

# DBSCAN & GMM

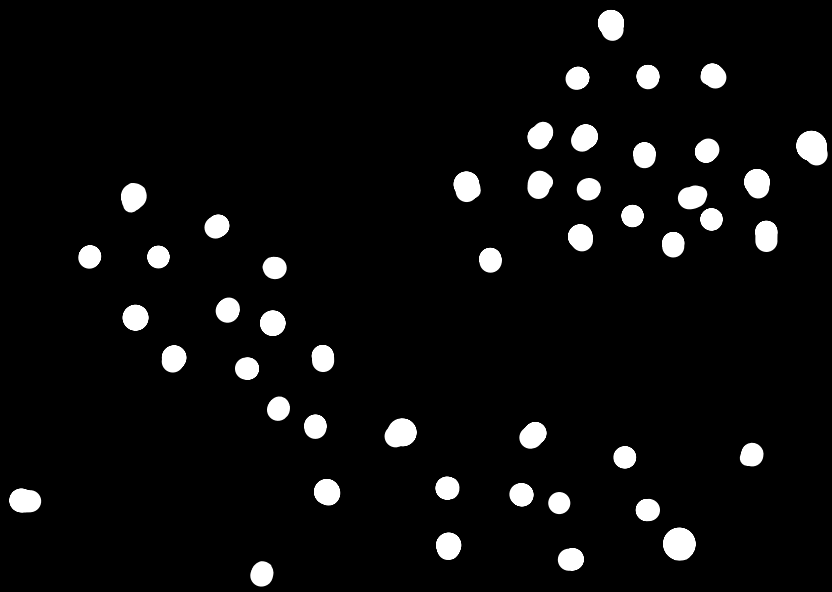
## [Clustering]

→ 2 new ideas !!

→ Comparison of algorithms

# DBSCAN → 3rd big idea

Density based spatial clustering application  
with noise.



Q: How many clusters do you see?

