# PRINCIPAL COMPONENT ANALYSIS (PCA)

→ Numerical data

Data Matrix $\quad A$
$\qquad N_{obsev} \times N_{features}$

$$= \begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \\ \circ \end{pmatrix} \quad \begin{matrix} obs\,1 \\ \vdots \\ obs\,N \end{matrix}$$

$$f_1 \quad \cdots \quad \cdots \quad f_M$$

M-dimensional data

$\downarrow$ PCA

d-dimensional projection $\quad (d \leq\leq M)$

Ex: $M = 2$

$(f_1, f_2)$

$\downarrow$ PCA

$d = 1$



$f_2$

$\vec{v_1}$ : 1st eigenvector $\rightarrow$ direction of largest variability $\rightarrow$ 1st principal component

$\vec{v_2}$ : eigenvector 2nd principal component

projection of data into 1st PC

$\lambda_1 > \lambda_2$

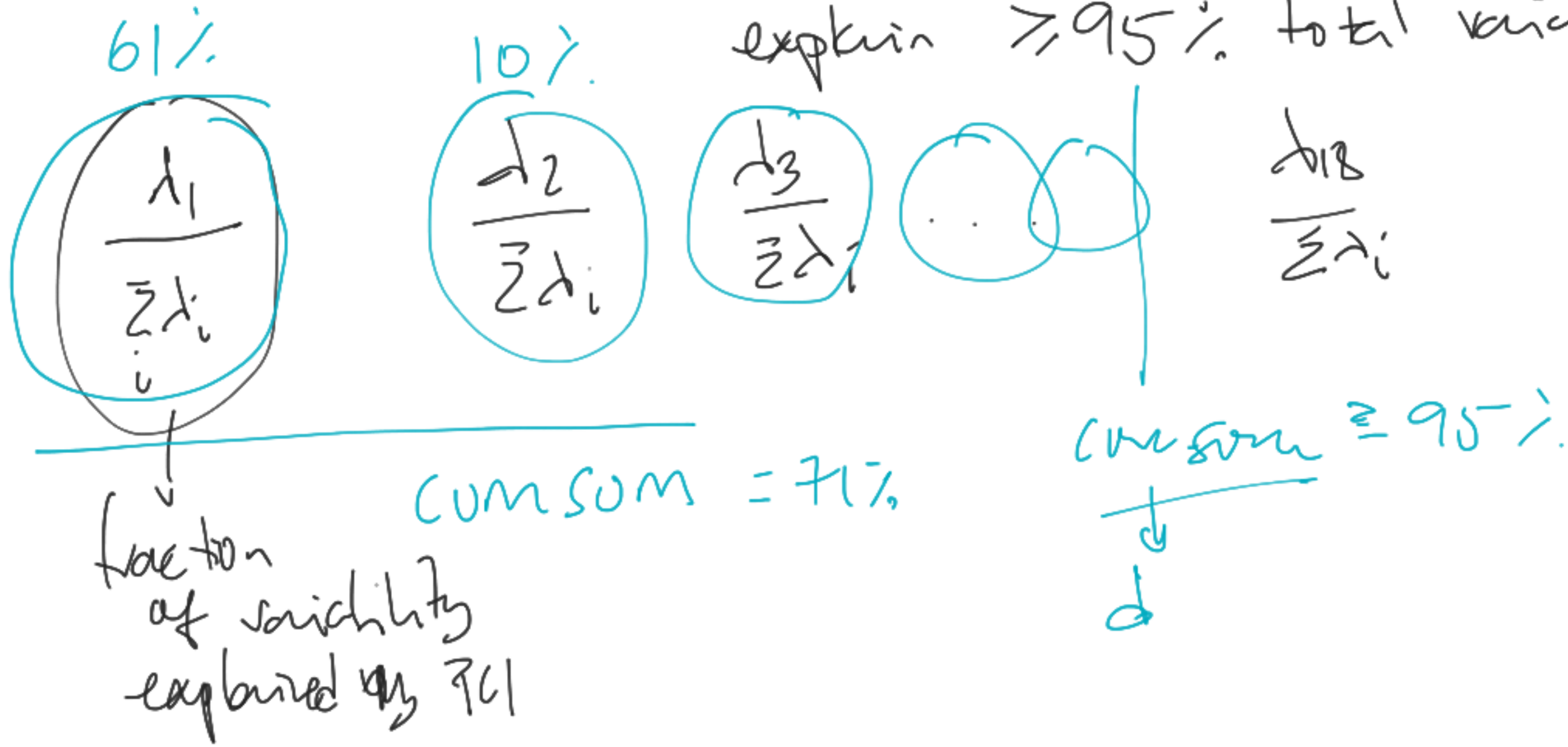$f_1$

directions of the PC's

diagonalize the data covariance matrix $C$

$\vec{v_1}, \lambda_1$   $\vec{v_2}, \lambda_2$   variability explained by each PC

