

Supervised ← given label

Unsupervised ← not given

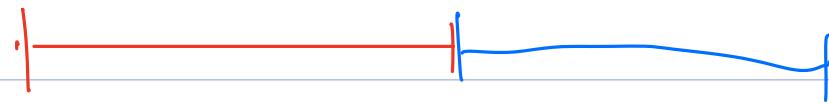
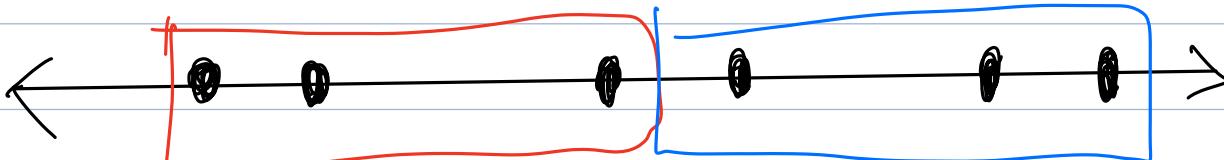
K-Means Clustering

↑
number

$K=2$



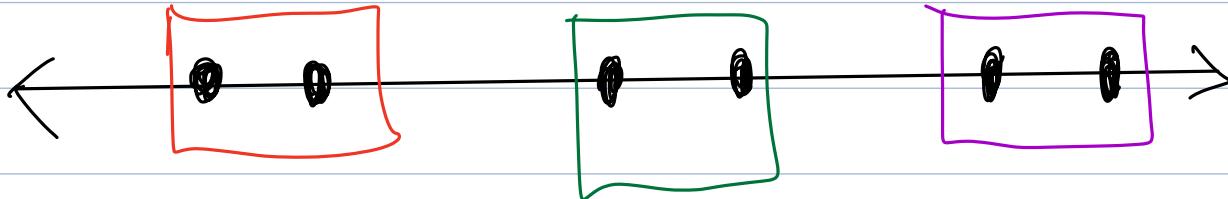
$K=2$



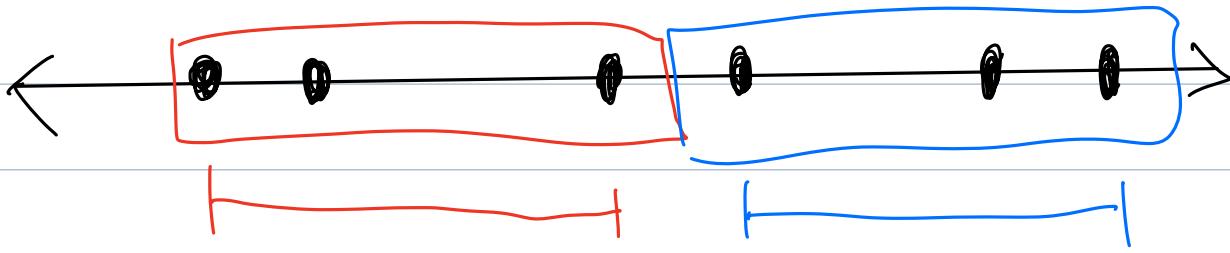
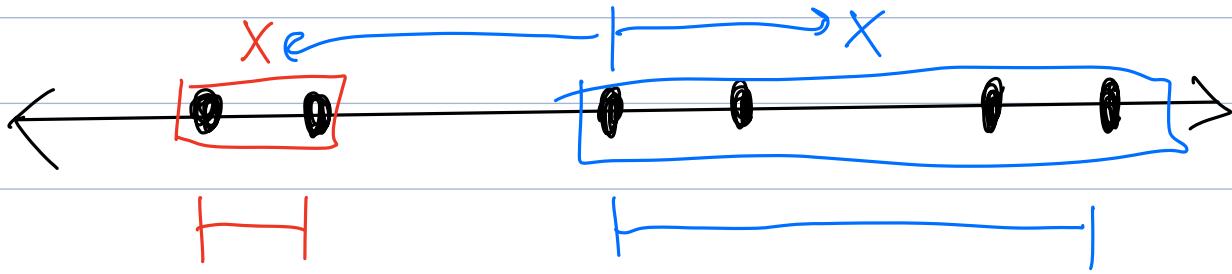
$k=3$

