T1.1 → Intro data analysis & ML

T1.2 → CLUSTERING ALGORITHMS

$\quad\quad\quad\quad$ ↳ k-means

$\quad\quad\quad\quad$ ↳ Hierarchical clustering

$\quad\quad\quad\quad$ ↳ Gaussian Mixture Models (GMM)

$\quad\quad\quad\quad\quad\quad$ • Model selection

$\quad\quad\quad\quad\quad\quad$ • Parameter estimation in
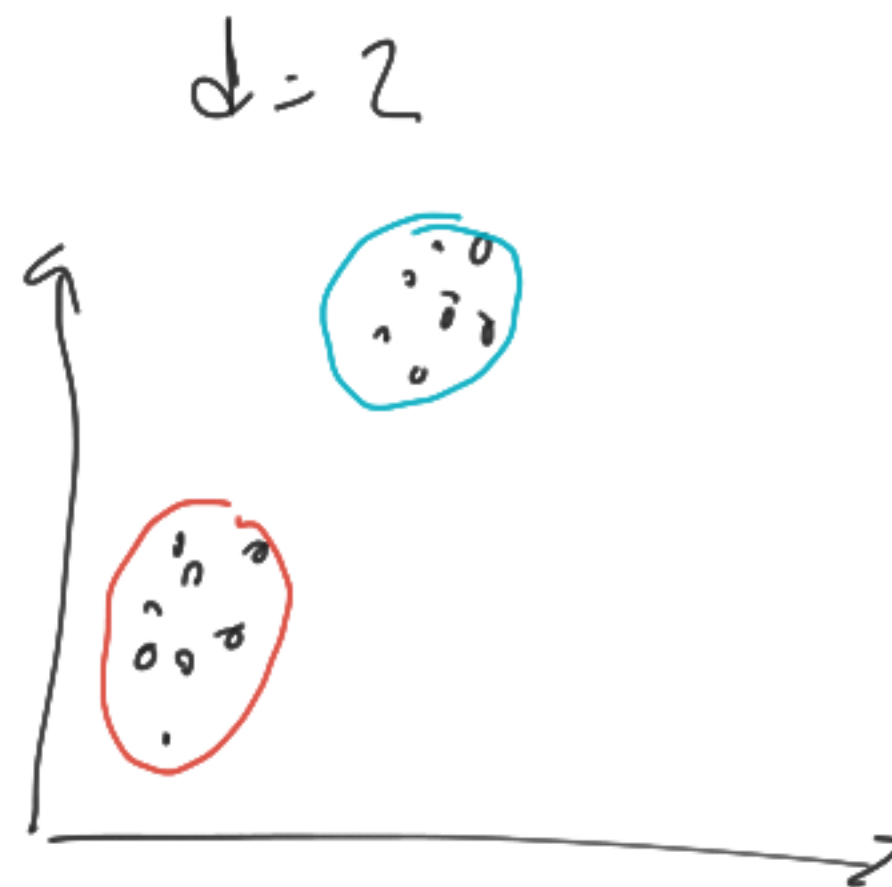$\quad\quad\quad\quad\quad\quad\quad\quad\quad$ probabilistic models.

# CLUSTERING

d-dimensional feature space $\vec{x} \in \mathbb{R}^d$

From a set of observations

$$\{ \vec{x}_1, \vec{x}_2, \vec{x}_3, \cdots, \vec{x}_N \}$$

## UNSUPERVISED METHODS !

$d = 2$

DATA ML $\longrightarrow$ CLASS MEMBERSHIP OF EACH OBSERVATION

data matrix
_____
observations × features

$W_j \quad j = 1 \ldots NOBS$

$W_1 = $ 'FCB supporter'

$W_2 = $ 'other team'

$W_3 = $ 'other team''



obs

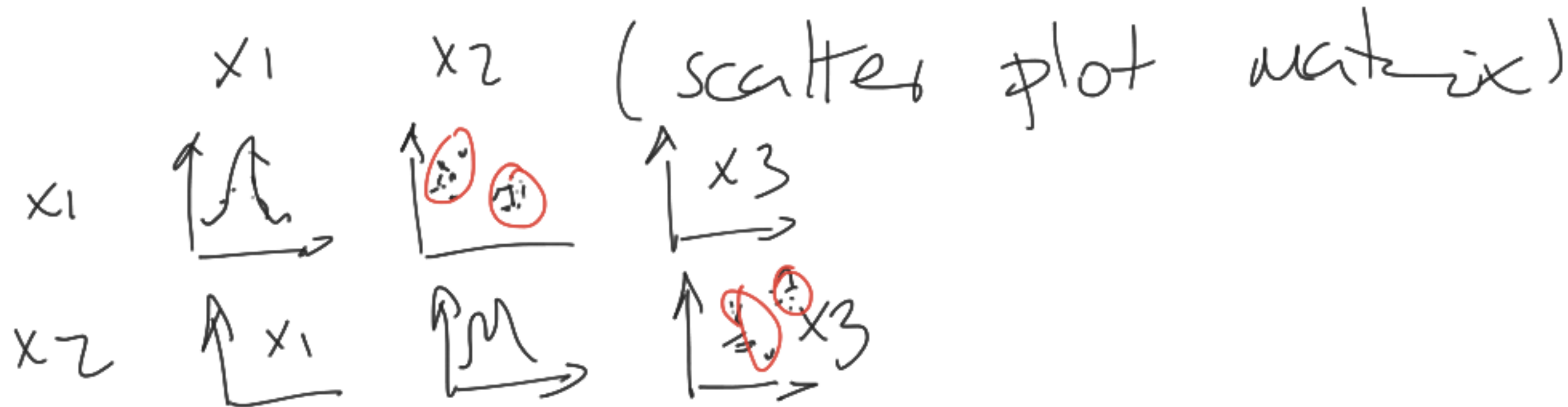features / variable / attributes.