a) 1

b) 2

c) 3

d) 4

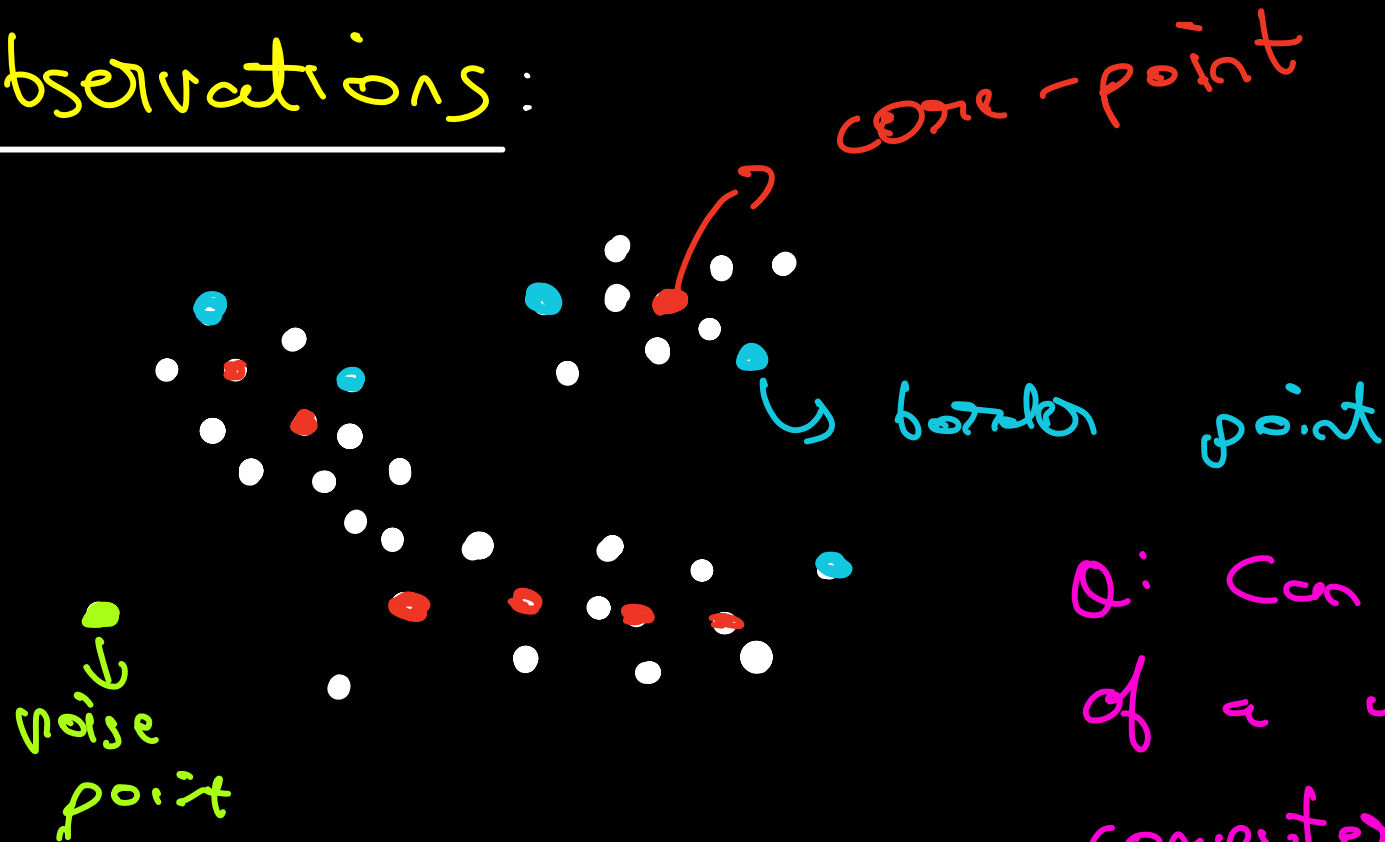
Q: Do you think KMeans will work?  
↳ No

Idea:

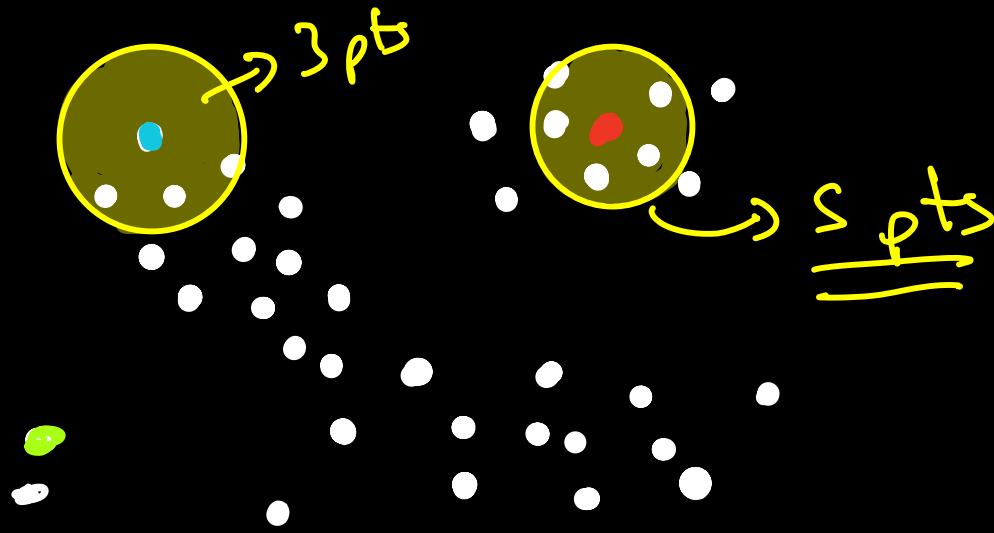
→ if a point is surrounded by

many other points its in the cluster!!

