

# Intro to ML

October 20<sup>th</sup>, 2021

CHAPTER 6:

# Dimensionality Reduction

# Why Reduce Dimensionality?

- Reduces time complexity: Less computation
- Reduces space complexity: Fewer parameters
- Saves the cost of observing the feature
- Simpler models are more robust on small datasets
- More interpretable; simpler explanation
- Data visualization (structure, groups, outliers, etc) if plotted in 2 or 3 dimensions


# Feature Selection vs Extraction

- **Feature selection:** Choosing  $k < d$  important features, ignoring the remaining  $d - k$

## Subset selection algorithms

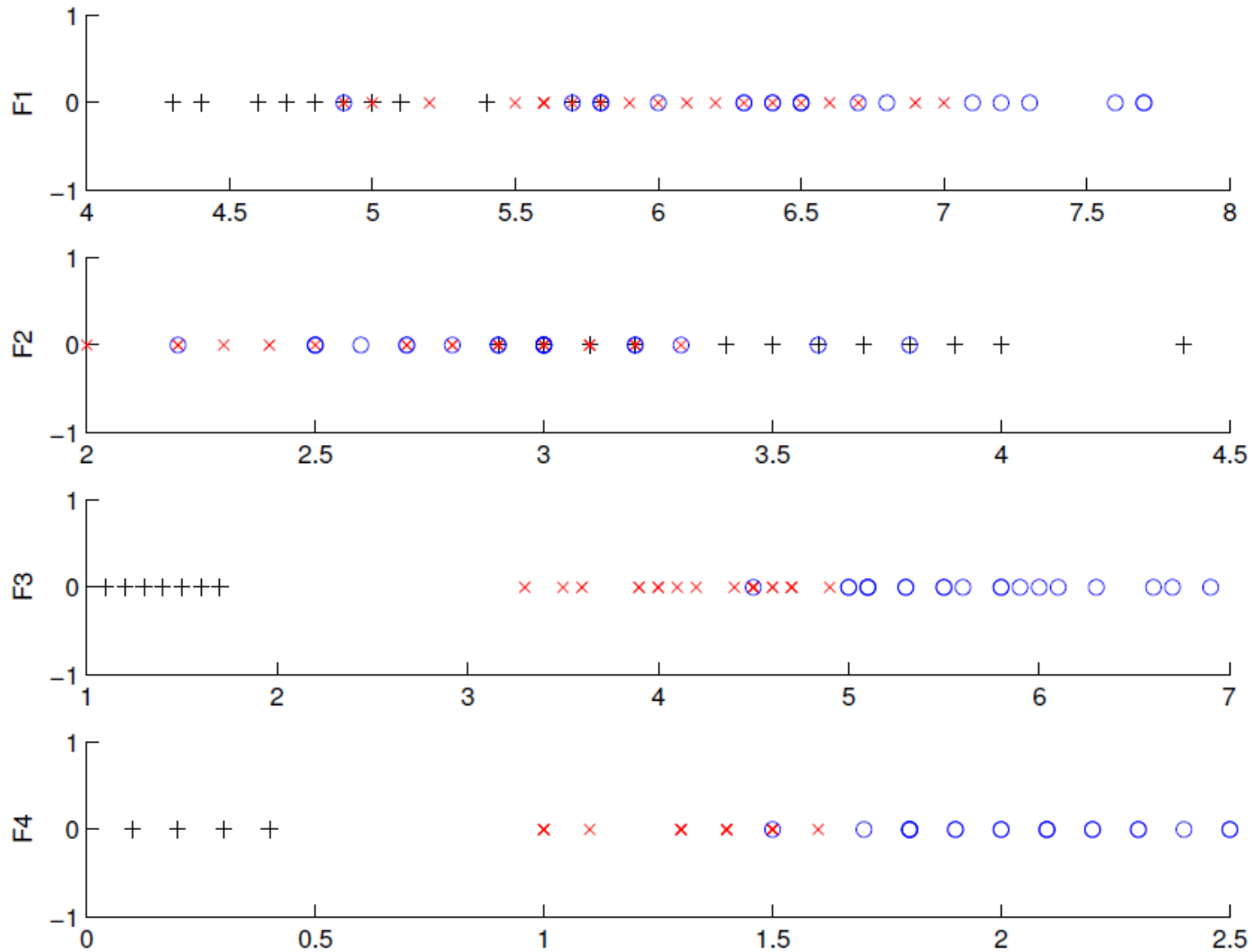
- **Feature extraction:** Project the original  $x_i, i = 1, \dots, d$  dimensions to new  $k < d$  dimensions,  $z_j, j = 1, \dots, k$