$W = \{ $ 'FCB supporter', 'other team' $\}$

data matrix $\longrightarrow$ unsupervised techniques

$$\left\{ \begin{array}{l} \text{data matrix} \\ \text{class - label vector} \\ \{\omega_1, \omega_2, \ldots \omega_{NOBS}\} \end{array} \right. \longrightarrow$$ supervised techniques

$\downarrow$

TRAIN / FITTIN / LEARNING

(ML)

k-means $\rightarrow$ require pre-define #clusters to be found in data
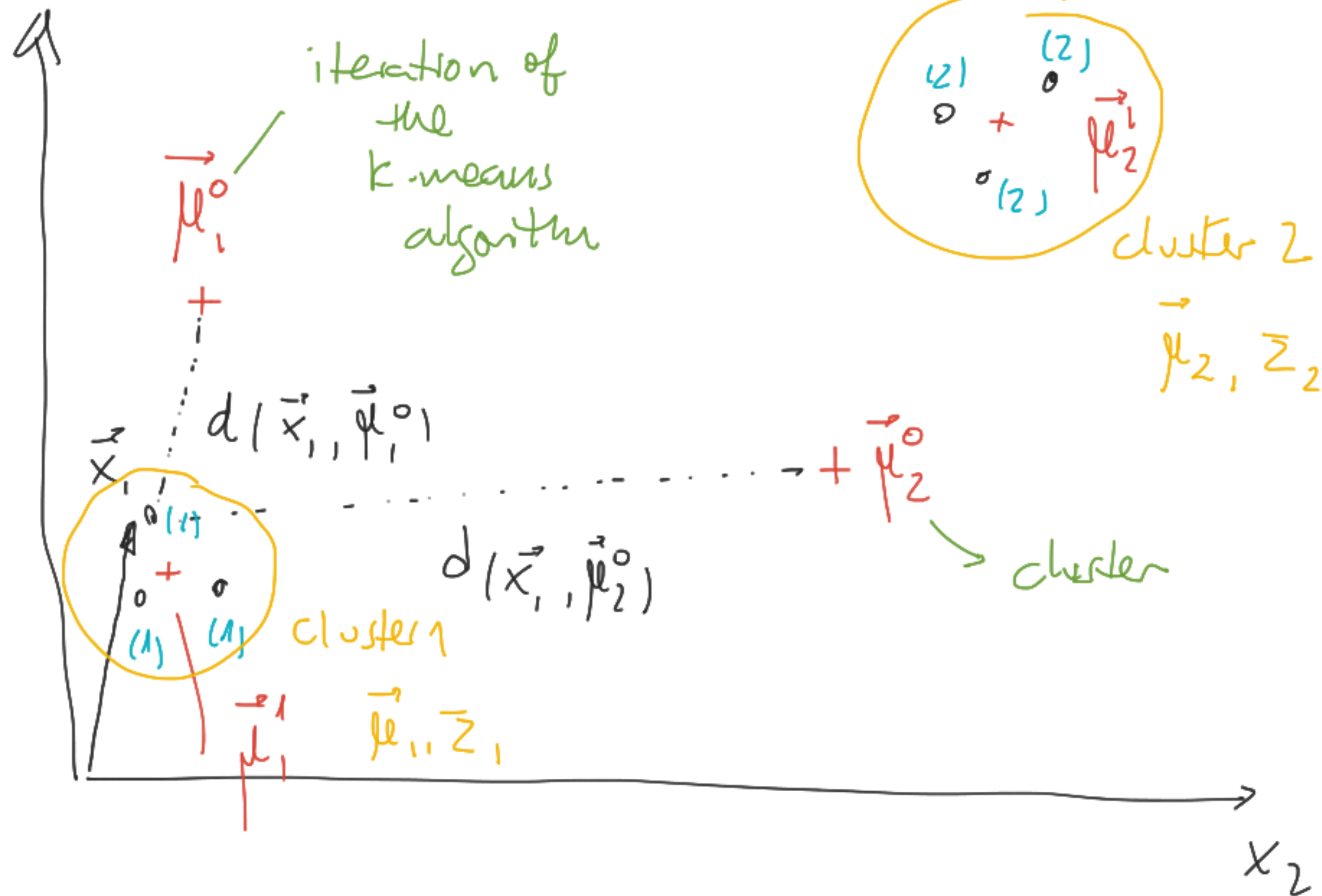
hierarchical

2d $\rightarrow$ easy

$>$2d $\rightarrow$ plot data using a pairplot

x1    x2    (scatter plot matrix)

k-means $x_1$

$d = 2$

$\vec{\mu}_1, \vec{\mu}_2$ ?