## Observations:



Q: Can you think of a way for a computer to find these??



- Draw a circle of radius  $\epsilon$
- Count # pts in circle
- if #pts  $>$  min pts
  - ↳ core pt
  - else
    - ↳ non-core pt.

→ if any non-core pt is inside circle of any core pt, then → border pt

else

→ noise pt

→ animation

→ categorise each pt into

→ core

→ border

→ noise



Join them based on  
neighbours, don't join  
2 separate borders,

Pros:

→ works with arbitrary shapes

→ No need to decide 'K'

Cons:

→ Does not work well with sparse  
points (high dim)

→ needs entire data set for  
inference

→ code

→ Time complexity :  $O(n^2)$   
↓  
need to calc distances  
of all points w.r.t all

Q. Which of the following algos  
can be used to detect outliers

a) Kmeans

b) Hierarchical

c) DBSCAN

d) Clustering algos can't be used for  
outliers.

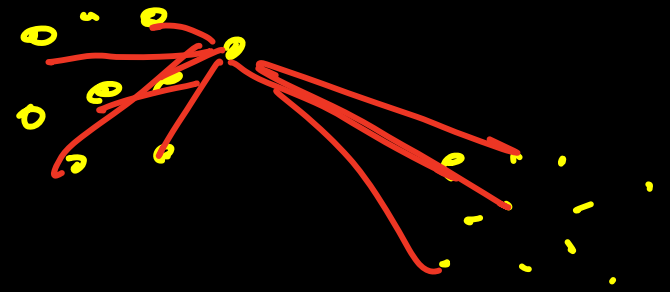
## Deciding epsilon

One way to estimate for far away clusters is:

→ calc distances between each point

→ plot a histogram of those distances

→ You may get 2 peaks, eps in b/w these peaks.



eps -> starting pt

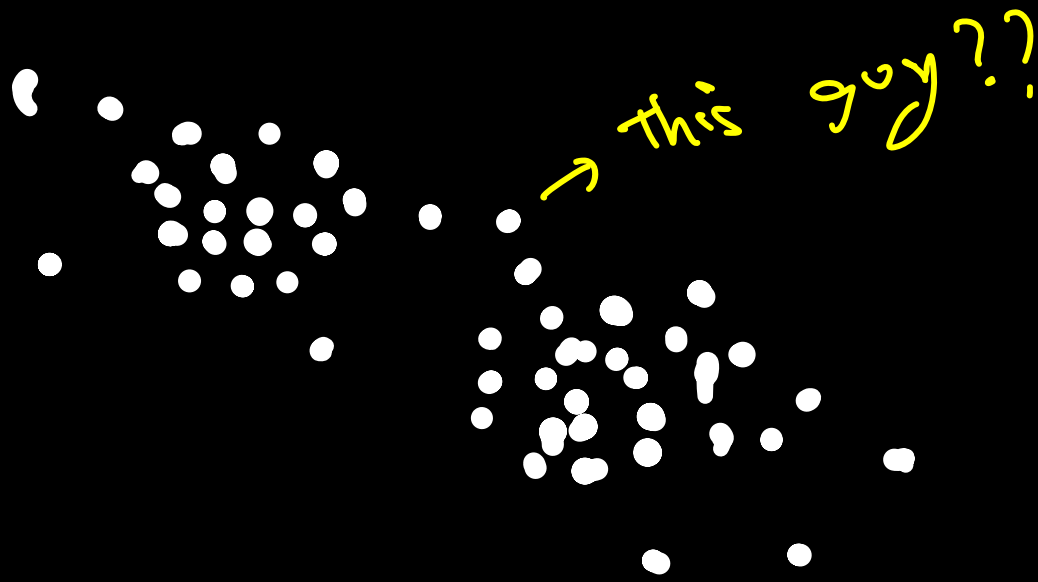
# Gaussian Mixture Models

Soft Clustering  $\rightarrow$  4<sup>th</sup> big idea!