# Subset Selection

- There are  $2^d$  subsets of  $d$  features
  - Exhaustive search is not possible
- **Forward search:** Add the best feature at each step
  - Set of features  $F$  initially  $\emptyset$ .
  - At each iteration, find the best new feature
$$j = \operatorname{argmin}_i E ( F \cup x_i )$$


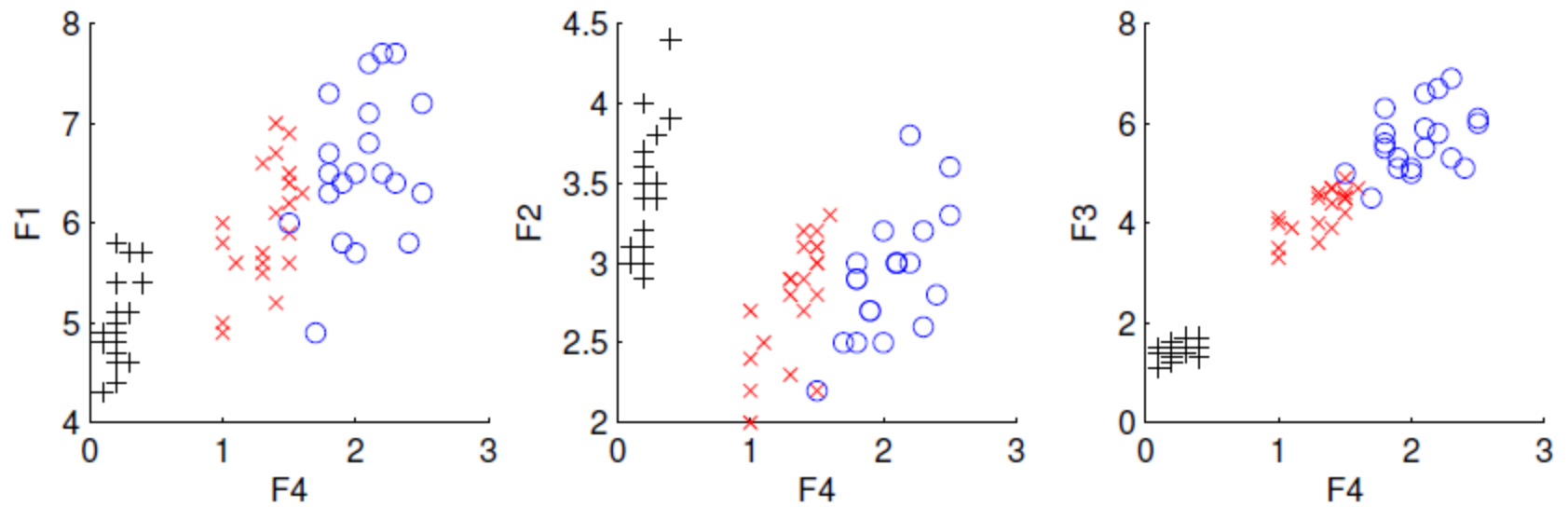
Use train, validation, test setup
  - Add  $x_j$  to  $F$  if  $E ( F \cup x_j ) < E ( F )$
- **Backward search:** Start with all features and remove one at a time, if possible.
- Hill-climbing  $O(d^2)$  algorithm
- Floating search (Add  $k$ , remove  $l$ ,  $k > l$ )

## Iris data: Single feature



Chosen

Iris data: Add one more feature to F4



Chosen

# Principal Components Analysis

- Find a low-dimensional space such that when  $\mathbf{x}$  is projected there, information loss is minimized.
- The projection of  $\mathbf{x}$  on the direction of  $\mathbf{w}$  is:  $z = \mathbf{w}^T \mathbf{x}$
- Find  $\mathbf{w}$  such that  $\text{Var}(z)$  is maximized information

$$\begin{aligned}\text{Var}(z) &= \text{Var}(\mathbf{w}^T \mathbf{x}) = E[(\mathbf{w}^T \mathbf{x} - \mathbf{w}^T \boldsymbol{\mu})^2] \\ &= E[(\mathbf{w}^T \mathbf{x} - \mathbf{w}^T \boldsymbol{\mu})(\mathbf{w}^T \mathbf{x} - \mathbf{w}^T \boldsymbol{\mu})] \\ &= E[\mathbf{w}^T (\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T \mathbf{w}] \\ &= \mathbf{w}^T E[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T] \mathbf{w} = \mathbf{w}^T \boldsymbol{\Sigma} \mathbf{w}\end{aligned}$$

where  $\text{Var}(\mathbf{x}) = E[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T] = \boldsymbol{\Sigma}$



