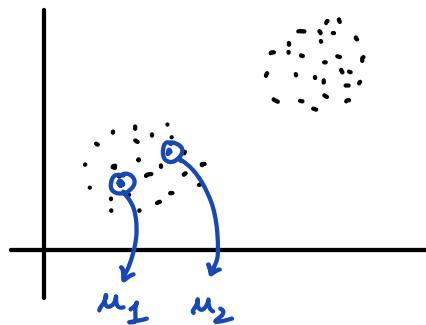


## CENTROID BASED CLUSTERING:

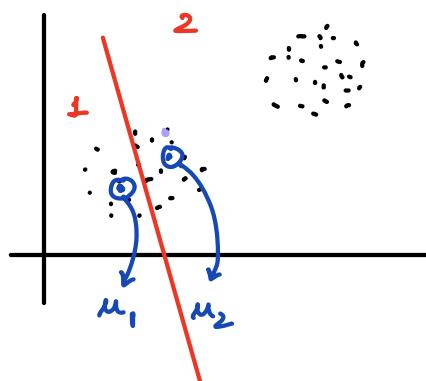
### Illustration



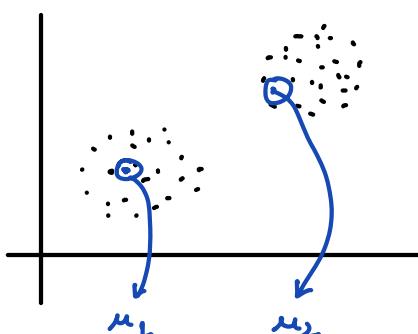
Initialization

### Iteration 1

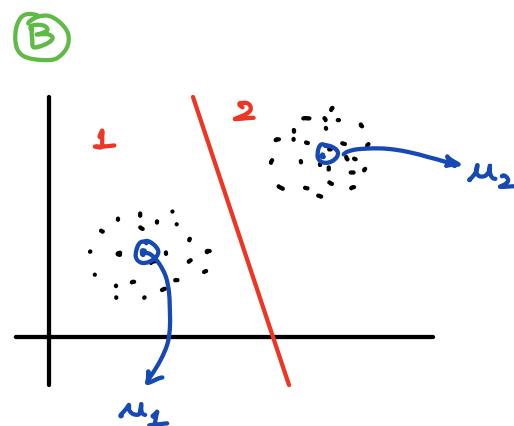
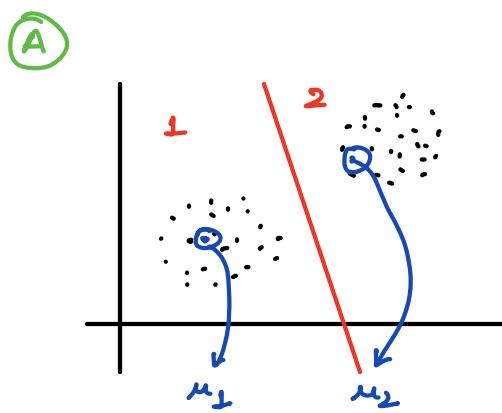
- A) find cluster assignment  
(assign each point to closest center)



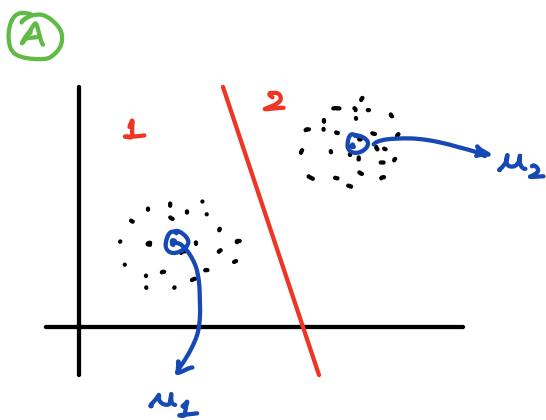
- B) Recompute  $\mu$



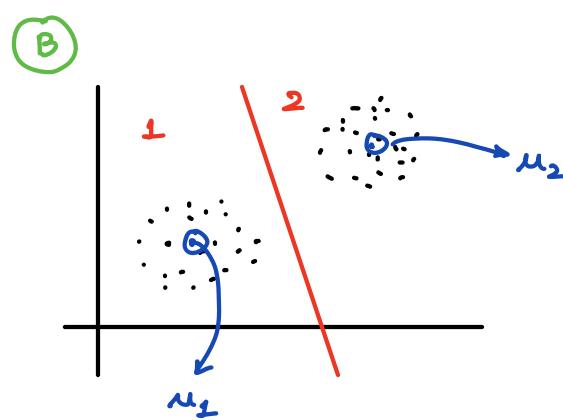
## Iteration 2 :



Iteration 3 : Algo converged (Assignment and  $\mu$  does not change)



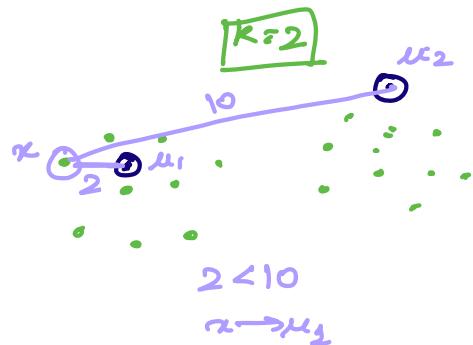
E Step



M Step

### Code (E Step)

- go to every data point  $\rightarrow m$
  - [ calculate the distance of this  $\rightarrow k \cdot n$  data point from  $\mu$  of every cluster
  - then find out the minimum distance
  - Assigns data point to the cluster with min distance
- $\hookrightarrow m(mk + k + 1) = O(mnk)$



### Code (M Step)

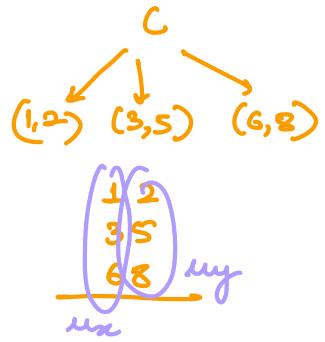
- Iterate over all clusters

Iterate over all points present in the cluster and take the mean

total  
 $m$ : no. of datapoints



$$\frac{m}{2} \cdot n + \frac{m}{4} \cdot n + \frac{m}{4} \cdot n = m \cdot n$$



Time Complexity:

$$x^{(i)} \in \mathbb{R}^n$$

$K$ : no. of clusters

$m$ : # points

$$O(m \cdot k \cdot n + n \cdot m)$$

↓ E Step      ↓ M Step

for a particular cluster:  $n \cdot m_k$

↓  
no. of points in  $k$  cluster

$$\sum_{k=1}^K n m_k = n m_1 + n m_2 + \dots + n m_K = n \cdot m$$

Maths:

$$x = \{x^1, x^2, x^3, \dots, x^m\}$$

$$u = \{u_1, u_2, \dots, u_K\}$$

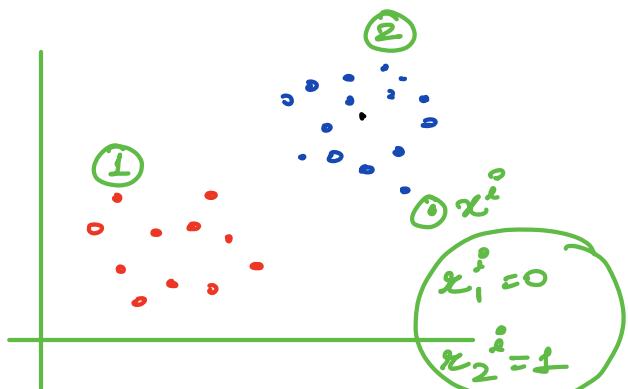
$K$ =clusters

Loss/Inertia

$$J = \sum_{i=1}^m \sum_{k=1}^K x_k^{(i)} (x^{(i)} - u_k)^2$$

all datapoints      all clusters

distance



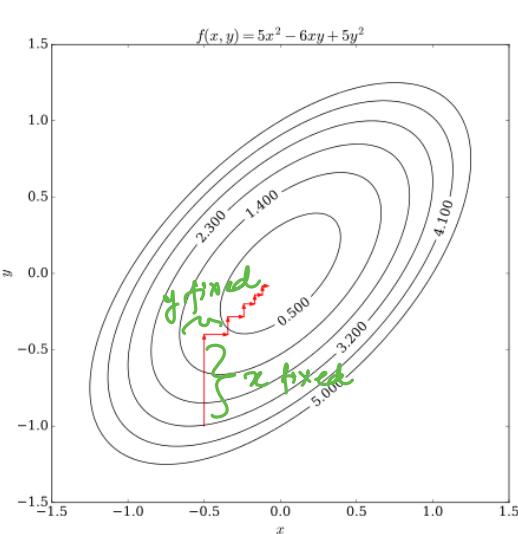
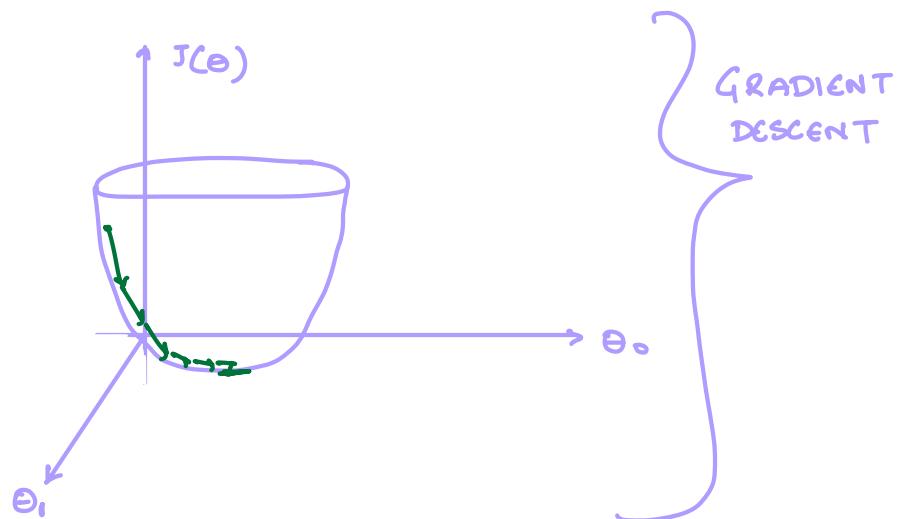
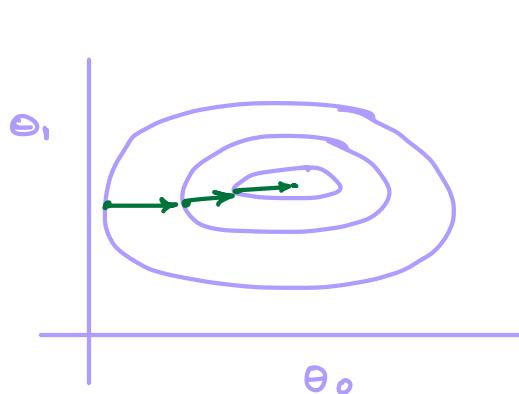
} compute the distance of every point from every cluster.

$$r_k^{(i)} = \begin{cases} 1 & \text{if } i\text{th point } \in k\text{th cluster} \\ 0 & \text{otherwise} \end{cases}$$

We want to learn  $r_k^{(i)}$  and  $\theta$

This is done by using COORDINATE DESCENT.

↓  
optimize one set of variable  
at a time while keeping  
others fixed.



Minimize your loss w.r.t only  
one variable at a time.

$\min_{\theta} (10, 20, 30, 40, 5, 50)$

for what value  
of  $f'(x^{(i)} - \mu_j)^2$  is  
minimum

$= \underline{\underline{S}}$

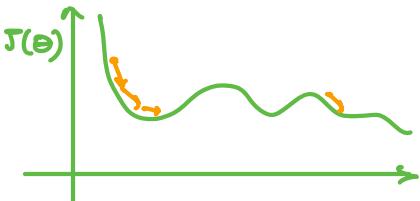
$\arg \min_{\theta} (10, 20, 30, 40, 5, 50)$

$$r_k^{(i)} = \begin{cases} 1 & \text{if } k = \underset{j}{\operatorname{argmin}} \|x^{(i)} - \mu_j\|^2 \\ 0 & \text{otherwise} \end{cases}$$

40, 5, 50  
3, 4, 5

= 4

$$J = \sum_{i=1}^m \sum_{k=1}^K r_k^{(i)} (x^{(i)} - \mu_k)^2$$



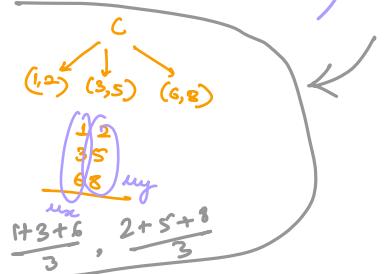
K-Means loss fn is non convex.

$$2 \sum_{i=1}^m r_k^{(i)} (x^{(i)} - \mu_k) = 0$$

$$\sum_{i=1}^m r_k^{(i)} x^{(i)} - \mu_k \sum_{i=1}^m r_k^{(i)} = 0$$

points w/ k-th cluster.

(Cluster center)



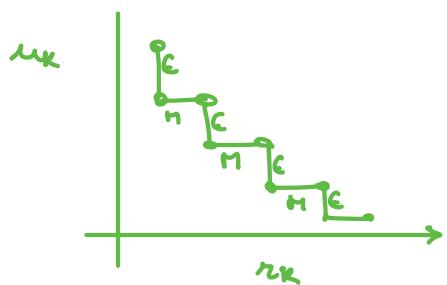
$$\mu_k = \frac{\sum_{i=1}^m r_k^{(i)} x^{(i)}}{\sum_{i=1}^m r_k^{(i)}}$$

M step

no of points in k-th cluster

M Step: find  $\mu_k$  keeping  $r_k^{(i)}$  fixed

E Step: find  $r_k^{(i)}$  keeping  $\mu_k$  as fixed.



it will converge to local minima.

Algo:

initialize cluster centroids  $\mu_1, \mu_2, \dots, \mu_k$   $\mu \in \mathbb{R}^n$

- Repeat until convergence:

- for every data point  $i$ , set

$$r_k^{(i)} = \begin{cases} 1 & \text{if } k = \underset{j}{\operatorname{arg\min}} \|x^{(i)} - \mu_j\|^2 \\ 0 & \text{otherwise} \end{cases}$$

] E Step

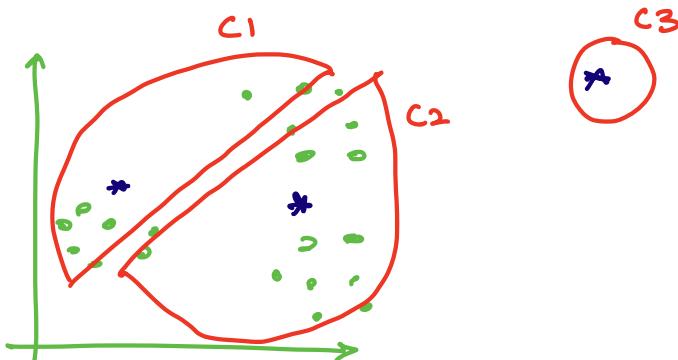
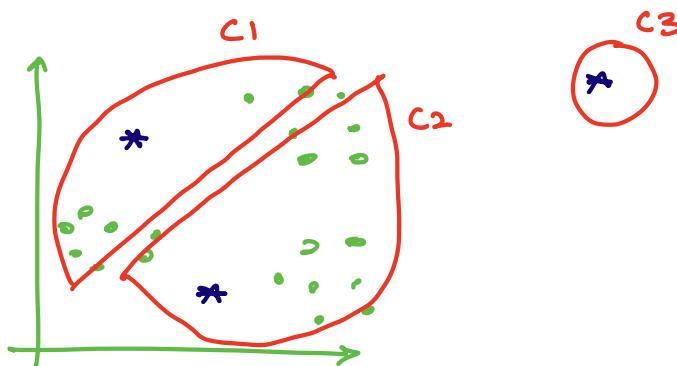
- for each  $k$ , set  $\mu_k = \frac{\sum_{i=1}^m r_k^{(i)} x^{(i)}}{\sum_{i=1}^m r_k^{(i)}}$

] M Step

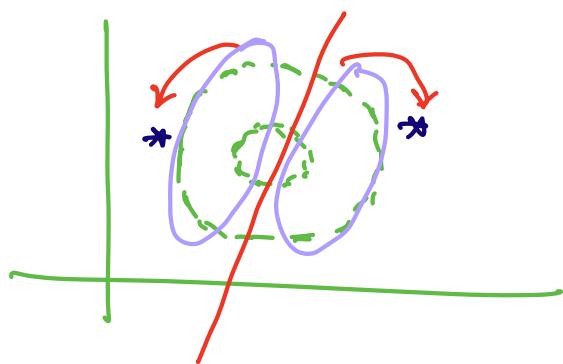
when k-means does not work ?