Variation

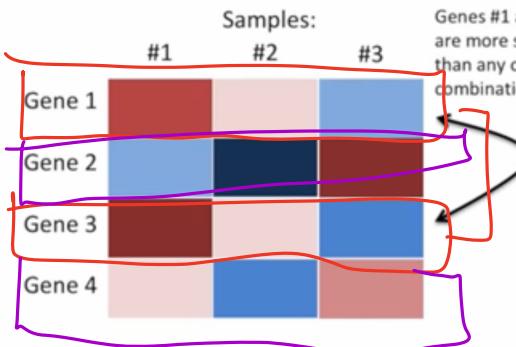
different k



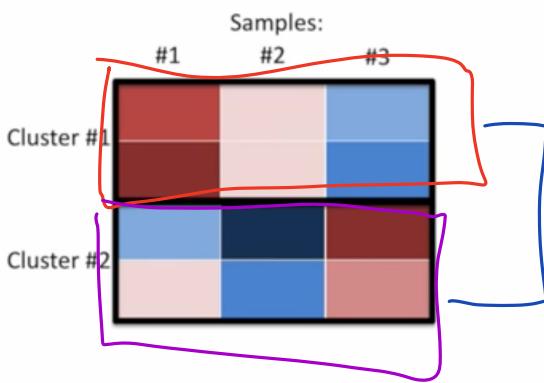
Random initial points (centre)



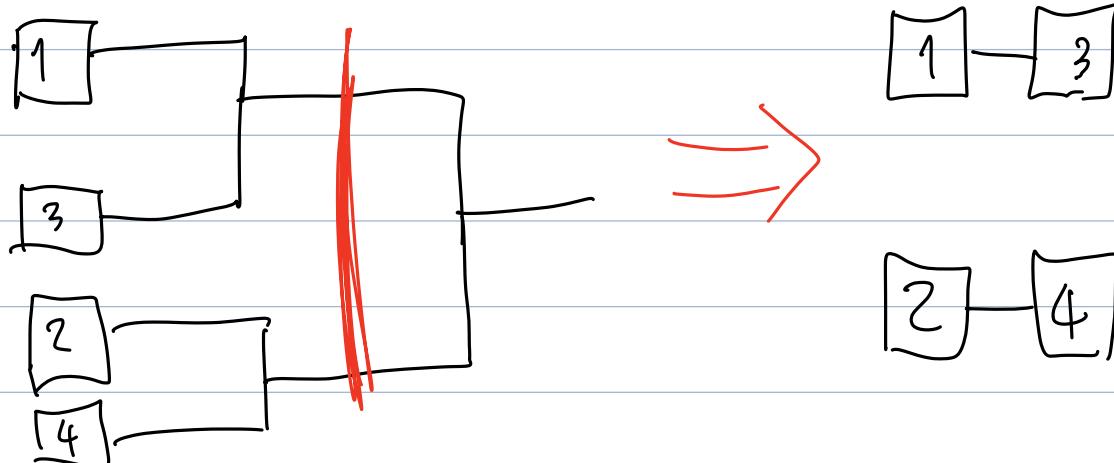
Hierachical Clustering

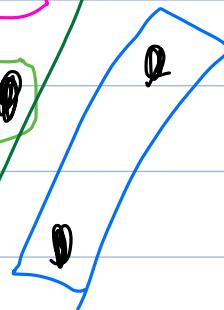
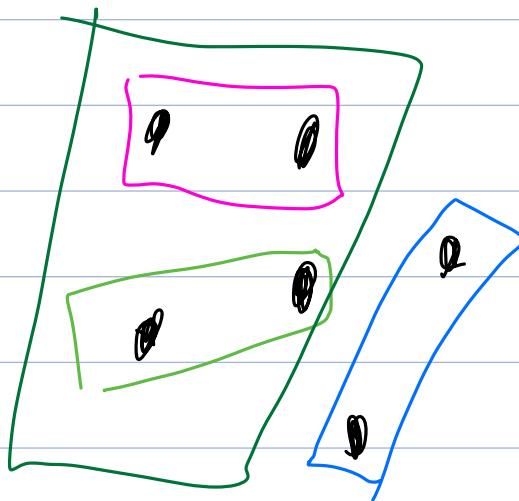
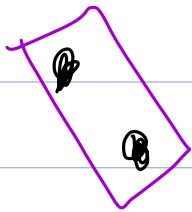


- Conceptually...
- 1) Figure out which gene is most similar to gene #1.
 - 2) Figure out which genes is most similar to gene #2... (and then #3 and then #4).
 - 3) Of the different combinations, figure out which two genes are the most similar. Merge them into a cluster.

