

Lecture 11:

CLUSTERING

as an
example of an
inference /
probabilistic
modelling task

x

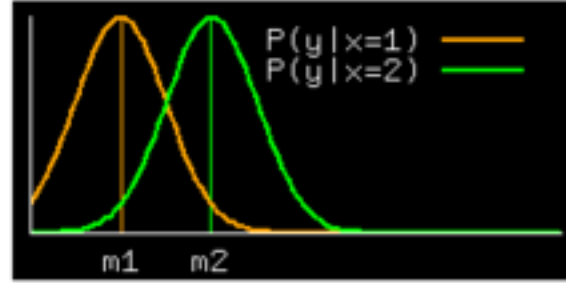
1

2

A channel has binary input $x \in \{1, 2\}$ and real output y .

$$P(x = 1) = \pi_1$$

$$P(y | x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y - m_x)^2}{2\sigma^2}}$$

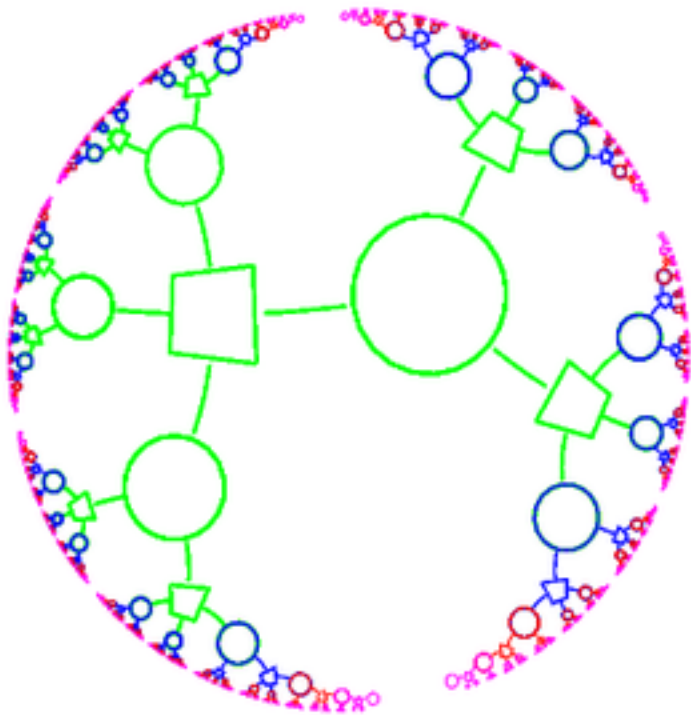


You see y . What is x ?

Write $P(x = 1 | y)$ in the form

$$P(x = 1 | y) = \frac{1}{1 + \exp(-\text{something elegant})}$$

Information theory, pattern recognition, and neural networks



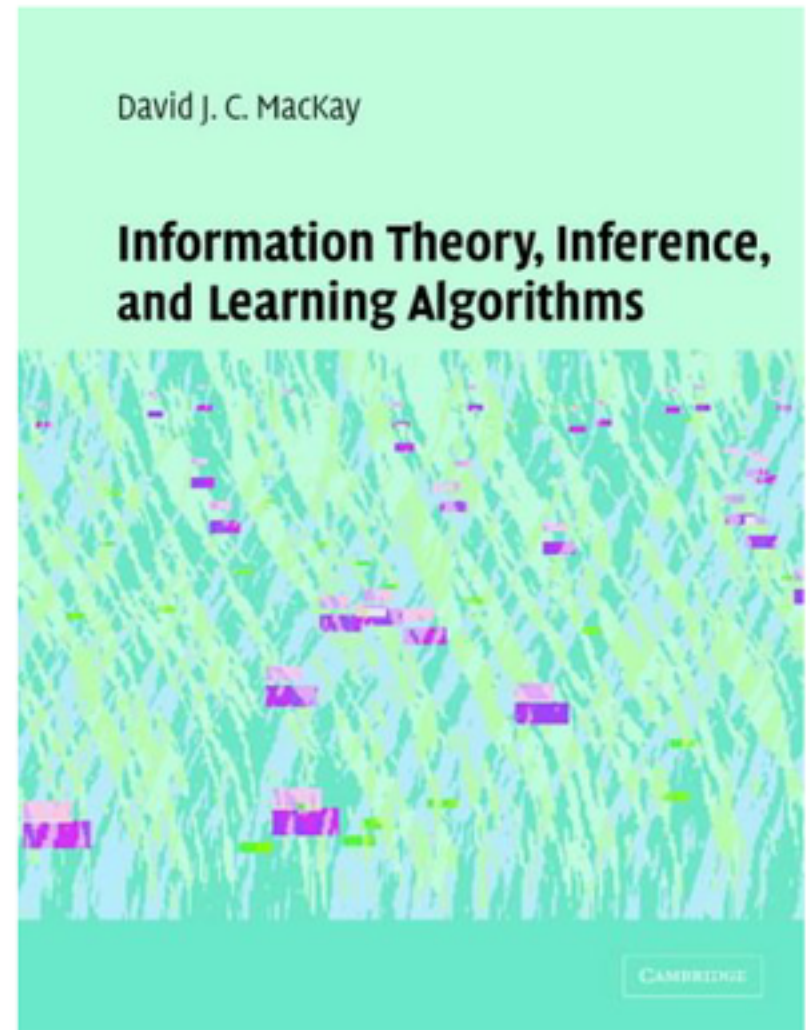
- 1 Noisy-channel coding
- Source coding (Data compression)
 - 2 Information content, entropy
 - 3 Typicality and the source coding theorem
 - 4 Symbol codes
 - 5 Symbol codes and Arithmetic coding
- Noisy-channel coding
 - 6 Information measures for noisy channels
 - 7 Capacity of a noisy channel
 - 8 The noisy-channel coding theorem
- Inference + probabilistic methods
 - 9-10 Inference
 - 11 Clustering
 - 12-13 Monte Carlo methods
 - 14 Variational methods

www.inference.phy.cam.ac.uk/itprnn/

www.inference.phy.cam.ac.uk/itila/

The course

www.inference.phy.cam.ac.uk/itprnn/



The book

www.inference.phy.cam.ac.uk/itila/

Monte Carlo methods

Simple Monte Carlo methods

- Importance sampling
- Rejection sampling

Markov-chain Monte Carlo methods

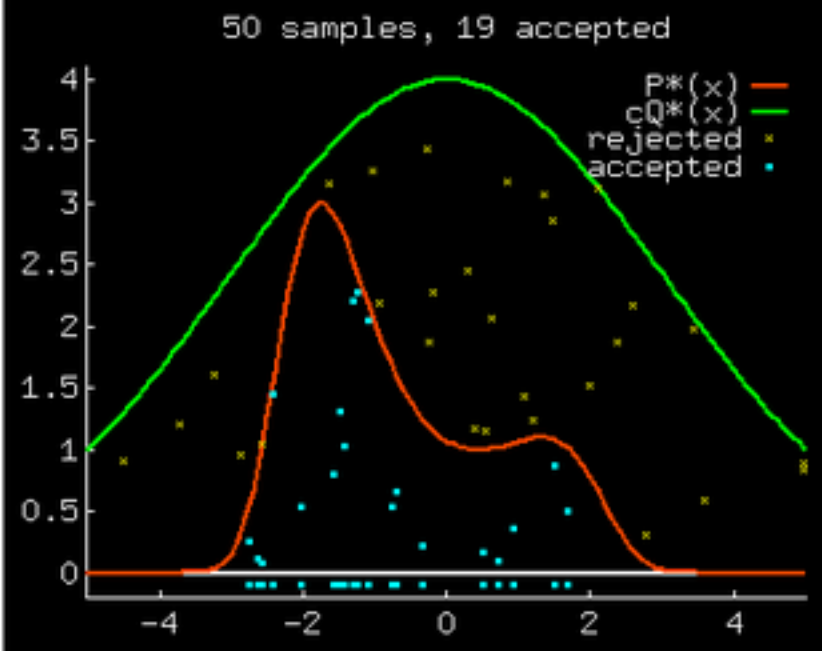
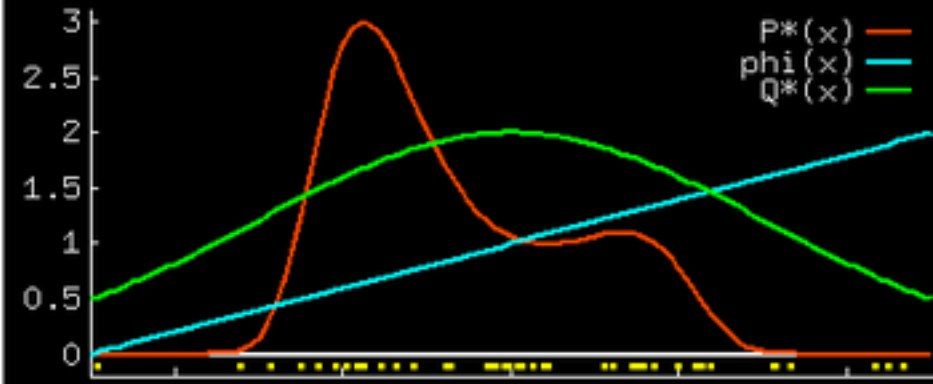
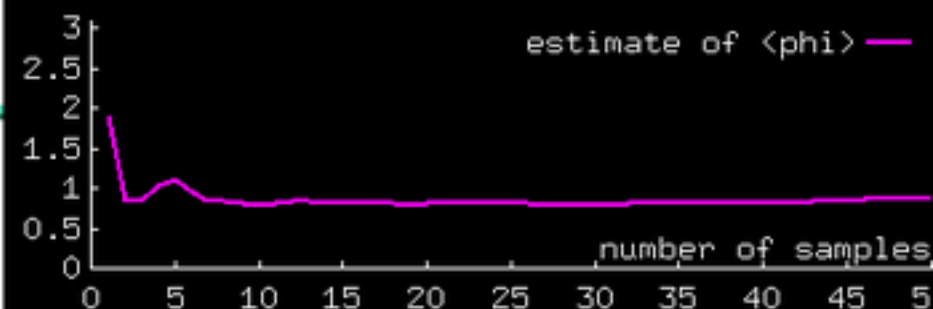
- Metropolis method
- Gibbs sampling
- Slice sampling

Reducing random-walk behaviour

- Hamiltonian Monte Carlo
- Overrelaxation

Exact sampling

[itp/exact/rc](#) RUNME



Lecture 11

Clustering

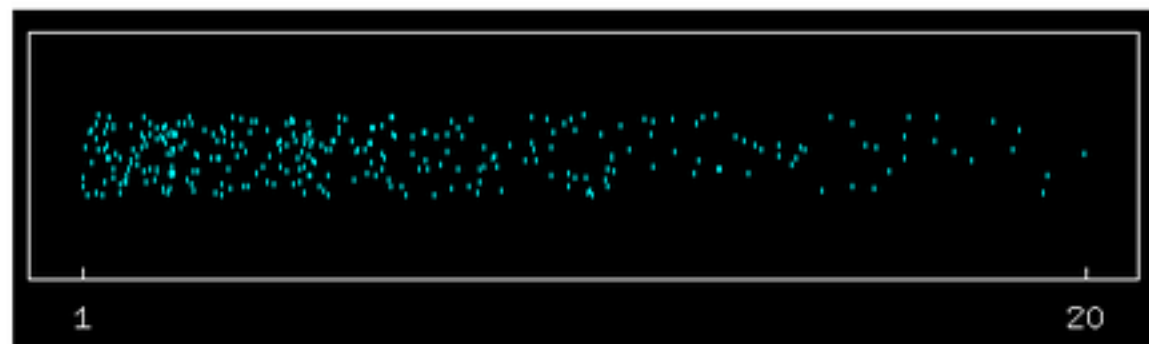
Lecture 12

Monte Carlo methods

Lecture 13

Advanced Monte Carlo methods

Inferring parameters in science experiments

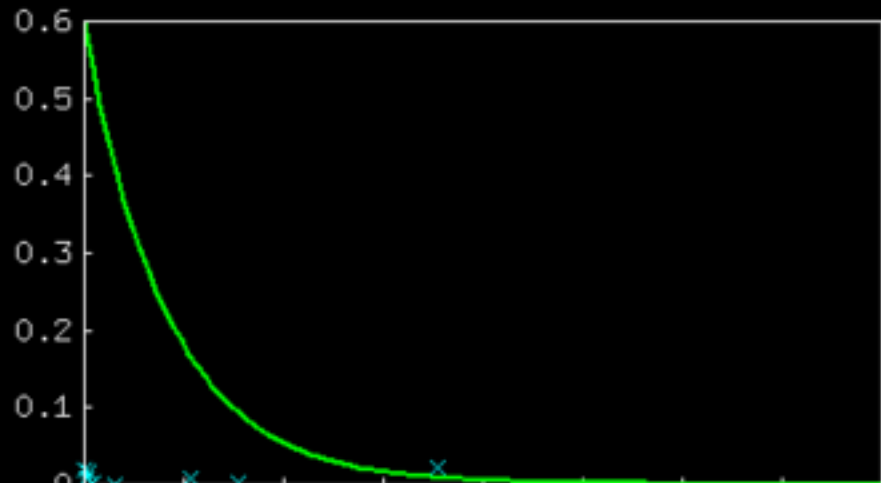


A source emits particles; they decay at locations exponentially distributed with lengthscale λ . Decays can be observed only if the location falls in a window $(x_{\min}, x_{\max}) = (a, b) = (1, 20)$.