- random initialization (empty)
- non linearly separable  $\odot$
- one cluster is large as compared to other.

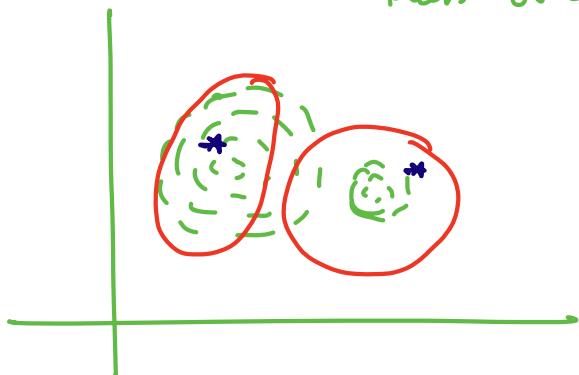
Random Initialization



No linearly Separable



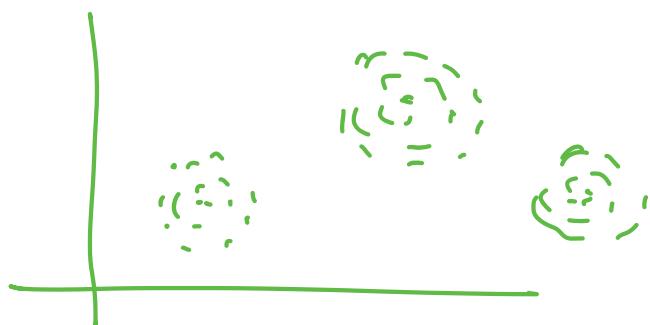
one cluster is larger than other



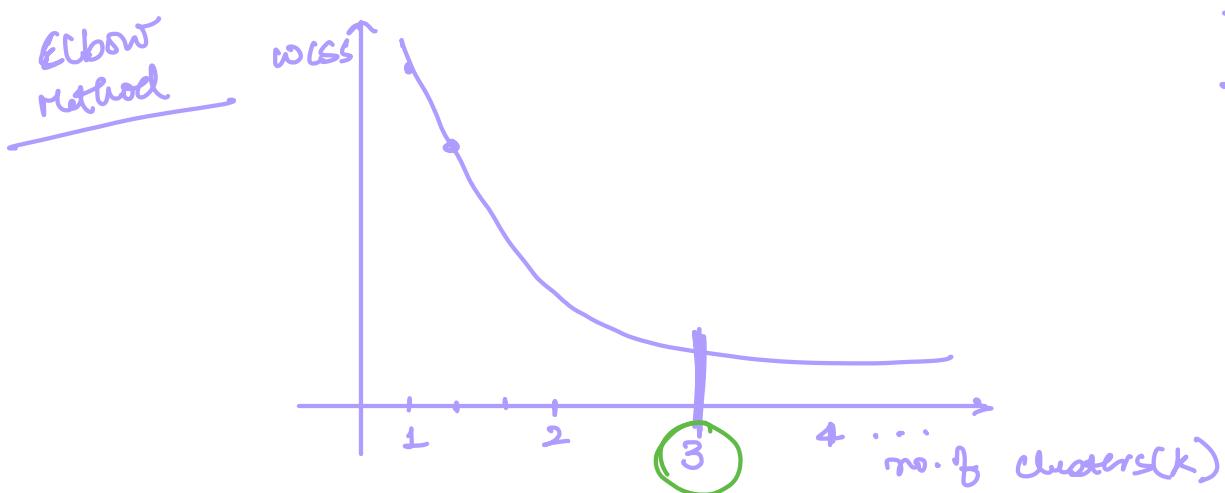
$k$  = no of clusters.  
↳ hyperparameter

$k=1$   
2  
3  
⋮  
⋮

Within Cluster Sum of Squares (WCSS)

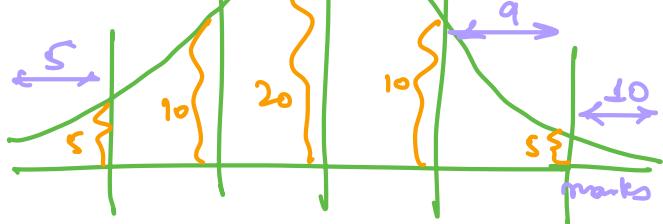
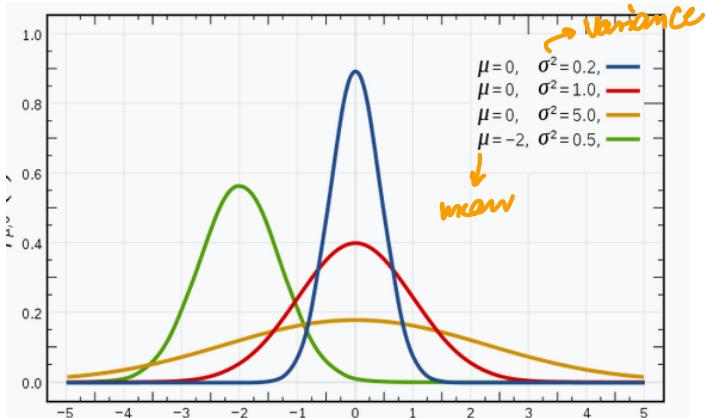


$$WCSS = \sum_{i \text{ in } C1} \text{distance}(x^{(i)}, u_1) + \sum_{i \text{ in } C2} \text{distance}(x^{(i)}, u_2) \dots$$



# Gaussian Mixture Models

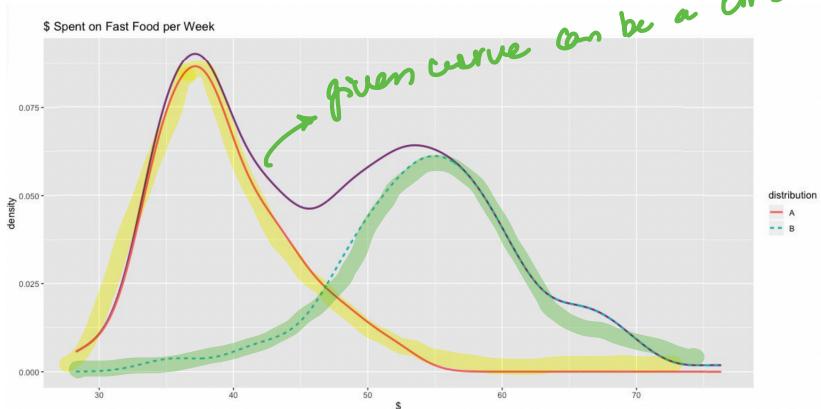
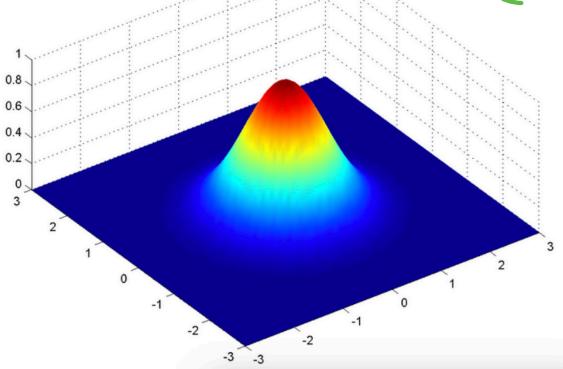
Normal Distribution



$$f(x | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

$$\sigma^2 = \frac{\sum (z_i - \bar{z})^2}{n-1}$$

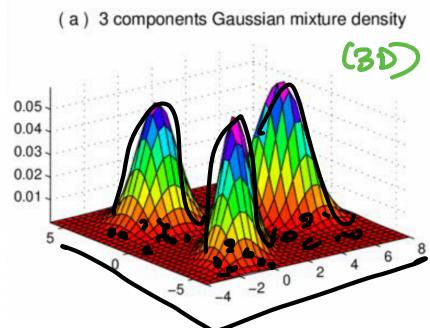
(3D)



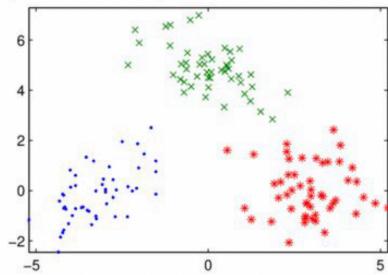
can be a combination of 2 different normal distributions' curves.

## Multivariate Normal Distributions

(2D)



(b) Data from 3 components Gaussian mixture density



- Hard Assignment
  - All Variances the Same
  - Roughly same # of datapoints
- K means**
- $E_{CK}$   $\theta_{CK}$   $0 \text{ or } 1$
- cluster center / mean / centroid

- Soft (probabilistic) Assignment
  - Variances can be different
  - Explicitly model # of data points
- GMM**
- $p(x)$   $p_1$   $p_2$   $p_3$
- $\mu$ ,  $\sigma^2$

## K-Means Algorithm

1. Choose  $k$  random points to be cluster centers
2. For each data point, assign it to the cluster whose center is closest
3. Using these assignments, recalculate the centers  $\rightarrow M$
4. Repeat 2 and 3 until either:
  - a. Cluster membership does not change
  - b. Centers change only a tiny amount

$$\sum_{i=1}^m \sum_{k=1}^K x_k^{(i)} \quad (\text{K-means})$$

$$\sum_{n=1}^N \sum_{k=1}^K x_{nk} \quad (\text{GMM})$$

## Gaussian Mixture Model (EM Algorithm)

1. Choose  $k$  random points to be cluster centers (or estimate using k-means...etc)
2. For each data point, calculate the probability of belonging to each cluster
3. Using these probability weights, recalculate the means + variances (and weights)
4. Repeat 2 and 3 until distributions converge.

$x \rightarrow c_1 ? p_1$   
 $\vdots$   
 $c_2 ? p_2$   
 $\vdots$   
 $c_3 ? p_3$



$\hat{x}_1, \hat{x}_2, \hat{x}_3$

Data is a weighted combination of various distributions.

$$p(x) = \sum_{k=1}^K w_k p_k(x) \rightarrow \text{Normal Distributions}$$

$$\text{Var}(x, x) = \text{variance}(x) = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}$$

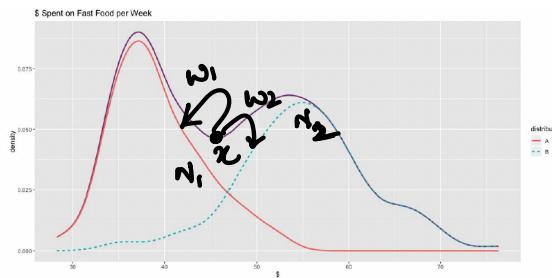
$$\text{Cov}(x, y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{n-1}$$

$$p(x) = \sum_{k=1}^K w_k p_k(x)$$

In Gaussian, it is normal distribution.

$$p(x) = \sum_{k=1}^K w_k \mathcal{N}(x | \mu_k, \Sigma_k)$$

mean  
Covariance



$$p(x) = \sum_{k=1}^K w_k \mathcal{N}(x|\mu_k, \Sigma_k)$$

Ref : <https://www.youtube.com/watch?v=sAjOZ9hkn84>

$$p(x) = \sum_{k=1}^K w_k \mathcal{N}(x|\mu_k, \Sigma_k)$$

Probability of being in group k

Likelihood of seeing x in group k

### Posterior Probabilities

$$p(\text{cluster } k|x) = \frac{w_k \mathcal{N}(x|\mu_k, \Sigma_k)}{\sum_{j=1}^K w_j \mathcal{N}(x|\mu_j, \Sigma_j)}$$

Prior probability of being in cluster k

Likelihood of seeing x in cluster k

Posterior probability of being in cluster k

$r_k^{(i)}$

### Maximum Likelihood Estimation

$$p(\mathbf{X}|\mathbf{w}, \mu, \Sigma) = p(x_1, x_2, \dots, x_n|\mathbf{w}, \mu, \Sigma) =$$

$$\prod_{n=1}^N \sum_{k=1}^K w_k \mathcal{N}(x_n|\mu_k, \Sigma_k)$$

$$\log(p(\mathbf{X}|\mathbf{w}, \mu, \Sigma)) = \sum_{n=1}^N \log \left\{ \sum_{k=1}^K w_k \mathcal{N}(x_n|\mu_k, \Sigma_k) \right\}$$

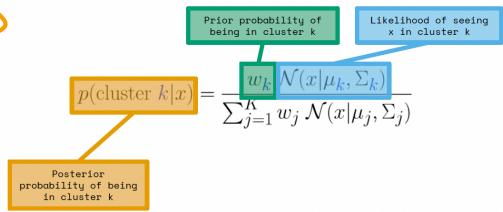
Goal: choose  $w$ ,  $\mu$ ,  $\Sigma$  that maximize the log likelihood

mean covariance

## Formulas (E-Step)

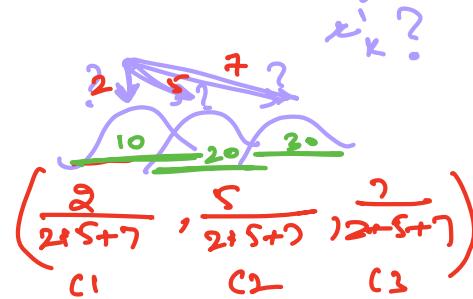
$$r_{nk} = \frac{w_k N(x_n | \mu_k, \Sigma_k)}{\sum_j w_j N(x_n | \mu_j, \Sigma_j)}$$

**Responsibilities** are the posterior probability of a data point being in cluster  $k$



$$\omega_1 = \frac{10}{10+20+30}$$

GMM  
E Step:



## Formulas (M-Step)

$$\mu_k = \frac{1}{N_k} \sum_{n=1}^N r_{nk} x_n$$

$$N_k = \sum_{n=1}^N r_{nk}$$

$$\Sigma_k = \frac{1}{N_k} \sum_{n=1}^N r_{nk} (x_n - \mu_k)(x_n - \mu_k)^T$$

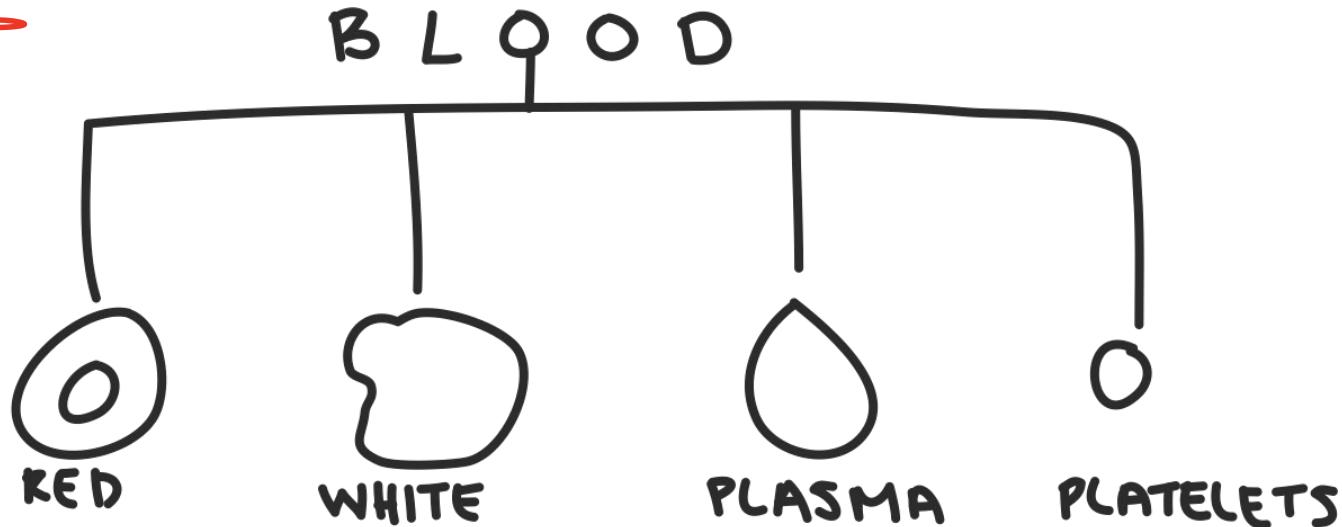
$$\omega_k = \frac{N_k}{N}$$

- GMM does **soft assignment**, every data point belongs to every cluster with some probability
- Data points that are more likely to be in a cluster have **more influence** over its parameters
- GMM uses the EM algorithm to iteratively update the cluster distributions. It first assigning a **responsibility** to each data point (**E-step**), and then using them to calculate **weighted means** and **variances** for each cluster (**M-step**)
- Responsibilities measure the **probability of a data point being in each cluster** (technically the **posterior probability**).  $\omega_k$
- Responsibilities contain information about **how common a cluster is** as well as the **likelihood of a data point belonging to that cluster**  $N(\mu_k, \Sigma_k)$ .