iteration of
the
k-means
algoritm

$\vec{\mu}_1^0$

$d(\vec{x}_1, \vec{\mu}_1^0)$

$\vec{x}_1$

(1)

$d(\vec{x}_1, \vec{\mu}_2^0)$

$+ \vec{\mu}_2^0$

cluster

(1)    (1)    cluster 1

$\vec{\mu}_1^1$    $\vec{\mu}_1, \Sigma_1$

$\vec{\mu}_1^0$

(2)    (2)

(2)    $+$    $\vec{\mu}_2^1$
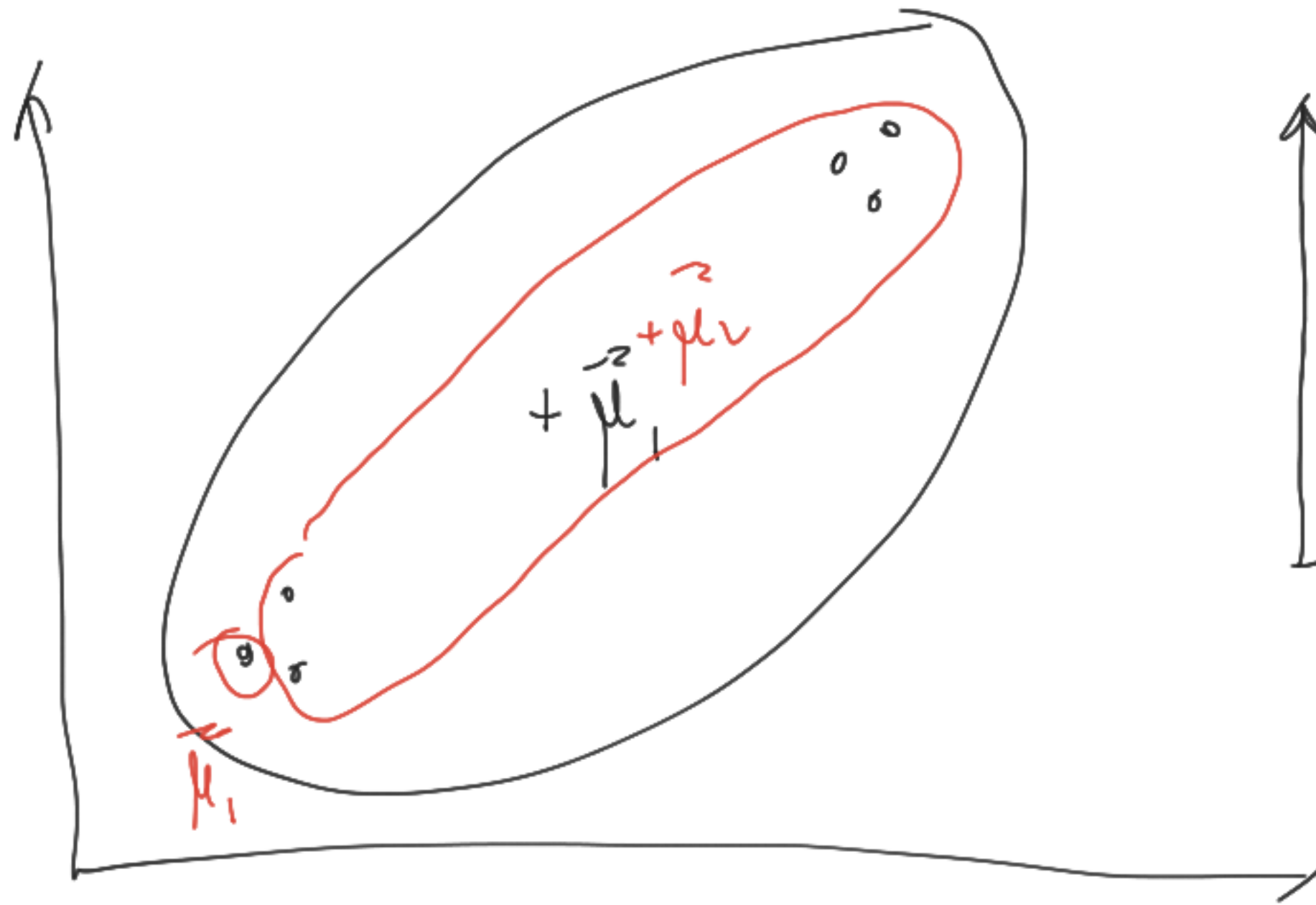
(2)

cluster 2

$\vec{\mu}_2, \Sigma_2$

$x_2$

hierarchical : (Agglomerative)