- First projection: Maximize  $\text{Var}(z_1)$  subject to  $||\mathbf{w}_1||=1$

$$\max_{\mathbf{w}_1} \mathbf{w}_1^T \Sigma \mathbf{w}_1 - \alpha (\mathbf{w}_1^T \mathbf{w}_1 - 1) \quad \leftarrow \text{A Lagrange problem}$$

$\Sigma \mathbf{w}_1 = \alpha \mathbf{w}_1$  that is,  $\mathbf{w}_1$  is an eigenvector of  $\Sigma$

Choose the one with the largest eigenvalue for  $\text{Var}(z)$  to be max

- Second principal component: Max  $\text{Var}(z_2)$ , s.t.,  $||\mathbf{w}_2||=1$  and orthogonal to  $\mathbf{w}_1$

$$\max_{\mathbf{w}_2} \mathbf{w}_2^T \Sigma \mathbf{w}_2 - \alpha (\mathbf{w}_2^T \mathbf{w}_2 - 1) - \beta (\mathbf{w}_2^T \mathbf{w}_1 - 0)$$

$\Sigma \mathbf{w}_2 = \alpha \mathbf{w}_2$  that is,  $\mathbf{w}_2$  is another eigenvector of  $\Sigma$  with the second largest eigenvalues and so on

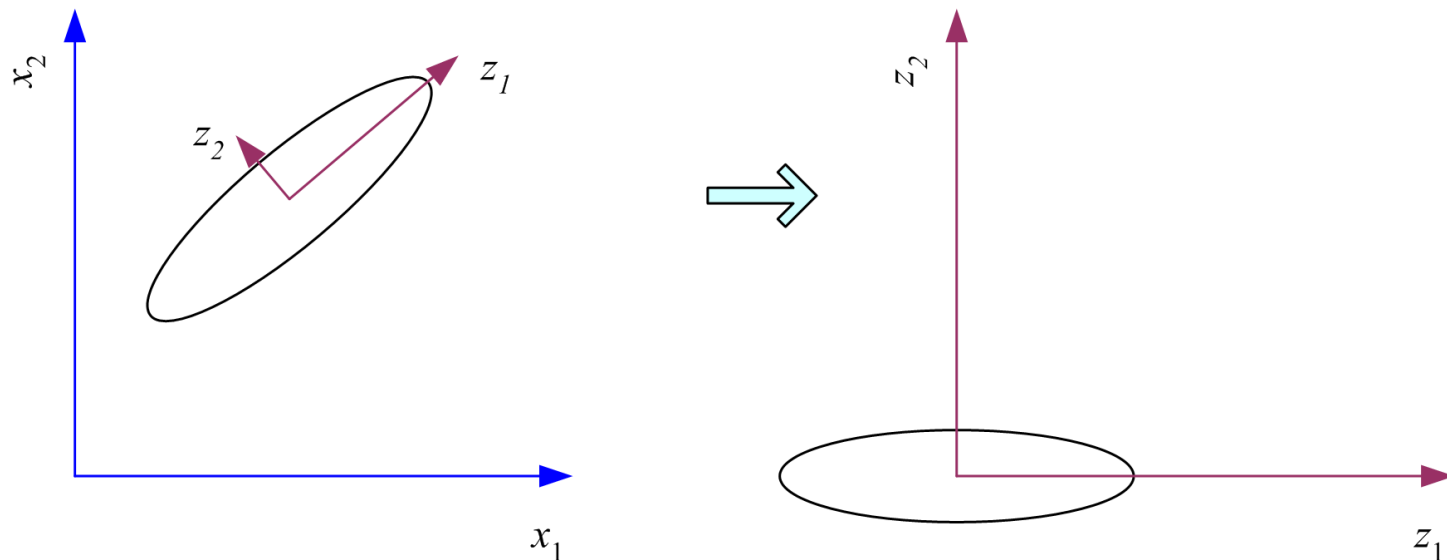
- $\Sigma$  is symmetric  $\rightarrow$  two eigenvectors orthogonal
- $\Sigma$  is positive definite  $\rightarrow$  all eigenvalues are positive
- $\Sigma$  is singular  $\rightarrow$  it has a lower rank (dimension require  $k < d$ )

