

Gaussian Mixture Model (GMM)

For every sample x in the training set,

$$P(x|\theta) = \sum_{k=1}^K \pi_k \cdot G(x|\mu_k, \Sigma_k)$$

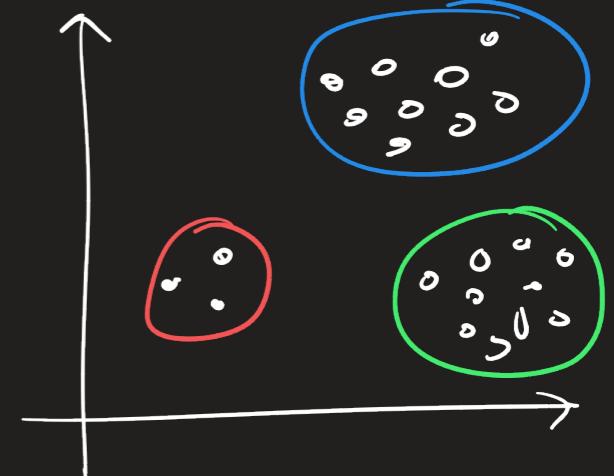
$$\theta = \left\{ \pi_k, \mu_k, \Sigma_k \right\}_{k=1}^K$$

$K \equiv \#$ Gaussian components (hyperparameter)

weight of each component
OR responsibility
of k^{th} Gaussian
component in
representing
the entire
mixture.

such that $\sum_{k=1}^K \pi_k = 1$

and $0 \leq \pi_k \leq 1$

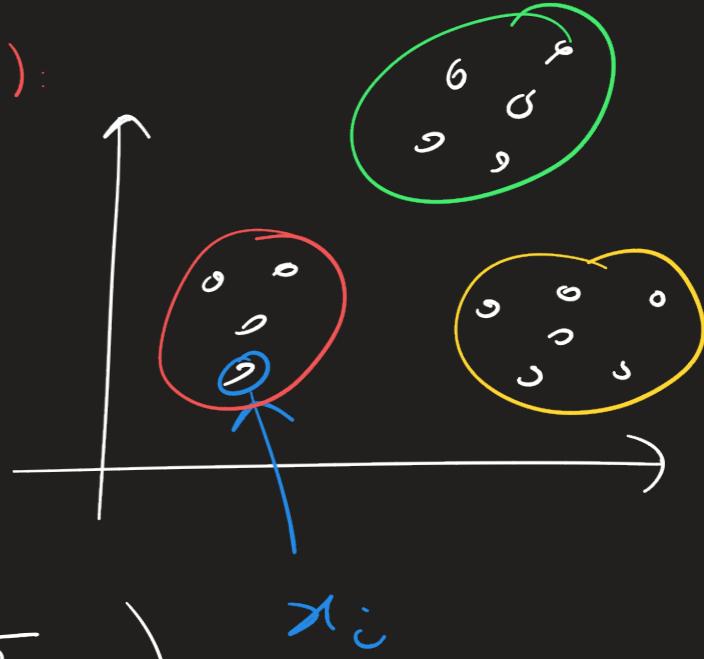


DATA SET: $\{x_i\}_{i=1}^N$, i.i.d.

OBSERVED DATA Likelihood (as a mixture model):

$$\mathcal{L} = \prod_{i=1}^N P(x_i | \theta)$$

$$= \prod_{i=1}^N \sum_{k=1}^K \pi_k G(x_i | \mu_k, \Sigma_k)$$



Set of parameters:

$$\theta = \{\pi_k, \mu_k, \Sigma_k\}_{k=1}^K$$

Optimization problem:

$$\arg \max_{\theta} \mathcal{L}$$

② Applying \log to L^o :

Log-likelihood:

$$L = \ln L^o \\ = \sum_{i=1}^N \ln \left(\sum_{k=1}^K \pi_k \cdot G(x_i | \mu_k, \Sigma_k) \right)$$