pdist $(\vec{\mu_i}, \vec{\mu_j})$

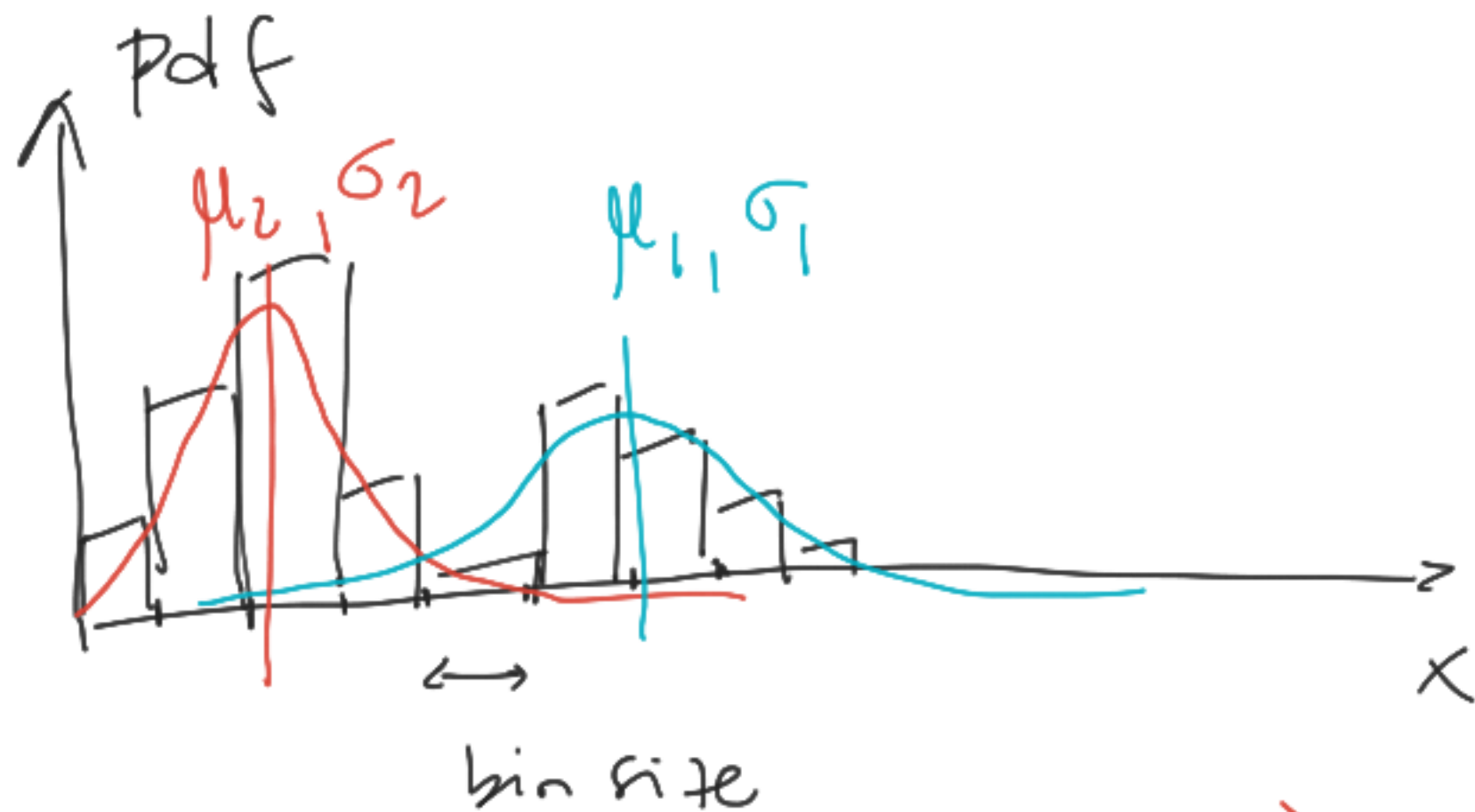K = 2 clusters

# GMM $\quad \vec{x} \in \mathbb{R}^d$

$$p(\vec{x}) = \sum_{j=1}^{G} \pi_j \cdot N(\vec{\mu}_j, \bar{\Sigma}_j)$$
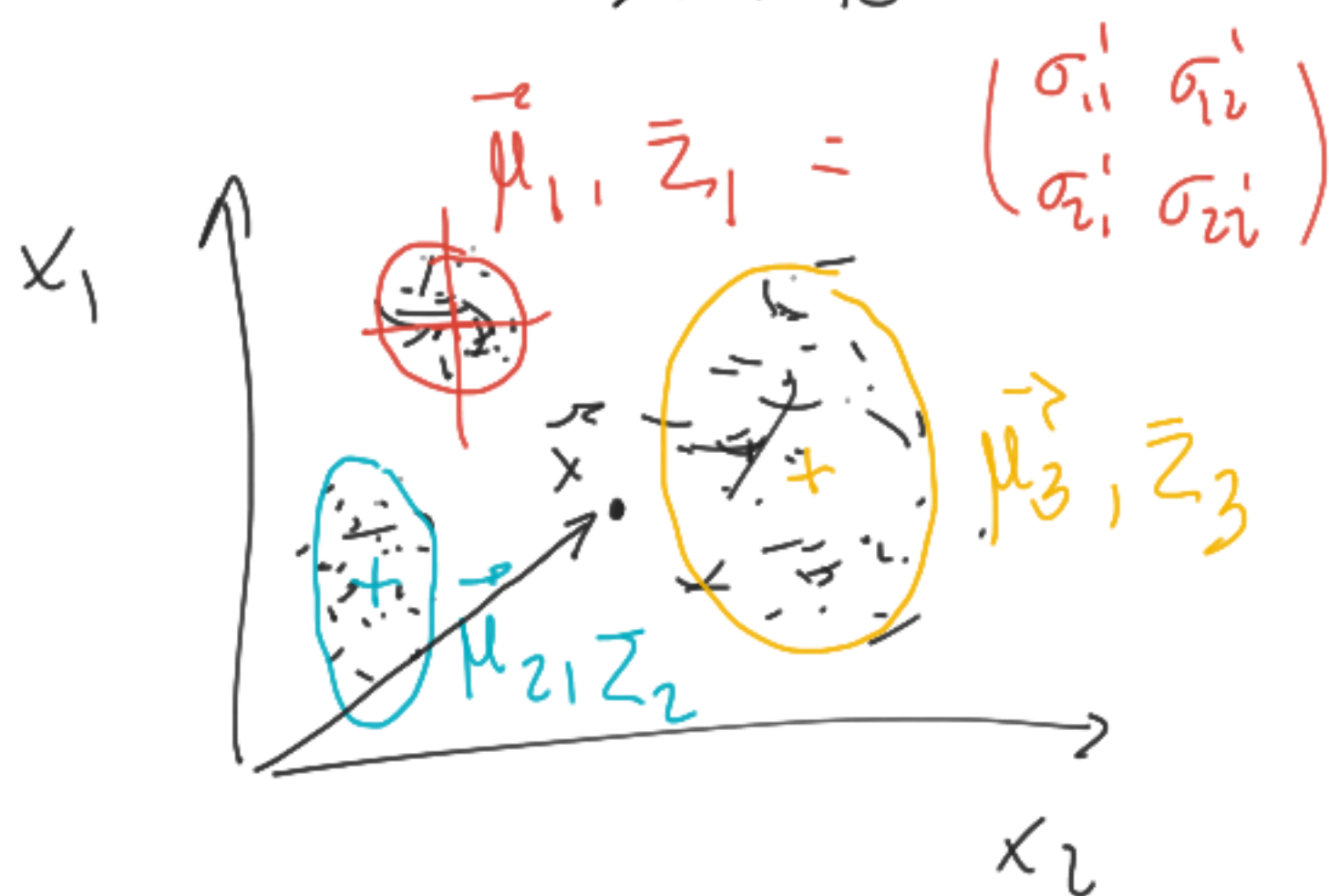
weights

mean
of
gaussian $j$

$(1 \times d)$

covariance matrix
of gaussian $j$

$(d \times d)$

Ex: $d = 1$ (univariate), obser. # $G = 2$ gaussians



bin size

$\vec{\mu}_1, \bar{Z}_1 = \begin{pmatrix} \sigma_{11} & \sigma_{12} \\ \sigma_{21} & \sigma_{21} \end{pmatrix}$

Ex. $d = 2$ $G = 3$

$p(\vec{x}) = \sum_{j=1}^{3} \pi_j \cdot N(\vec{\mu}_j, \bar{Z}_j)$

$\boxed{\text{GMM}}$ -2 pre-define # clusters → # gaussians

<span style="color:red"># parameters (GMM)</span>

| | | |
|---|---|---|
| $k=1$ | $\vec{\mu}_1, \bar{\Sigma}_1$ | $d + d \times d$ |

$\vec{\mu}_1$

$\bar{\Sigma}_1$ $\quad \bar{\mu}_1, \vec{\mu}_2 \quad \bar{\Sigma}_1, \bar{\Sigma}_2 \quad \pi_1$

$k=2 \qquad \vec{\mu}_1, \bar{\Sigma}_1, \vec{\mu}_2, \bar{\Sigma}_2, \pi_1 \quad (\pi_1 + \pi_2 = 1) \qquad \overbrace{2d} + \overbrace{2d \cdot d} + 1$