# What PCA does

$$\mathbf{z} = \mathbf{W}^T(\mathbf{x} - \mathbf{m})$$

where the columns of  $\mathbf{W}$  are the eigenvectors of  $\Sigma$  and  $\mathbf{m}$  is sample mean

Centers the data at the origin and rotates the axes



# How to choose k ?

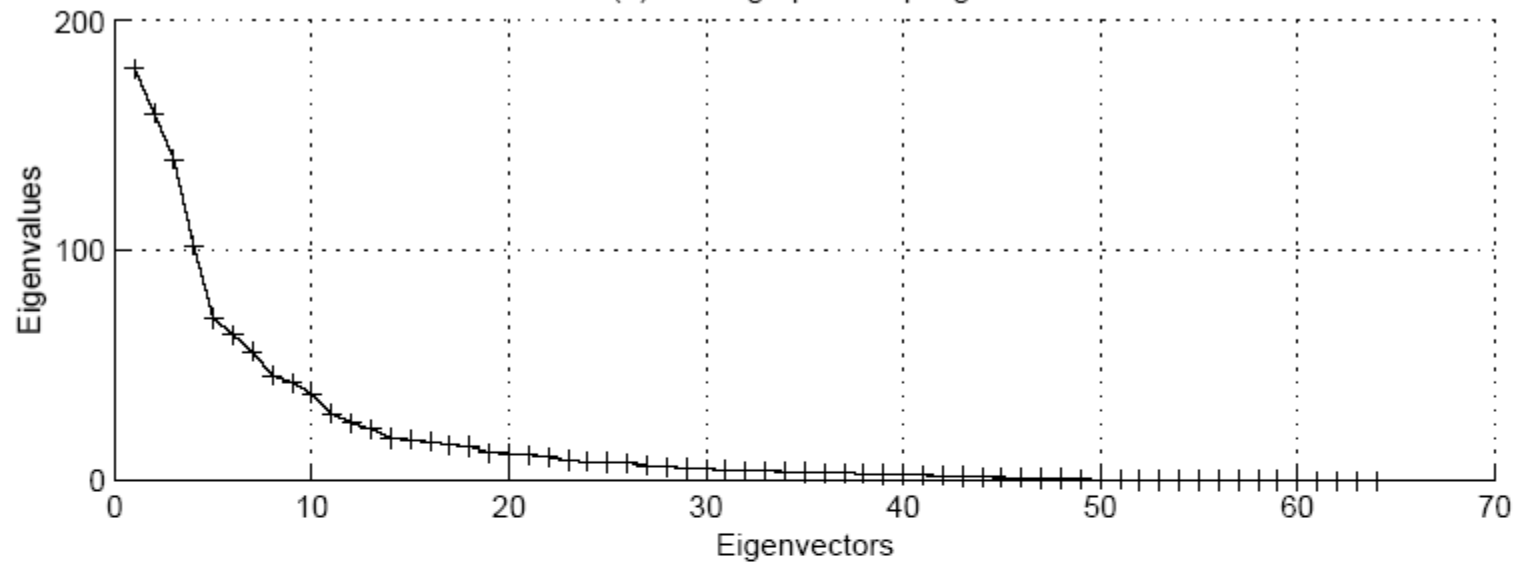
- Proportion of Variance (PoV) explained

$$\frac{\lambda_1 + \lambda_2 + \dots + \lambda_k}{\lambda_1 + \lambda_2 + \dots + \lambda_k + \dots + \lambda_d}$$

