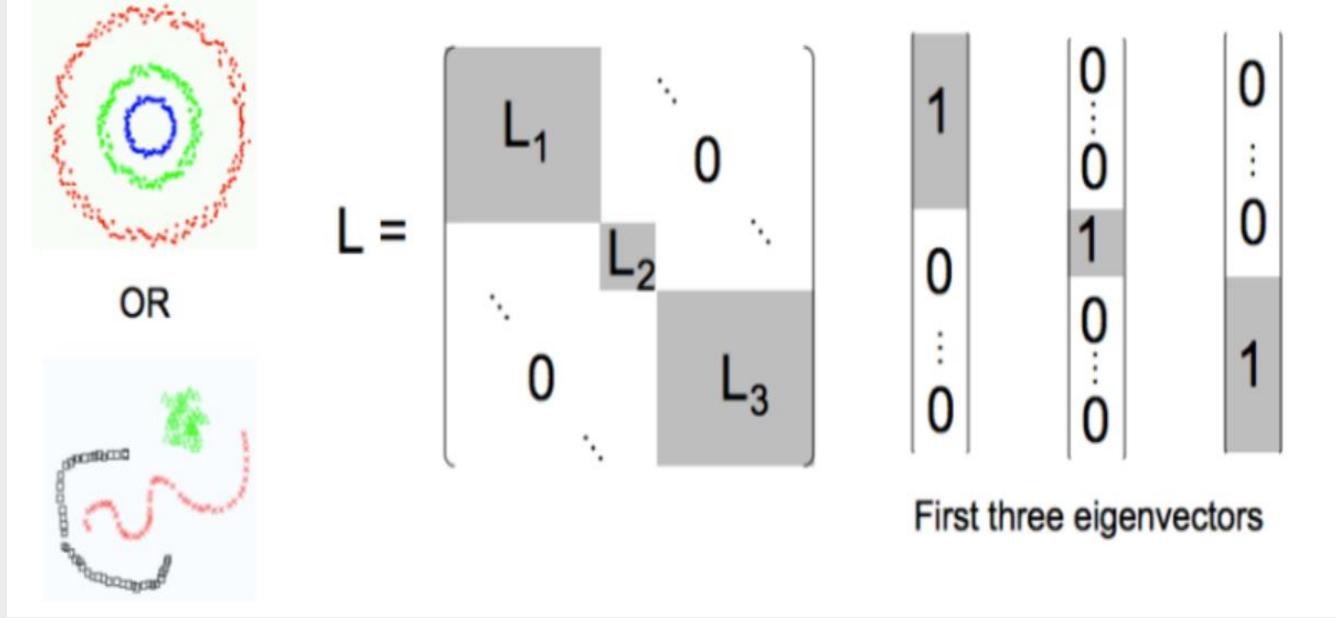
$$D_{ij} = \text{diag}[N]$$



Laplacian Matrix

$$L = D - A$$

Simplify L



3. Find N smallest eigenvalues (diagonalize it)
4. Project everything into N-dim space
5. Run k-means!

Affinity Matrix

How *close* are the data points i and j with a hard cut-off
Typically doesn't use a Gaussian distance metric

$$A_{ij} = \exp \left[-\epsilon |x_i - x_j|^2 \right]$$

Degree Matrix

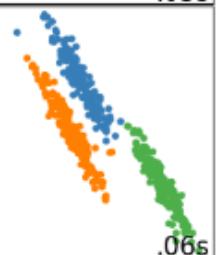
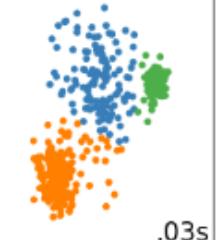
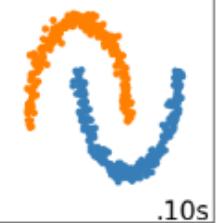
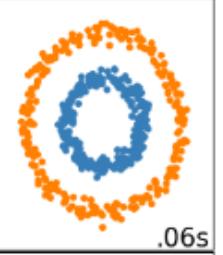
How many neighbors each point has