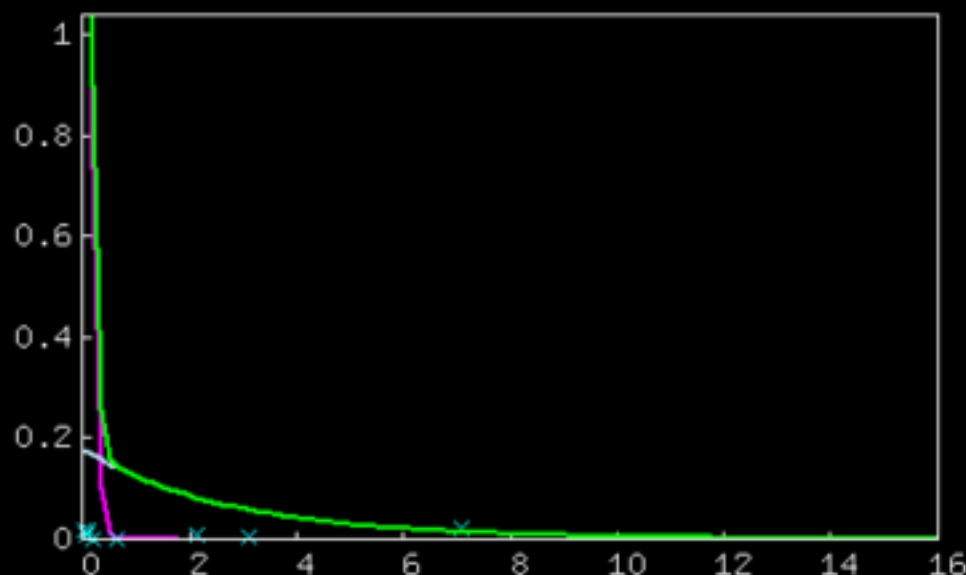
N decays are observed at locations $\{x_1, \dots, x_N\}$.
What is λ ?

+ What about the possibility that the distribution is a mixture of two exponentials?

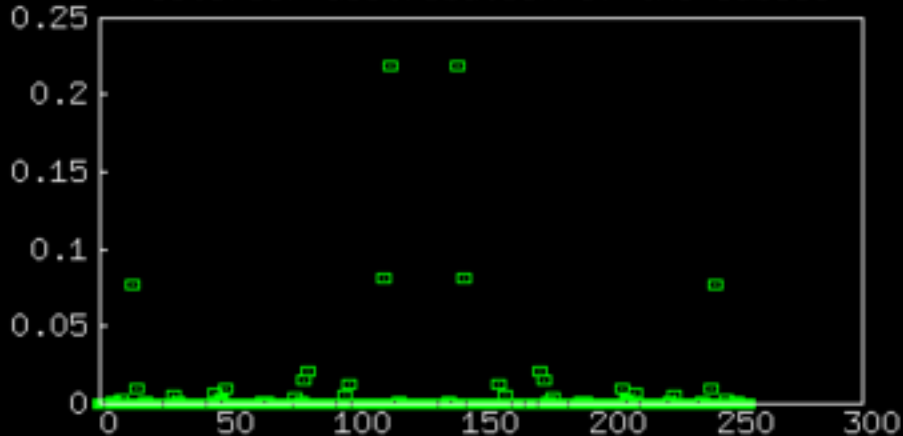
maximum likelihood distribution



ML model with two tau's

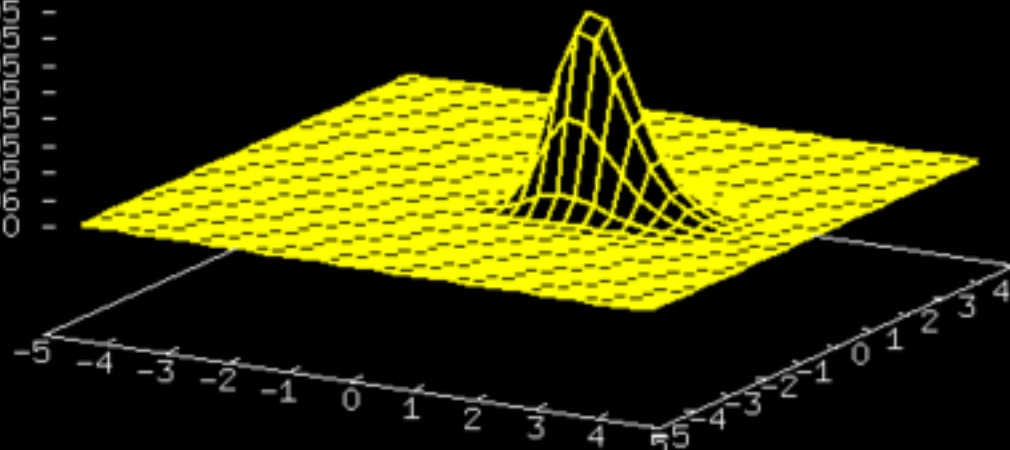


Posterior distribution of the labels



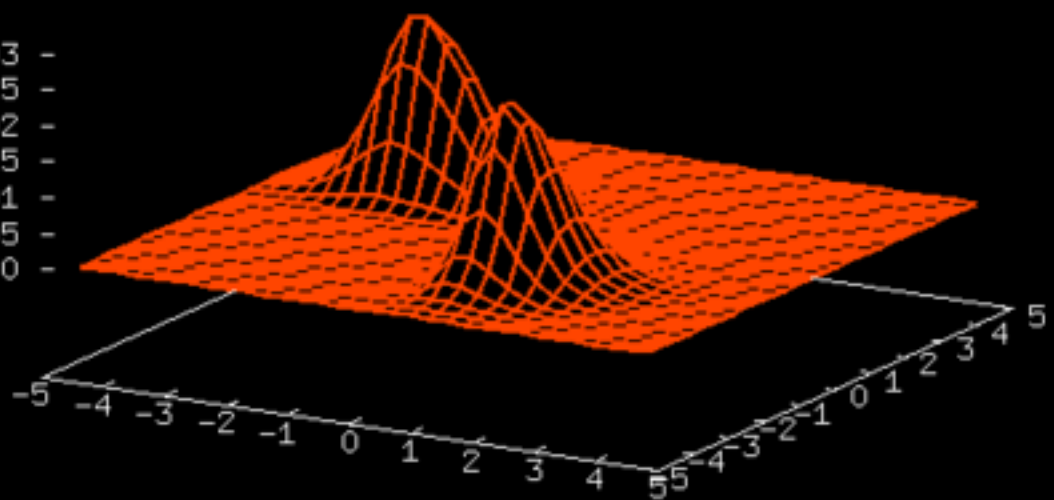
$P(\text{Data}|\tau_1, \tau_2, c_1..c_N = 00000111)$

4e-05 -
 3.5e-05 -
 3e-05 -
 2.5e-05 -
 2e-05 -
 1.5e-05 -
 1e-05 -
 5e-06 -
 0 -



Marginal likelihood of tau1 and tau2 - $P(\text{Data}|\tau_1, \tau_2)$

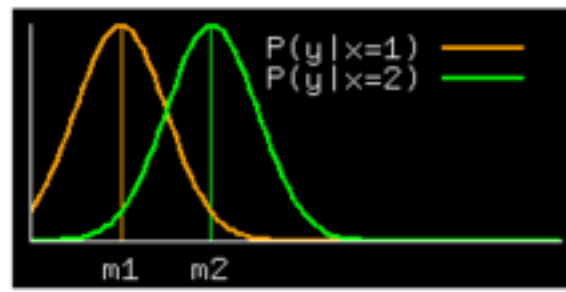
.0003 -
 .00025 -
 .0002 -
 .00015 -
 .0001 -
 5e-05 -
 0 -



A channel has binary input $x \in \{1, 2\}$ and real output y .

$$P(x = 1) = \pi_1$$

$$P(y | x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y - m_x)^2}{2\sigma^2}}$$

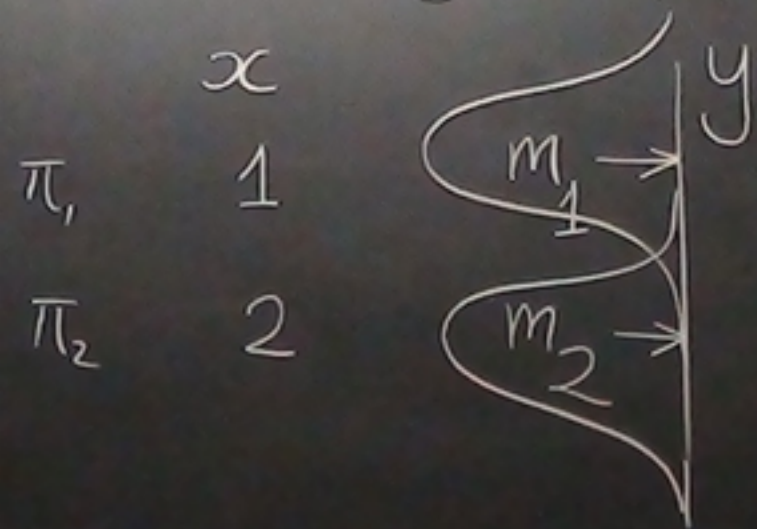


You see y . What is x ?

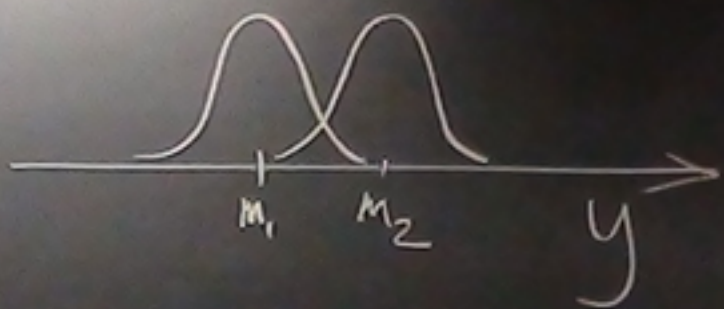
Write $P(x = 1 | y)$ in the form

$$P(x = 1 | y) = \frac{1}{1 + \exp(-\text{something elegant})}$$

The Gaussian channel (with two inputs)



$$P(x=1 | y) =$$

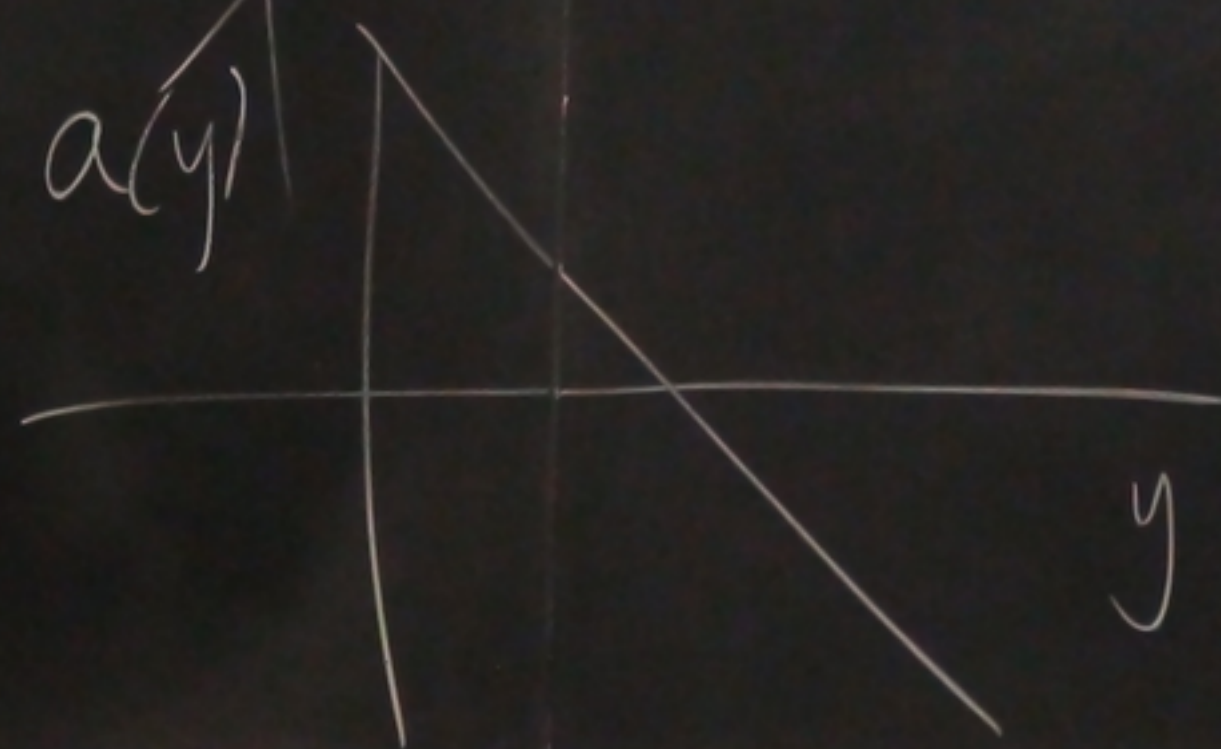
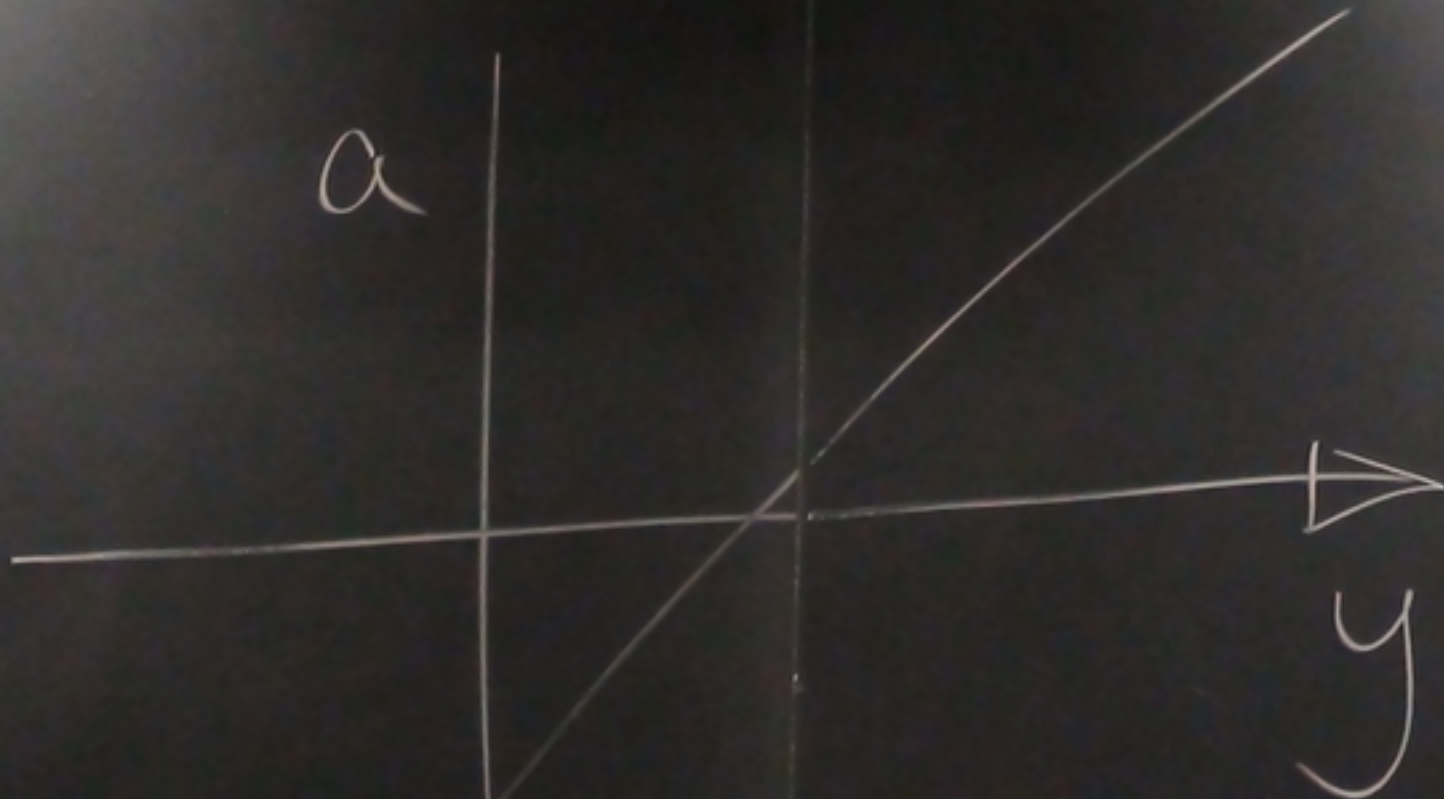