Problem: With classification<sup>n</sup> algos I could get probabilities. How do I get probability of a pt belonging to a cluster?

Q: Any ideas??





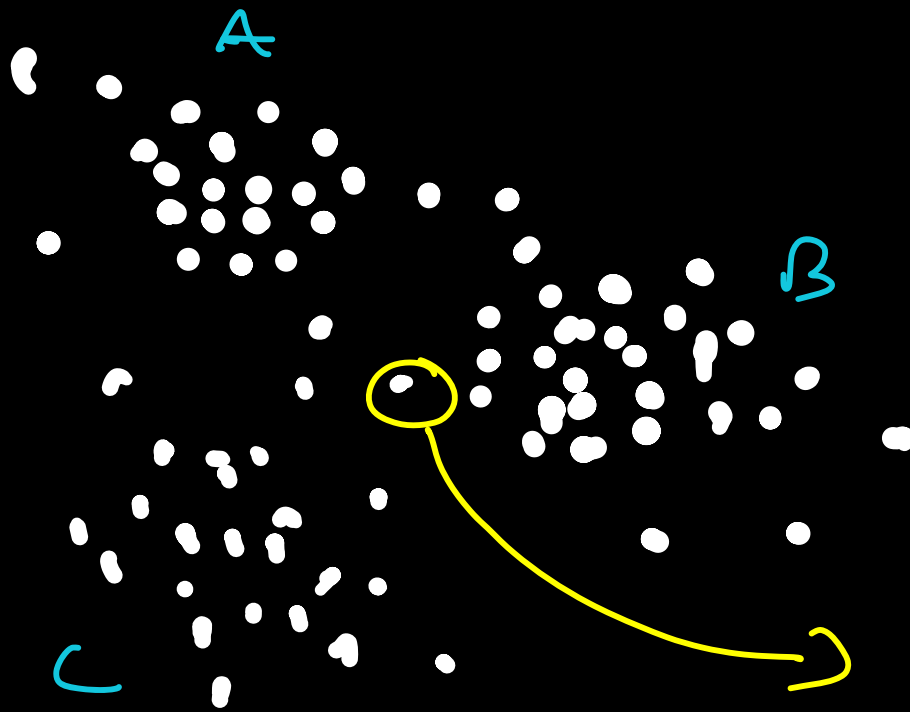
Q! How many clusters do you see?

→ In business we can make multiple policies:

→ Eg: A Rich / Premium — More ads

B Medium — discounts + ads

C Discount lovers — More discounts



closest to B

then to C

then to A

Q: So what % of ads and discount  
do I give to this guy?