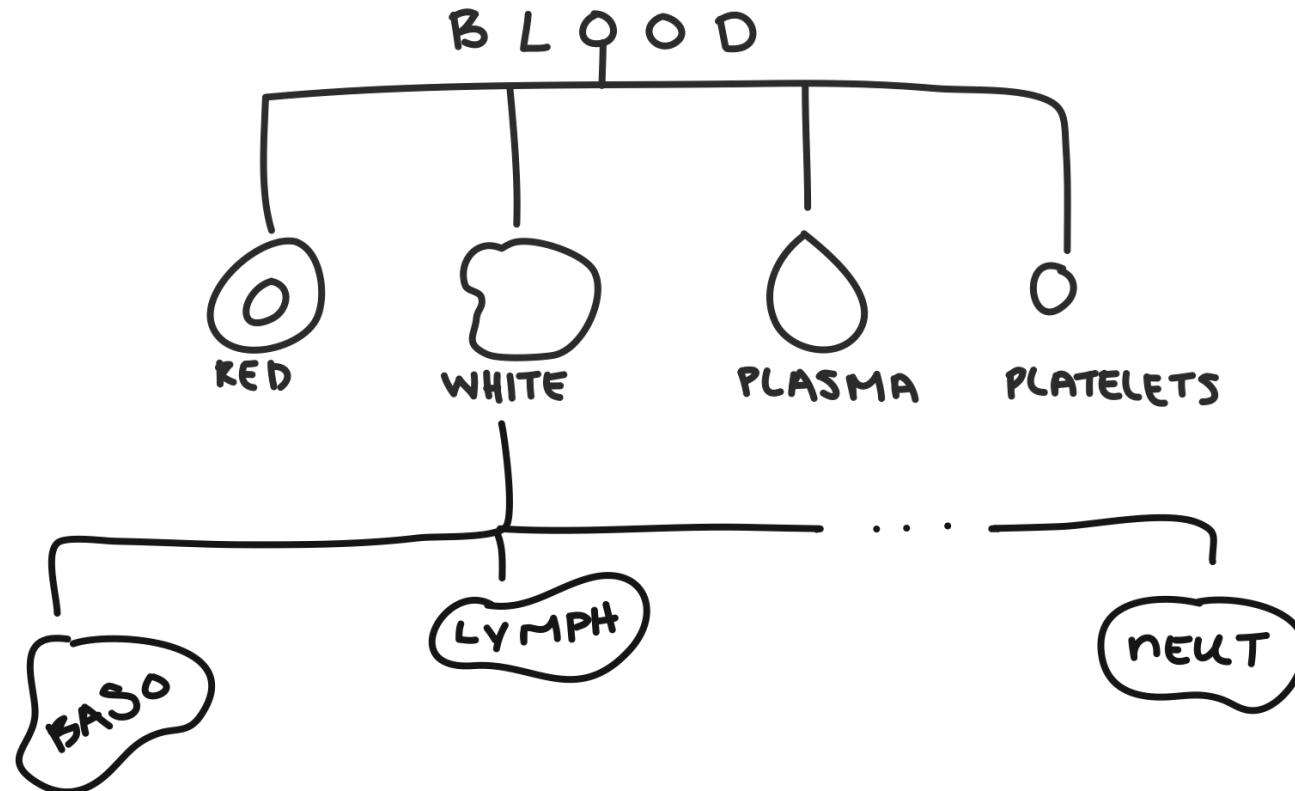
3 GMM

Hierarchical Based Clustering → Agglomerative (Bottom up)  
HAC → Divisive (Top Down)