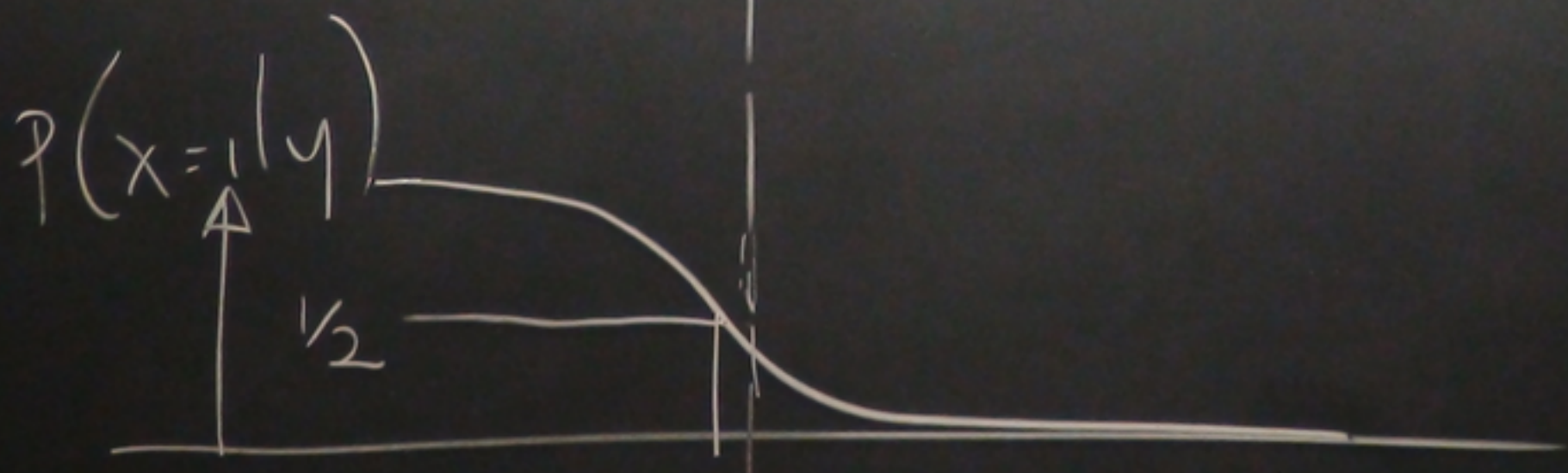
$$P(y|x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y-m_x)^2}{2\sigma^2}}$$

$$= \frac{P(y|x=1) \pi_1}{\sum_{k=1}^2 P(y|x=k) \pi_k} = \frac{1}{1 + e^{-a(y)}}$$

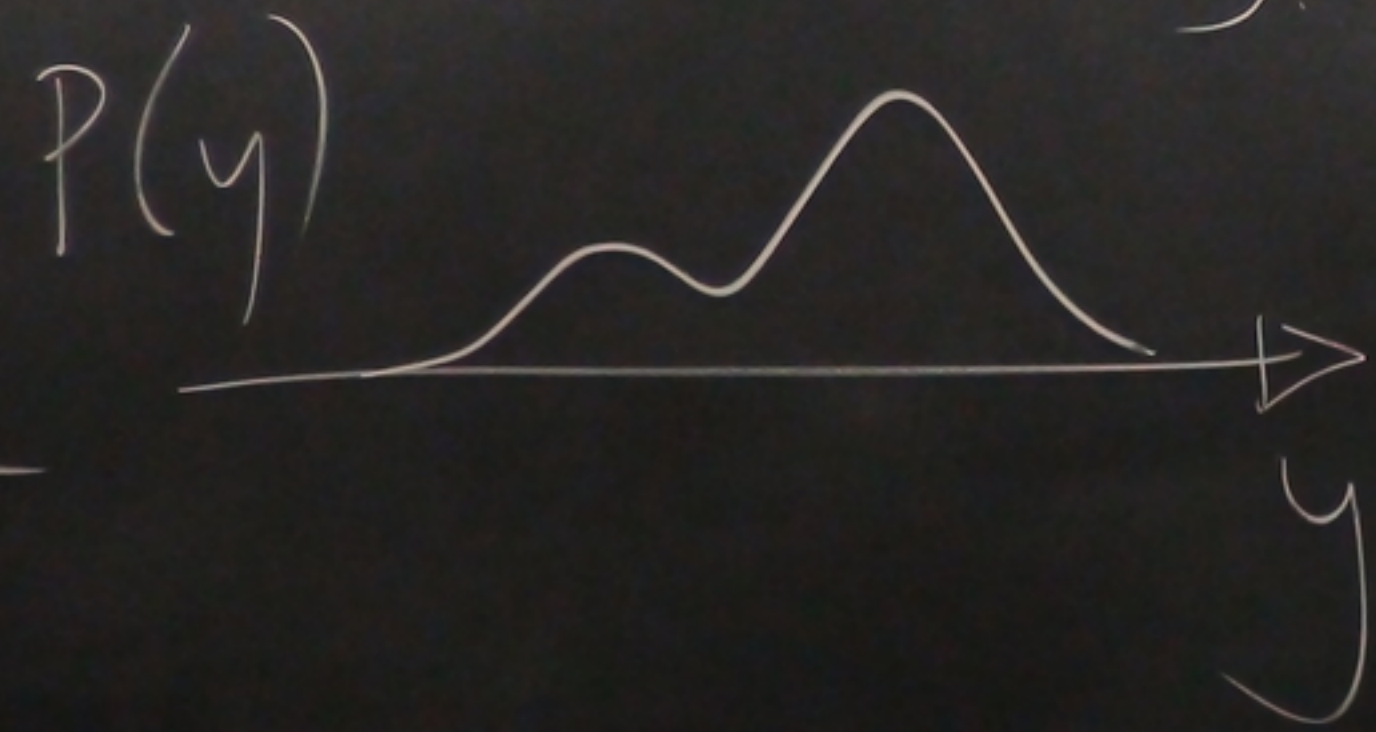
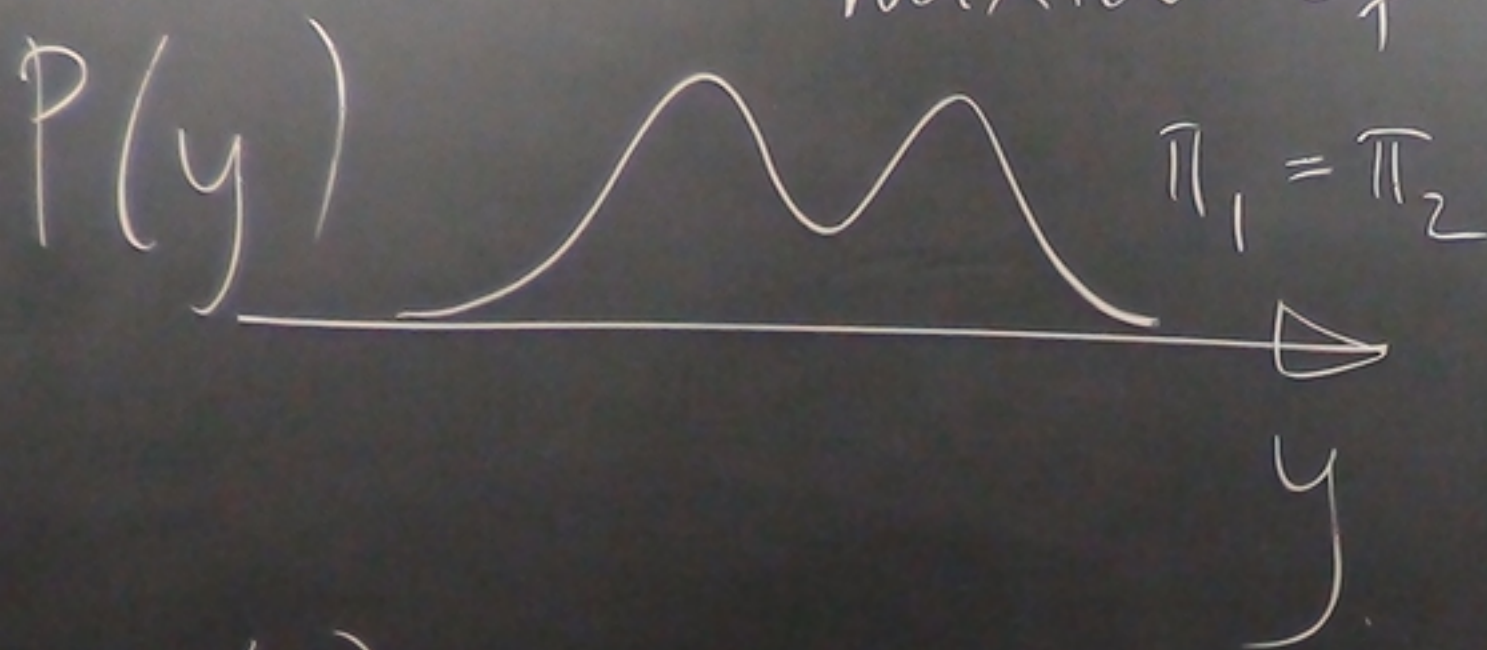
$$a(y) = \frac{(y - m_2)^2}{2\sigma^2} - \frac{(y - m_1)^2}{2\sigma^2} + \log \frac{\pi_1}{\pi_2}$$

$$= \frac{1}{\sigma^2} \left(y - \frac{m_1 + m_2}{2} \right) (m_1 - m_2) + \log \frac{\pi_1}{\pi_2}$$

