

# Dimensionality Reduction



आई आई टी हैदराबाद  
IIT Hyderabad

# ML Problems

	<i>Supervised Learning</i>	<i>Unsupervised Learning</i>
<i>Discrete</i>	classification or categorization	clustering
<i>Continuous</i>	regression	dimensionality reduction

# What is Dimensionality Reduction

- Refers to the mapping of the original high-dimensional data onto a lower-dimensional space.
- Criterion for feature reduction can be different for different problems.
  - Unsupervised setting: minimize the information loss
  - Supervised setting: maximize the class discrimination
- Given a set of data points of  $p$  variables the linear transformation (projection)

$$\boxed{G^T} \boxed{X} = \boxed{Y}$$

# Dimensionality Reduction vs Feature Selection

- **Feature reduction**

- All original features are used
- The transformed features are linear combinations of the original features (in case of linear DR).

- **Feature selection**

- Only a subset of the original features are used.

# Why DR?

- Most machine learning techniques may not be effective for high-dimensional data
  - **Curse of Dimensionality**
    - Query accuracy and efficiency degrade rapidly as the dimension increases
    - Lower space and time complexity
    - Visualization, Data compression, Noise/irrelevant feature removal
- The **intrinsic** dimension may be small
  - For example, the number of genes responsible for a certain type of disease may be small

# High-dimensional data are strange

- Consider the hypersphere of radius  $r$  on a space of dimension  $d$

$$\mathcal{S} = \left\{ \mathbf{x} \mid \sum_{i=1}^d x_i^2 \leq r^2 \right\}$$

- Its volume is

$$V_d(r) = \frac{r^d \pi^{\frac{d}{2}}}{\Gamma\left(\frac{d}{2} + 1\right)}$$

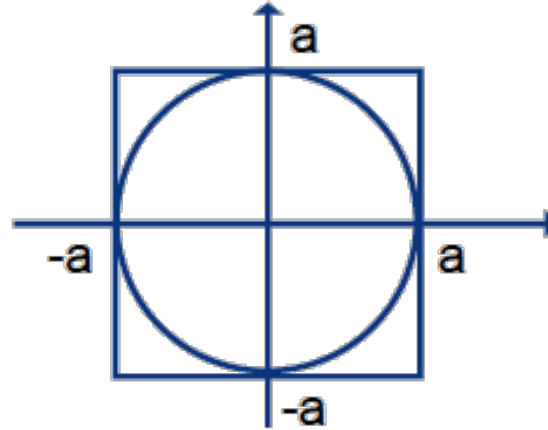
- Where  $\Gamma(n)$  is the Gamma function

$$\Gamma(n) = \int_0^{\infty} e^{-x} x^{n-1} dx$$

Source: N Vasconcelos

# Hypercube vs Hypersphere

- Consider the hyper-cube  $[-a, a]^d$  and the inscribed hyper-sphere, i.e.



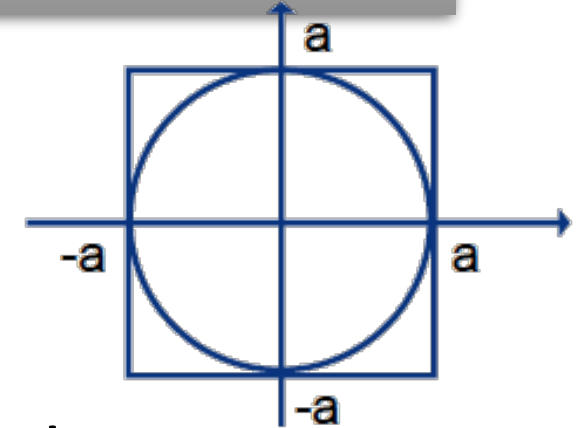
- What does your intuition tell you about the relative sizes of these two objects?
  - Volume of sphere  $\sim$  volume of cube
  - Volume of sphere  $\gg$  volume of cube
  - Volume of sphere  $\ll$  volume of cube