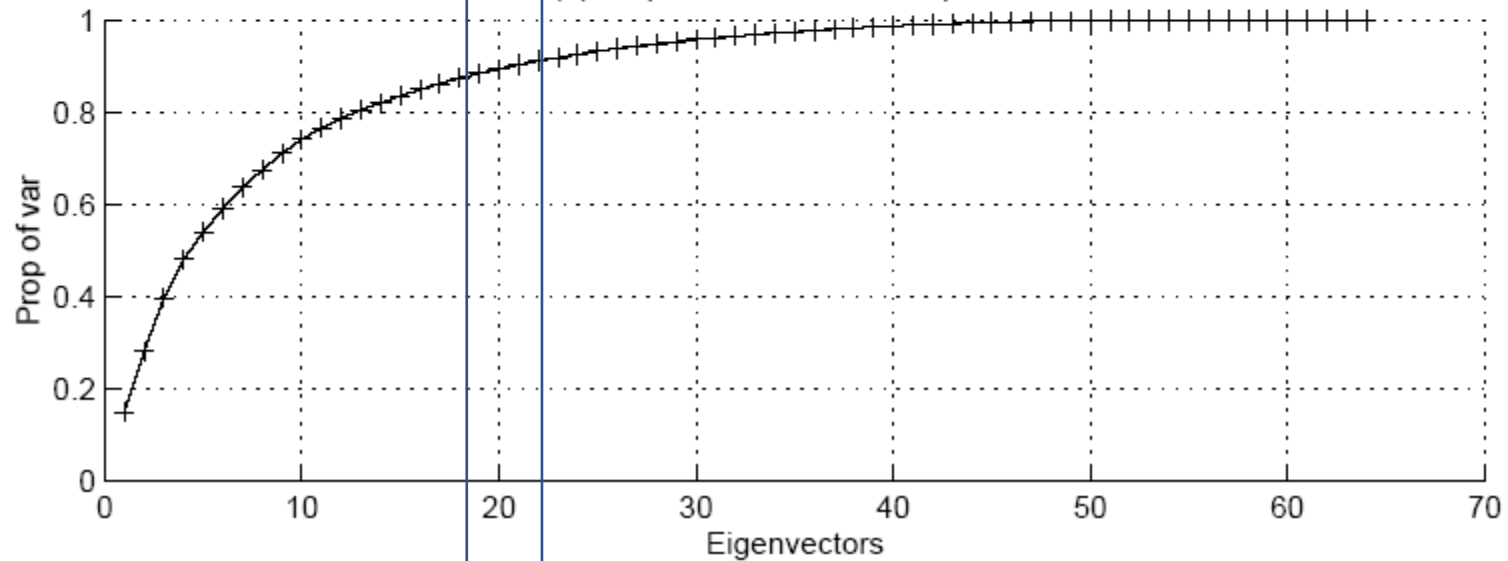
when  $\lambda_i$  (eigenvalues) are sorted in descending order

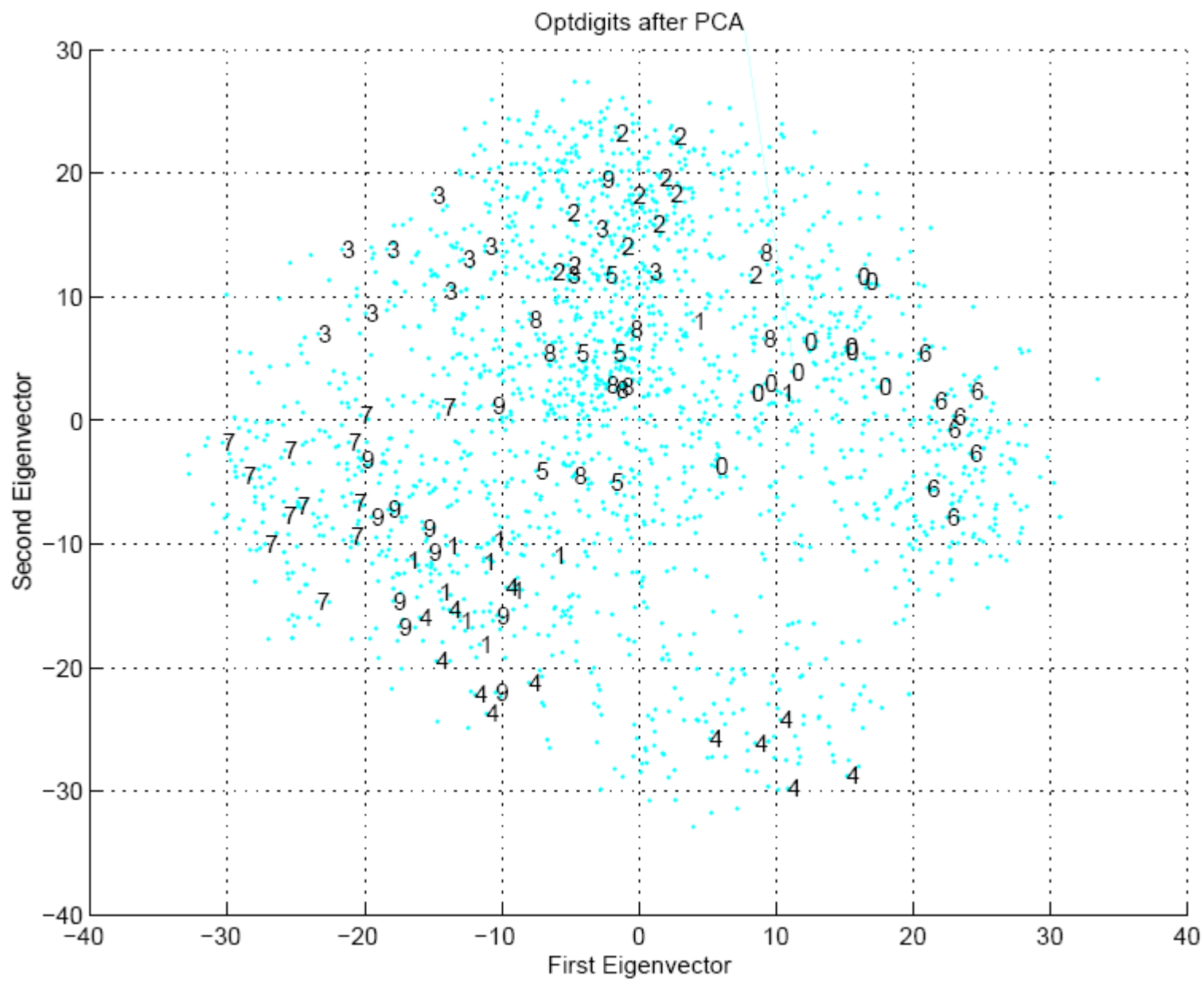
- Typically, stop at  $\text{PoV} > 0.9$
- Scree graph plots of PoV vs  $k$ , stop at “elbow”
  - *adding another eigenfactor does not add much variances*