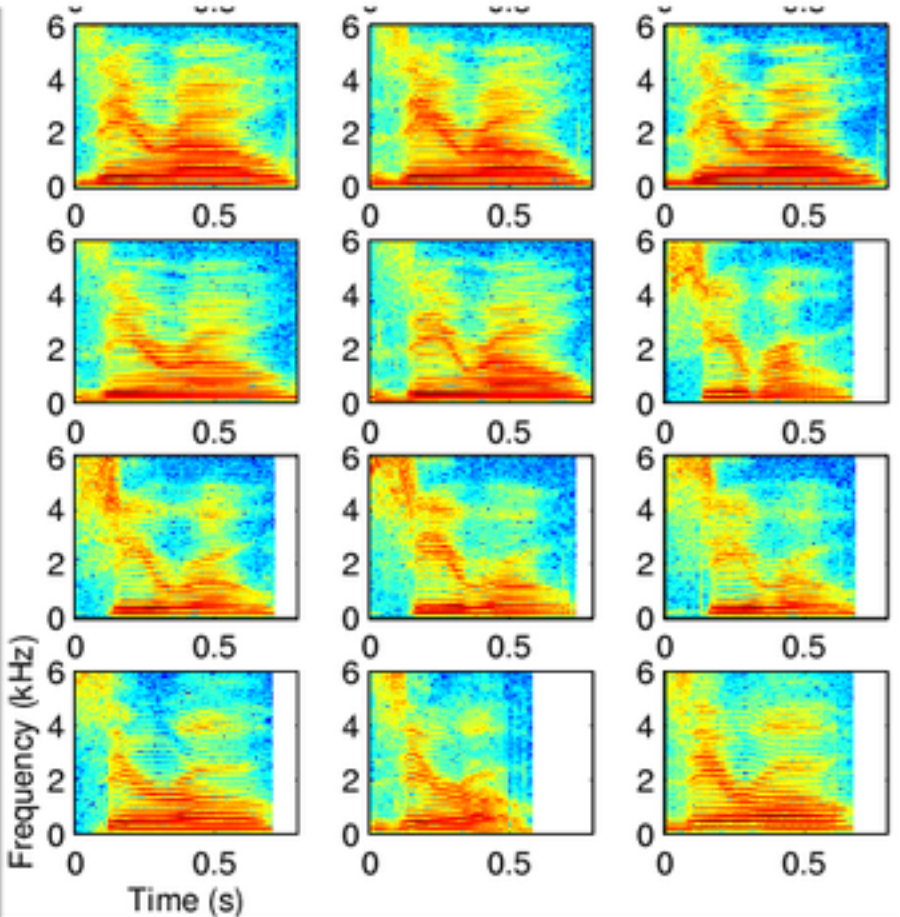
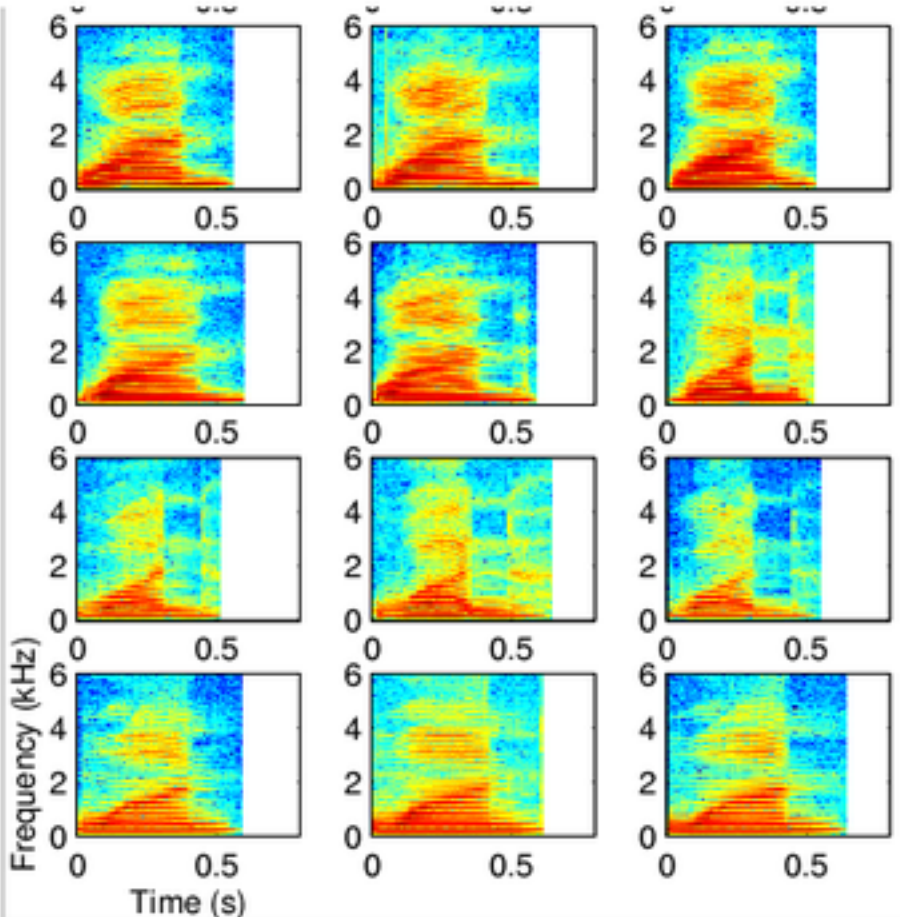