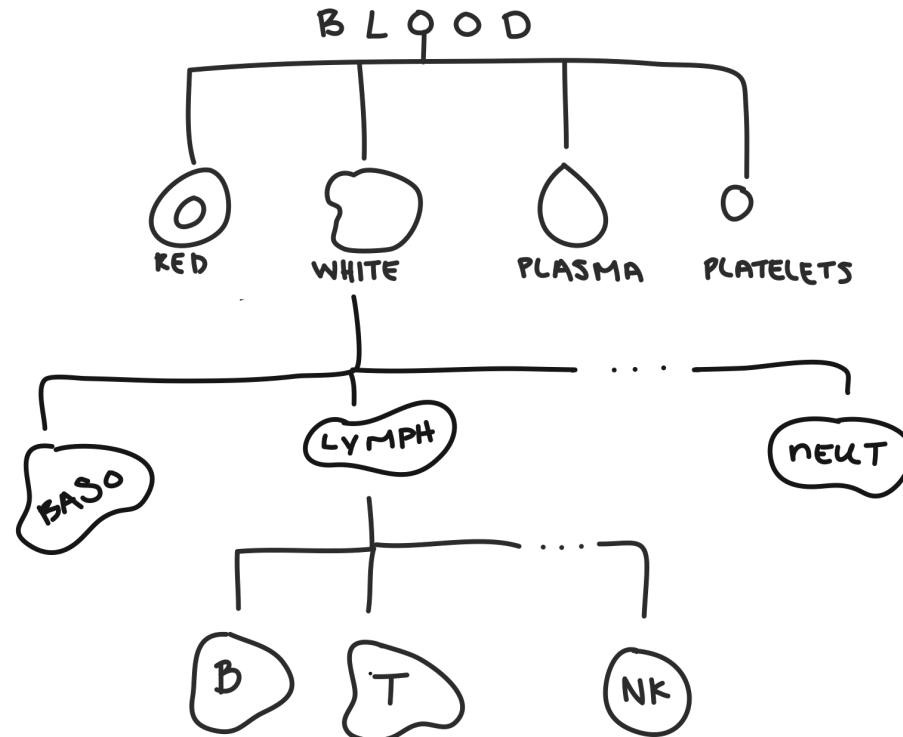
Divisive



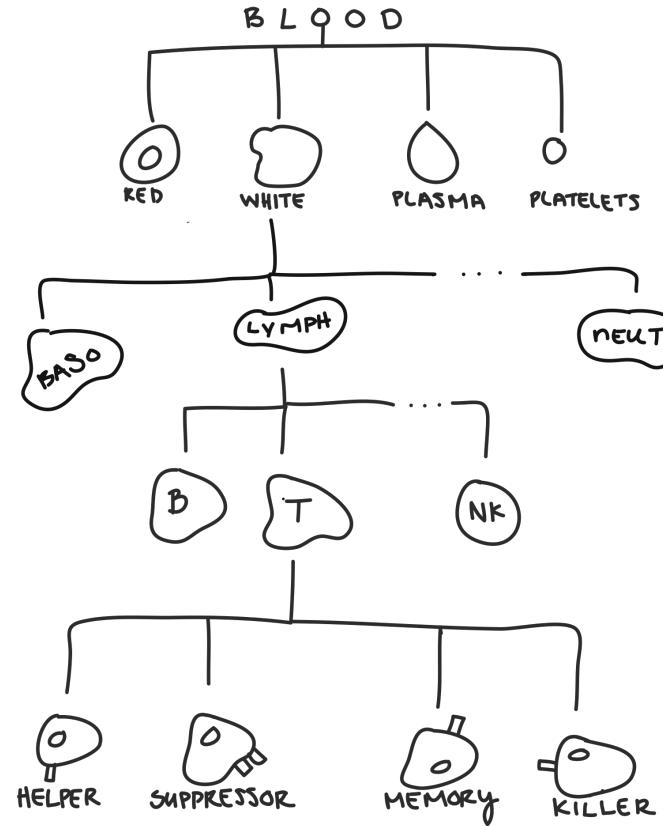
# HAC



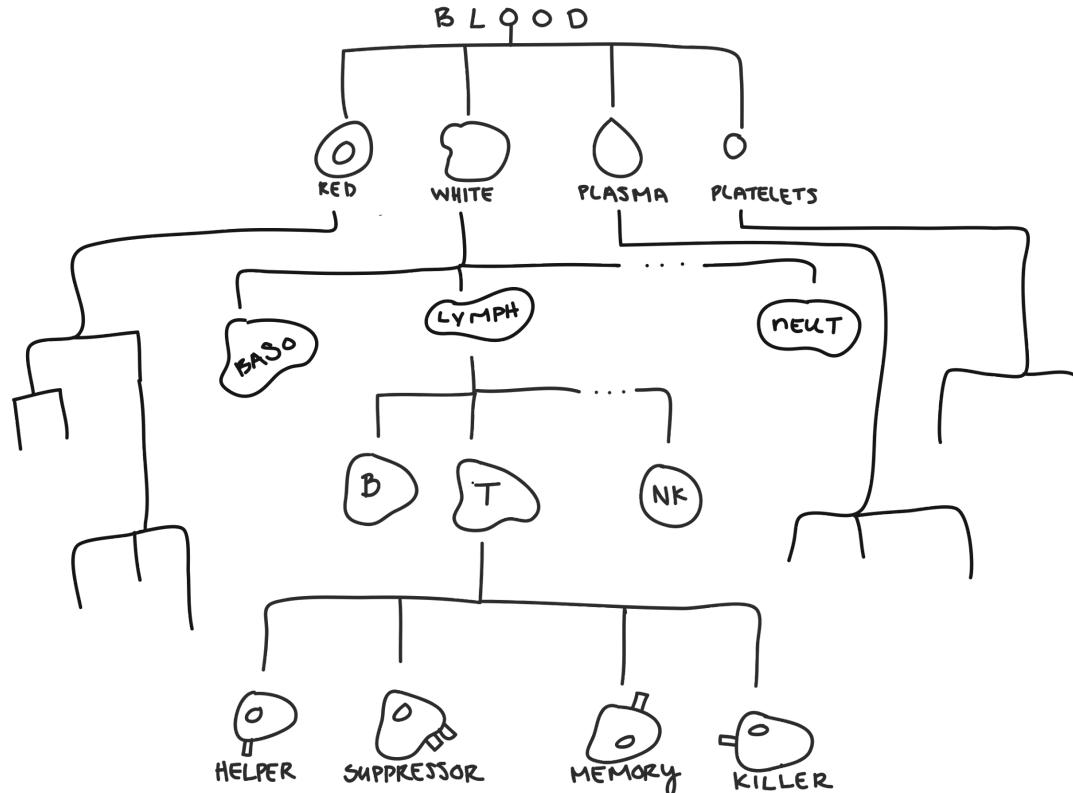
# HAC



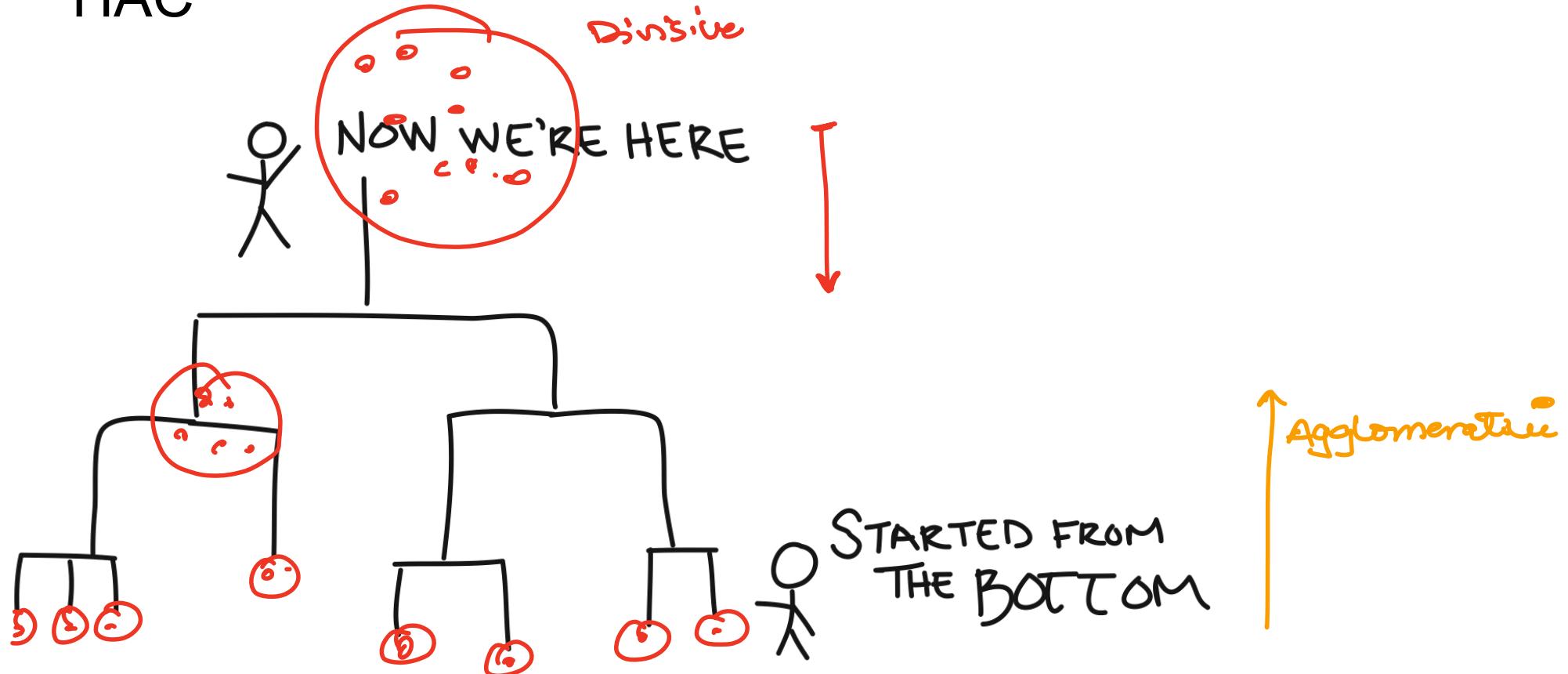
# HAC



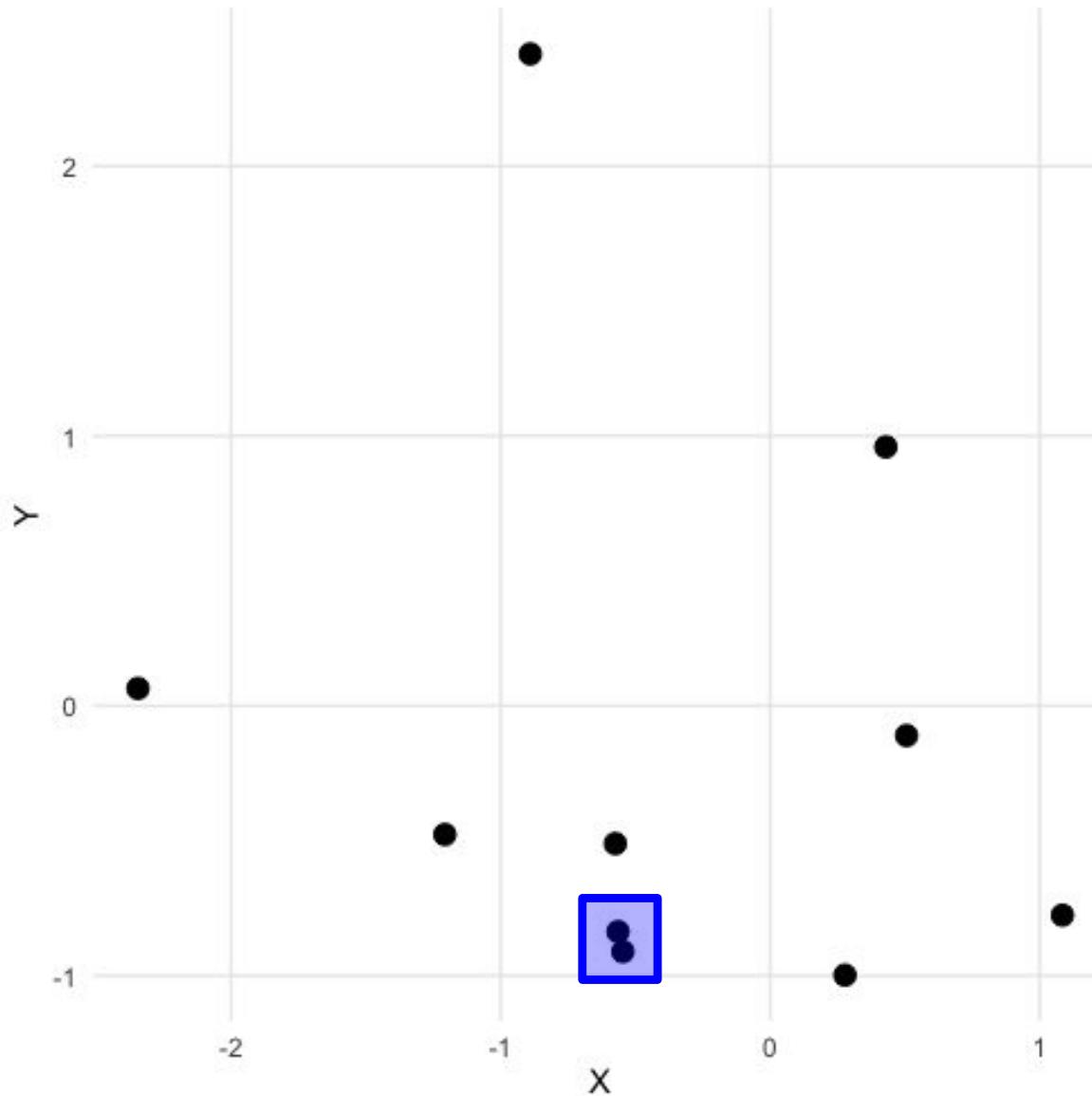
# HAC



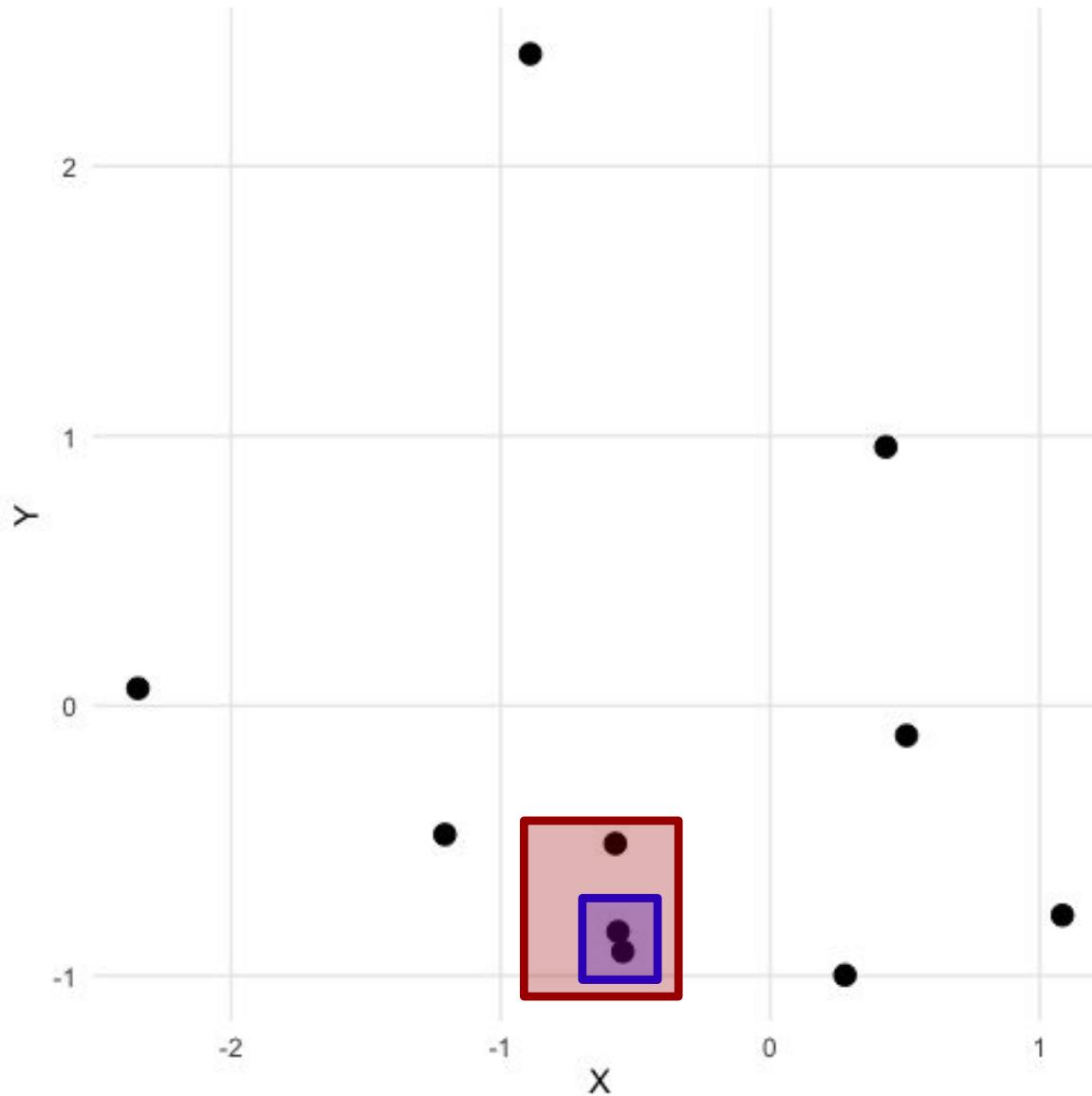
HAC



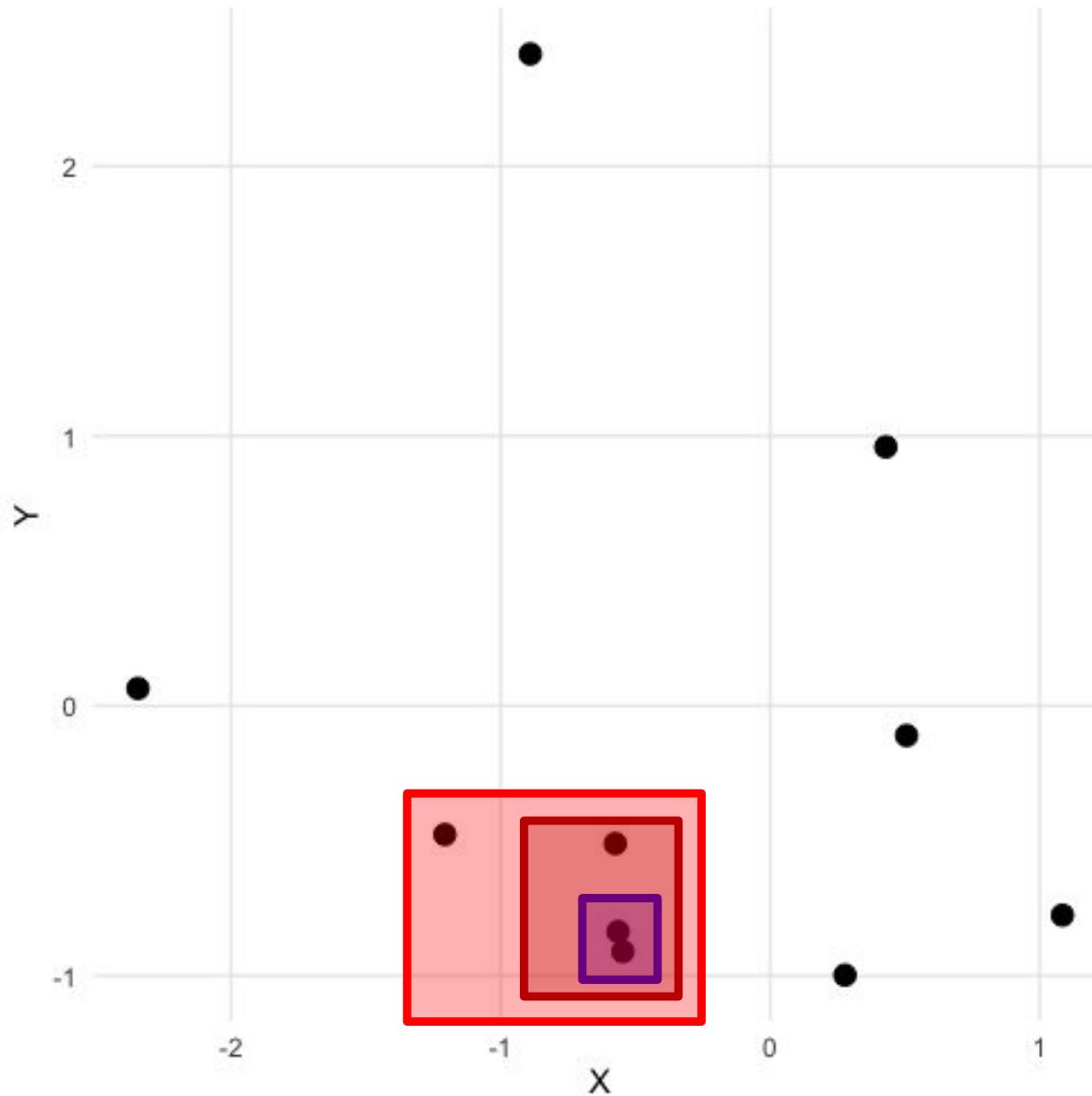
# Algorithm



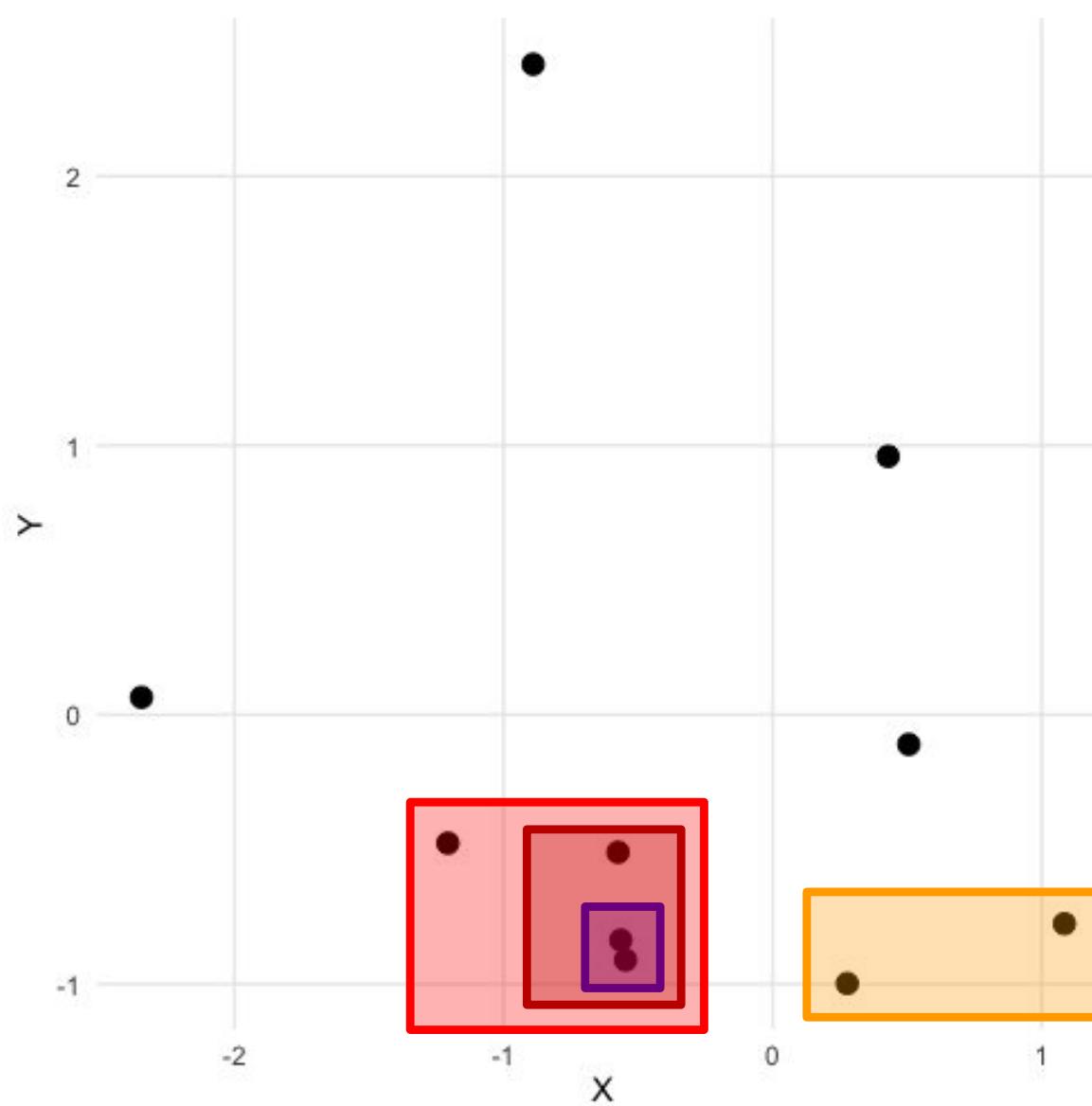