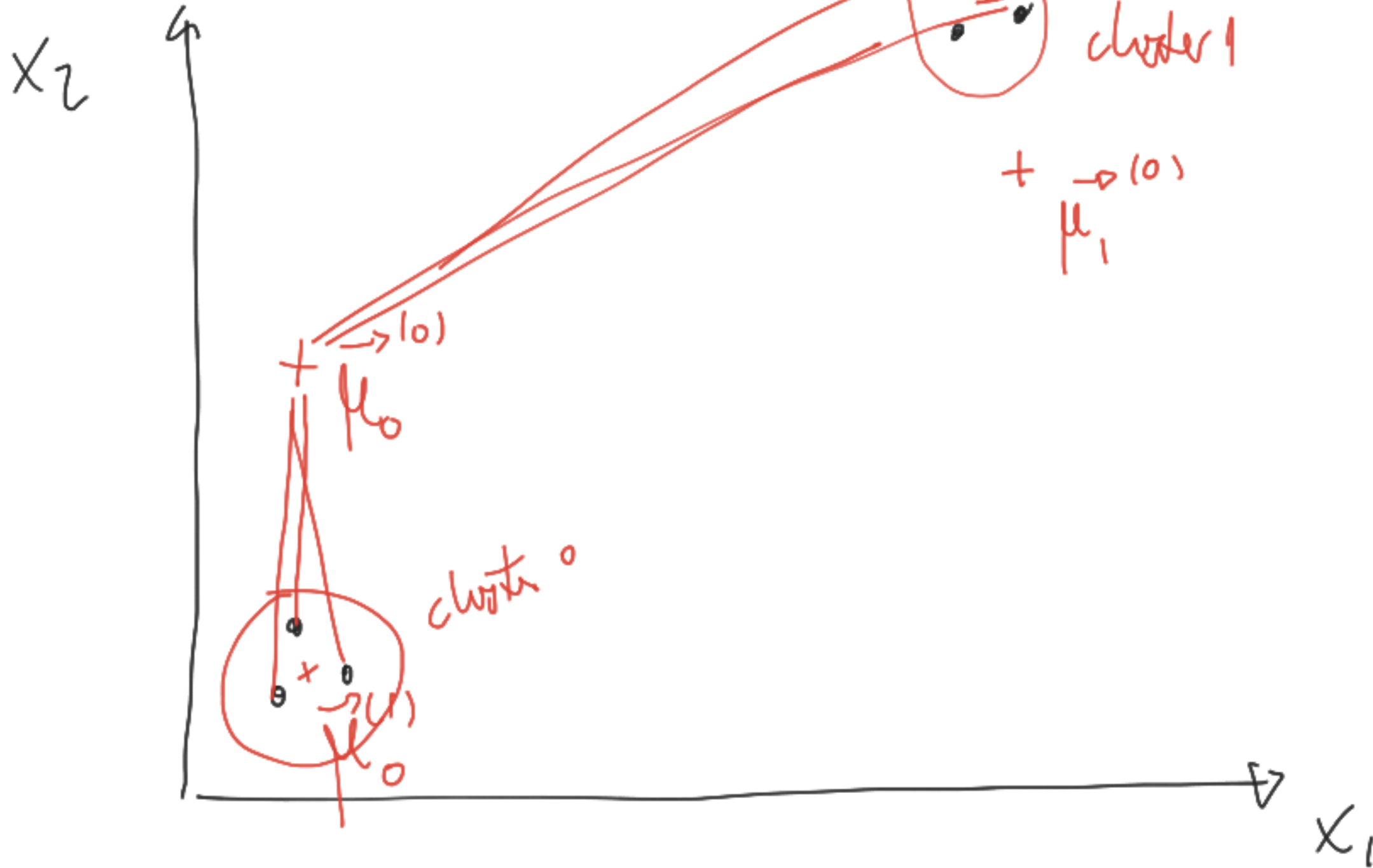$$M = 18 \xrightarrow{\text{PCA}} \vec{v}_i, \lambda_i, \quad i = 1 \ldots 18$$

How do I know how many to retain? (d)