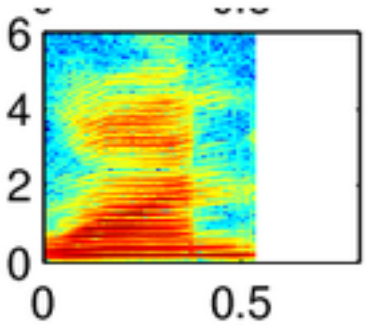
mixture of 2 gaussians



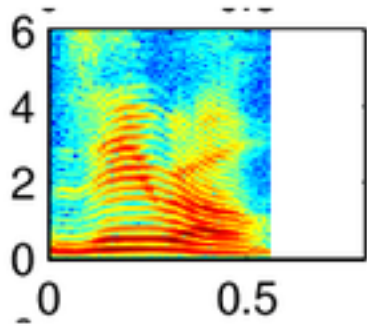
Clustering



One

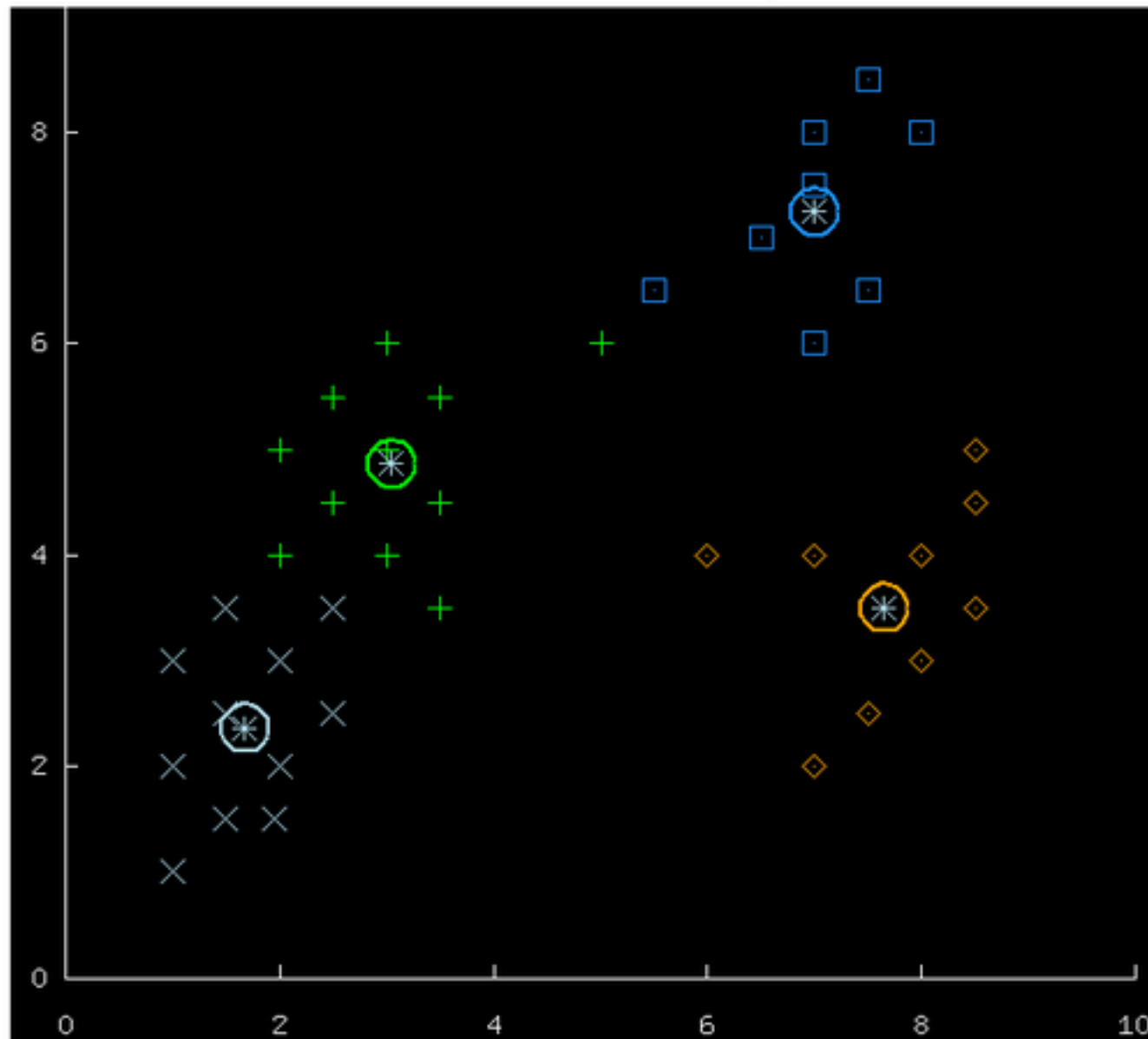


Zero



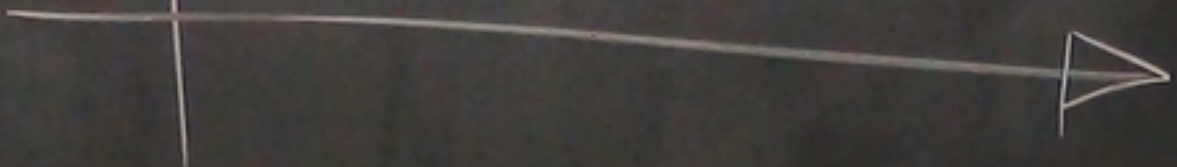
Clustering

K means clustering



ellipticity

x_2

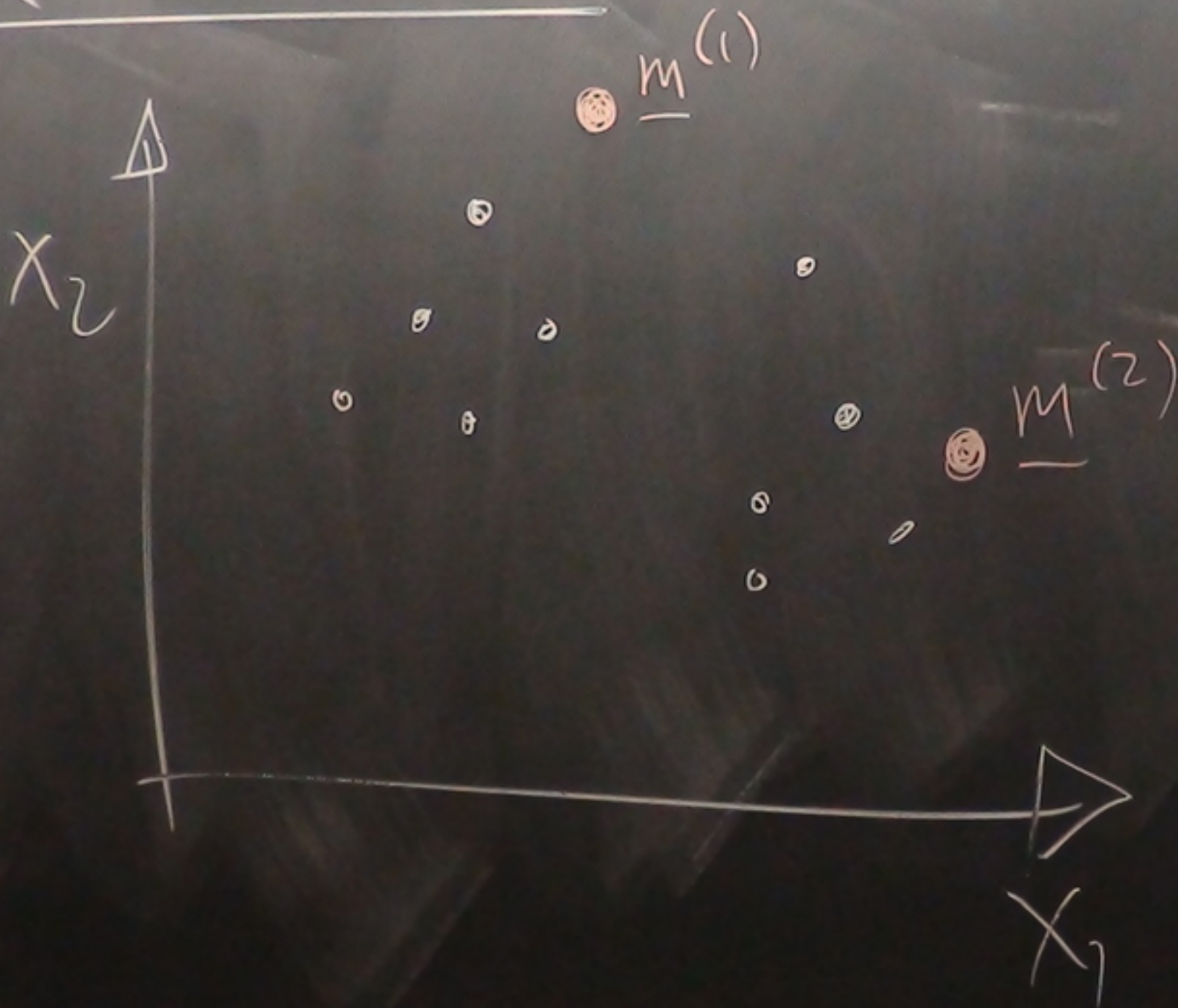


x_1

reflectance

at 600nm

K MEANS



1 Initialize K means $\left\{ \underline{m}^{(k)} \right\}_{k=1}^K$ at random.

2 Assign $\underline{x}^{(n)}$ to nearest \underline{m}

Define $d = \frac{1}{2} \left(m^{(k)} - x^{(n)} \right)^2$

$d(m^{(k)}, x^{(n)})$

K MEANS



1. Initialize

2. Assign

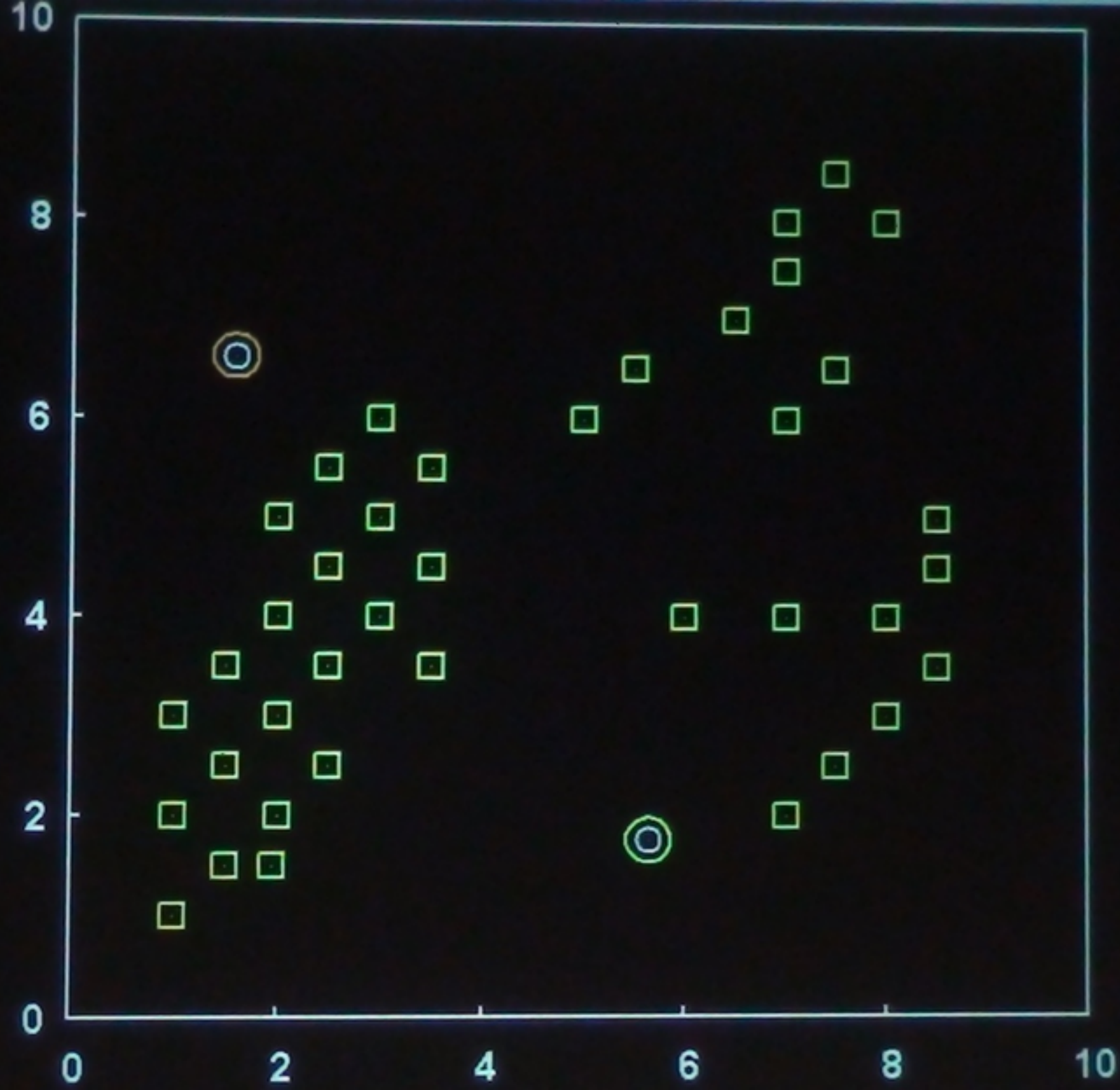
Assign $x^{(n)}$ to nearest \underline{m}

Define $d = \frac{1}{2} (m^{(k)} - x^{(n)})^2$

$d(m^{(k)}, x^{(n)})$

Responsibility $r_k^{(n)} = \begin{cases} 1 & \text{if } m^{(k)} \text{ is the closest mean to } x^{(n)} \\ 0 & \text{otherwise} \end{cases}$

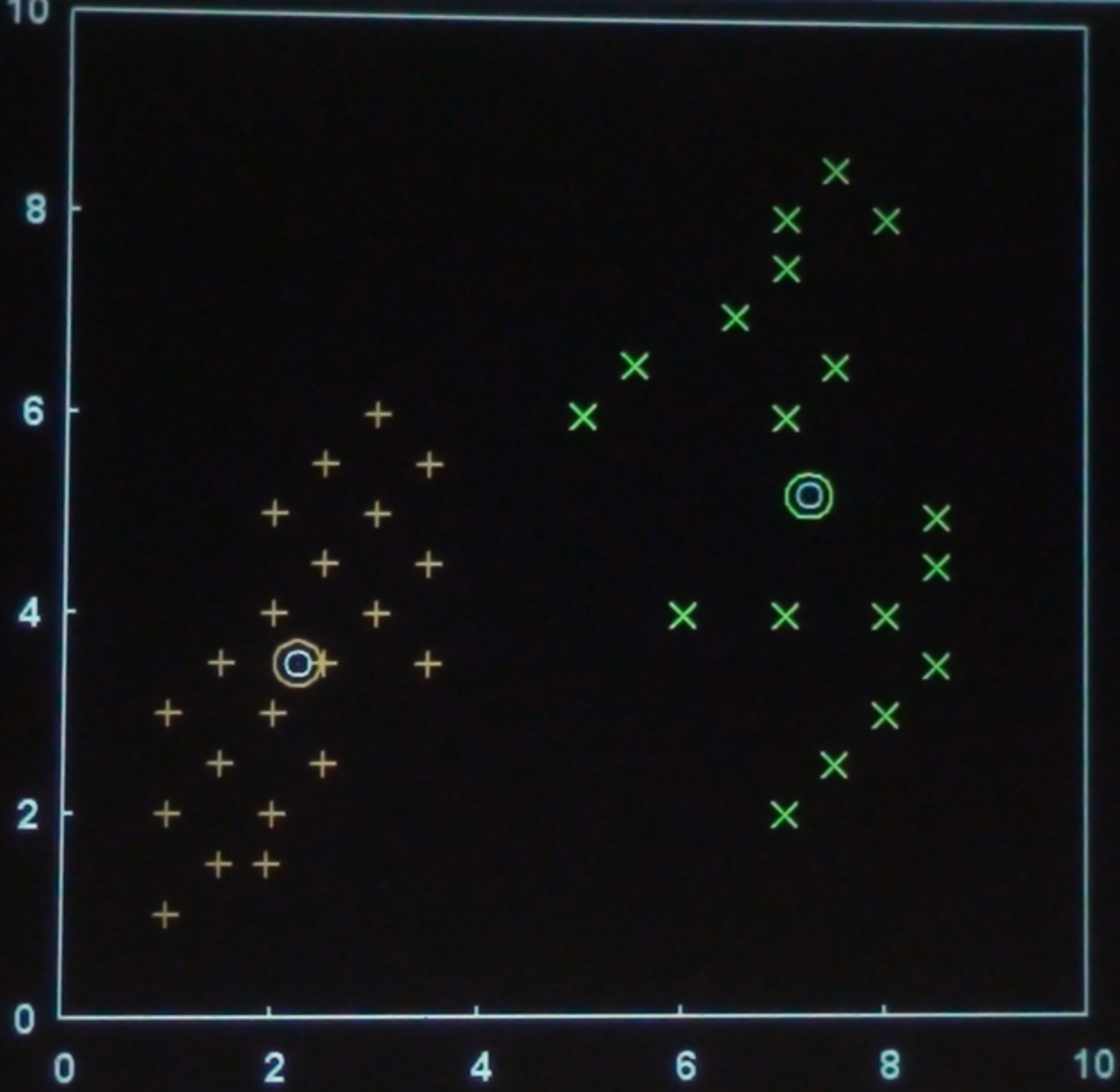
ie if k minimises $d(m^{(k)}, x^{(n)})$
wrt k'