# Algorithm



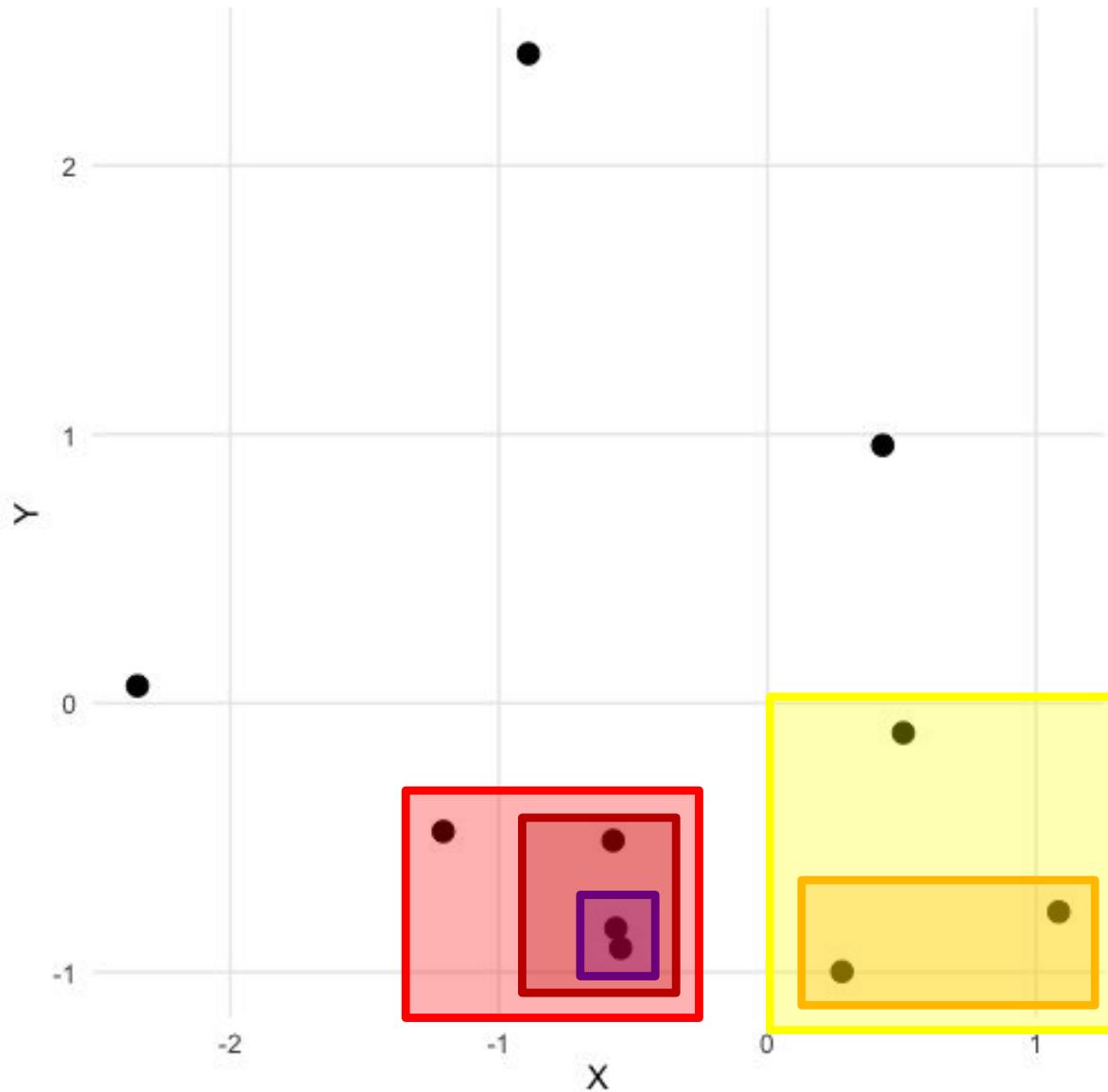
# Algorithm



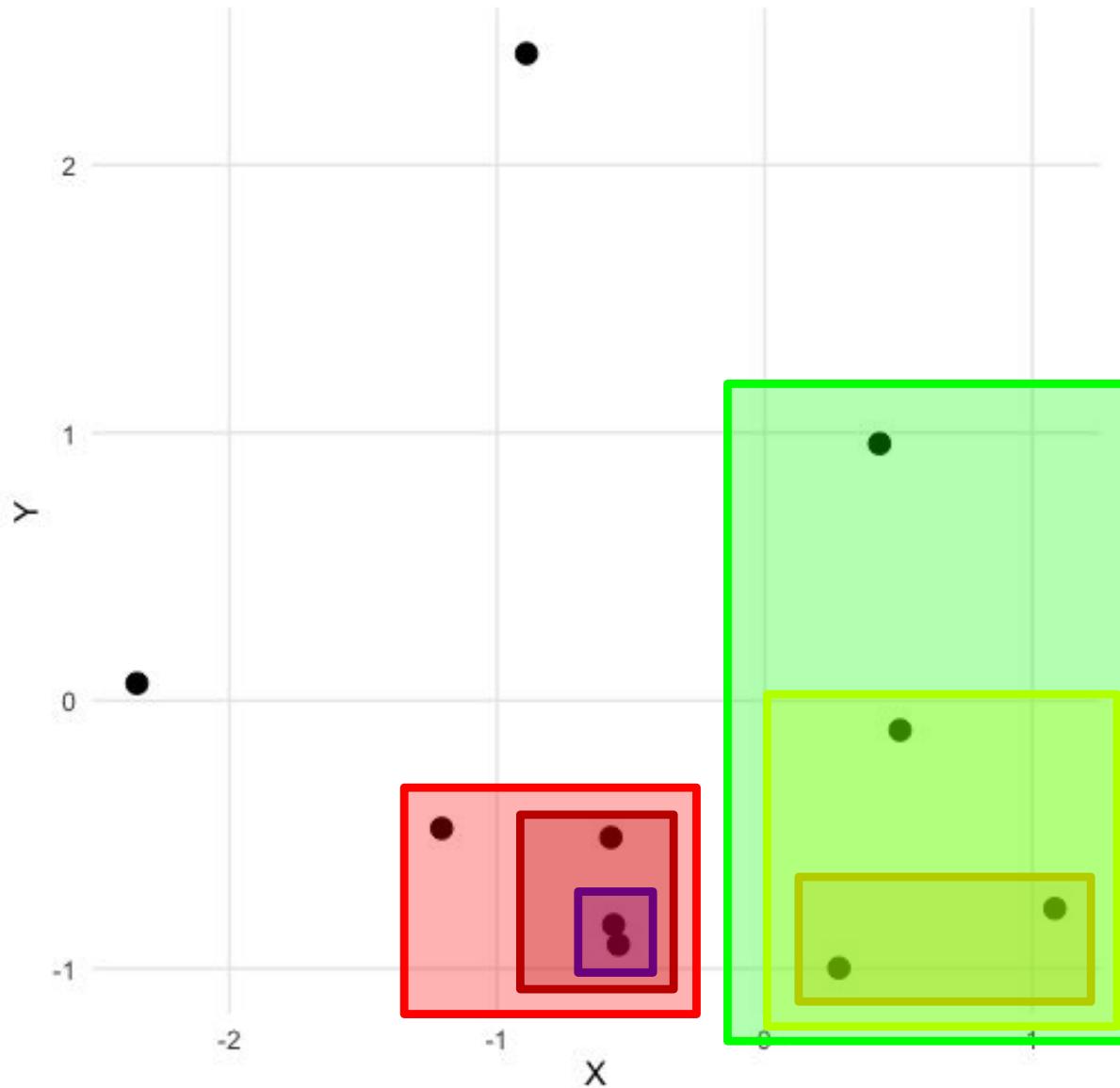
# Algorithm



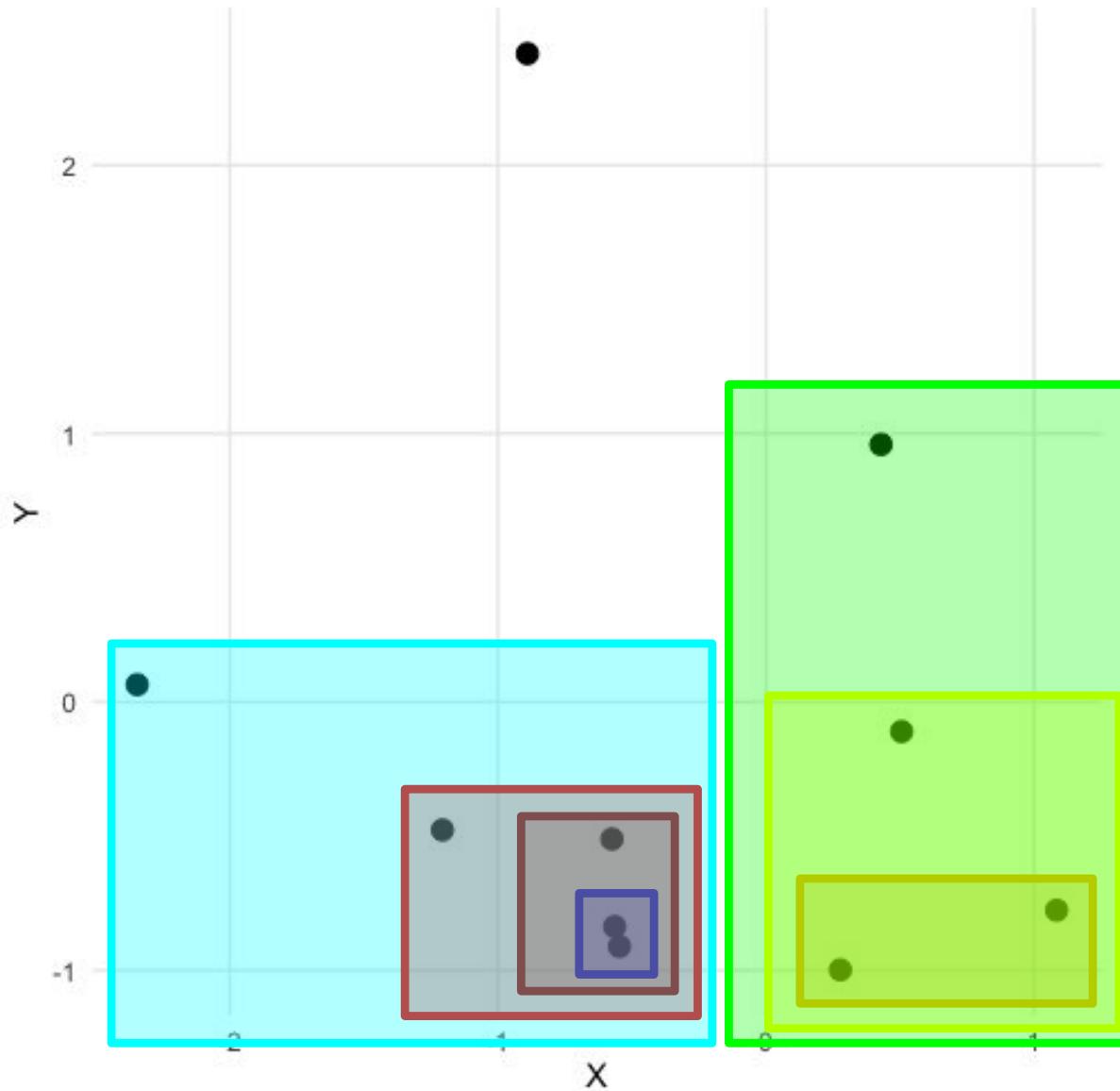
# Algorithm



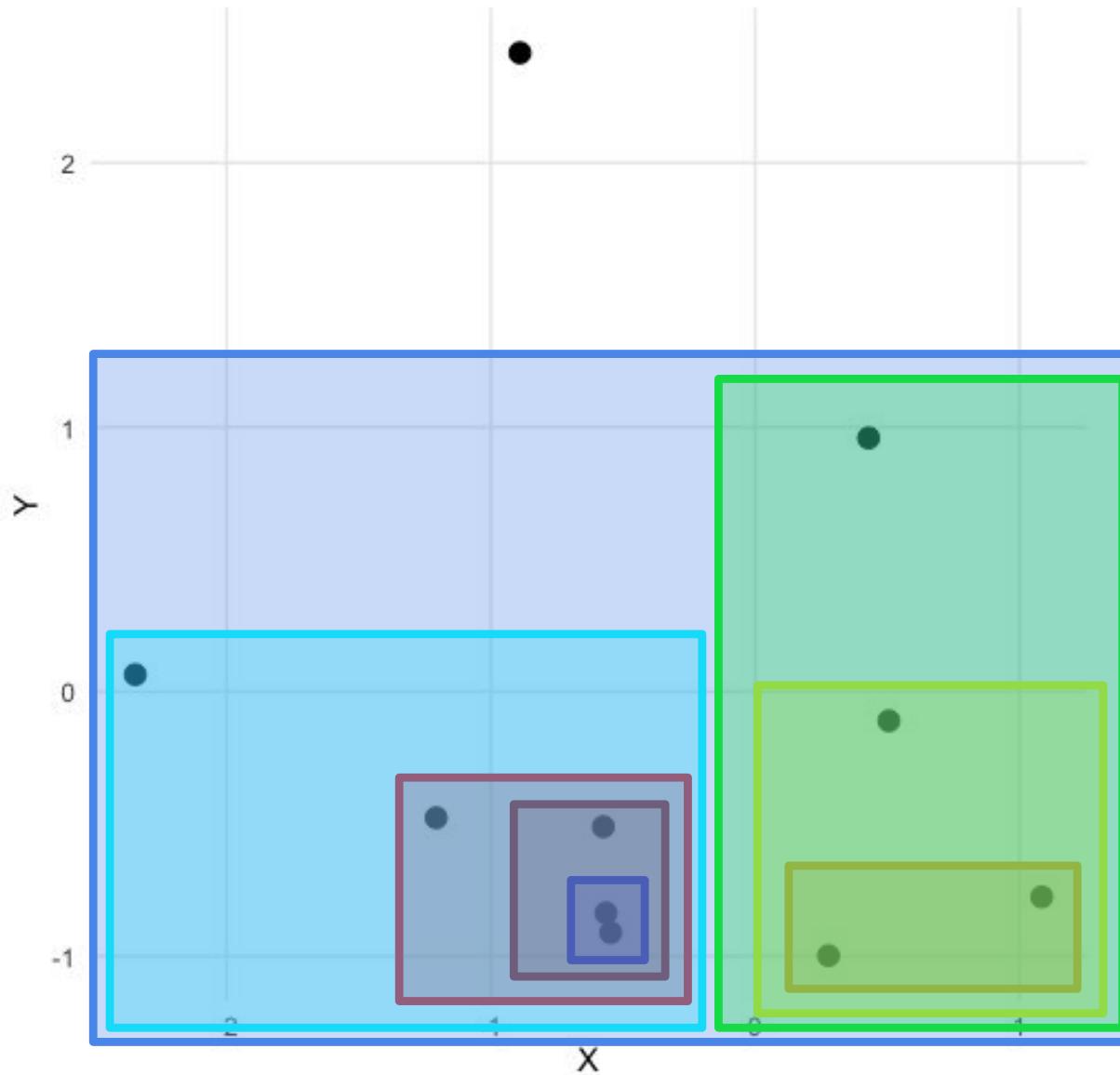
# Algorithm



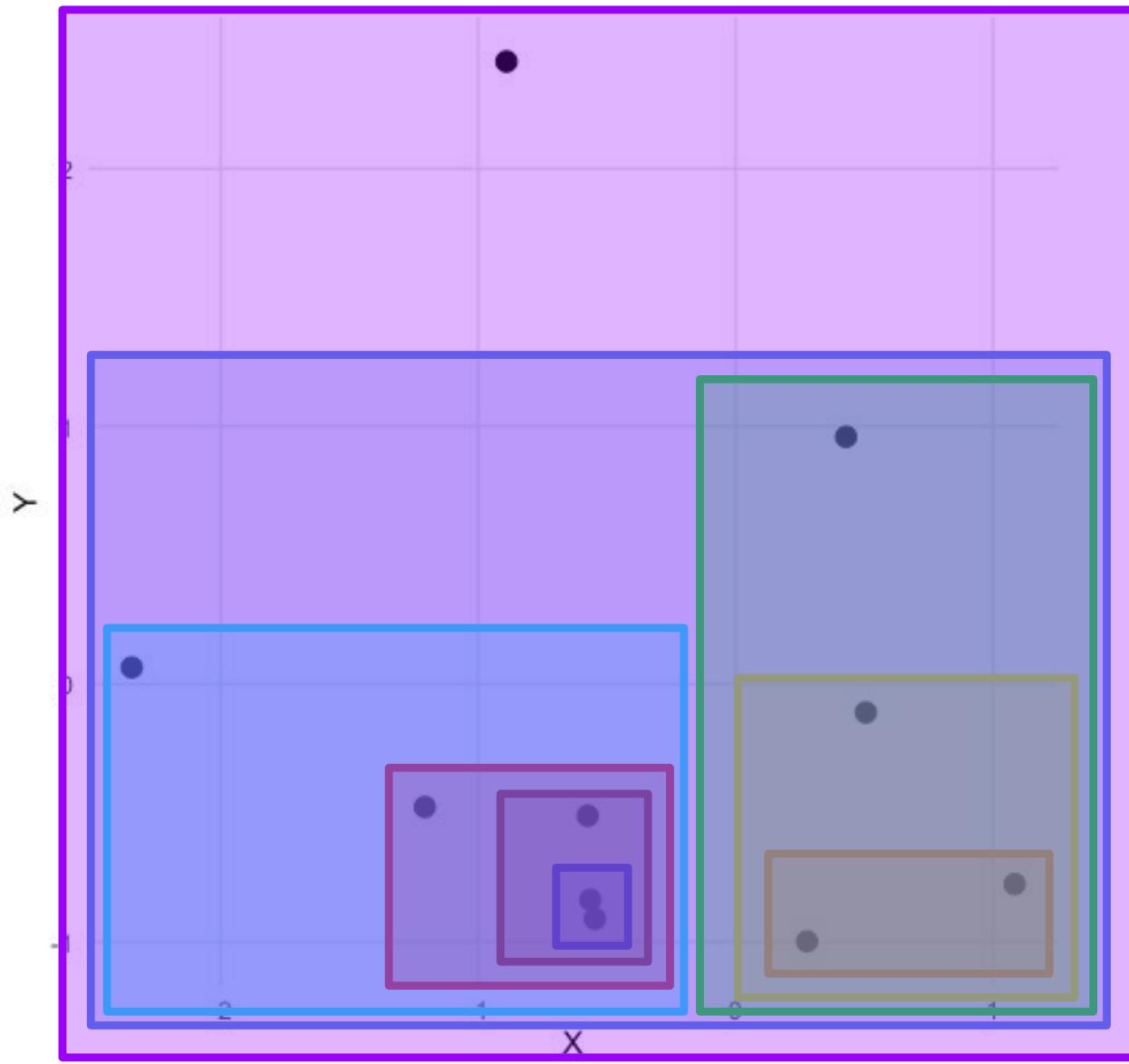
# Algorithm



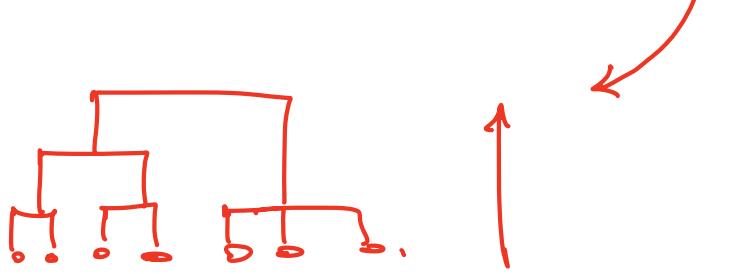
# Algorithm



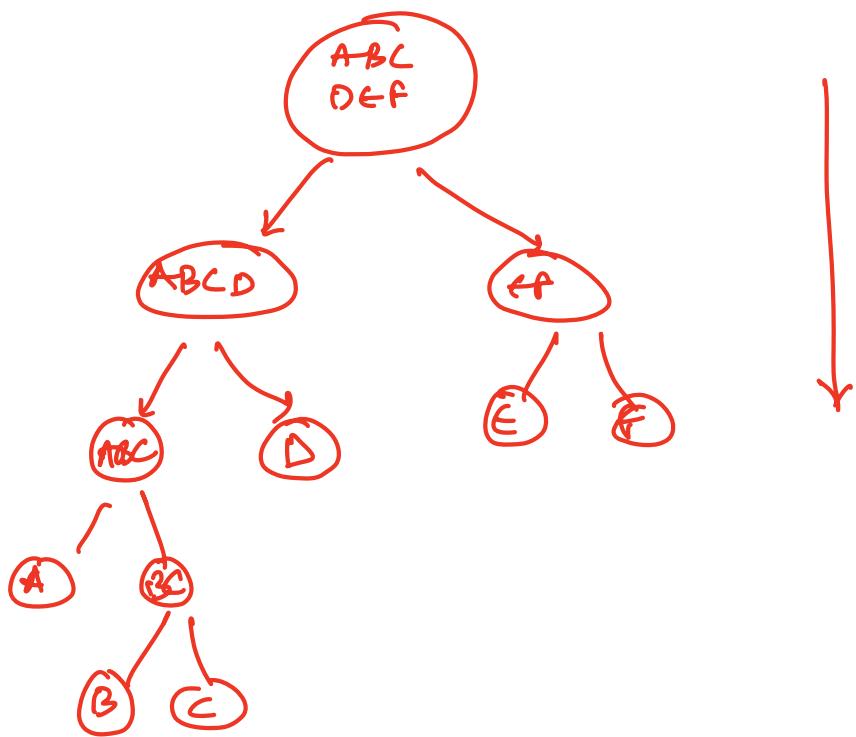
# Algorithm



Aggregation?  $\rightarrow$  BU