Model complexity

$\vdots$

<span style="color:red">PARSIMONY INDICES</span>
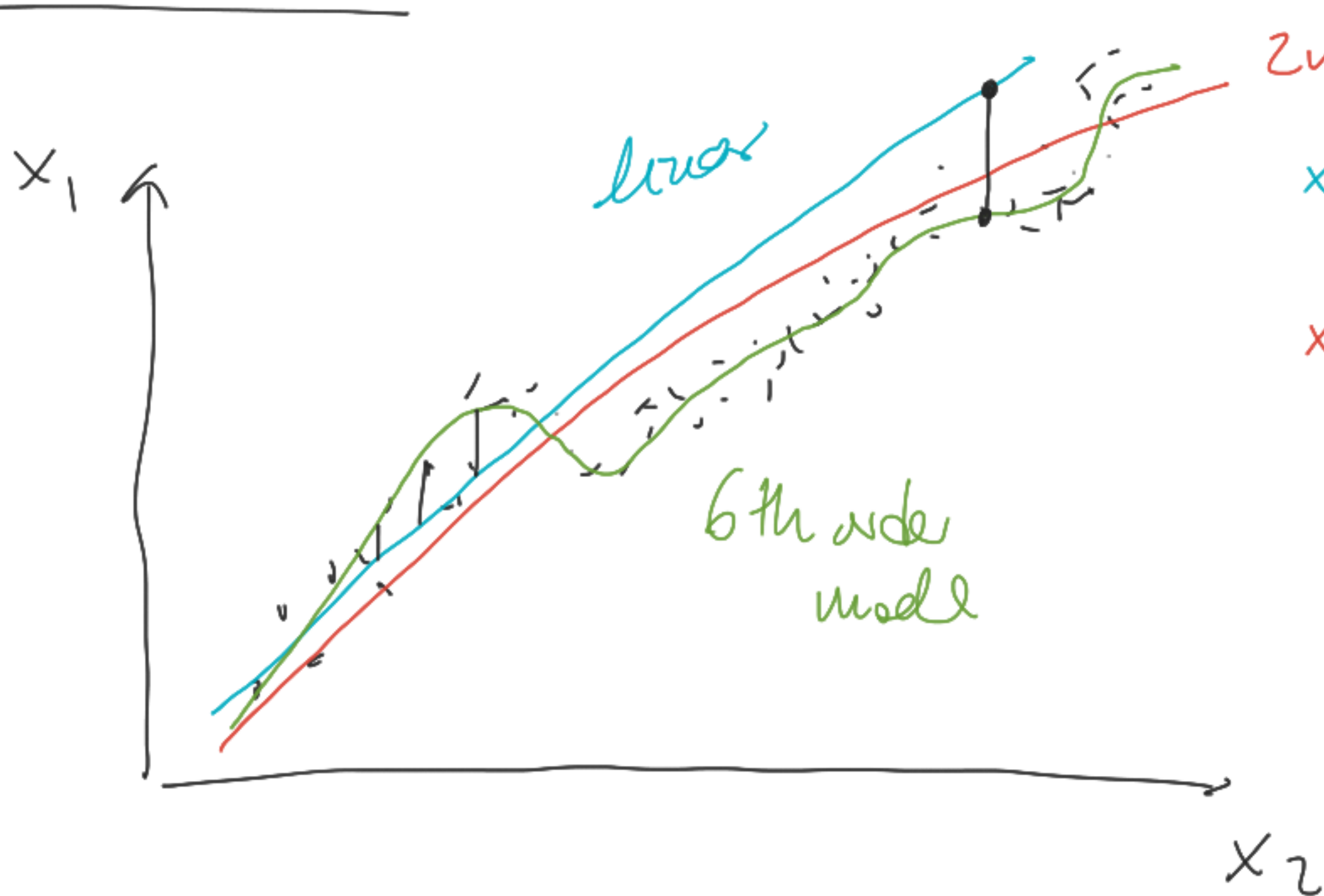
You need several observations to fit a GMM Model

PARSIMONY $=$ (RESIDUALS) MODEL (ERROR) MSE PERFORMANCE vs MODEL COMPLEXITY #parameters
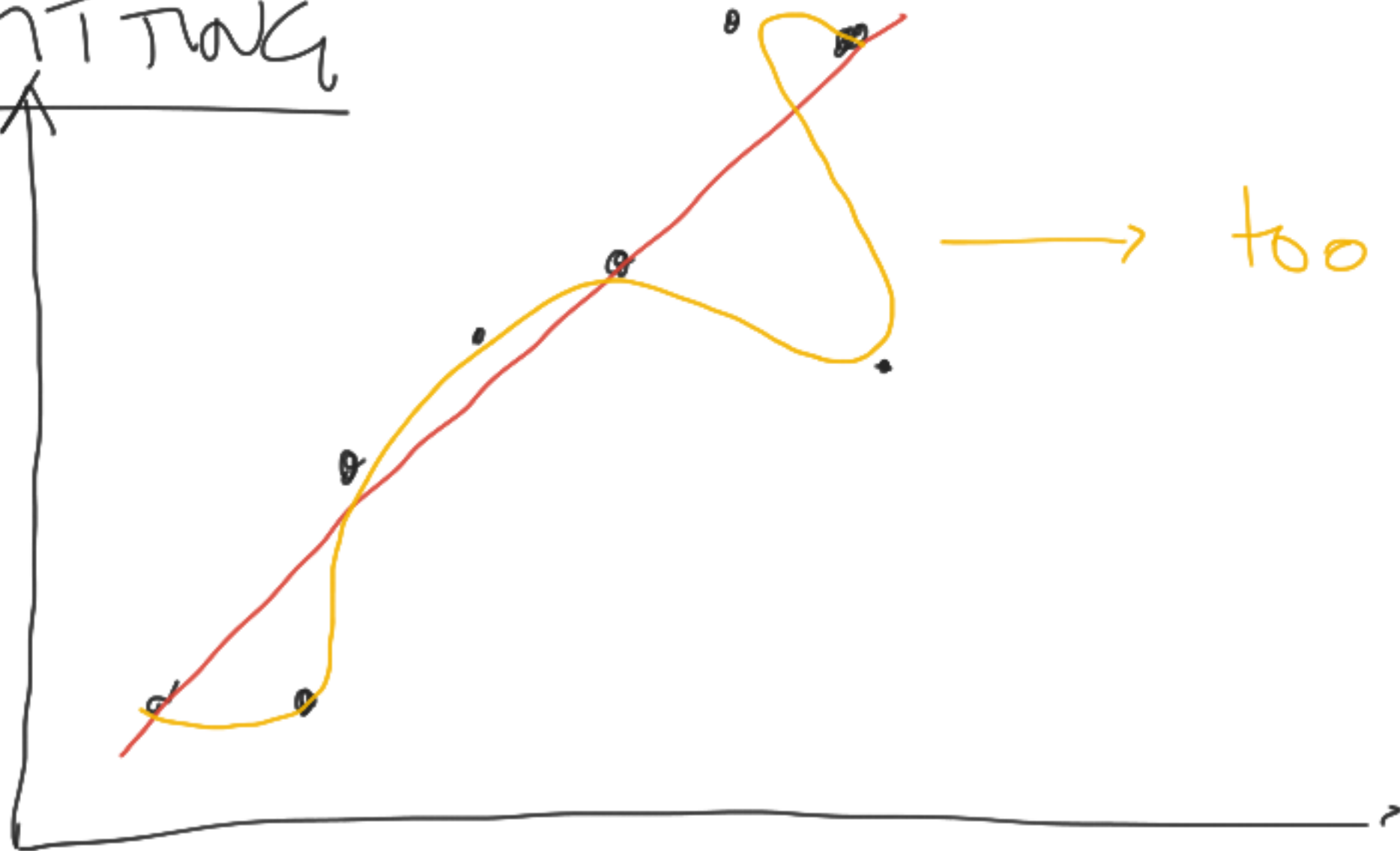