Source: N Vasconcelos

# Hypercube vs Hypersphere

- Let's compute the answer

$$f_d = \frac{Vol(sphere)}{Vol(cube)} = \frac{\frac{a^d \pi^{\frac{d}{2}}}{\Gamma(\frac{d}{2} + 1)}}{(2a)^d} = \frac{\pi^{\frac{d}{2}}}{2^d \Gamma(\frac{d}{2} + 1)}$$



- Sequence that does not depend on a, just on the dimension d!

d	1	2	3	4	5	6	7
f <sub>d</sub>	1	.785	.524	.308	.164	.08	.037

- It goes to zero, and goes to zero fast!

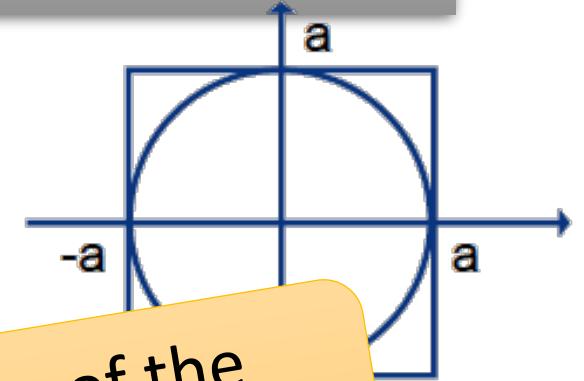
Source: N Vasconcelos



# Hypercube vs Hypersphere

- Let's compute the answer

$$f_d = \frac{Vol(sphere)}{Vol(cube)} = \frac{\frac{a^d \pi^{\frac{d}{2}}}{\Gamma(\frac{d}{2} + 1)}}{(2a)^d} = \frac{\pi^{\frac{d}{2}}}{2^d \Gamma(\frac{d}{2} + 1)}$$



- Sequence that does not depend on  $d$ !

As the dimension of the space increases, the volume of the sphere is much smaller (infinitesimal) than that of the cube!

5	6	7
.785	.524	.308
.164	.08	.037

- It goes to zero, and goes to zero fast!