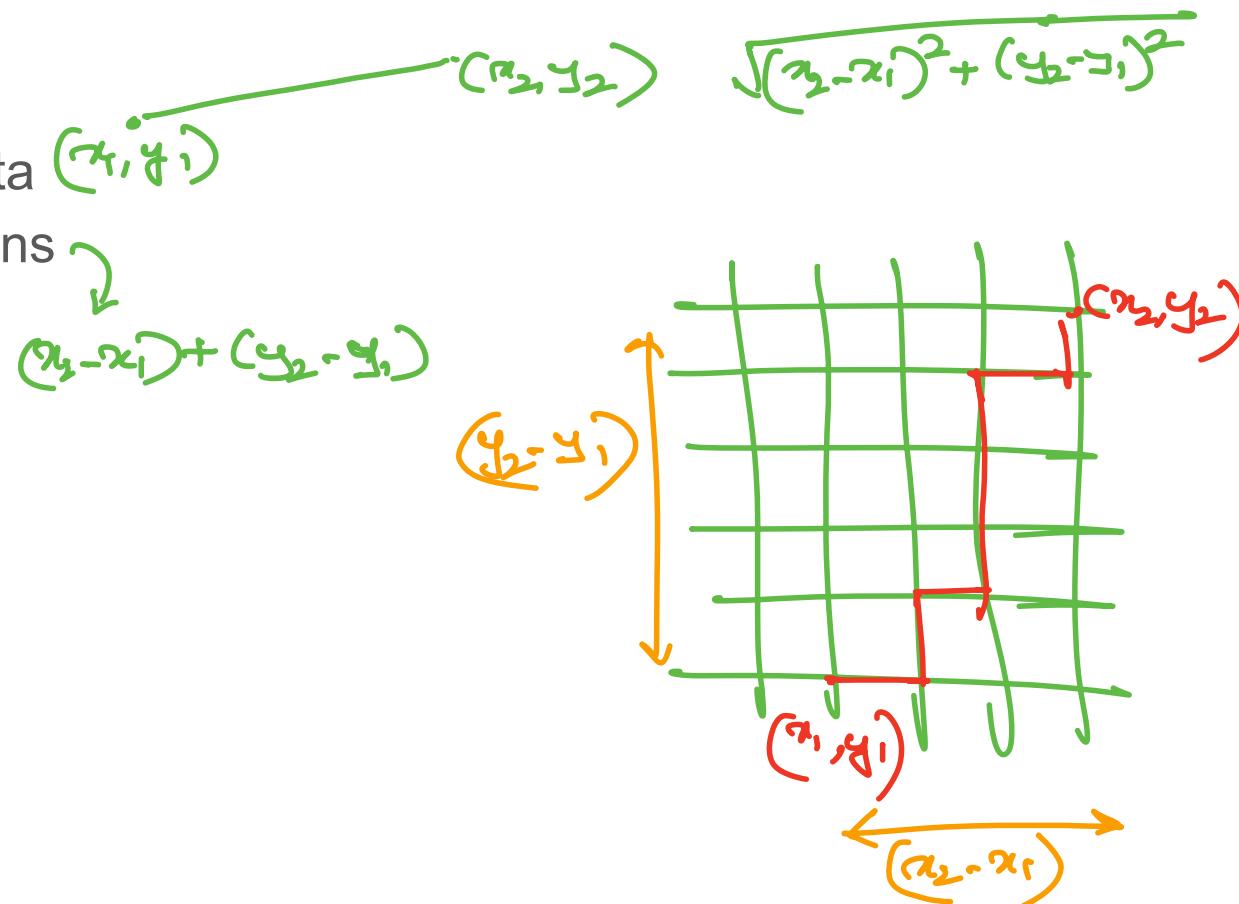
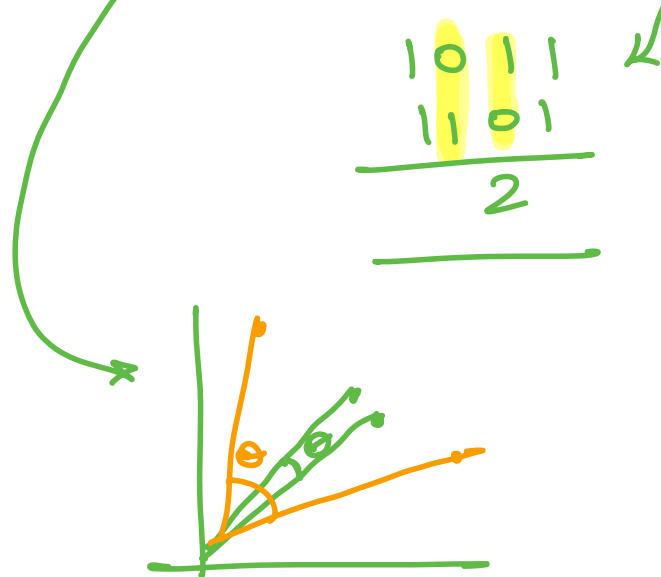


Divide  $\rightarrow$   
(Top Down)



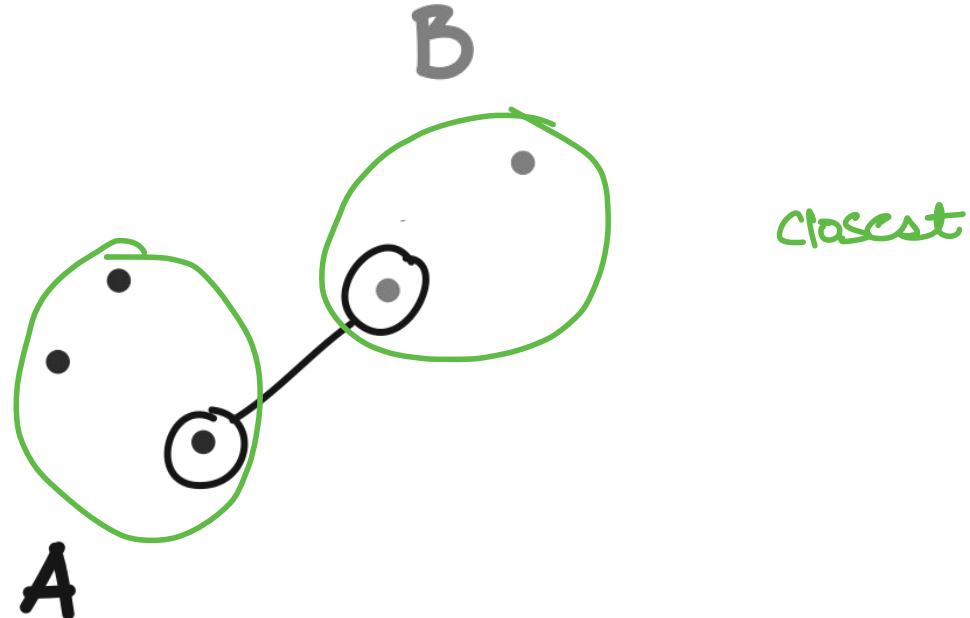
# Distance Metrics

- **Euclidean:** Continuous Data
- **Manhattan:** High Dimensions
- **Hamming:** Categories
- **Cosine:** Word Counts



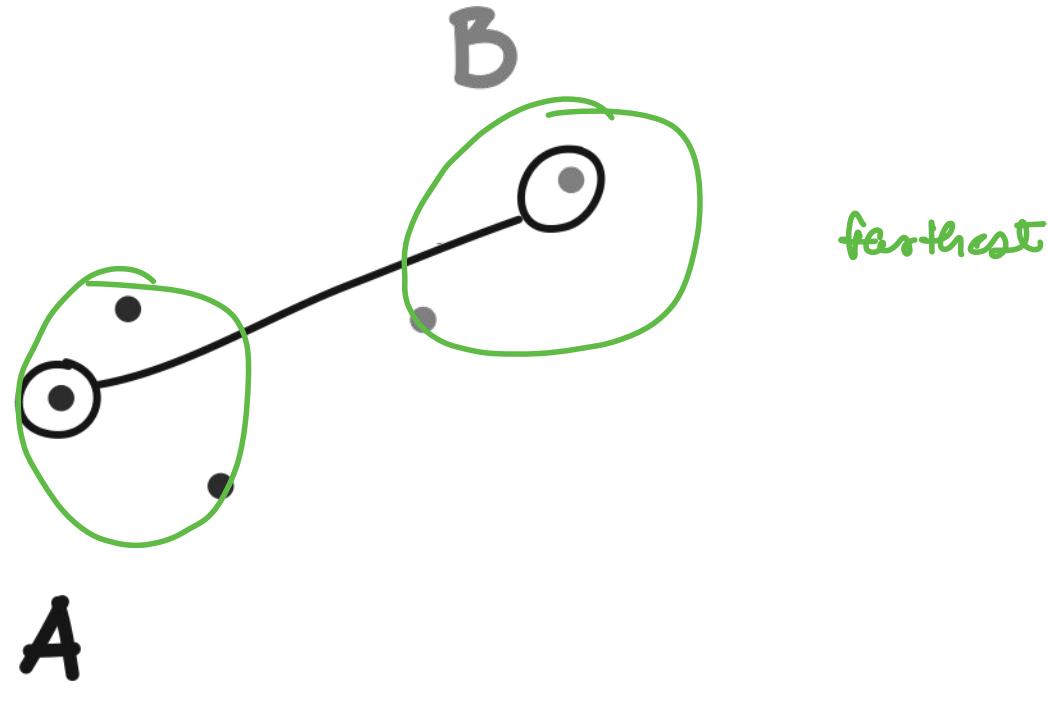
# Linkage Criteria

(Distance b/w 2 clusters)