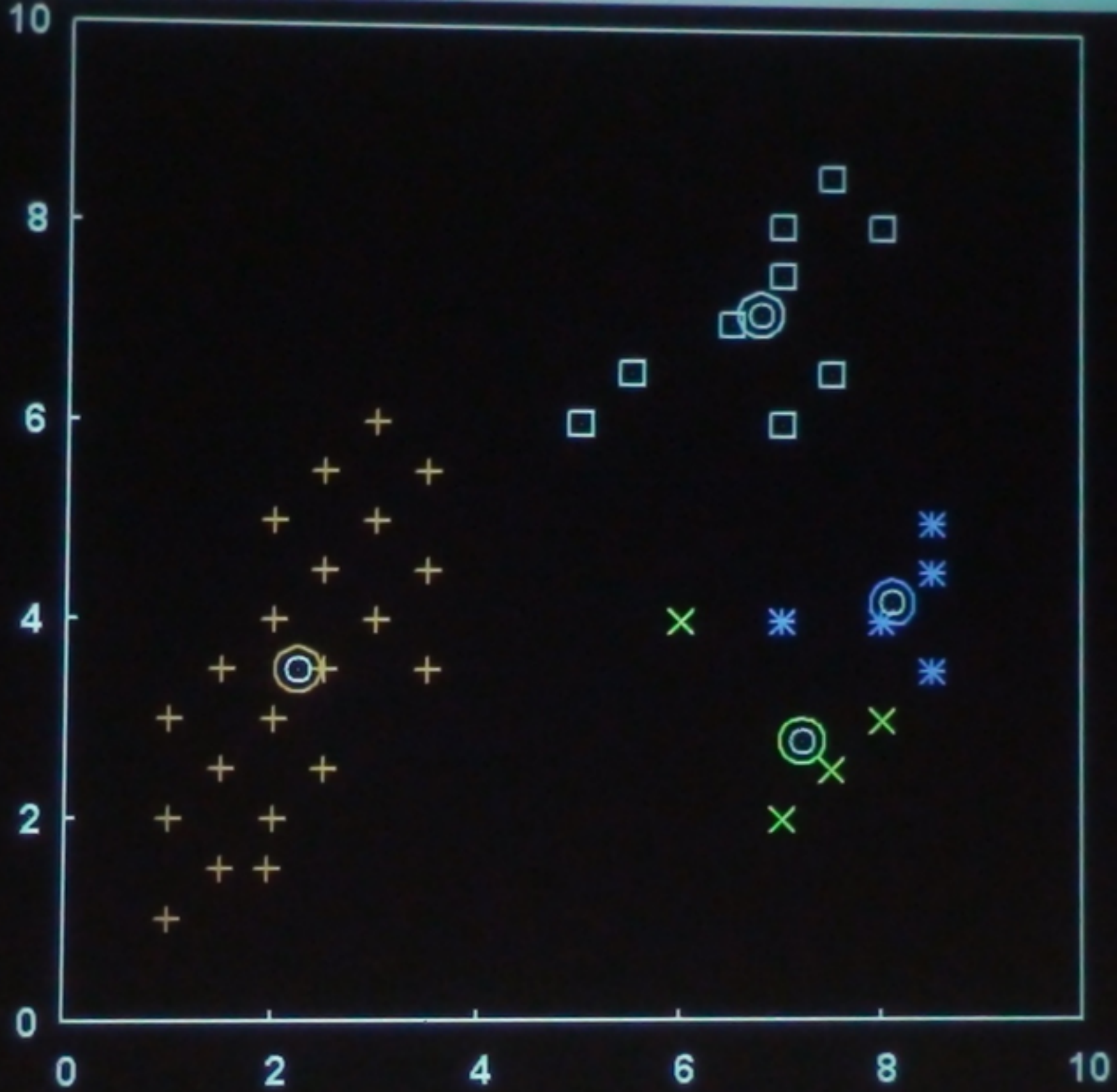


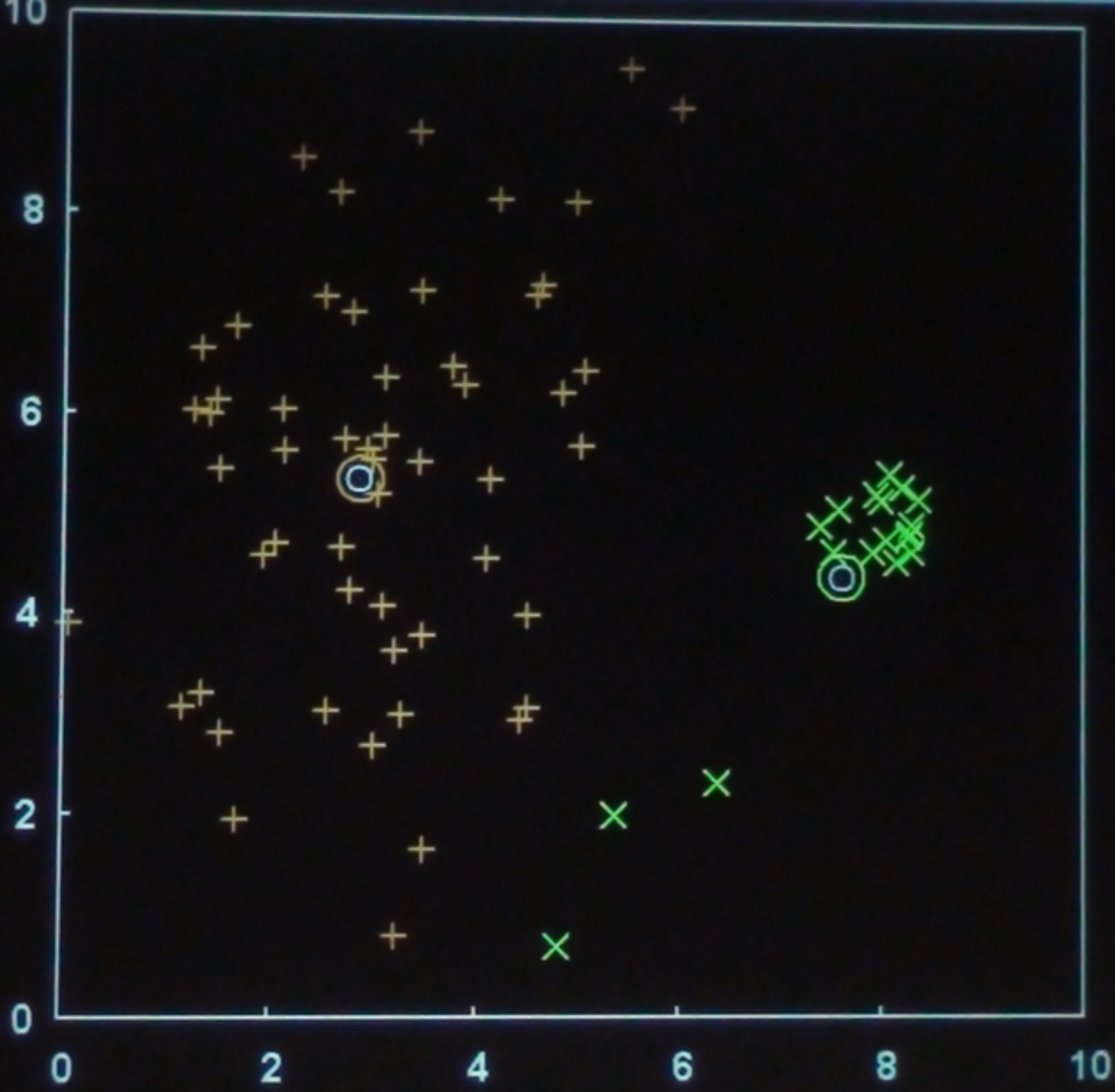
if k minimises $d(m^{(k)}, x^{(n)})$
wrt k'

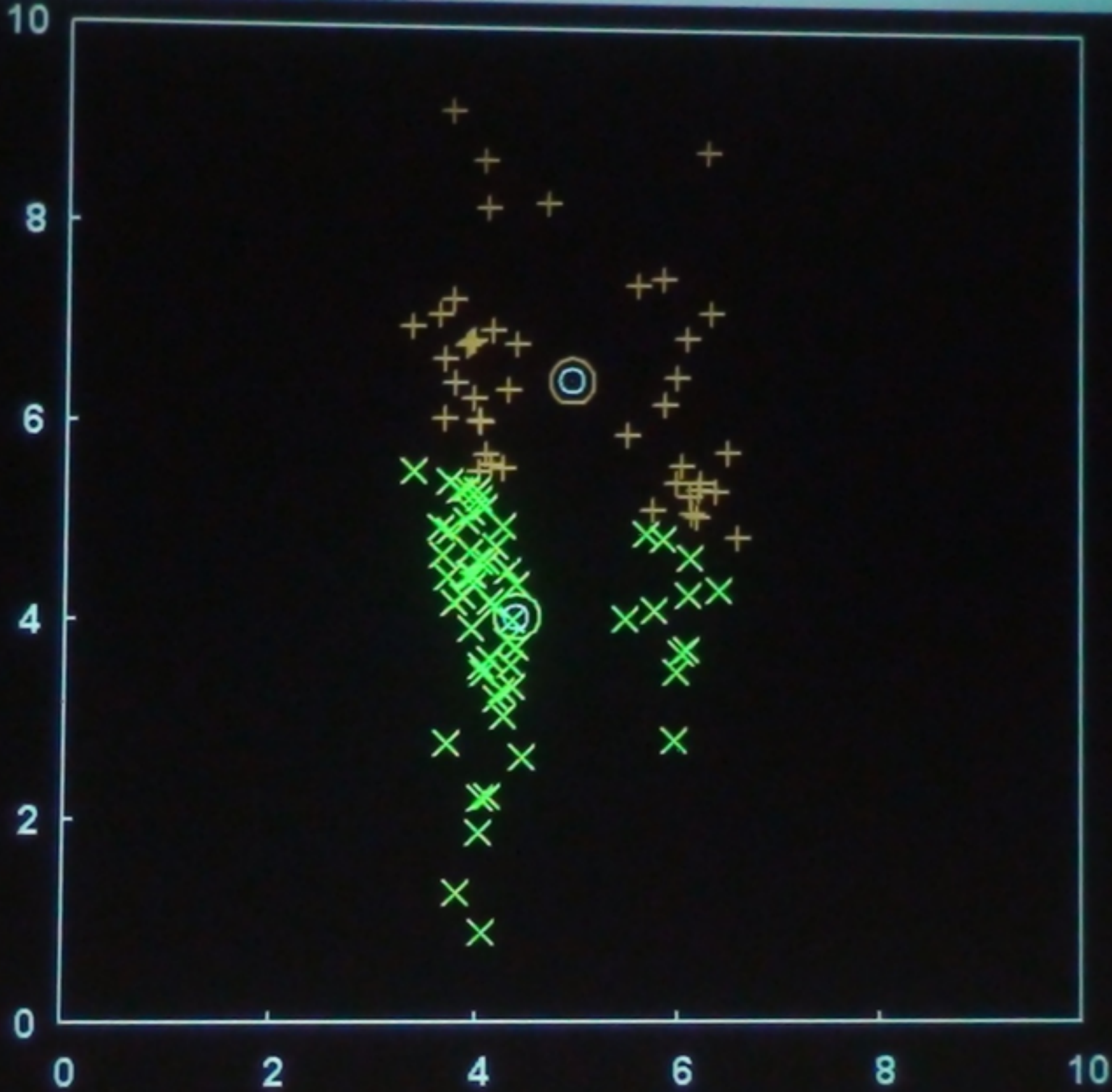
3 Update

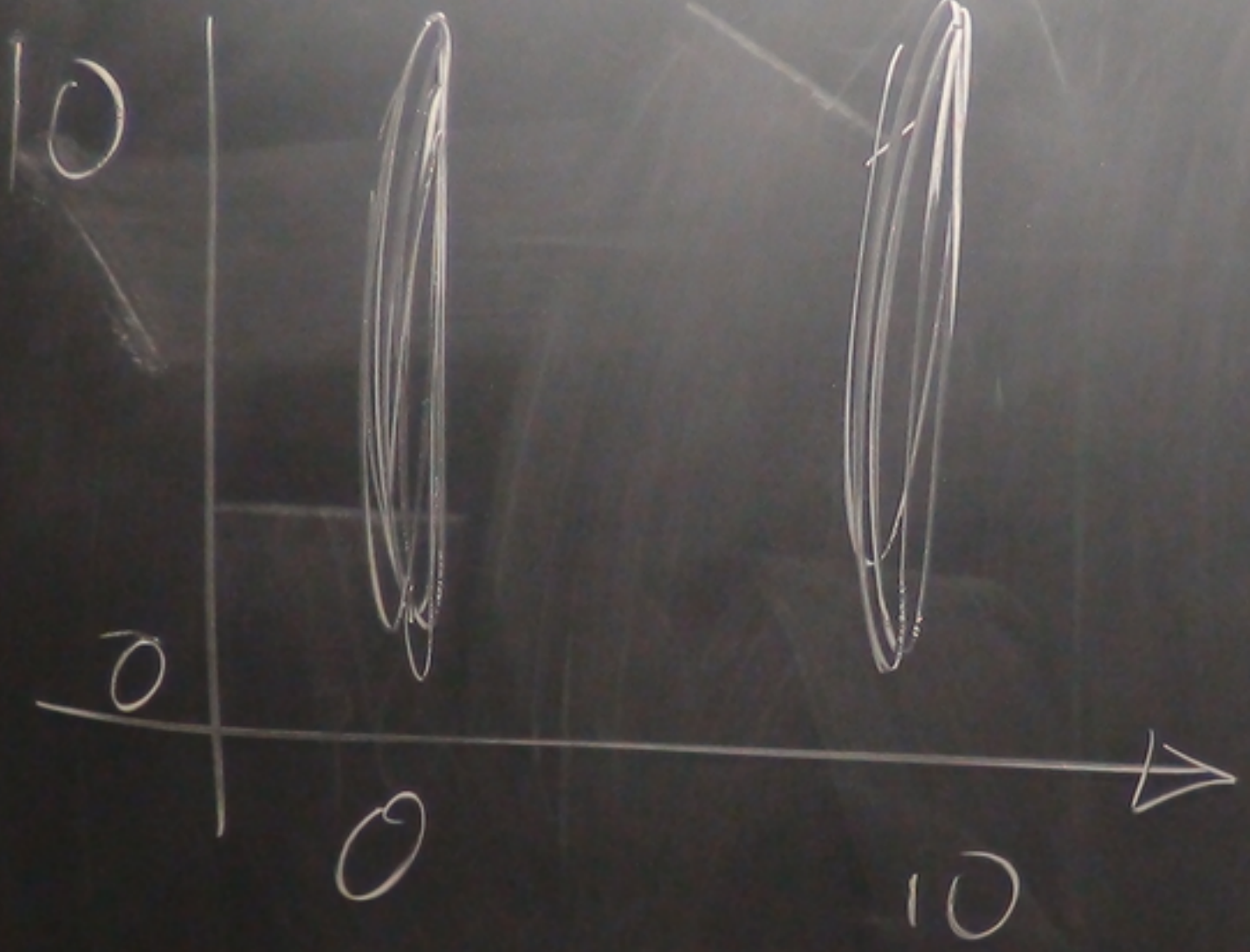
$$\underline{m}^{(k)} = \frac{\sum_{n=1}^N r_k^{(n)} x^{(n)}}{\sum_{n=1}^N r_k^{(n)}}$$



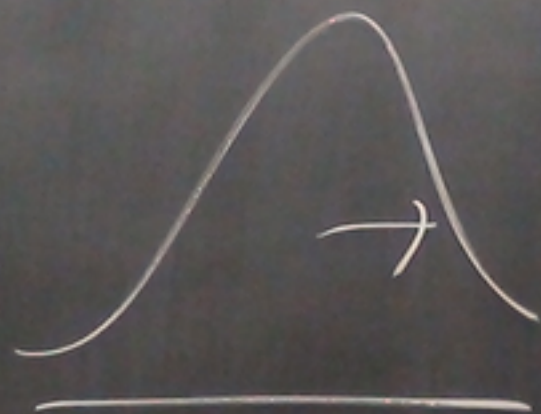




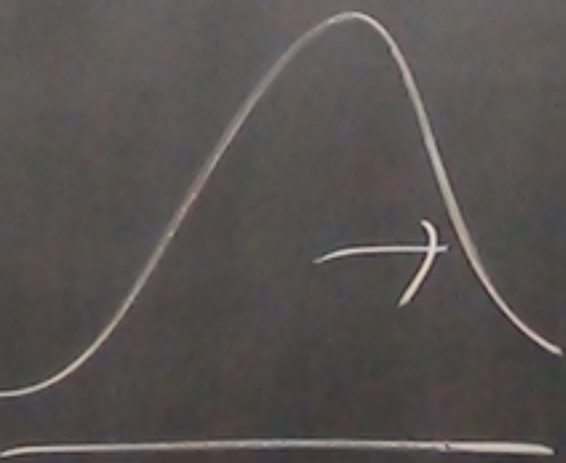




Interpretation



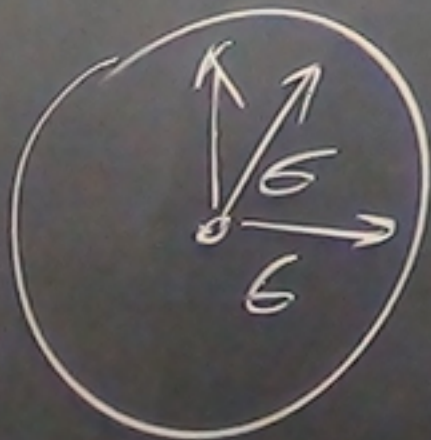
$$e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$



e

$$\frac{1}{Z} \left(\begin{matrix} \cdot \\ \cdot \\ \cdot \end{matrix} \right)$$

\ll



spherical

K Gaussians

all same σ

$$\pi_1 = \pi_2 \dots \pi_K$$

k means is
a MAP*

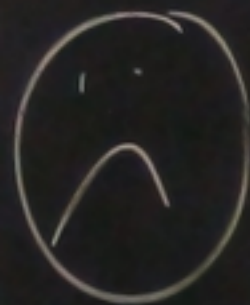


algm

that adjusts " R_n "

& $\{m_k\}$

to maximize posterior prob.