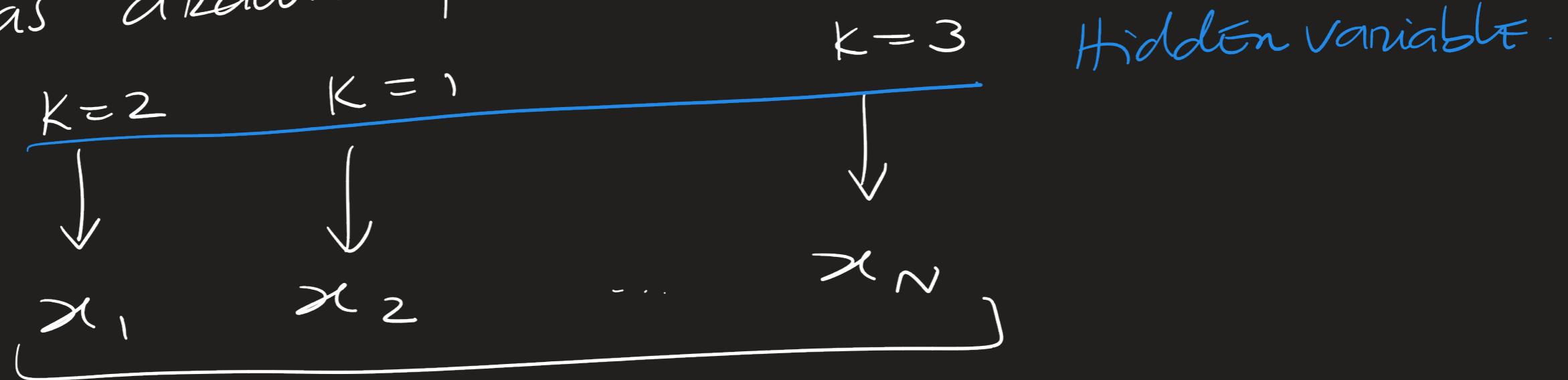
intractable density function

$$\frac{\partial L}{\partial \mu_k} = 0 \dots \text{No analytical solution} \dots$$

Expectation - Maximization (EM)

Algorithm

It introduces a hidden latent variable, z , that encodes a label from which Gaussian component each sample x_i was drawn from.