$x_i \rightarrow$  50% B

20% A

30% C

$$= 0.5(Dis + Ads) + 0.2(Ads) + 0.3(Dis)$$

$$= 0.8(Dis) + 0.7(Ads)$$

$$= \frac{0.8}{0.7+0.8} (Discount) + \frac{0.7}{0.7+0.8} (Ads)$$

$$= 53\% \text{ Discount}$$

$$= 47\% \text{ Ads}$$

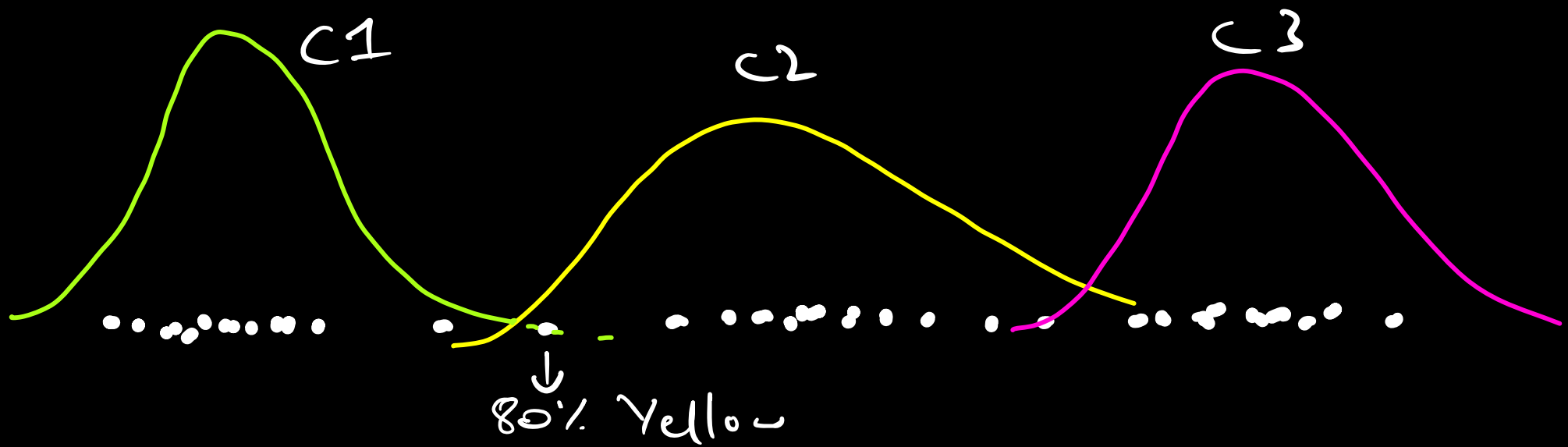
Whatever budget we have for customer, for this guy we should spend 53% on Discounts, 47% on Ads.

Idea: Use  $n$ -d gaussian dist to express clusters!!

Lets discuss this in 1-d first



Q: How many clusters?



19%. Green

1%. Pink

Q: What do you need for gaussian?

→  $\mu, \sigma$

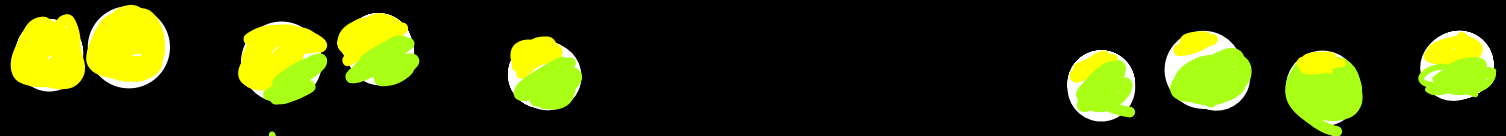
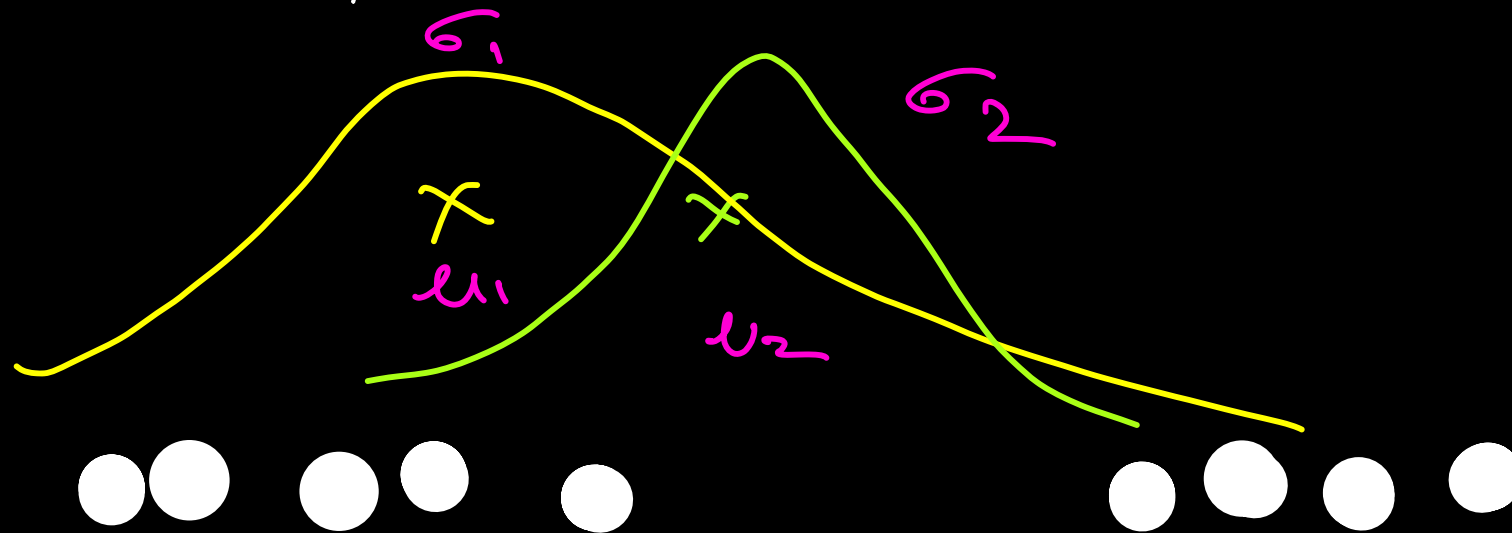
↳ I want 3 clusters:

$\mu_1$	$\mu_2$	$\mu_3$
$\sigma_1$	$\sigma_2$	$\sigma_3$

Algorithm:

Very similar to K-means

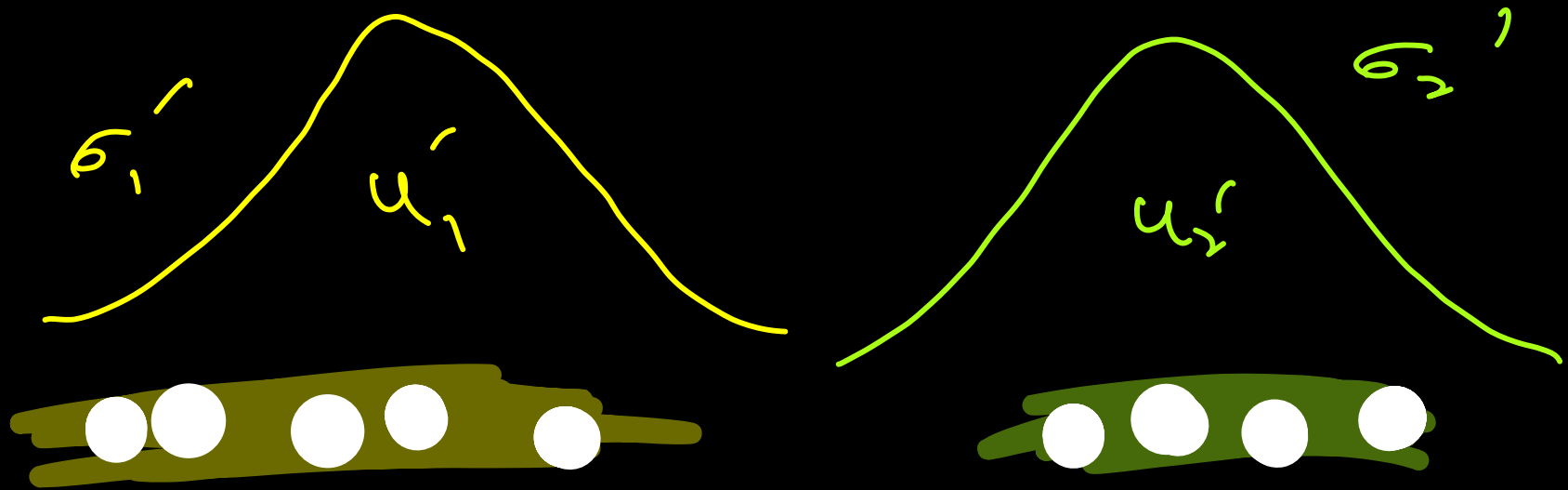
→ Random  $\mu, \sigma$  initialise



$$\mu_2' = \sum p_i(\text{G}) \cdot x_i$$