linear

2nd order

$x_1 = a x_2 + b$    2

$x_2 = a x_2^2 + b x_2 + \cdots$    3

$\vdots$

6th order model

avoid data overfitting

# OVERFITTING

$\longrightarrow$ too complex!

R      J

OVERFITTING $\iff$ GENERALIZATION

# GMM + MODEL SELECTION

### BIC
### PARSIMONY INDEX

fit GMM ($G=1$)

$BIC_1$    $-748$

DATA

fit GMM ($G=2$)

$BIC_2$   $-800$

fit GMM ($G=3$)

$BIC_3$   $-790$

fit GMM ($G=4$)

$BIC_4$   $-796$

$$BIC = -ERROR + COMPRESSION$$

$$-2\ln(\hat{\mathcal{L}}) + k\ln(n)$$

(ML)

#parameters

How many clusters?

$$\underset{j=1..4}{\arg\max} \{BIC_j\} \quad \text{2 clusters}$$

# PARAMETER ESTIMATION IN PROBABILISTIC MODELS

Probabilistic model:   $\vec{x}$ : data observation (d-dimensions)

$$\theta : \text{model parameters}$$

$$\downarrow$$

$$p(\vec{x}, \theta) \longrightarrow \text{joint pdf}$$

Ex:   $p(\vec{x} \mid \theta) = N(\vec{x} \mid \vec{\mu}, \vec{\Sigma})$

$$\theta : \text{model parameters}$$

$$p(\vec{x} \mid \theta) = \sum_{j=1}^{G} \pi_j \cdot N(\vec{\mu}_j, \vec{\Sigma}_j)$$

# PARAMETER ESTIMATION :-

sample of observations $\{\vec{x}_1^2, \vec{x}_2^{-2}, \ldots, \vec{x}_n\} \longrightarrow \theta^*$

SAMPLE ESTIMATE of model parameters

$\longrightarrow$ MAXIMUM LIKELIHOOD ESTIMATE (ML)

Gaussian multivariate (unknown $\vec{\mu}$) $\longrightarrow \theta = \vec{\mu}$ ( $\Sigma$ is known)

likelihood function $P(\vec{x} \mid \vec{\mu}) = N(\vec{x} \mid \vec{\mu}, \Sigma)$

define log-likelihood $\mathcal{L} = \ln P(\vec{x}|\vec{\mu})$     $N(\vec{x}|\vec{\mu},\bar{\bar{\Sigma}}) = \dfrac{1}{(2\pi)^{d}|\bar{\bar{\Sigma}}|^{+\frac{1}{2}}}$

$$\mathcal{L} = -\frac{1}{2}\ln\left[(2\pi)^{d}\cdot|\bar{\bar{\Sigma}}|\right] - \frac{1}{2}\underbrace{(\vec{x}-\vec{\mu})^{T}\cdot\bar{\bar{\Sigma}}^{-1}\cdot(\vec{x}-\vec{\mu})}$$

mahalanobis distance

( quadratic weighted with variance in each direction)

<u>Maximum likelihood :</u>

$$\partial_{\vec{\mu}}\ln P(\vec{x}|\vec{\mu}) = \partial_{\mu}\mathcal{L} =$$

$$= \bar{\bar{\Sigma}}^{-1}\cdot(\vec{x}-\vec{\mu}) \longrightarrow \vec{\mu}_{j} = \frac{1}{n}\sum_{i=1}^{n}x_{i}^{j}$$

# LATENT VARIABLE MODELS → hidden variable

GMM: $p(\vec{x}) = \sum\limits_{k=1}^{G} \pi_k \cdot N(\vec{x} | \vec{\mu_k}, \vec{\Sigma_k})$