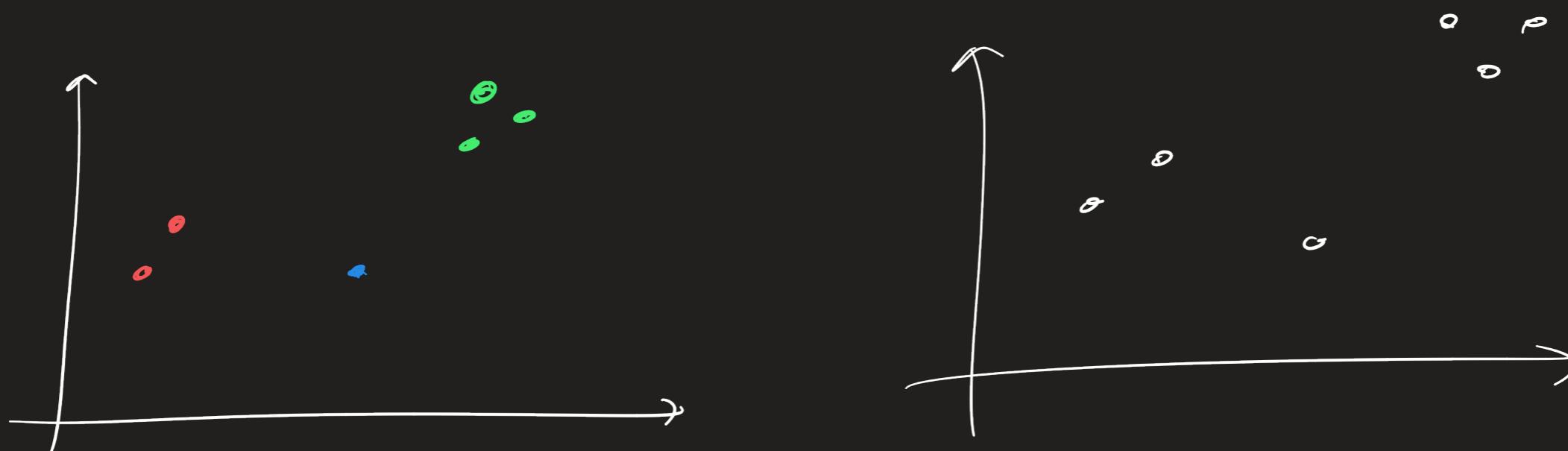
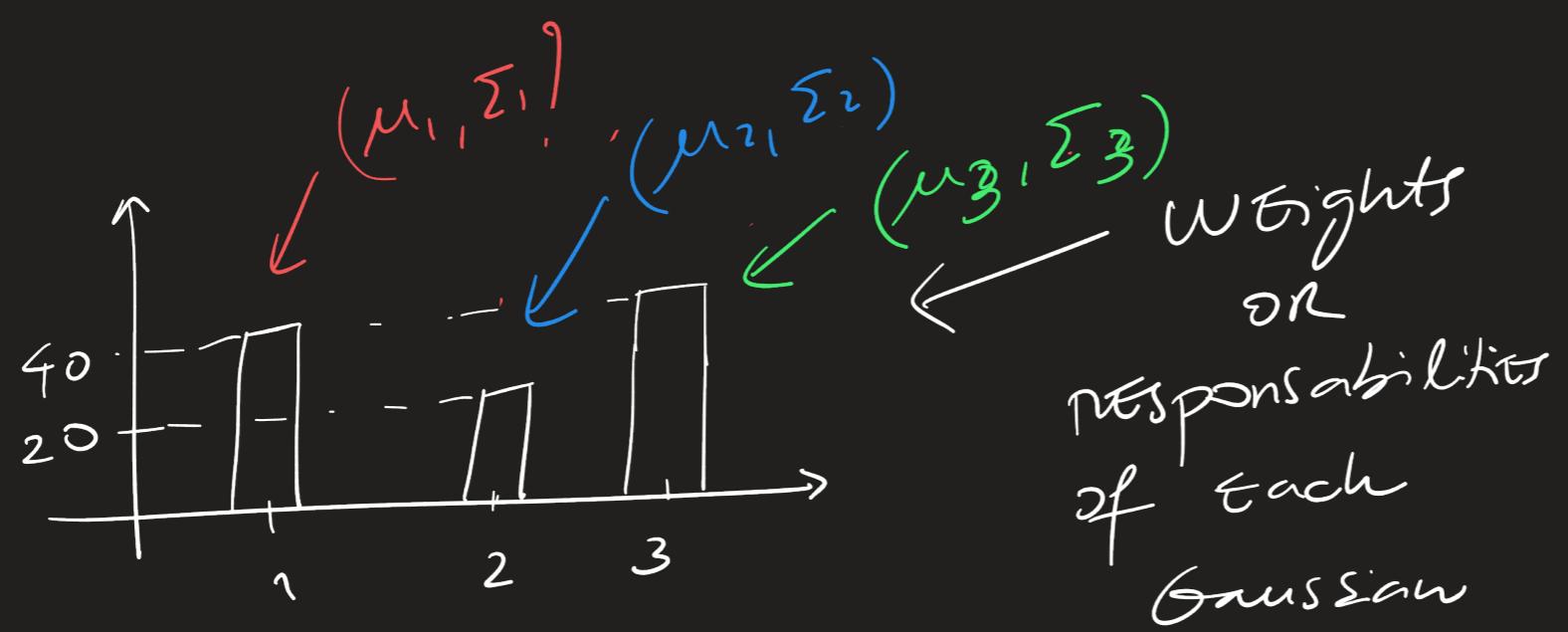


x_i will have a label z_i
which corresponds to the Gaussian
component it was drawn from.

$$z_i \in \{1, 2, \dots, K\}$$



Gaussian
Component
label.



$$\mathcal{L}^o = \prod_{i=1}^N \sum_{k=1}^K \pi_k \cdot G(x_i | \mu_k, \Sigma_k)$$

Hidden
latent
variable:

z_i = true und-seizable value
for the Gaussian component
that x_i was drawn

from

$$z_i \in \{1, 2, \dots, K\}$$

COMPLETE
DATA
Likelihood

$$\mathcal{L}^c = \prod_{i=1}^N \pi_{z_i} \cdot G(x_i | \mu_{z_i}, \Sigma_{z_i})$$

such that $\sum_{k=1}^K \pi_k = 1$

and $0 \leq \pi_k \leq 1$

EM Algorithm

- ① Initialize iteration $t = 0$.
- ② Initialize the parameter values
 $\theta^{(t)} = \{\pi_k^{(t)}, \mu_k^{(t)}, \sum_k^{(t)}\}_{k=1}^K$
- ③ Expectation (E)- step
Fix $\theta^{(t)}$, and find the labels
 $\{z_i^{(t)}\}_{i=1}^N$.
- ④ Maximization (M)- step.
Fix the labels $\{z_i^{(t)}\}_{i=1}^N$. Update the parameters $\theta^{(t+1)}$
- ⑤ go back to step ③ until convergence criteria is met.