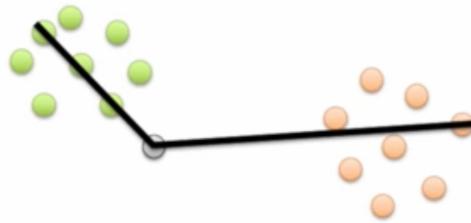


- Conceptually...
- 1) Figure out which gene is most similar to cluster #1.
 - 2) Figure out which genes is most similar to gene #2... (and then #4).
 - 3) Of the different combinations, figure out which two genes are the most similar. Merge them into a cluster.
 - 4) Go back to step 1.



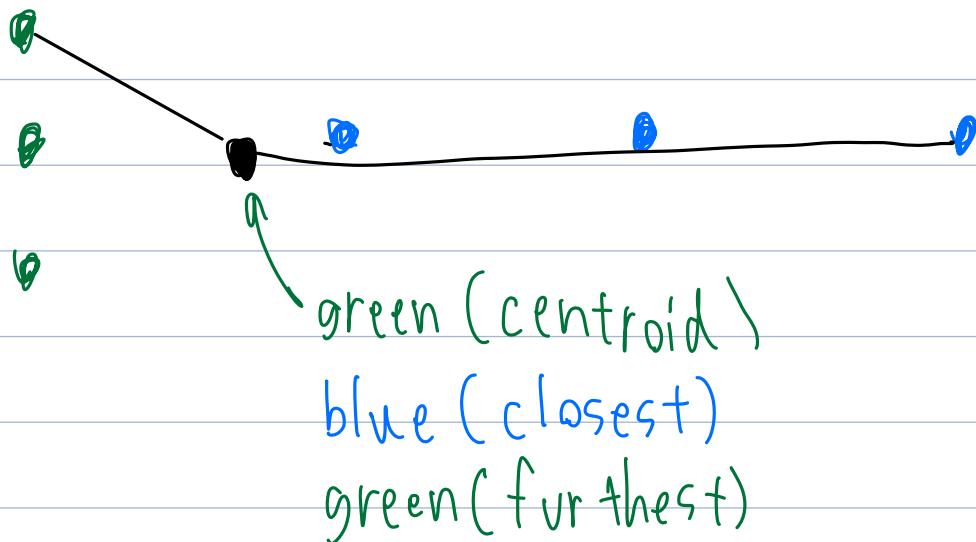


For the sake of visualizing how the different methods work, imagine our data was spread out on an X-Y plane.