**Single**

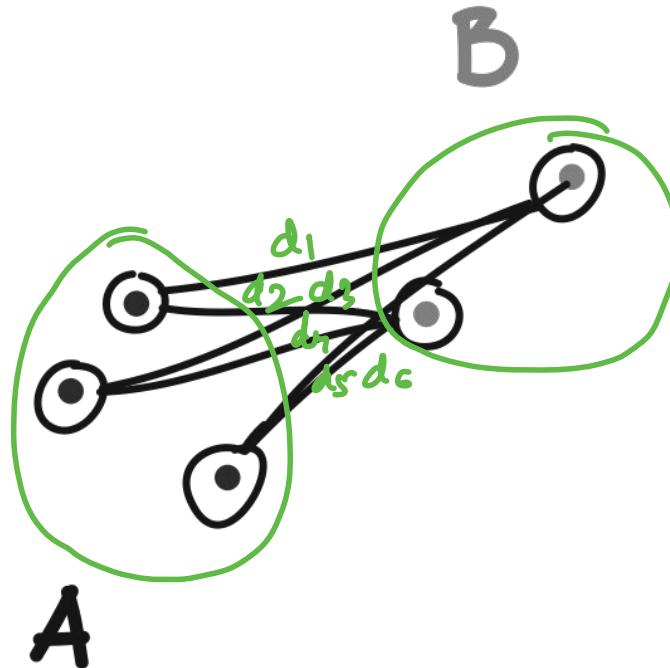
# Linkage Criteria



**Complete**

# Linkage Criteria

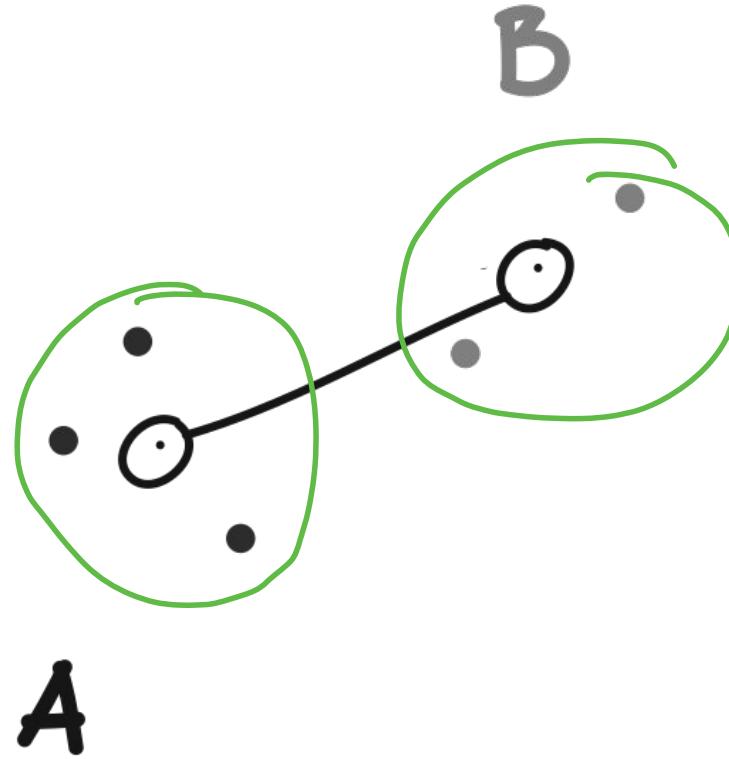
Average



$$\frac{d_1 + d_2 + d_3 + d_4 + d_5 + d_6}{6}$$

# Linkage Criteria

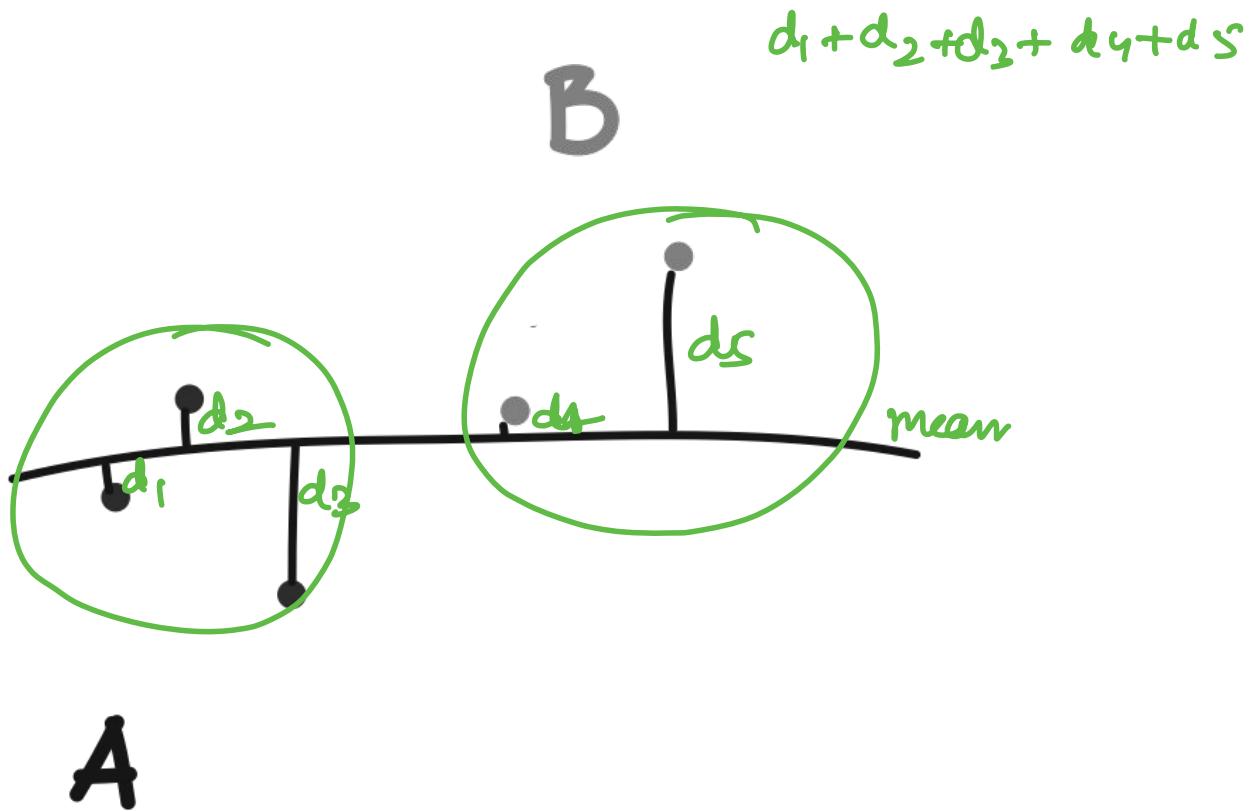
**Centroid**



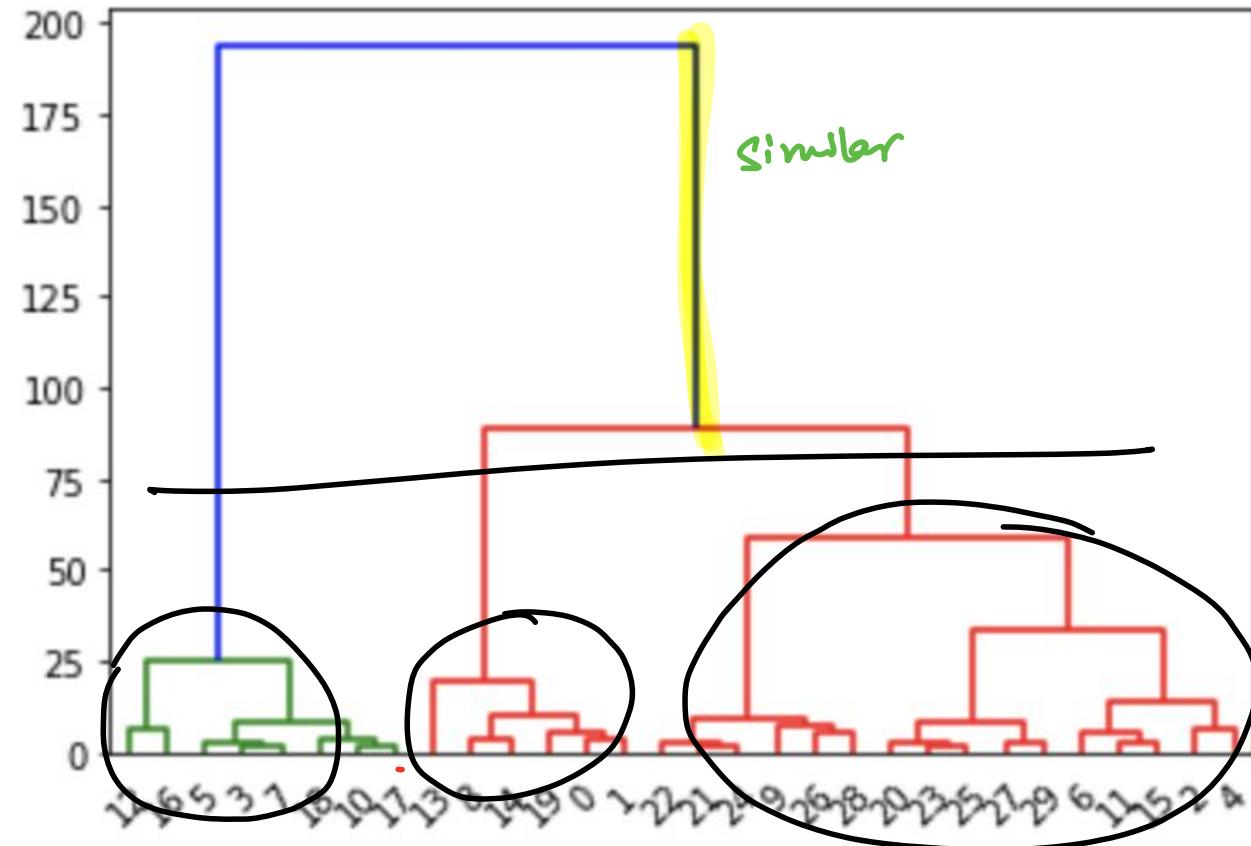
# Linkage Criteria

(Variance Based)

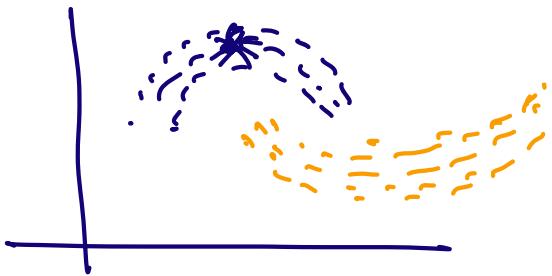
**Ward's**



# Reading a Dendrogram



## DBSCAN:



## feature selection:

so features  $\rightarrow$  all features are not relevant.

## feature selection?

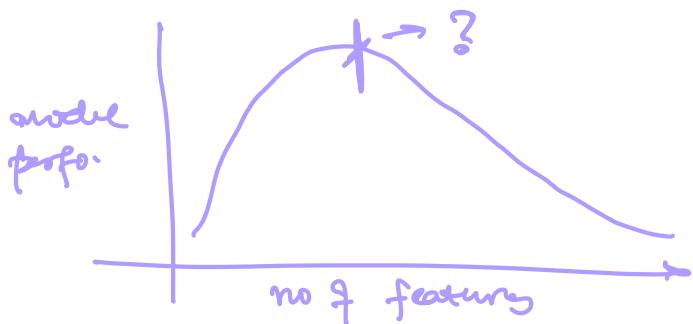
Curse of Dimensionality

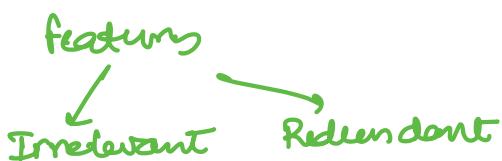
↪ lot of features, model performance  $\downarrow$

## Predict Salary

Humidity	Temp.	College	GPA	Branch	Exp.
----------	-------	---------	-----	--------	------

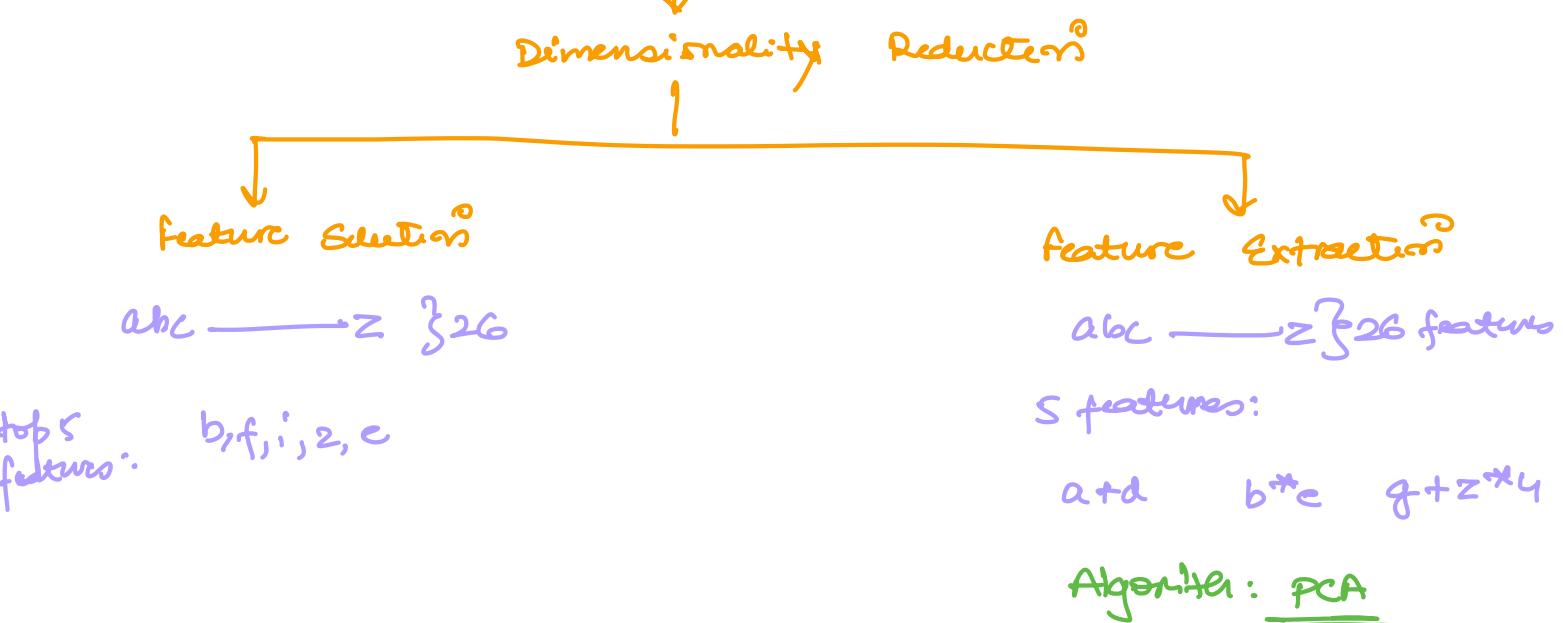

  
irrelevant X



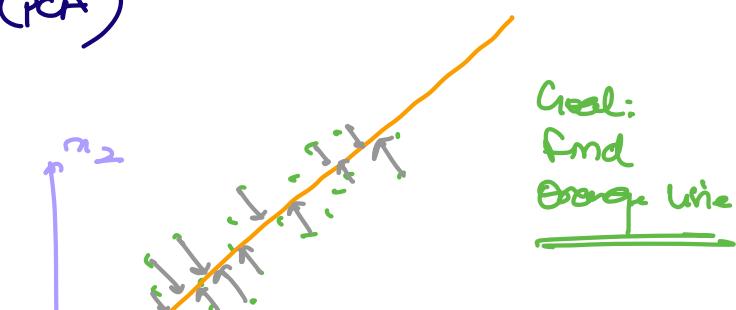
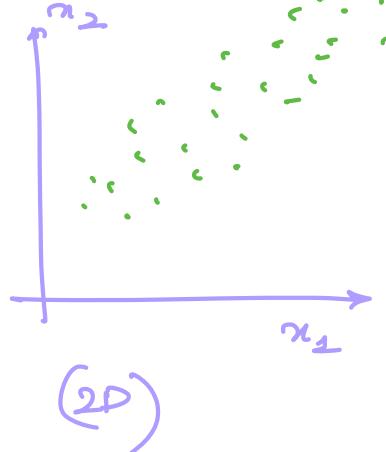

  
Irrelevant      Redundant

ht (cm)	wt (kg)	age	db
X	X		X

## Curse of Dimensionality



## Principal Component Analysis (PCA)



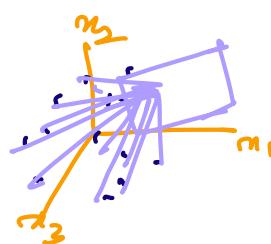
Goal:  
find  
Eigen Line

(2D)

(1D)