The optimization function for
(or objective)

EM algorithm is:

$$Q(\theta, \theta^{(t)}) = E_Z [\ln(\mathcal{L}^e) | X, \theta^{(t)}]$$

Expected value
over R.V. Z .

θ = learnable parameters (variables)

$\theta^{(t)}$ = numerical values for θ .

X = training data

Reminder

① X is a discrete R.V., $\xrightarrow{\text{P.M.F.}}$

$$E_x[X] = \sum_i x_i \cdot p_x(x_i)$$

② X is a continuous R.V., $\xrightarrow{\text{P.D.F.}}$

$$E_x[X] = \int_{-\infty}^{+\infty} x_i \cdot f_x(x_i) dx_i$$

③ X is discrete R.V.,

$$E_x[f(x)] = \sum_i f(x_i) \cdot p_x(x_i)$$

— // —
LOTUS (Law of the Unconscious Statistician)

$$\mathcal{L}^c = \prod_{i=1}^N \pi_{z_i} \cdot G(x_i | \mu_{z_i}, \Sigma_{z_i})$$

$$Q(\theta, \theta^{(t)}) = \mathbb{E}_z [\ln \mathcal{L}^c | X, \theta^{(t)}]$$

$$= \sum_{z_i=1}^K \frac{\ln(\mathcal{L}^c)}{\downarrow} \cdot \frac{P_z(z_i | X, \theta^{(t)})}{\downarrow}$$

optimize this
 term in the
 M-step. optimizing for
 in E-step
 (Finds the z_i 's)
 (Finds the $\{\mu_k, \Sigma_k, \pi_k\}_{k=1}^K$)

$$= \sum_{i=1}^N \sum_{z_i=1}^K \left[\ln(\pi_{z_i}) - \frac{d}{2} \ln(2\pi) - \frac{d}{2} \ln(|\Sigma_{z_i}|) - \frac{1}{2} (x_i - \mu_{z_i})^\top \Sigma_{z_i}^{-1} (x_i - \mu_{z_i}) \right] \cdot P(z_i | X, \theta^{(t)})$$

$$= \sum_{i=1}^N \sum_{k=1}^K \left[\ln(\pi_k) - \frac{d}{2} \ln(2\pi) - \frac{d}{2} \ln(|\Sigma_k|) - \frac{1}{2} (x_i - \mu_k)^\top \Sigma_k^{-1} (x_i - \mu_k) \right] \cdot P(z_i = k | X, \theta^{(t)})$$

E - STEP:

Fix the current value for parameters, $\theta^{(t)}$.

Estimate the following: $G(x_i | \mu_k, \Sigma_k) \xrightarrow{\pi_k}$

$$P(z_i = k | x_i, \theta^{(t)}) = \frac{P(x_i | z_i = k, \theta^{(t)}) \cdot P(z_i | \theta^{(t)})}{P(x_i | \theta^{(t)})}$$

$$= \frac{G(x_i | \mu_k^{(t)}, \Sigma_k^{(t)}) \cdot \pi_k^{(t)}}{\sum_{j=1}^K G(x_i | \mu_j^{(t)}, \Sigma_j^{(t)}) \cdot \pi_j^{(t)}} = U_{ik} \quad \begin{matrix} \text{membership of} \\ \text{sample } x_i \text{ in} \\ \text{component } k. \end{matrix}$$

In computation, this can be written as
a matrix, U , the membership matrix.
 $N \times K$

$$U = \begin{bmatrix} & 1 & 2 & \dots & K \\ 1 & 0.1 & 0.4 & \dots & 0.5 \\ 2 & 0.3 & 0.5 & \dots & 0.2 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ N & \vdots & \vdots & \ddots & \vdots \end{bmatrix} \quad \begin{matrix} \leftarrow \text{membership values} \\ \text{of sample } x_1 \text{ in} \\ \text{each gaussian} \\ \text{component} \\ \sum_{j=1}^K U_{ij} = 1, \forall i \end{matrix}$$

$\uparrow \quad \begin{matrix} \text{membership of all points} \\ \text{in gaussian comp. } K. \end{matrix}$

$$P(z_i = 1 | x_i, \theta^{(t)}) = U_{i1}$$

$$\sum_{i=1}^N U_{i1}$$

Find the column labeled
with largest probability

$$\rightarrow \begin{bmatrix} x_1 \in C_K \\ x_2 \in C_2 \\ \vdots \end{bmatrix} = \hat{z} = \begin{bmatrix} K \\ 2 \\ \vdots \\ 1 \end{bmatrix}_{N \times 1}$$

Component
labels for
all samples

M-STEP

Fix the z 's from E-step and
optimize for $\theta = \{\pi_k, \mu_k, \Sigma_k\}_{k=1}^K$

We can use MLE or MAP for parameter estimation.