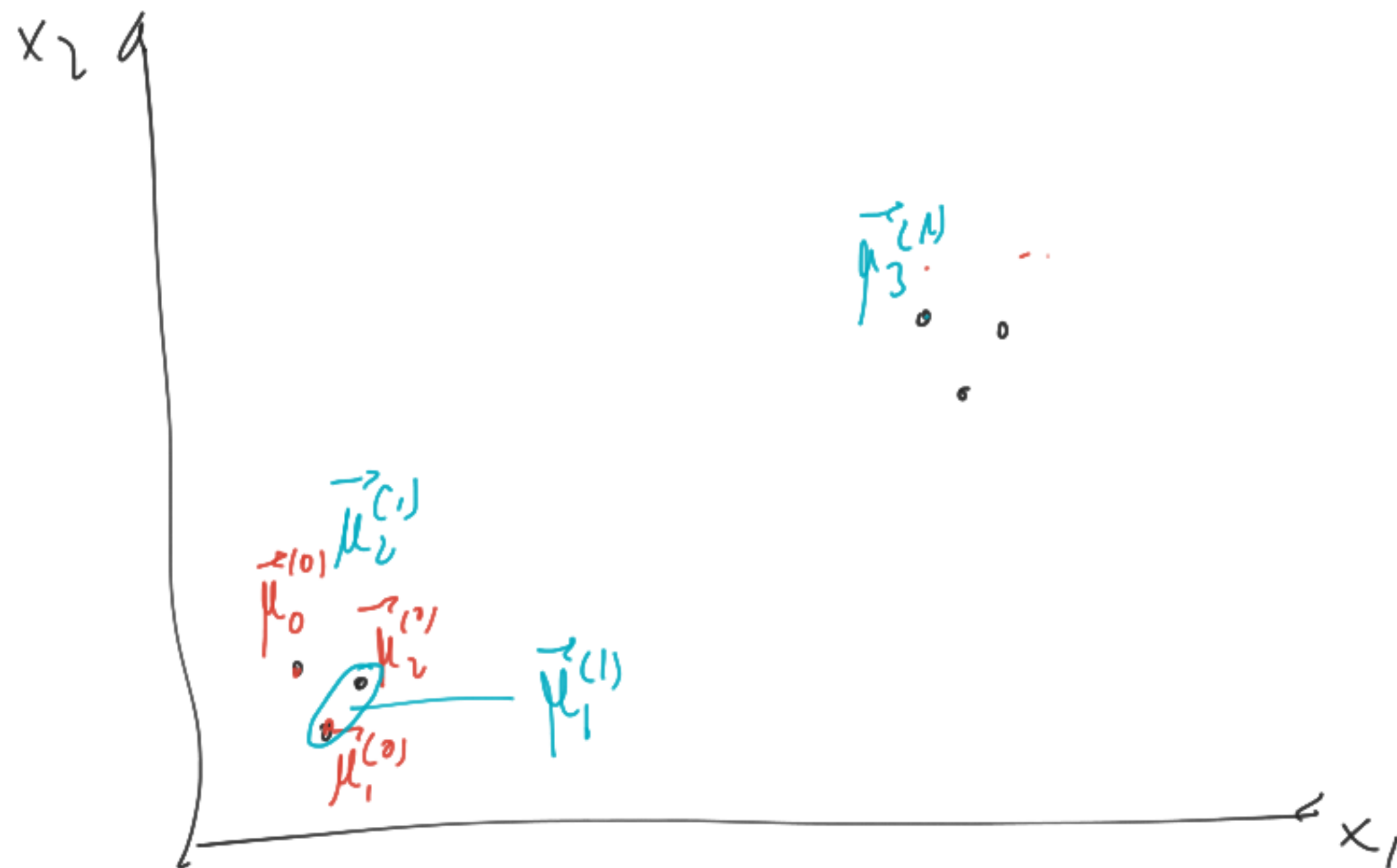rule: keep as many PC's as needed to explain $\geq 95\%$ total variability data