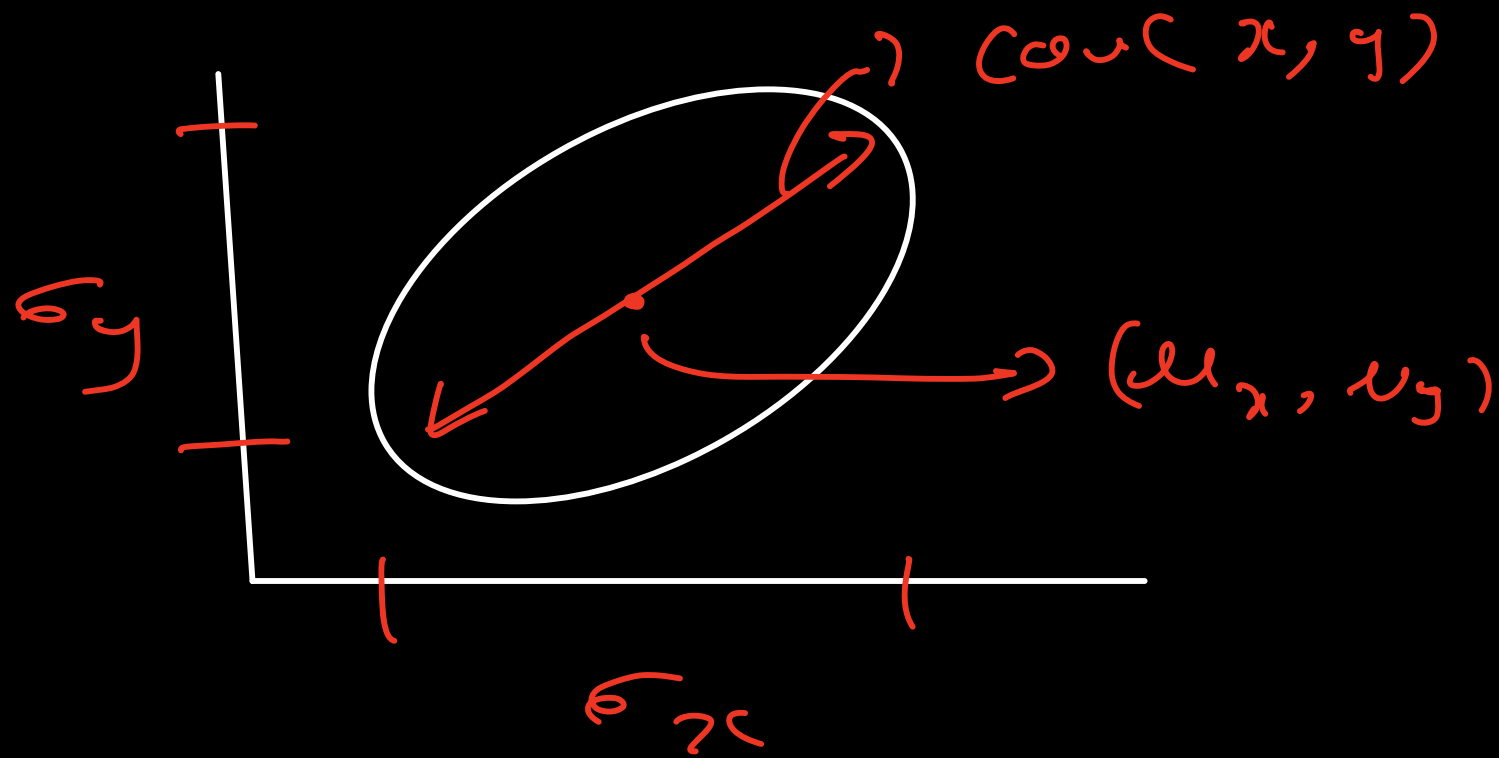
$$\mu_1' = \sum p_i(\text{Y}) \cdot x_i$$

Similarly we will calculate variance



After multiple updates, you will have tightly fitting gaussians.