$$D_{ij} = \text{diag}[N]$$

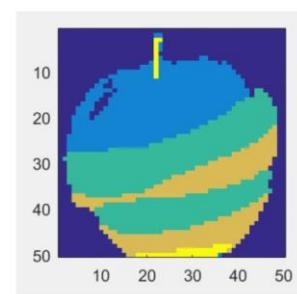
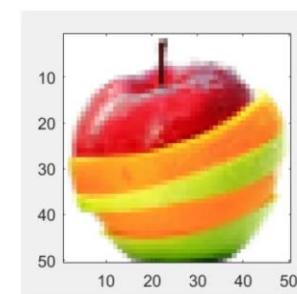
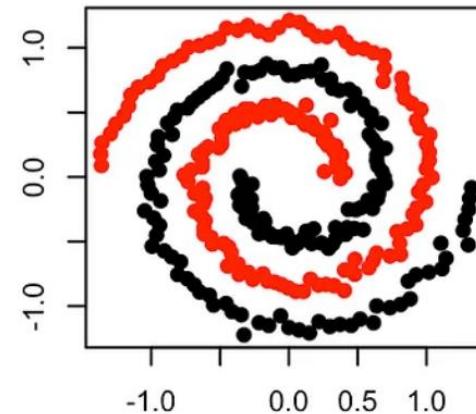
Laplacian Matrix