(a) Scree graph for Optdigits



(b) Proportion of variance explained





Visual display

# After PCA

- From  $d$ -dimension to  $k$ -dimension
  - $K < d$ , parametric discriminant function rely on fewer parameters (due to lesser features)
  - Since these features are projected to be uncorrelated (orthogonal is you assume Normal distribution)
    - The multivariate Gaussian's covariance matrix can now be assumed to be diagonal
- PCA is an unsupervised method of dimension reduction -> does not require the knowledge of label, just the data

# Another formulation of PCA

- Find a matrix  $W$  such that when we have  $z=W^T x$ , we will get  $\text{cov}(z) = D$ , where  $D$  is any diagonal matrix (uncorrelated  $z_i$ )
- Form a  $d \times d$  matrix  $C$ , whose  $i$ th column is the normalized eigenvectors  $c_i$  of  $S$ ,  $C^T C = I$

$$S = S C C^T = S(c_1, c_2, c_3 \dots c_d) C^T = (\lambda_1 c_1, \lambda_2 c_2 \dots \lambda_d c_d) C^T = C D C^T$$

Where  $D$  = diagonal matrix with eigenvalues on the diagonal

- We also call this spectral decomposition of  $S$
- $C^T S C = D$
- $z=W^T x$ ,  $\text{Cov}(z) = W^T S W$ , want this to be a diagonal matrix, just set  $W=C$

# Factor Analysis

- Find a small number of unobservable factors  $\mathbf{z}$ , which when combined generate  $\mathbf{x}$ :

$$x_i - \mu_i = v_{i1}z_1 + v_{i2}z_2 + \dots + v_{ik}z_k + \varepsilon_i$$

where  $z_j, j=1, \dots, k$  are the latent factors with

$$E[z_j]=0, \text{Var}(z_j)=1, \text{Cov}(z_i, z_j)=0, i \neq j,$$

$\varepsilon_i$  are the **noise sources**

$$E[\varepsilon_i]=0, \text{Var}(\varepsilon_i)=\psi_i, \text{Cov}(\varepsilon_i, \varepsilon_j)=0, i \neq j, \text{Cov}(\varepsilon_i, z_j)=0$$