$\{m^{(k)}\}_{k=1}^K$ at random. & set a "stiffness" β

$$r_k^{(n)} = \frac{e^{-\beta d(m^{(k)}, x^{(n)})}}{\sum_{k'} e^{-\beta d(m^{(k')}, x^{(n)})}}$$

"softmax"

Update $m^{(k)}$

=

$$\frac{\sum_{n=1}^N r_k^{(n)} X^{(n)}}{\sum_{n=1}^N r_k^{(n)}}$$

Figure 1

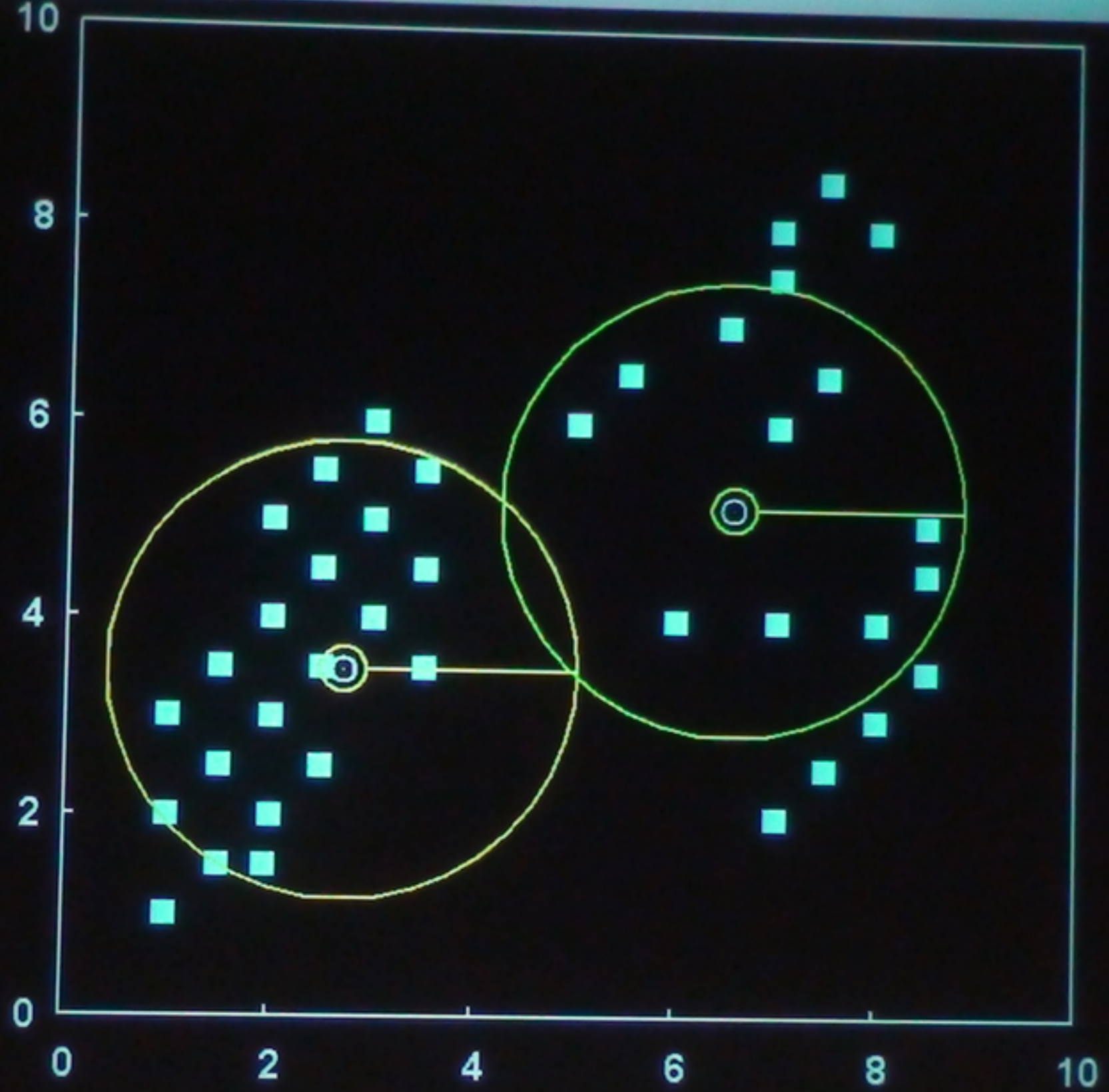
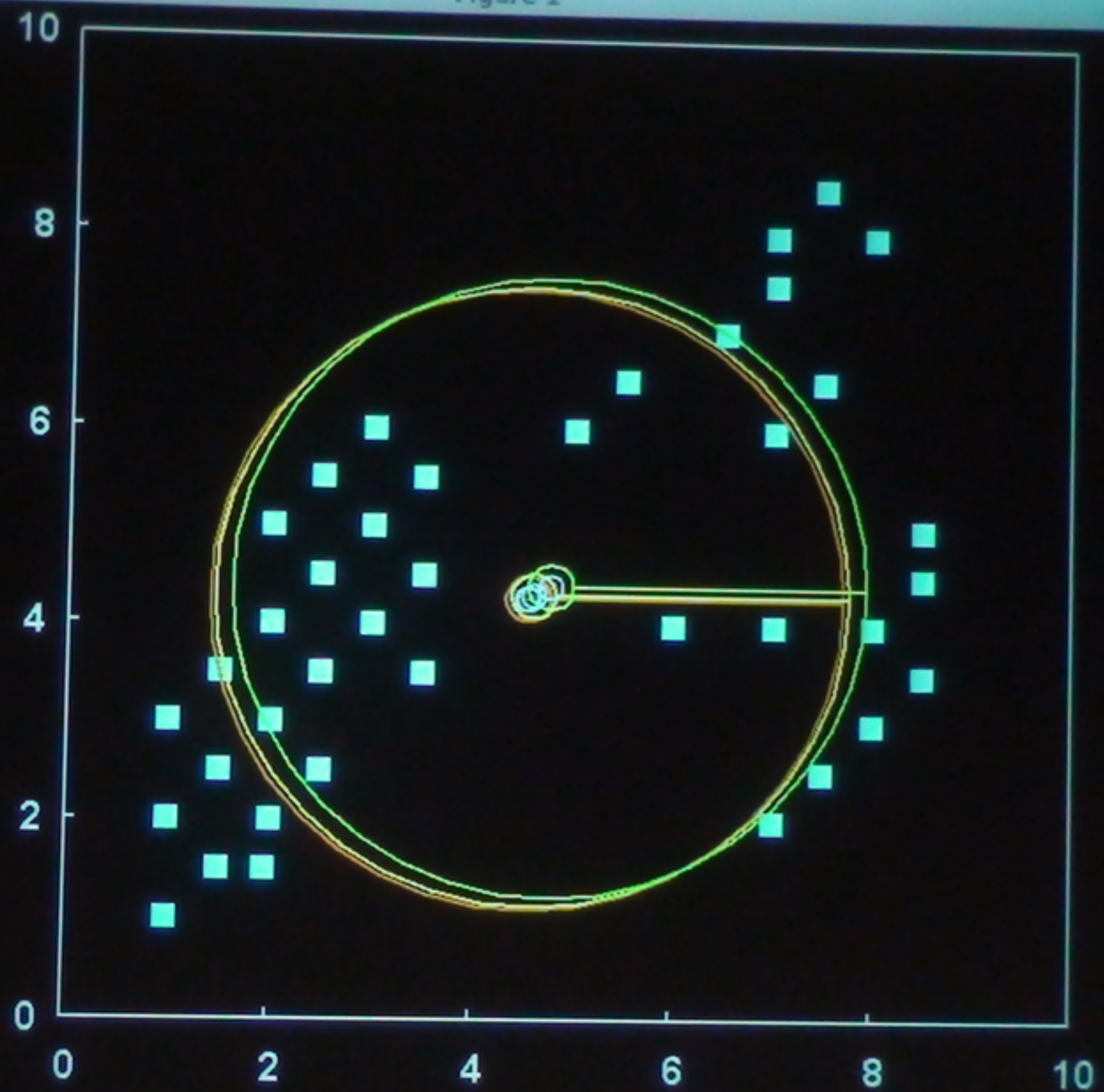


Figure 1



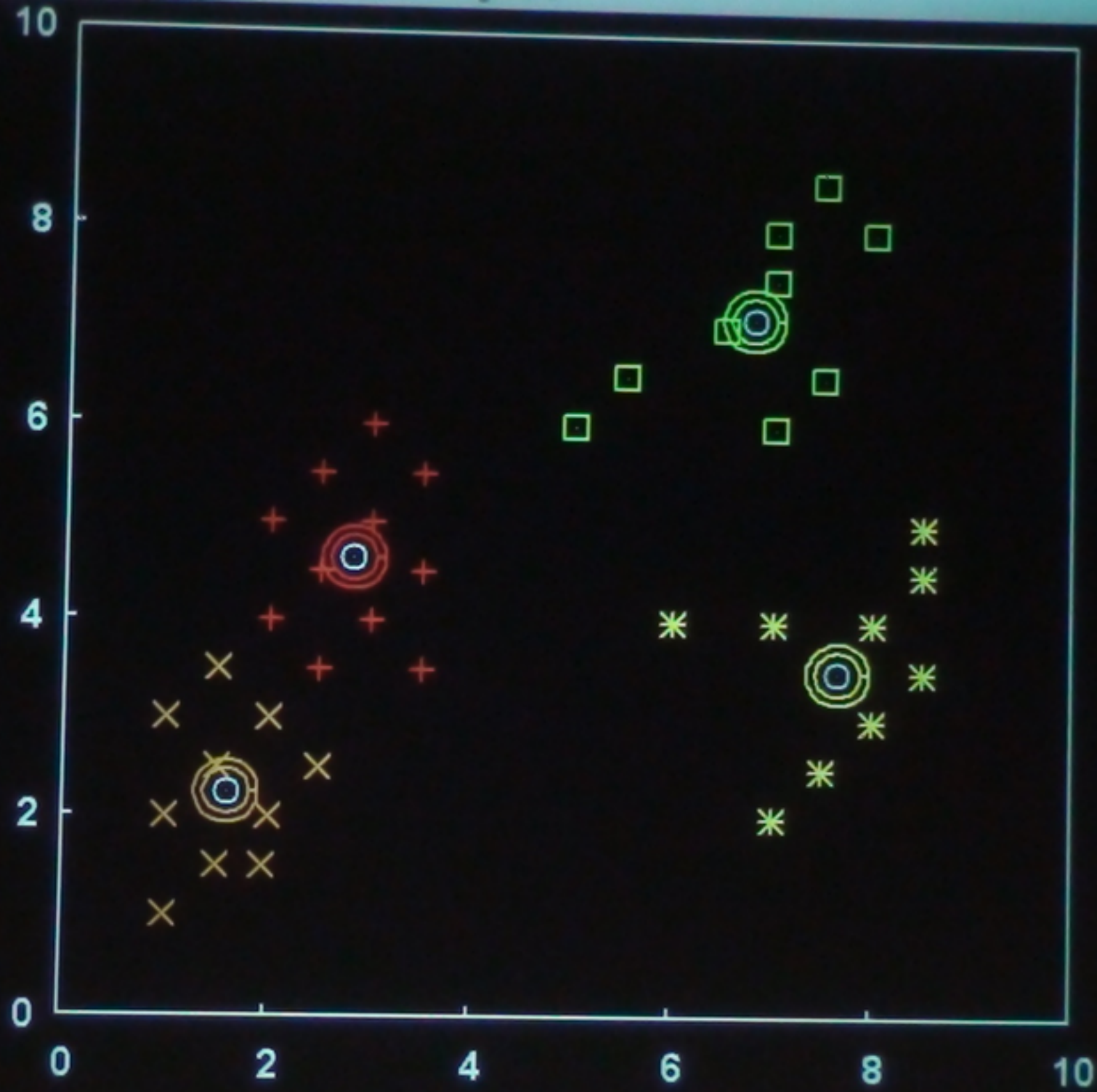
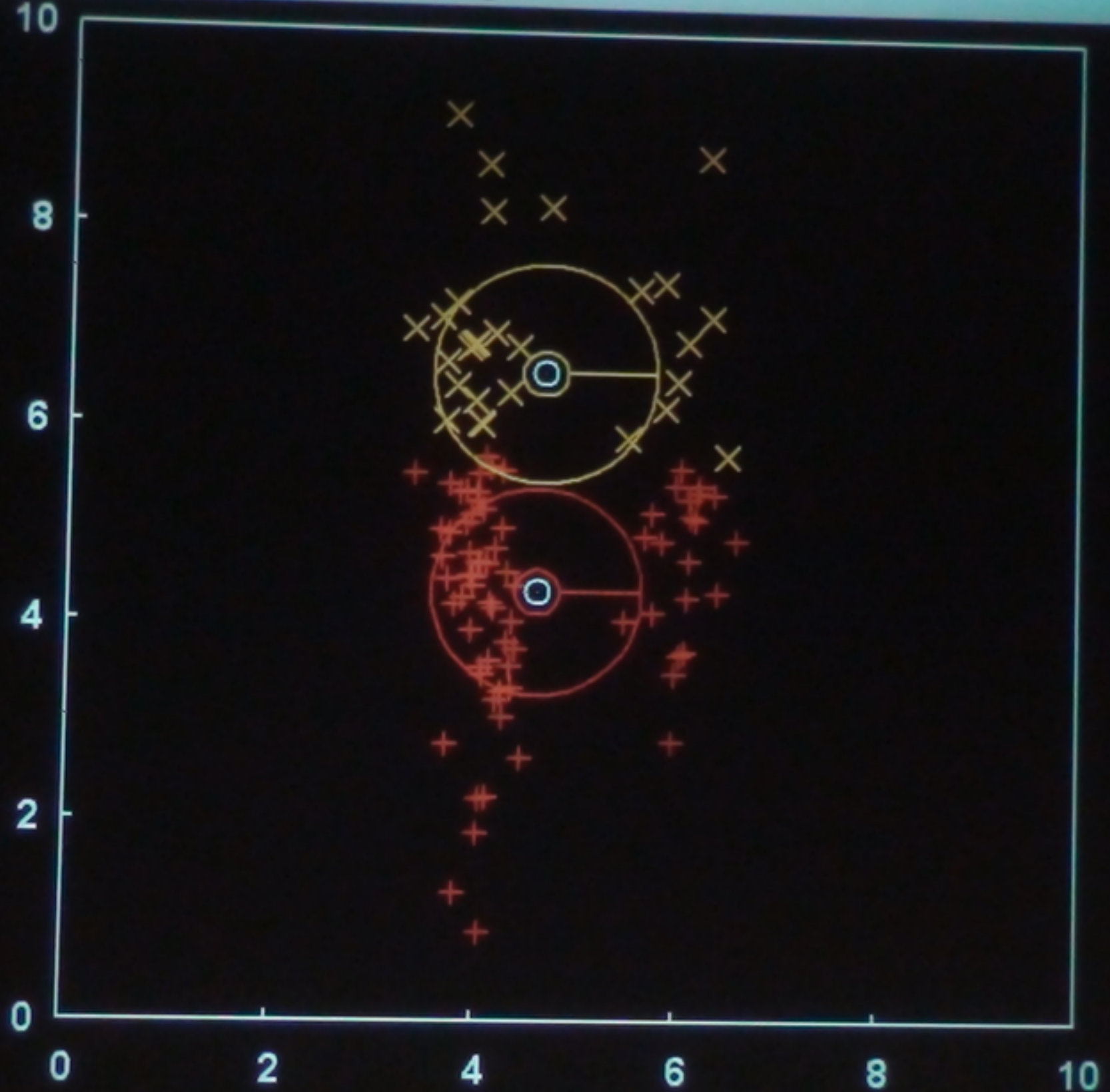
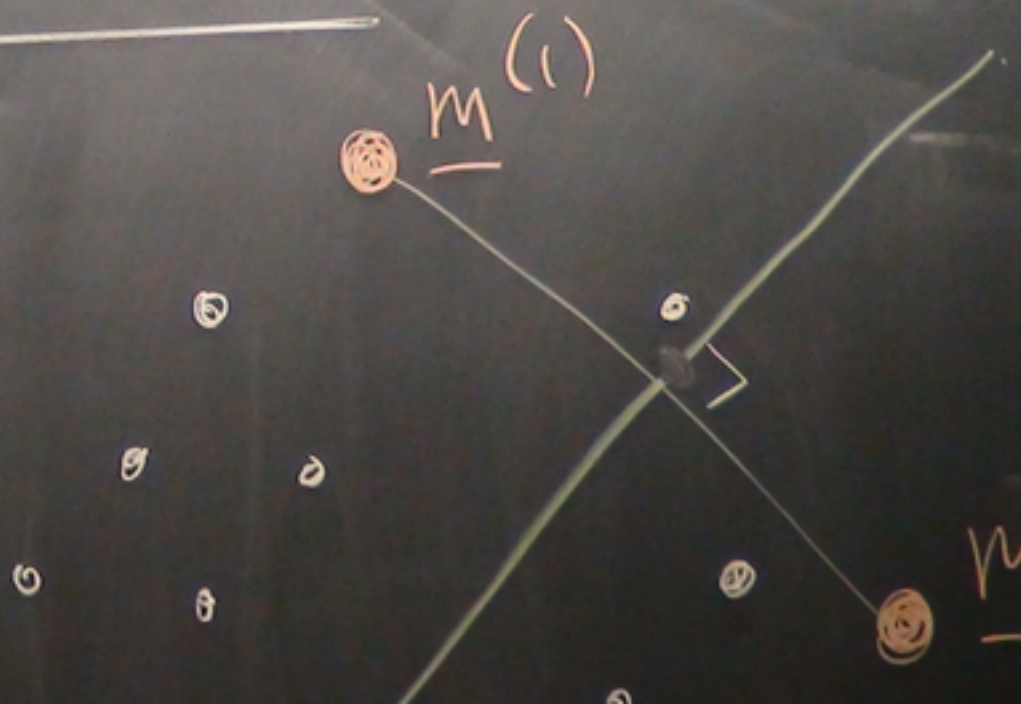
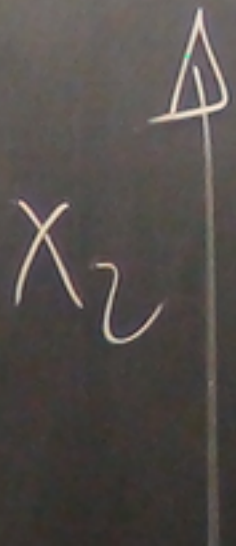