$$L = D - A$$

Spectral
Clustering



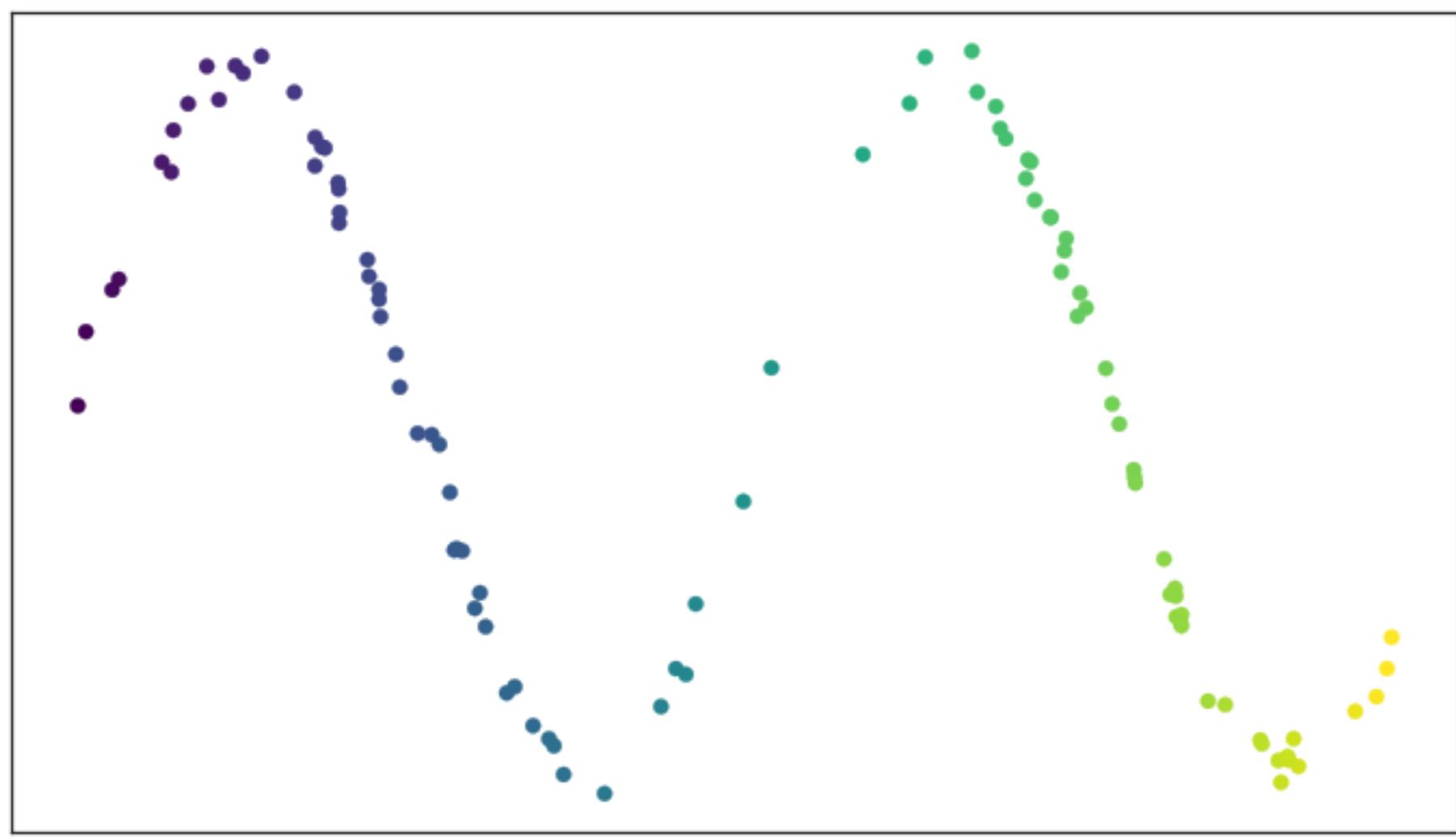
3

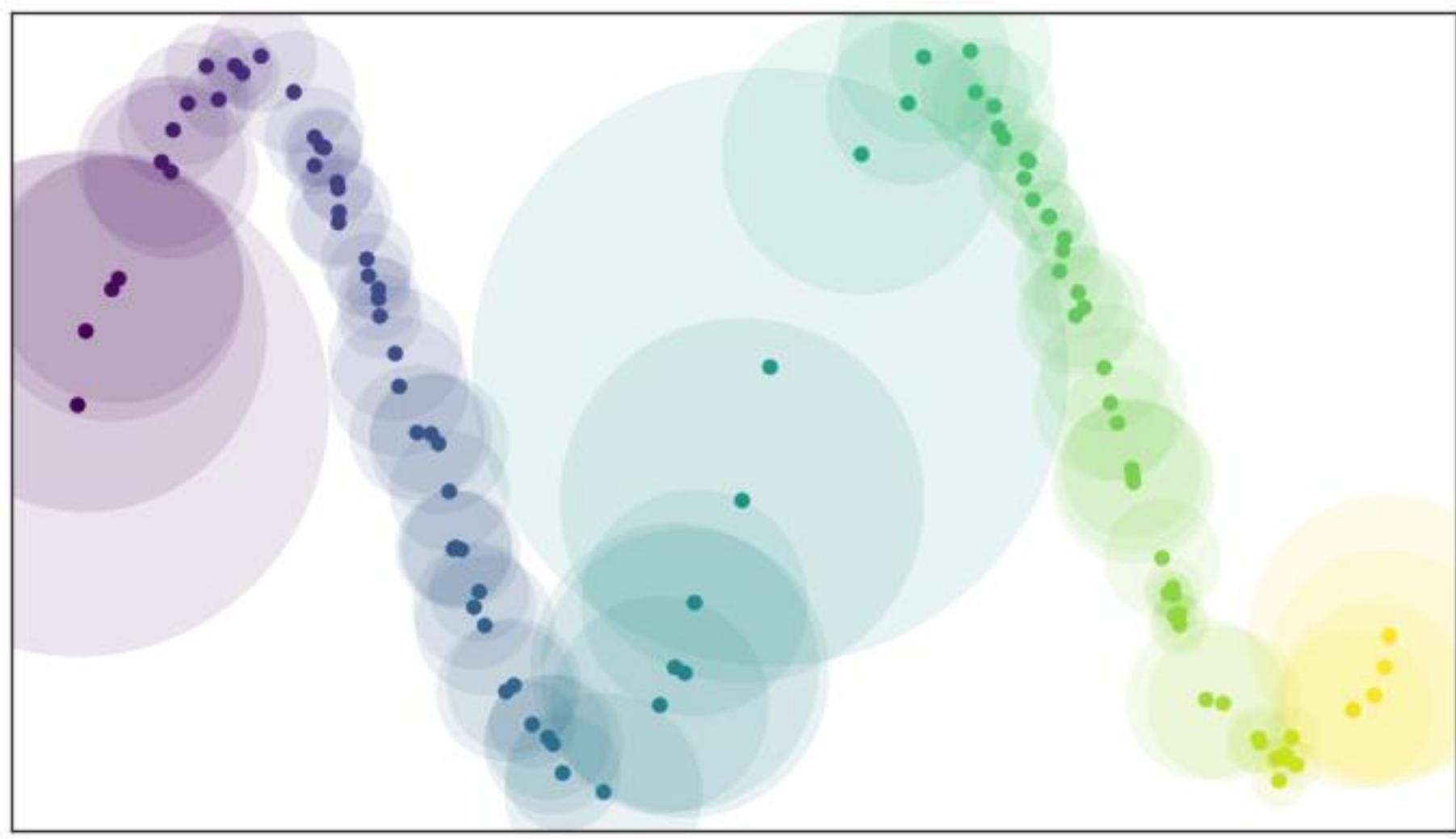


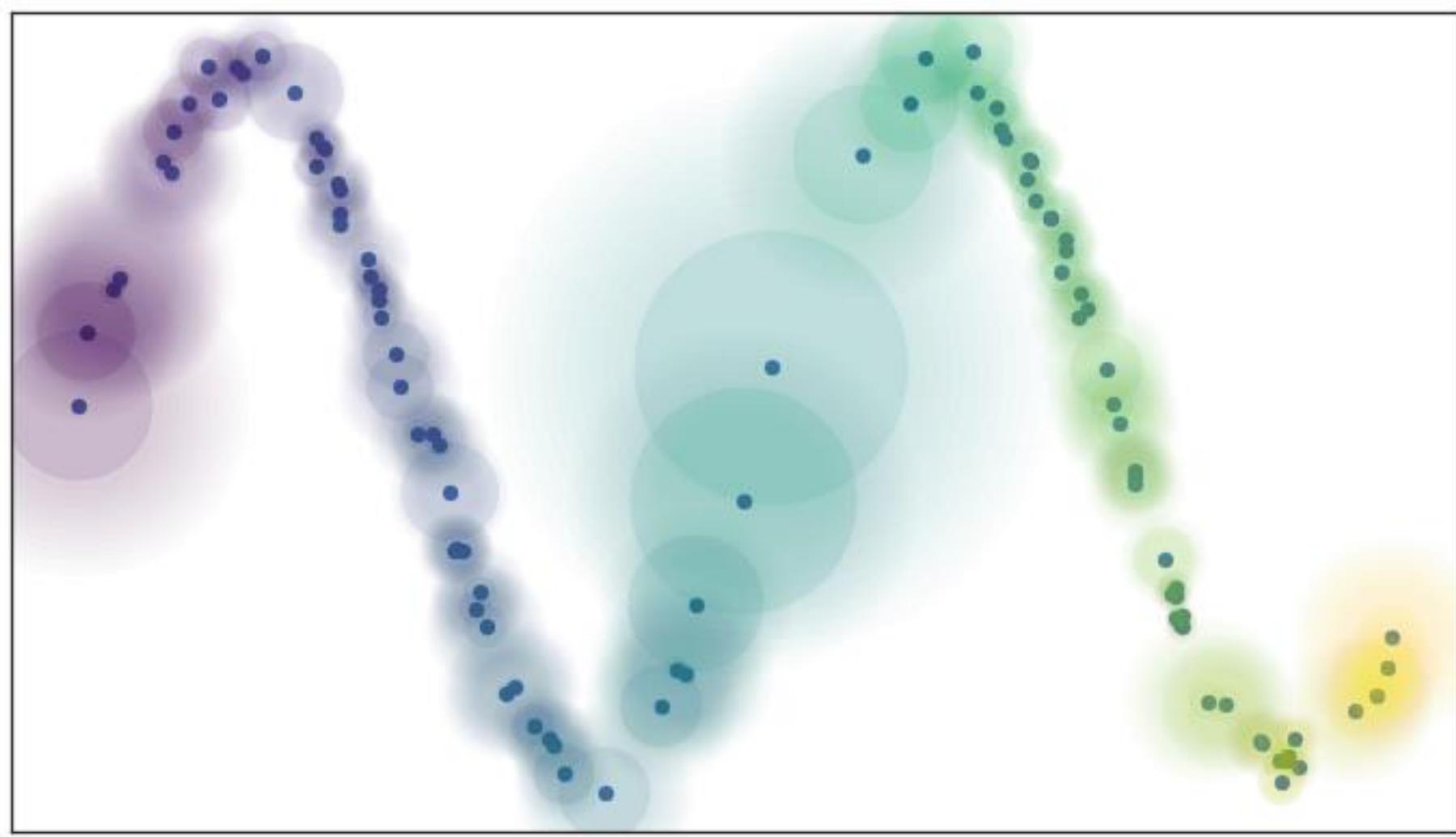
UMAPs

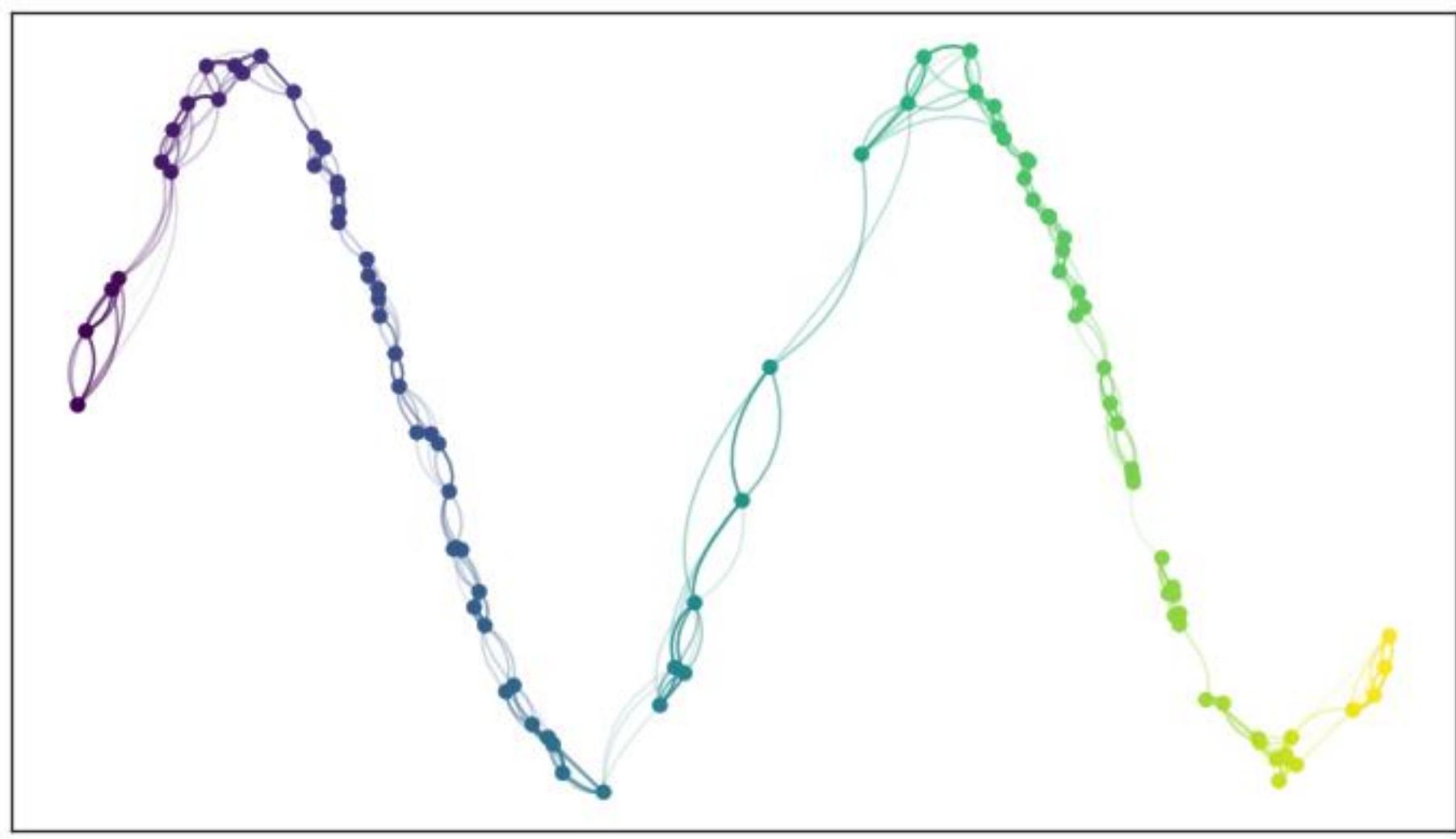