61%

$$\frac{\lambda_1}{\sum_i \lambda_i}$$
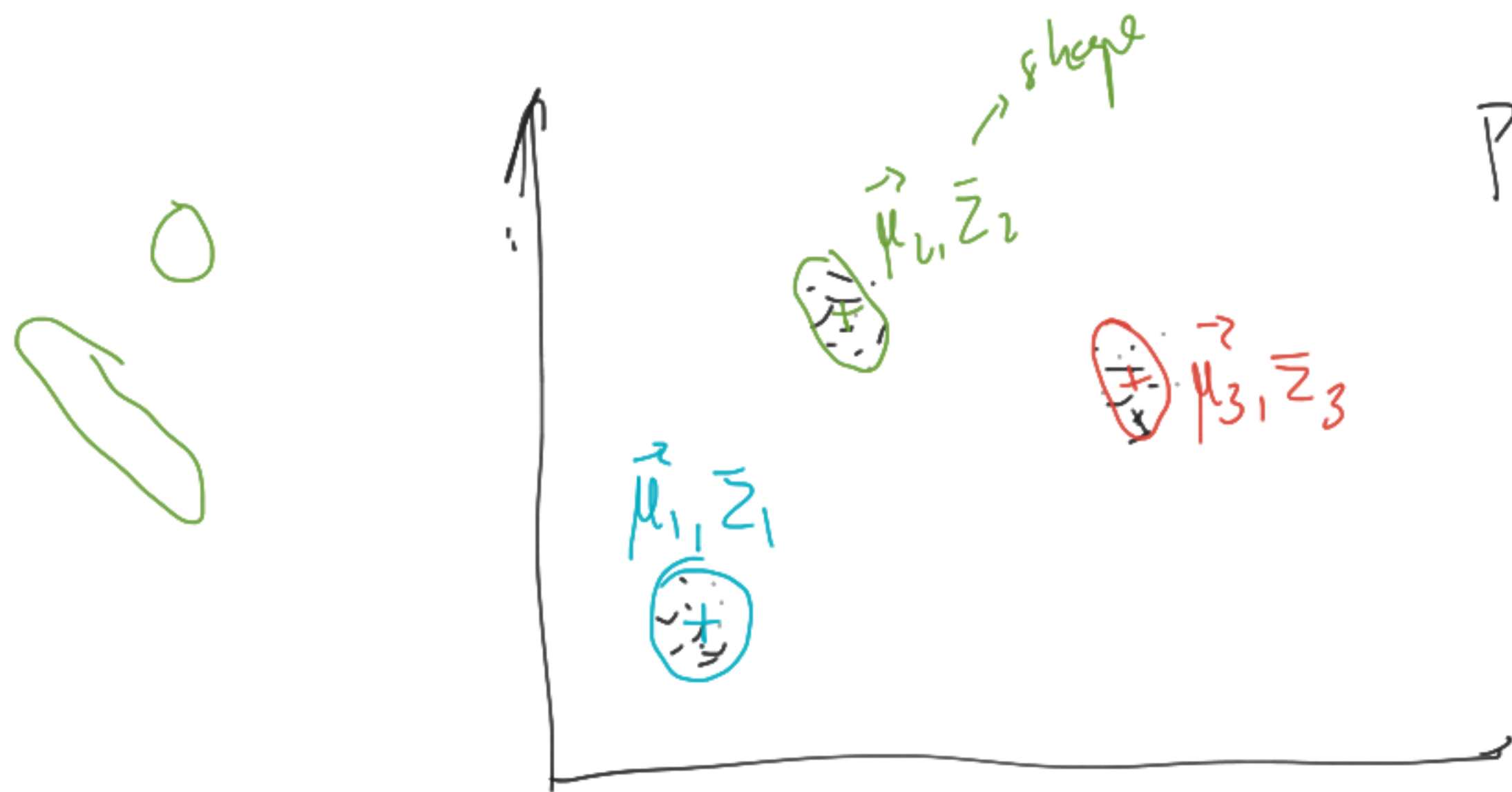
10%