### Applications:

- Data Compression

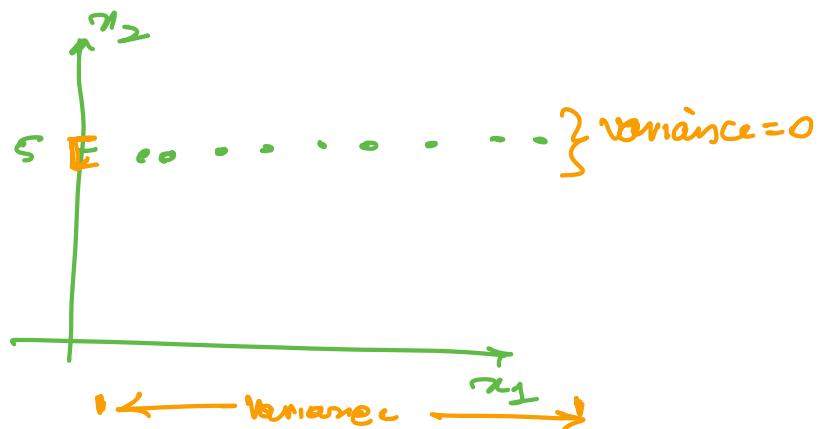


- Data Visualization

- Speed up Computation

Variance = Spread

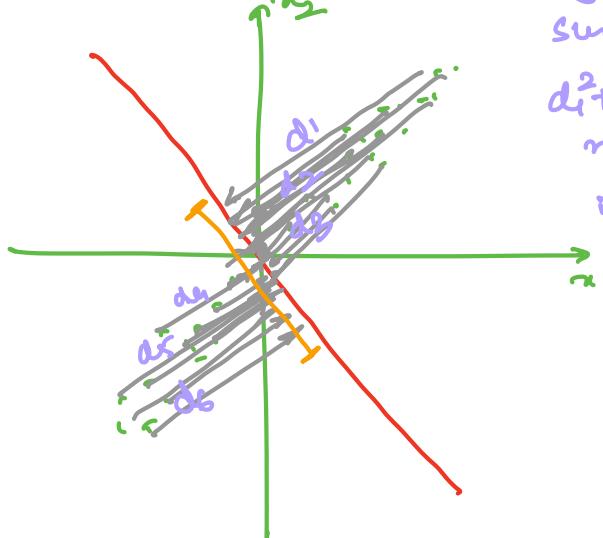
Eg:



ignore  $x_2$   
choose  $x_1$

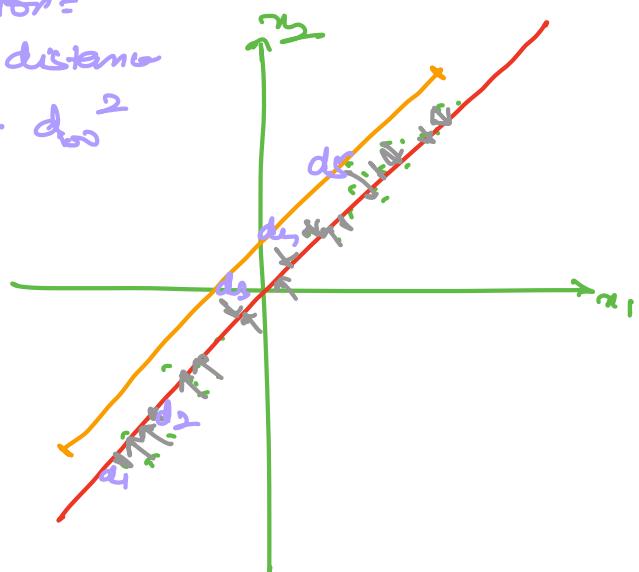
$$\text{Variance} = \frac{\sum (x - \bar{x})^2}{N}$$

Eg:



Projection Error =  
sum of  $d_i^2$   
 $d_1^2 + d_2^2 + \dots + d_m^2$

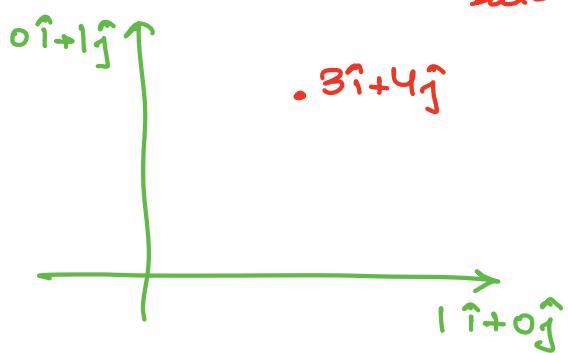
$$\sum_{i=1}^m d_i^2$$



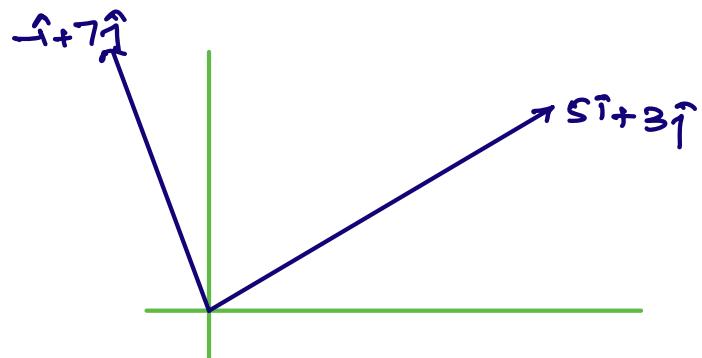
spread less  
Variance low  
(Projection error high)

spread high  
Variance high  
PE low

Vectors:



- original axis
- transformed axis
- data point



$$\text{Transformed} = 3(\text{Transformed } x) + 4(\text{Transformed } y)$$

$$\begin{bmatrix} 5 \\ 3 \end{bmatrix}$$

$$\begin{bmatrix} -1 \\ 7 \end{bmatrix}$$

$$= 3 \begin{bmatrix} 5 \\ 3 \end{bmatrix} + 4 \begin{bmatrix} -1 \\ 7 \end{bmatrix} = \begin{bmatrix} 15 \\ 9 \end{bmatrix} + \begin{bmatrix} -4 \\ 28 \end{bmatrix} = \begin{bmatrix} 11 \\ 37 \end{bmatrix}$$

Transformation  
Matrix

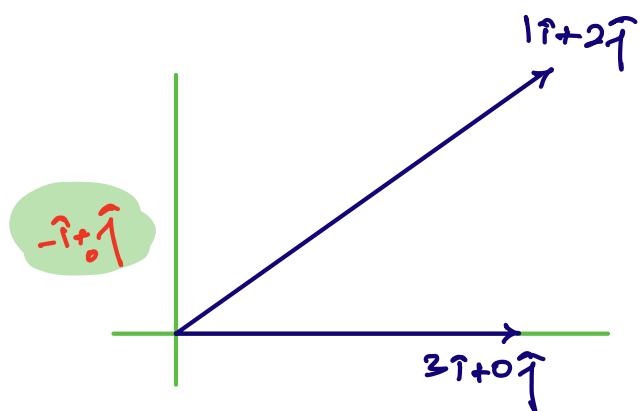
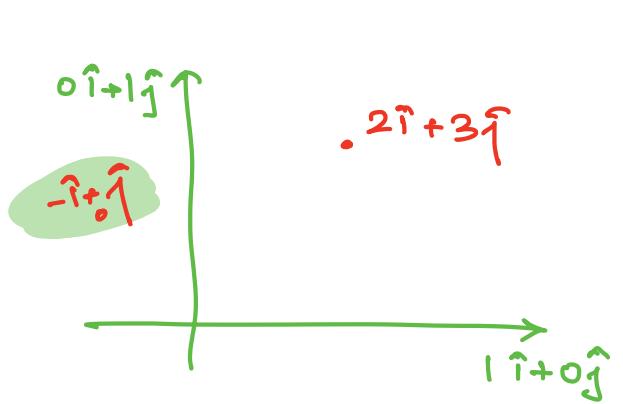
$$\begin{bmatrix} 5 \\ 3 \end{bmatrix} \quad \begin{bmatrix} -1 \\ 7 \end{bmatrix} \quad \begin{bmatrix} 3 \\ 4 \end{bmatrix}$$

new x      new y      coordinates  
vector

$$= \begin{bmatrix} 15-4 \\ 9+28 \end{bmatrix} = \begin{bmatrix} 11 \\ 37 \end{bmatrix}$$

$$11\hat{i} + 37\hat{j}$$

Eigenvectors & Eigenvalues



$$2\hat{i} + 3\hat{j} \quad \begin{bmatrix} 3 & 1 \\ 0 & 2 \end{bmatrix} \begin{bmatrix} 2 \\ 3 \end{bmatrix} = \begin{bmatrix} 6+3 \\ 0+6 \end{bmatrix} = \begin{bmatrix} 9 \\ 6 \end{bmatrix} = 3 \begin{bmatrix} 3 \\ 2 \end{bmatrix}$$

$$-i + j \quad \begin{bmatrix} 3 & 1 \\ 0 & 2 \end{bmatrix} \begin{bmatrix} -1 \\ 1 \end{bmatrix} = \begin{bmatrix} -3+1 \\ 2 \end{bmatrix} = \begin{bmatrix} -2 \\ 2 \end{bmatrix} = 2 \begin{bmatrix} -1 \\ 2 \end{bmatrix}$$

$A \cdot v = \lambda \cdot v$



$$A\vec{v} = \lambda \cdot I \cdot \vec{v}$$

$$(A - \lambda I) \vec{v} = 0$$

$$|A - \lambda I| = 0$$

$$\begin{bmatrix} 3 & 1 \\ 0 & 2 \end{bmatrix} - \begin{bmatrix} \lambda & 0 \\ 0 & \lambda \end{bmatrix} = 0$$

$$\begin{bmatrix} 3-\lambda & 1 \\ 0 & 2-\lambda \end{bmatrix} = 0$$

$$(2-\lambda)(2-\lambda) = 0$$

$$\underline{\lambda = 2, 3} \rightarrow \underline{\text{Eigen Values}}$$

$$I = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$\lambda I = \begin{bmatrix} \lambda & 0 \\ 0 & \lambda \end{bmatrix}$$

$\text{cov}(x_1, x_1)$	$\text{cov}(x_1, x_2)$
$\text{cov}(x_2, x_1)$	$\text{cov}(x_2, x_2)$

Covariance Matrix =

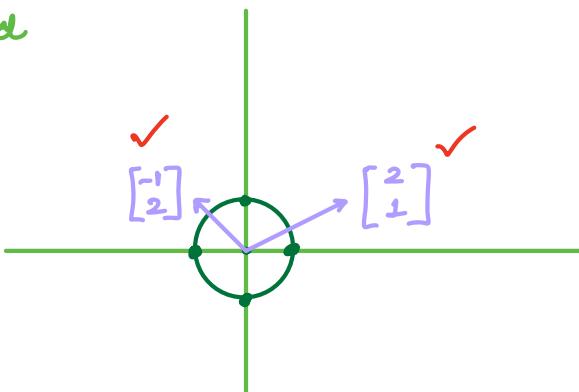
$$\begin{pmatrix} \text{cov}(x_1, x_1) & \text{cov}(x_1, x_2) \\ \text{cov}(x_2, x_1) & \text{cov}(x_2, x_2) \end{pmatrix} = \begin{pmatrix} 9 & 4 \\ 4 & 3 \end{pmatrix}$$

$$\text{cov}(x, y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{n-1}$$

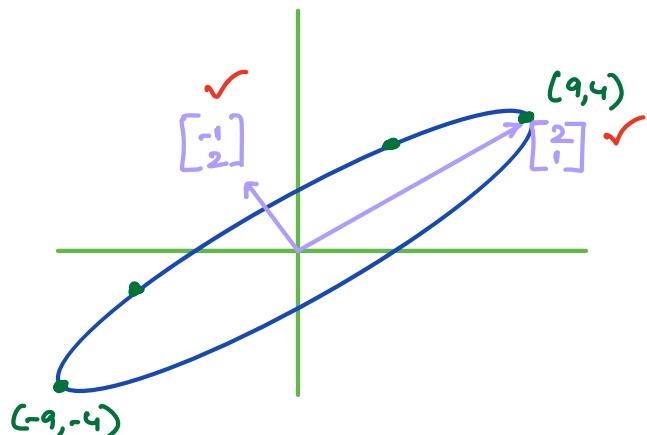
$$\begin{aligned} \text{cov}(x, x) &= \text{variance}(x) \\ &= \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1} \end{aligned}$$

- Symmetric
- Used as transformation matrix

old



new



$$\begin{pmatrix} 9 & 4 \\ 4 & 3 \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix} = x \begin{bmatrix} 9 \\ 4 \end{bmatrix} + y \begin{bmatrix} 4 \\ 3 \end{bmatrix} = \begin{bmatrix} 9x + 4y \\ 4x + 3y \end{bmatrix} \checkmark$$

→ Transformation matrix

$$(x, y) \rightarrow (9x + 4y, 4x + 3y)$$

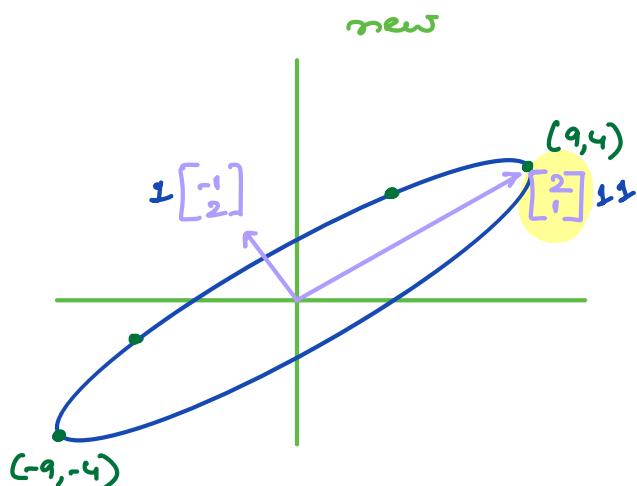
$$(0, 0) \rightarrow (0, 0)$$

circle

$$\left\{ \begin{array}{l} (1, 0) \rightarrow (9, 4) \\ (0, 1) \rightarrow (4, 3) \\ (-1, 0) \rightarrow (-9, -4) \\ (0, -1) \rightarrow (-4, -3) \end{array} \right. \quad \left. \begin{array}{l} \text{flat} \rightarrow \text{dimple} \end{array} \right\}$$

$$\left. \begin{array}{l} (2, 1) \rightarrow (22, 11) = 11(2, 1) \\ (-1, 2) \rightarrow -1(-1, 2) \end{array} \right. \quad \left. \begin{array}{l} \text{Eigen vector} \\ \text{Eigen value} \end{array} \right\}$$

In PCA, goal was to find new axis.



$\begin{bmatrix} -1 \\ 2 \end{bmatrix}$  and  $\begin{bmatrix} 2 \\ 1 \end{bmatrix}$  are going to be the new axis.

These are the lines having highest spread.

Highest variance will be given by Eigen vector.

EV1 →  $\frac{11}{\sqrt{11}}$  ] Eigen value 1  
EV2 →  $\frac{1}{\sqrt{11}}$  ] magnitude

Spread

We will project all points on Eigen vector 1 bcz its magnitude is larger.

### Steps:

1. Load Data

2. Standardization

$$\begin{aligned} \mu &= 0 & (\text{mean} = 0) \\ \sigma &= 1 & (\text{std} = 1) \end{aligned}$$

3. Covariance Matrix : how are 2 features related to each other.

n features  
4 4

	1	2	3	4
1				
2				
3				
4				

n x n  
4 x 4

4. Eigen Values & Eigen Vector

④ {  
evec1 → value 1  
evec2 → value 2  
evec3 → value 3  
evec4 → value 4

5. Larger eigen value means spread / variance is more along that vector.

evec1 → 50  
evec2 → 30  
evec3 → 60  
evec4 → 10

{ eigen values

Choose the axis along which spread is larger.  
Spread is given by eigen value.

Sort them acc. to eigen value

$$\begin{array}{l} \text{evec}_3 \rightarrow 60 \\ \text{evec}_1 \rightarrow 50 \\ \text{evec}_2 \rightarrow 30 \\ \text{evec}_4 \rightarrow 10 \end{array} \quad \left. \begin{array}{c} \text{Sorting} \\ \text{vectors} \end{array} \right\} \xrightarrow{k^2} \underline{k \text{ vectors.}}$$

6. Pick top  $k$  eigenvectors.

new dimension.

7. Projections

$$\begin{bmatrix} \checkmark \\ \text{matrix of data points} \\ \underbrace{\phantom{...}}_{n \text{ features}} \\ (m \times n) \end{bmatrix} \times \begin{bmatrix} \checkmark \\ \text{evec}_3 & \text{evec}_1 \\ | & | \\ \text{evec}_1 & \text{evec}_3 \end{bmatrix}_{(n \times k)} = \begin{bmatrix} \checkmark \\ \text{final points in } 2D \\ (m \times k) \end{bmatrix}$$

*new axis: max variance*