<https://umap-learn.readthedocs.io/en/latest/>

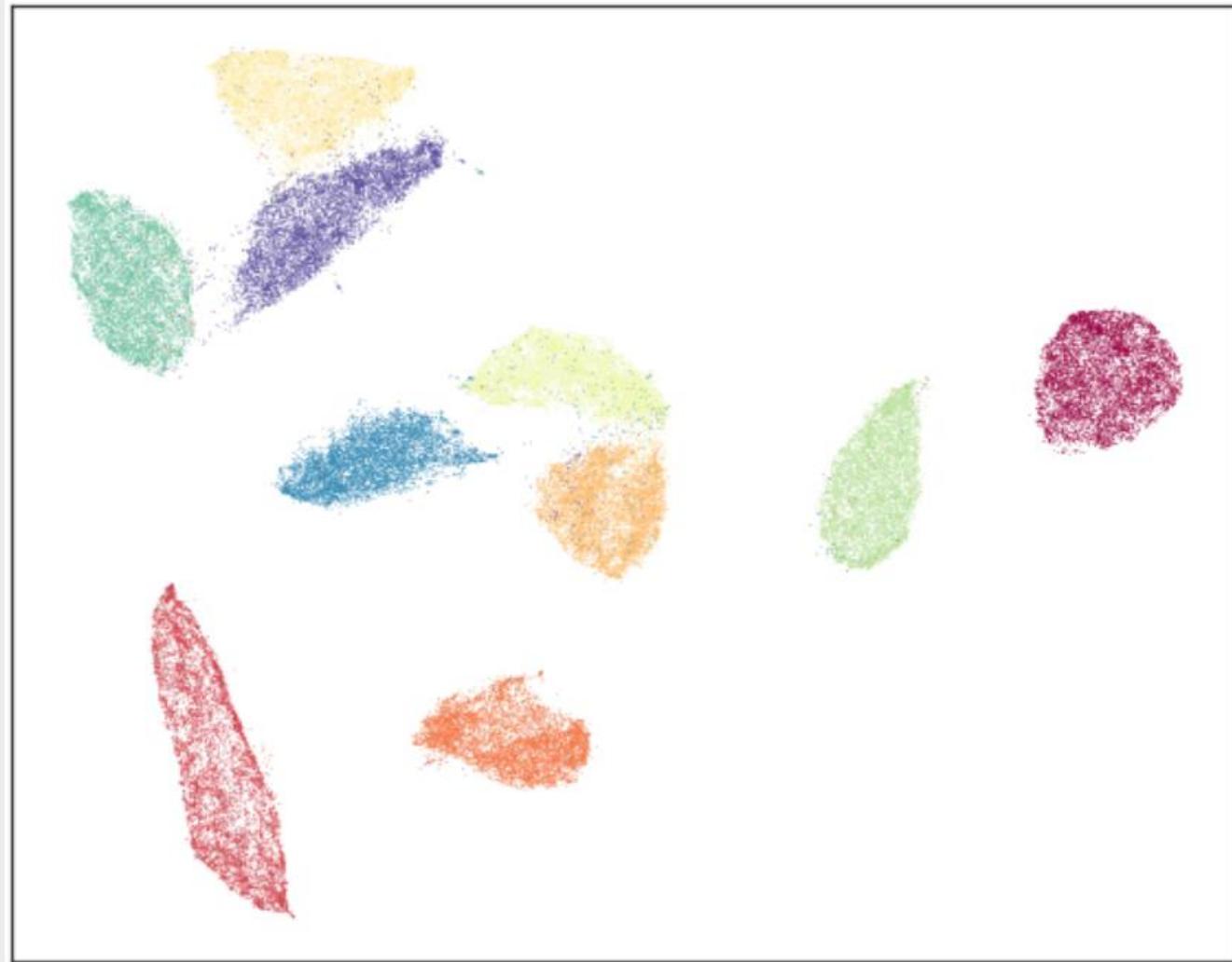
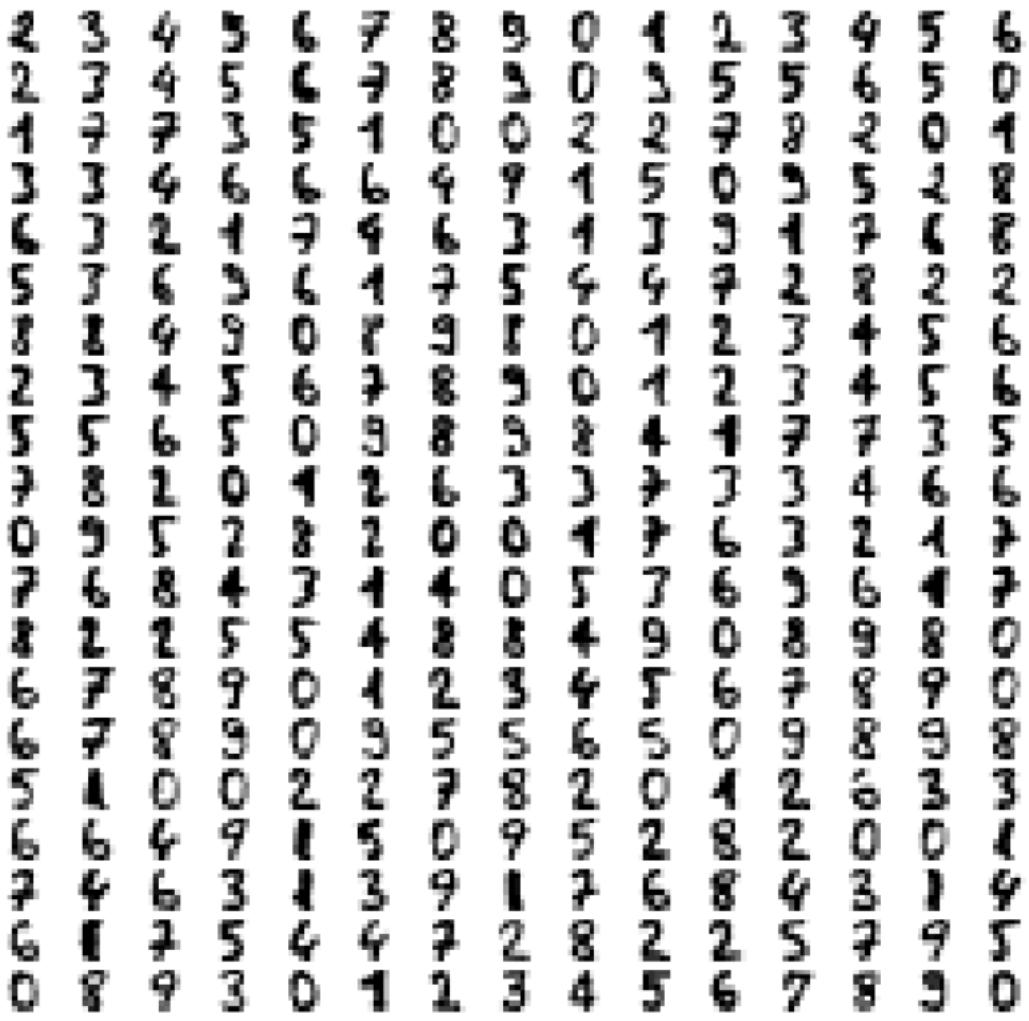








8x8 digits: 64-dimensional data



2-dimensional data w/ same graph

Final tip!

`sklearn.model_selection.KFold` ↗

```
class sklearn.model_selection.KFold(n_splits=5, *, shuffle=False, random_state=None)
```

[source]

K-Fold cross-validator.

Provides train/test indices to split data in train/test sets. Split dataset into k consecutive folds (without shuffling by default).

Each fold is then used once as a validation while the k - 1 remaining folds form the training set.

Read more in the [User Guide](#).

For visualisation of cross-validation behaviour and comparison between common scikit-learn split methods refer to [Visualizing cross-validation behavior in scikit-learn](#)