Figure 1



SOFT K MEANS II



& set a "stiffness" β

$$\sigma_i^{(k)} \\ i=1 \dots I$$

$$\beta d(m^{(k)}, x^{(n)})$$

"softmin"

$$\beta \equiv$$

$$e^{-\beta d(m^{(k)}, x^{(n)})}$$



$$\sigma_i^{(k)^2} =$$

$$\frac{\sum_k r_k^{(n)} \left(X_i^{(n)} - m_i^{(k)} \right)^2}{\sum_k r_k^{(n)}} \leftarrow$$

~~$$\sum_k r_k^{(n)}$$~~

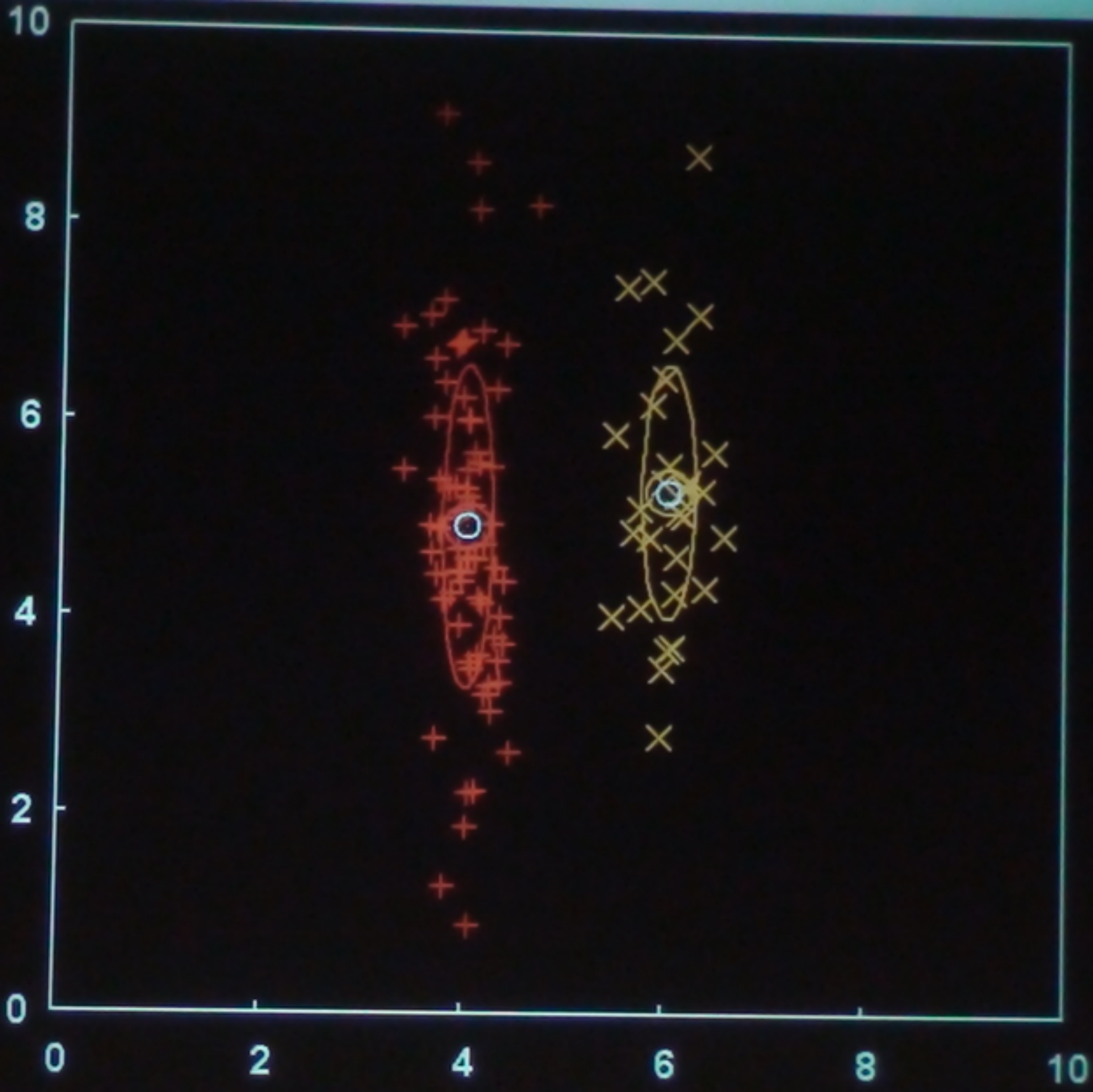
$$m^{(k)} =$$

$$\frac{\sum_{n=1}^N r_k^{(n)} X^{(n)}}{\sum_{n=1}^N r_k^{(n)}}$$

$$\sum_{n=1}^N r_k^{(n)}$$

$$\bar{\pi}_k = \frac{\sum r_k^{(n)}}{N}$$

0



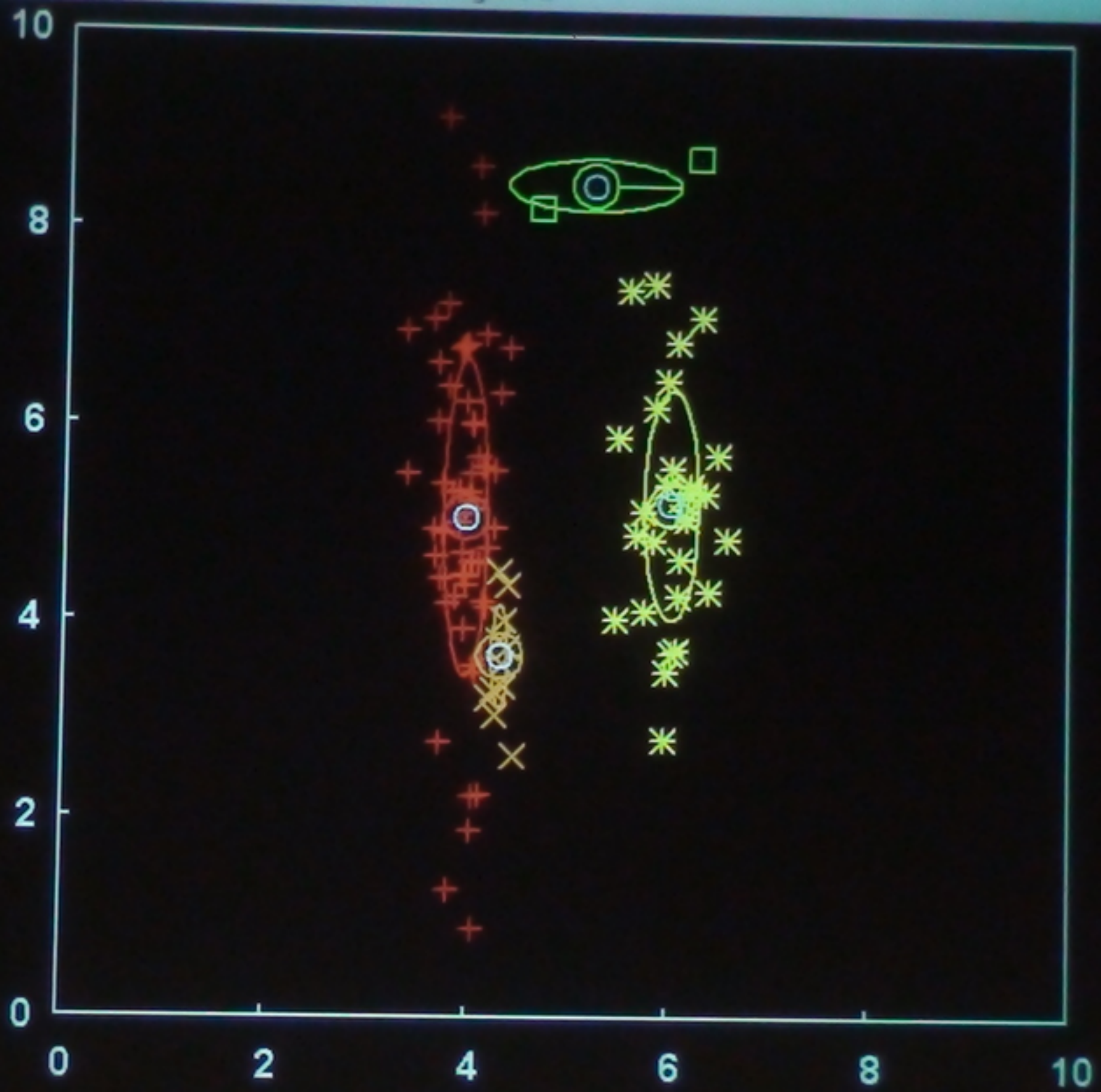


Figure 1