and  $v_{ij}$  are the **factor loadings**

Factors are:  
Unit normal and  
uncorrelated

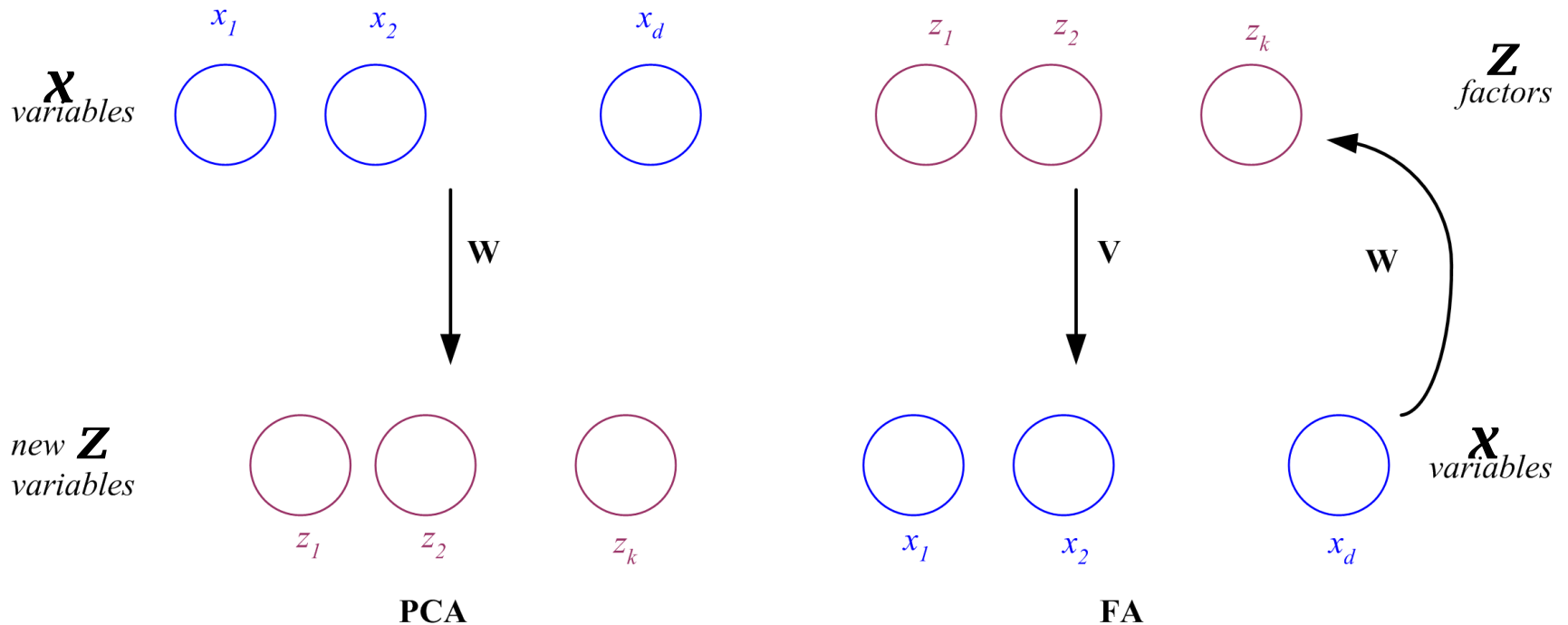


# PCA vs FA

- PCA From  $\mathbf{x}$  to  $\mathbf{z}$
- FA From  $\mathbf{z}$  to  $\mathbf{x}$