These illustrations will use MLE.

I'm going to further assume that

$$\Sigma_k = \sigma_k^2 \cdot I \quad (\text{isotropic covariance})$$

$$|\Sigma_k| = \underbrace{(\sigma_k^2)^d}_{\text{assuming } d \text{ features}}, \quad \Sigma_k^{-1} = (\sigma_k^2)^{-1} \cdot I$$

$$Q(\theta, \theta^{(t)}) = \sum_{k=1}^K \sum_{i=1}^N \left[\ln(\pi_k) - \frac{d}{2} \ln(2\pi) - \frac{d}{2} \ln(\sigma_k^2) - \frac{1}{2\sigma_k^2} (x_i - \mu_k)^T (x_i - \mu_k) \right] \cdot v_{ik}$$

where $v_{ik} = P(z_i = k | x_i, \theta^{(t)})$
 (which is fixed during M-step).

M-step:

1) starting w/ μ_k :

$$\frac{\partial Q}{\partial \mu_k} = 0 \Leftrightarrow \sum_{i=1}^N \frac{1}{\sigma_k^2} (x_i - \mu_k) \cdot v_{ik} = 0$$

$$\begin{aligned} &\Leftrightarrow \sum_{i=1}^N x_i \cdot v_{ik} = \sum_{i=1}^N \mu_k \cdot v_{ik} \\ &\Leftrightarrow \mu_k = \frac{\sum_{i=1}^N x_i \cdot v_{ik}}{\sum_{i=1}^N v_{ik}} \quad \text{weighted average.} \end{aligned}$$

2) w.r.t. σ_k^2 :

$$\frac{\partial Q}{\partial \sigma_k^2} = 0$$

$$\Leftrightarrow \sum_{i=1}^N \left[-\frac{d}{2} \cdot \frac{1}{\sigma_k^2} \cdot \frac{2}{(2\sigma_k^2)^2} (\mathbf{x}_i - \boldsymbol{\mu}_k)^T (\mathbf{x}_i - \boldsymbol{\mu}_k) \right] \cdot v_{ik} = 0$$

$$\Leftrightarrow \sigma_k^2 = \frac{\sum_{i=1}^N (\mathbf{x}_i - \boldsymbol{\mu}_k)^T (\mathbf{x}_i - \boldsymbol{\mu}_k) \cdot v_{ik}}{d \sum_{i=1}^N v_{ik}}$$