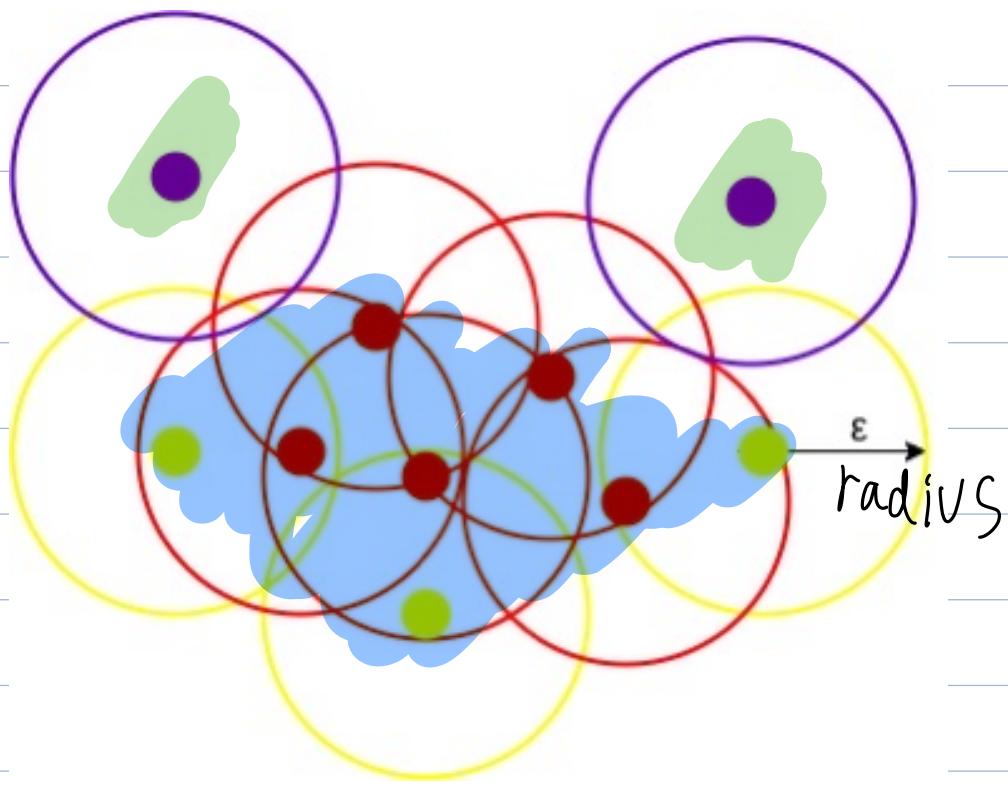


We can compare that point to...

- 1) The average of each cluster (this is called the “centroid”)
- 2) The closest point in each cluster (this is called “single-linkage”)
- 3) The furthest point in each cluster (this is called “complete-linkage”)



DBSCAN



- 1) Core
- 2) border
- 3) Noise

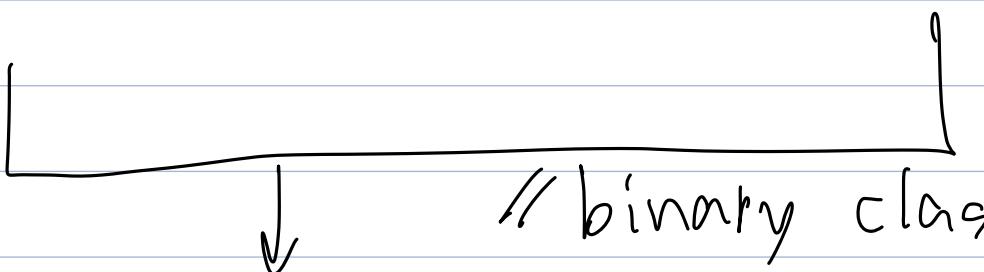
Specify
1) radius (ϵ) 2
2) Min points 4

4 If there are $\text{points}(\text{data}) \geq \text{MIN_POINTS}$: assign as "Core"

4 If there are $\text{points}(\text{data}) \leq \text{MIN_POINTS}$: assign as "border"

If $\text{points}(\text{data}) = 1$: assign as "Noise"

0 1 2 3



don't include when
scaling & one-hot encoding

Rules of thumb: Always OHE

but if data is in binary
then OHE is optional

$k_list = [1, 3, 5, 100]$

for k in k-list:

print(k) → 1

3

5

100

range(start, end)

↳ range(0, 10)

list(range(0, 10))

↳ [0, 1, 2, 3, ..., 10]

for i in range(0, 10)

print(i) → 0 1

print(k-list[i]) → 1 3

train - raw

target - train - raw

model

testing prediction:

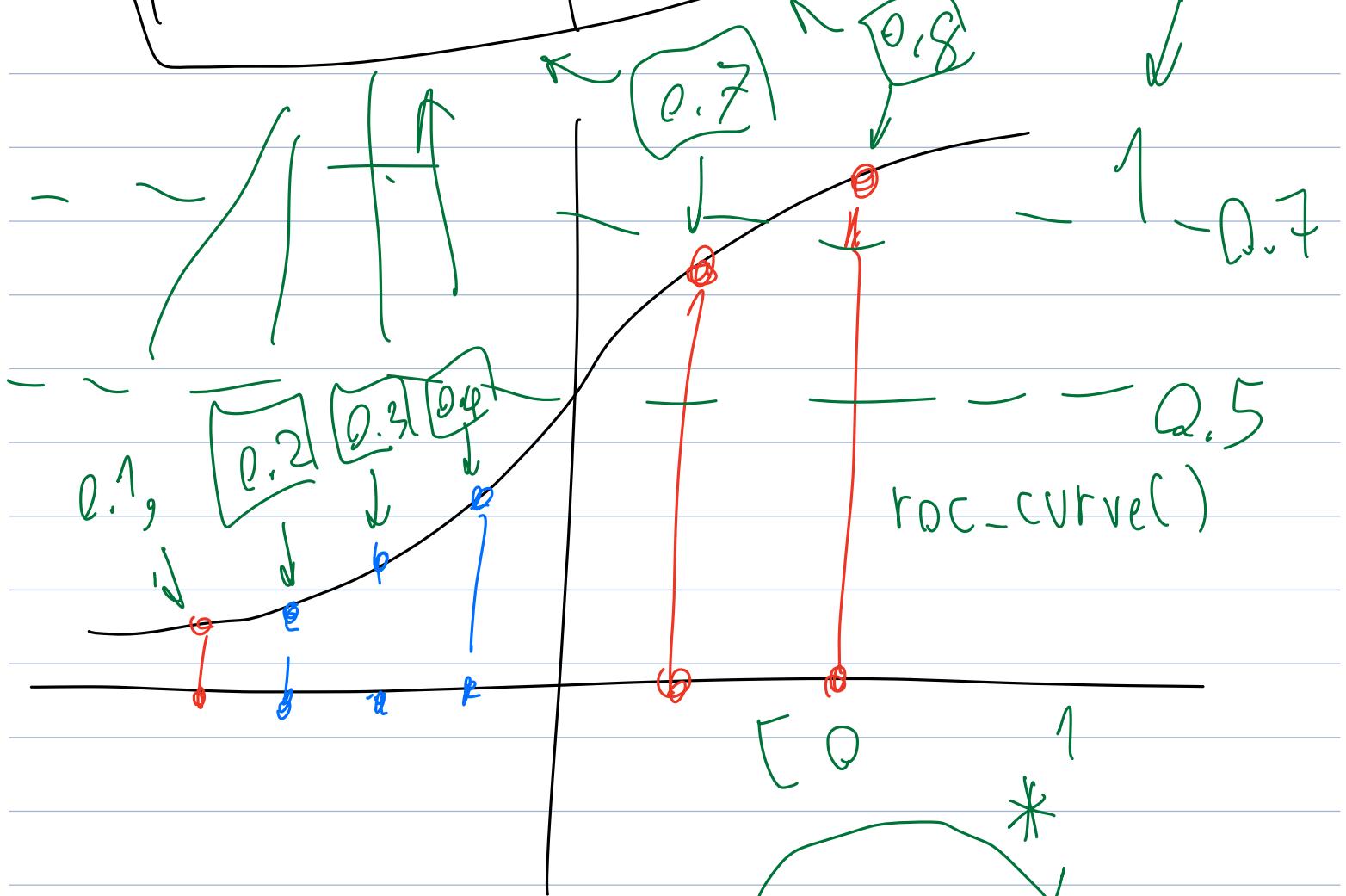
build the model

$X_train, X_test, Y_train, Y_test =$

`train-test-split()`

`predicted-proba()`

`predict()`



`predict_proba()[:, 1]` → `[0.2, 0.8]`
 probab of belongs to class 1
 probab of belongs to class 0

`num_pipeline.`
`fit_transform(heart_numerical_data)`

numerical

imp
AVG
SCA

raw_data

Categorical

Categorical - encode

full_pipeline.fit_transform(raw_data)
[data]

compare

train, test, target, test-target = trainSplit()

penalty: 'none'
'l1'

LinearRegression()

lin-reg

fit(train)

→ lin-reg.fit(train, target)

→ predicted = lin-reg.predict(test)

→ print - 4 - ... (predicted), ...)