LATENT
VARIABLE : $\vec{z} = (z_1, z_2, \ldots z_G)$    $z_j \in \{0, 1\}$
(clustering
variable)

one-hot encoding

$\vec{x}$ belongs to cluster 3 → $\vec{z} = (0, 0, 1, 0, 0 \ldots, 0)$

$$P(Z_K = 1) = \pi_K$$

joint distribution $\quad P(\vec{x}, \vec{z}) = P(\vec{x}|\vec{z}) P(\vec{z})$

$\uparrow$

latent

conditional probability $\quad P(\vec{z}|\vec{x}) \rightarrow$ posterior

$$\underset{\text{posterior}}{P(\vec{z}|\vec{x})} = \frac{\overset{\text{likelihood}}{P(\vec{x}|\vec{z})} \overset{\text{prior}}{P(\vec{z})}}{\underset{\text{normaliz. constant}}{P(\vec{x})}}$$

# EXPECTATION - MAXIMIZATION ALGORITHM: (EM -algorithm)

Iterative procedure to get the ML estimate of parameters in probabilitsic models with latent variables.
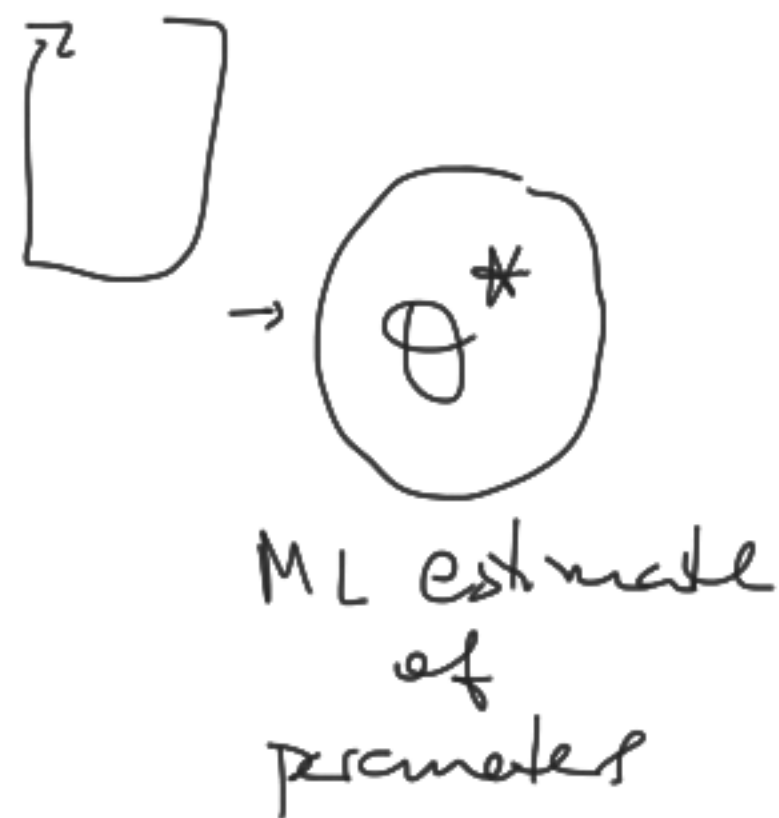
GMM + EM $\longrightarrow$ start from initial value for the parameters

$$\underline{\theta}^0 = \{ \pi_j^0, \vec{\mu}_j^0, \bar{\bar{\Sigma}}_j^0 \}$$

① E-STEP: Evaluate posterior $P(\vec{z} | \vec{x}, \underline{\theta}^0)$

② M-STEP: $\theta_{new} = \underset{\theta}{\arg\max} \; Q(\theta, \theta_0)$

expect. $Q(\theta, \theta_0) = \sum_{\vec{z}} P(\vec{z} | \vec{x}, \theta^0) \cdot \ln P(\vec{x}, \vec{z} | \theta)$

$\begin{bmatrix} \vec{z} \end{bmatrix} \longrightarrow \boxed{\theta^*}$

ML estimate
of
parameters

# fitting a GMM:

set
of
observations

$$\{ \vec{x}_1, \vec{x}_2, \cdots \vec{x}_n \}$$

$\downarrow$

EM-algorithm

$\downarrow$

$$\pi_j^*, \vec{\mu}_j^*, \vec{\Sigma}_j^* \qquad j = 1 \cdots G$$

$\downarrow$

$$p(\vec{x}) = \sum_{j=1}^{G} \pi_j^* \; N(\vec{x} \mid \vec{\mu}_j, \vec{\Sigma}_j) \longrightarrow$$

$$\boxed{\vec{z}}$$

assign each
obs.
to a cluster.

# GMM - options

→ Randomized initial values for $\vartheta$

→ Choice for the structure of the covariance matrix:

$$\bar{\Sigma} = \begin{pmatrix} \sigma & 0 \\ 0 & \sigma \end{pmatrix}$$

$$\bar{\Sigma} = \begin{pmatrix} \sigma_{11} & \sigma_{12} \\ \sigma_{21} & \sigma_{22} \end{pmatrix}$$

full-covariance

$$\bar{\Sigma} = \begin{pmatrix} \sigma_1 & 0 \\ 0 & \sigma_2 \end{pmatrix}$$

diagonal

tied

$$\bar{\Sigma}_1 \quad \bar{\Sigma}_2$$

$$\bar{\Sigma}$$

$x_1$

$x_2$