2D Gaussians!



# params to update = 5

↓  
Same algo

$\mu_x, \mu_y, \sigma_x, \sigma_y, \text{Cov}$

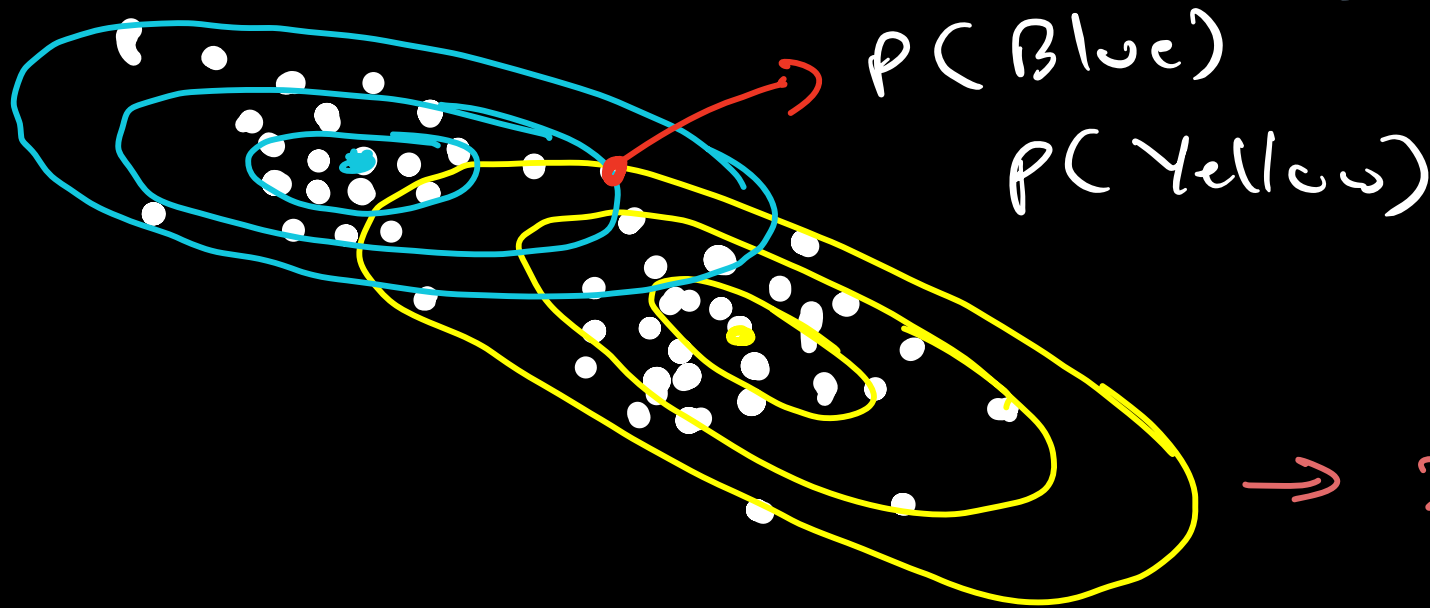


How to get these probabilities?



Gaussian Distributions

→ Modeling choice,  
you could create  
a variation with  
another dist.



$P(\text{Blue})$

$P(\text{Yellow})$

→ 2D Gaussian

→ animations

→ code

Pros and cons are similar to KMeans

Results are also very similar 1