$$\frac{\lambda_2}{\sum \lambda_i}$$

$$\frac{\lambda_3}{\sum \lambda_i}$$

$\cdots$

$$\frac{\lambda_{18}}{\sum \lambda_i}$$

fraction of variability explained by PC1

cumsum = 71%

cumsum $\geq 95\%$

$$\downarrow$$

$$d$$

CLUSTERING : K-means

$k : \# clusters \rightarrow 2$



$X_2$

$X_1$

$\vec{\mu}_1^{(1)}$

cluster 1

$+ \vec{\mu}_1^{(0)}$

$\vec{\mu}_0^{(0)}$

$\mu_0$

cluster 0

$\vec{\mu}_0^{(1)}$

0

# Agglomerative Clustering

# GMM — model statistics

$$P(\vec{x}) = \sum_{i=1}^{G} \pi_i \cdot N(\vec{\mu_i}, \Sigma_i)$$

shape

$\vec{\mu_2}, \Sigma_2$

$\vec{\mu_3}, \Sigma_3$

$\vec{\mu_1}, \Sigma_1$

# Model selection → Parsimony criterion: error vs complexity

y ↑

linear model $y = a + bx$
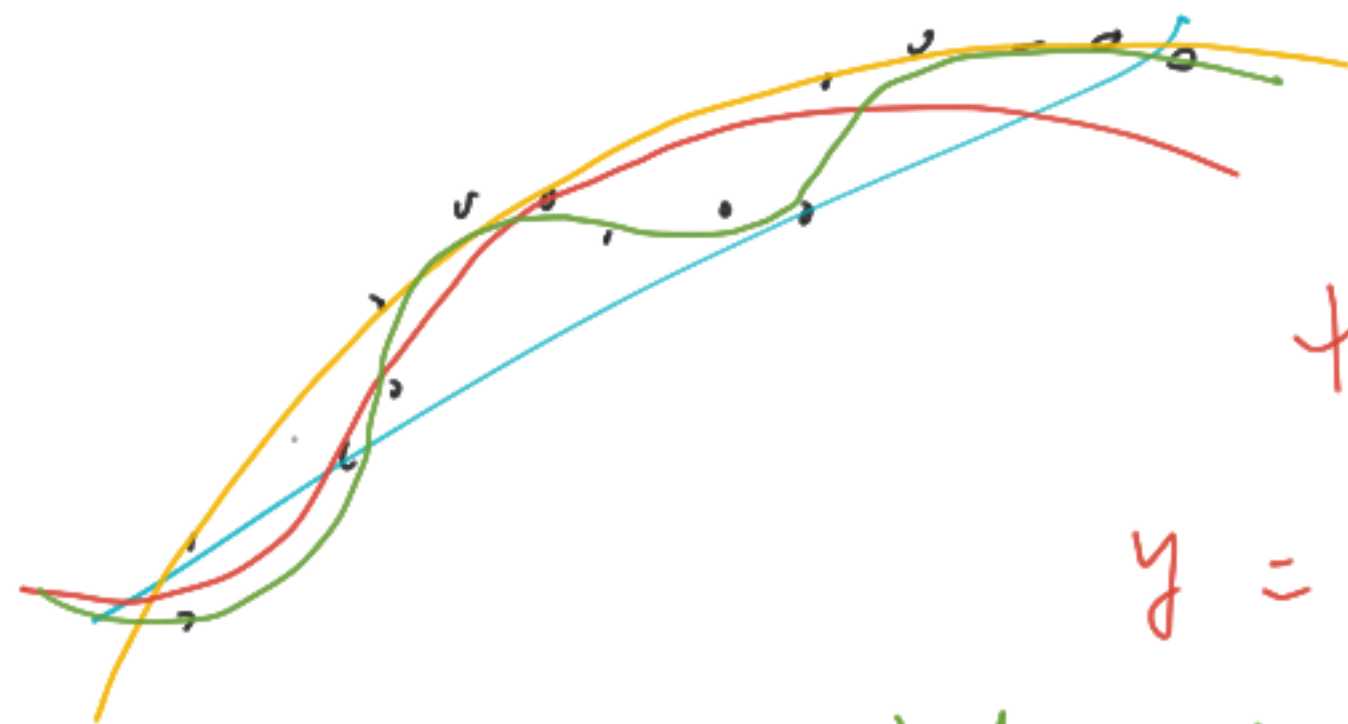
quadratic
$$y = a + bx + cx^2$$

third order model

$$y = a + \cdots$$

4th order

error → sum of
residuals
(MSE)

complexity → # parameters

→ x

| MODEL | MSE | #parameters | PARSIMONY INDEX |
|-------|-----|-------------|-----------------|
| 1st | 10 | 2 | 70 |
| 2nd | 1.6 | 3 | 60 → Best! |
| 3rd | 0.2 | 4 | 63 |
| 4th | 0.01 | 5 | 68 |

# SUPERVISED CLASSIFICATION

data

class
labels

$A_{Nobs \times Mfeatures}$

$\vec{w}$

$Nobs \times 1$

TRAIN MODEL



$f_1 \cdots f_M$     $\vec{w}$

obs1

70%

$w_1$
$w_2$
$\vdots$
$w_N$

obs N

30%

Ground truth

TEST MODEL

TRAIN A LDA classifier: training data

$P(\vec{x}|w=1) = N(\vec{\mu}_1, \bar{\Sigma}_1)$

- class $w = 1$
- class $w = 2$

$\vec{\mu}_1$    $d_1$    $\vec{x}_{new}$

$\bar{\Sigma}_1$    $d_2$    $\vec{\mu}_2$

$P(\vec{x}|w=2) = N(\vec{\mu}_2, \bar{\Sigma}_2)$

$\bar{\Sigma}_2$

sample
estimates
of
$\vec{\mu}_1, \vec{\mu}_2$
$\bar{\Sigma}_1, \bar{\Sigma}_1$

from
the training
samples.

KNN