≡ weighted

variance

$$\Sigma_k = \sigma_k^2 \cdot I$$

3) w.r.t. π_k

Must obey the constraints on π_k

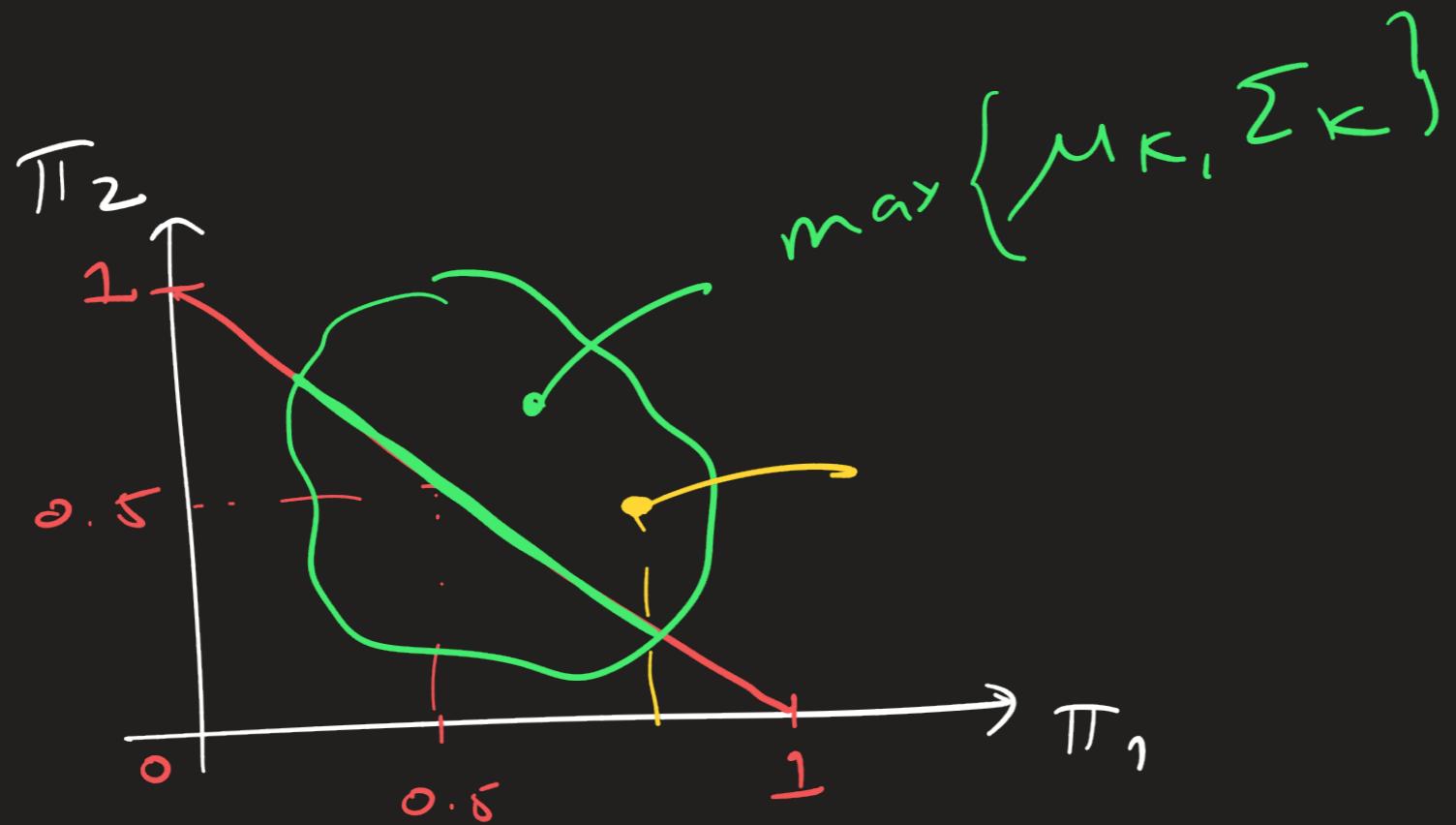
1) $\sum_{k=1}^K \pi_k = 1 \times \Leftrightarrow 1 - \sum_{k=1}^K \pi_k = 0$

2) $0 \leq \pi_k \leq 1$

To satisfy (1), we add it to the optimization function - Lagrangian
optimization problem.

$$Q_{\pi}(\theta, \theta^{(t)}) = Q(\theta, \theta^{(t)}) + \lambda \left(1 - \sum_{k=1}^K \pi_k \right)$$

Lagrangian
multiplier



APPENDIX E
on the
Bishop
textbook

$$\frac{\partial Q\pi}{\partial \pi_K} = 0$$

$$\Leftrightarrow \sum_{i=1}^N \frac{1}{\pi_K} \cdot v_{ik} - \lambda = 0$$

$$\Leftrightarrow \boxed{\pi_K = \frac{\sum_{i=1}^N v_{ik}}{\lambda}}$$

$$\boxed{\pi_K = \frac{1}{N} \sum_{i=1}^N v_{ik}}$$

Sum over K ,

$$\sum_{K=1}^K \pi_K = 1$$

$$\sum_{K=1}^K \left(\sum_{i=1}^N v_{ik} \cdot \frac{1}{\lambda} \right) = 1$$

$$\Leftrightarrow \sum_{i=1}^N \frac{1}{\lambda} \sum_{K=1}^K v_{ik} = 1$$

$$\Leftrightarrow \sum_{i=1}^N \frac{1}{\lambda} = 1 \quad \Leftrightarrow \frac{1}{\lambda} \cdot N = 1 \quad (\Rightarrow \lambda = N)$$

average
responsability
value for
component K .