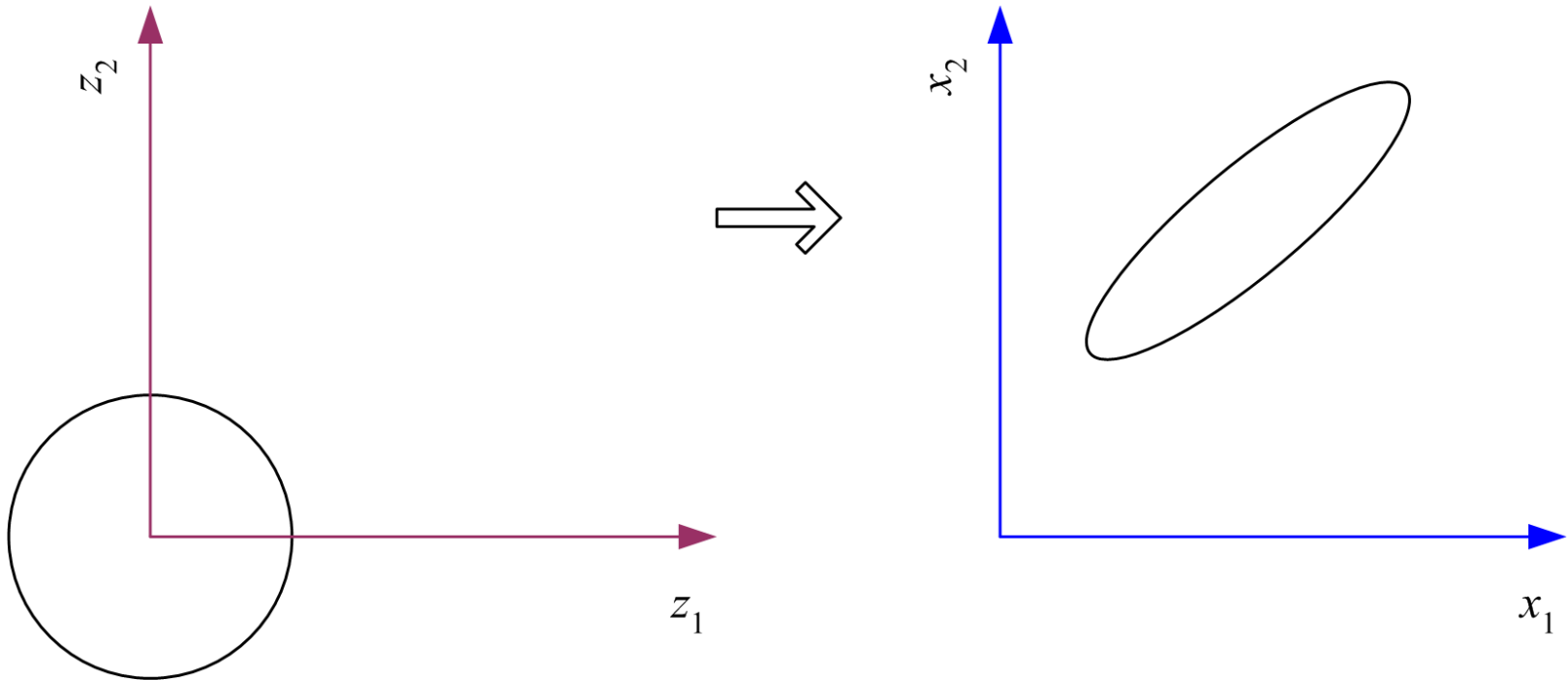
$$\mathbf{z} = \mathbf{W}^T(\mathbf{x} - \boldsymbol{\mu})$$

$$\mathbf{x} - \boldsymbol{\mu} = \mathbf{V}\mathbf{z} + \boldsymbol{\varepsilon}$$



# Factor Analysis

- In FA, factors  $z_j$  are stretched, rotated and translated to generate  $\mathbf{x}$



- $X_i - \mu_i = v_{i1}z_1 + v_{i2}z_2 + \dots + v_{ik}z_k + \varepsilon_i$
- $Z$

$$\text{Var}(X_i) = v_{i1}^2 + v_{i2}^2 + \dots + \psi_i$$

- Variance of  $X_i$  attributed to common factor 1-k and noise

$$\begin{aligned}\Sigma &= \text{cov}(X) = \text{cov}(Vz + \varepsilon) \\ &= \text{cov}(Vz) + \text{cov}(\varepsilon) \\ &= V \boxed{\text{cov}(z)} V^T + \psi \quad \text{Identity matrix} \\ &= VV^T + \boxed{\psi} \quad \text{Diagonal matrix}\end{aligned}$$

Assume there are two 'hidden factors', and two observable features

$$\text{Cov}(x_1, x_2) = v_{11}v_{21} + v_{12}v_{22}$$

- $\text{Cov}(x_1, x_2) = v_{11}v_{21} + v_{12}v_{22}$
- If  $x_1$   $x_2$  have high covariance due to the first 'factor'
  - The  $v_{11}v_{21}$  will be high or  $v_{12}v_{22}$  will be high
  - Either way the above terms will be high
- If they depend on different factors, one term high one term low, this summation will be small

$$\text{Cov}(x_1, z_2) = \text{Cov}(v_{12}z_2, z_2) = v_{12}$$

- $\text{Cov}(x, z) = V$
- Loading represent correlation between variables and the factors

• S estimate  $\Sigma$

• Factor analysis basically solves the following equation

$$\underline{S = VV^T + \Psi}$$

PCA: FIND THE  
EIGENVECTOR OF S

•  $z = W^T x$ ,  $\text{Cov}(z) = W^T S W$ , want this to be a diagonal matrix, just set  $W = C$