Source: N Vasconcelos

# Hypercube vs Hypersphere

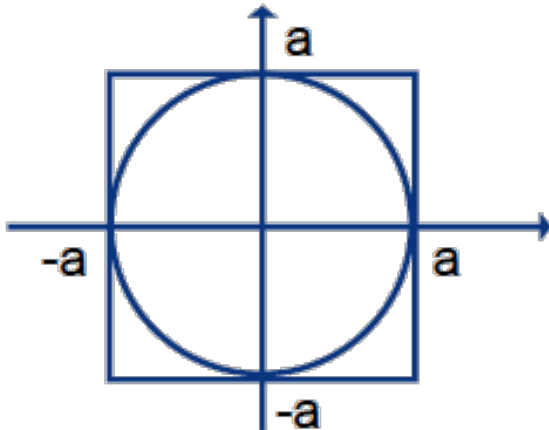
Actually not very surprising

1.  $d = 1$



Volume is the same

2.  $d = 2$

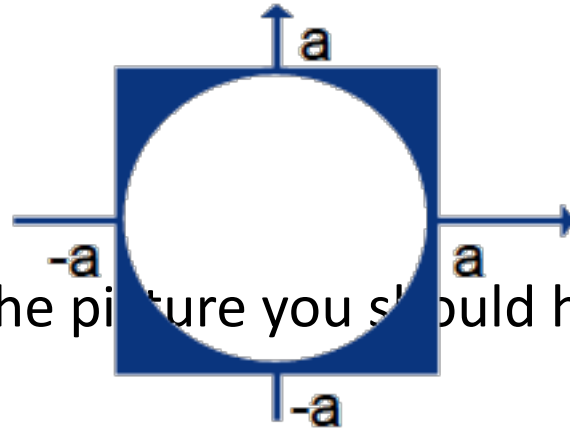


Volume of sphere is already smaller

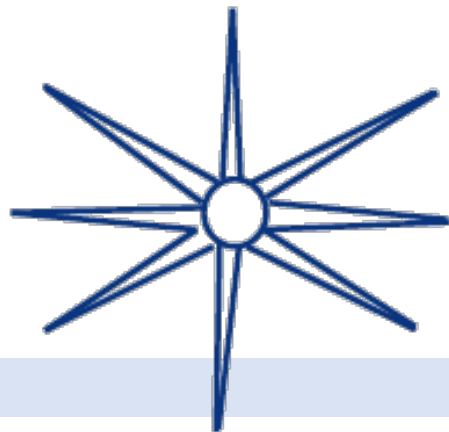
Source: N Vasconcelos

# Hypercube vs Hypersphere

- As the dimension increases the volume of the shaded corners becomes larger



- In high dimensions the picture you should have in mind is



all the volume of the cube  
is in these spikes!

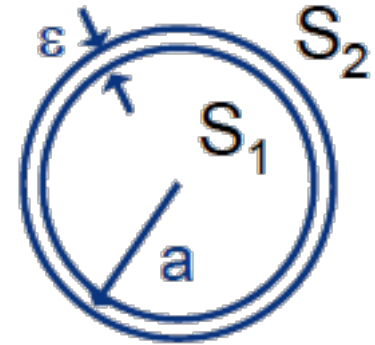
Source: N Vasconcelos



# The Curse of Dimensionality

- Consider the crust of unit hypersphere of thickness  $\epsilon$
- Let's compute the ratio of volumes

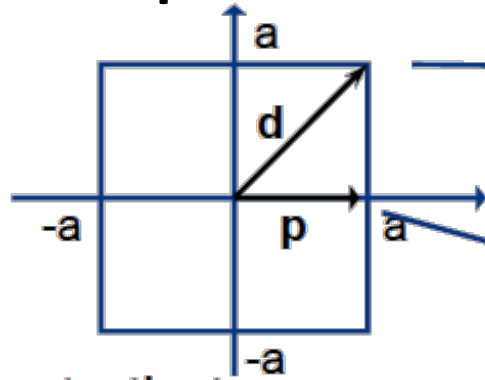
$$\frac{Vol(S_1)}{Vol(S_2)} = \frac{\frac{(a-\epsilon)^d \pi^{\frac{d}{2}}}{\Gamma(\frac{d}{2}+1)}}{\frac{a^d \pi^{\frac{d}{2}}}{\Gamma(\frac{d}{2}+1)}} = \frac{a^d \left(1 - \frac{\epsilon}{a}\right)^d}{a^d} = \left(1 - \frac{\epsilon}{a}\right)^d$$



- No matter how small  $\epsilon$  is, ratio goes to zero as  $d$  increases i. e. **“all the volume is in the crust!”**

# We can check mathematically

- Consider  $\mathbf{d}$  and  $\mathbf{p}$



$$\mathbf{d} = (a, a, \dots, a) \in \mathcal{R}^d$$

$$\mathbf{p} = (a, 0, \dots, 0) \in \mathcal{R}^d$$

- Note that

$$\frac{\|\mathbf{d}\|^2}{\|\mathbf{p}\|^2} = \frac{da^2}{a^2} = d \rightarrow \infty \quad \cos\theta = \frac{\mathbf{d}^T \mathbf{p}}{\sqrt{\|\mathbf{d}\|^2 \|\mathbf{p}\|^2}} \quad \lim_{d \rightarrow \infty} \frac{\text{dist}_{\max} - \text{dist}_{\min}}{\text{dist}_{\min}} \rightarrow 0$$
$$= \frac{a^2}{\sqrt{da^2 a^2}} = \frac{1}{\sqrt{d}} \rightarrow 0$$

- $\mathbf{d}$  orthogonal to  $\mathbf{p}$  as  $d$
- increases and infinitely larger!!!

Source: N Vasconcelos

# DR Methods

- **Linear**

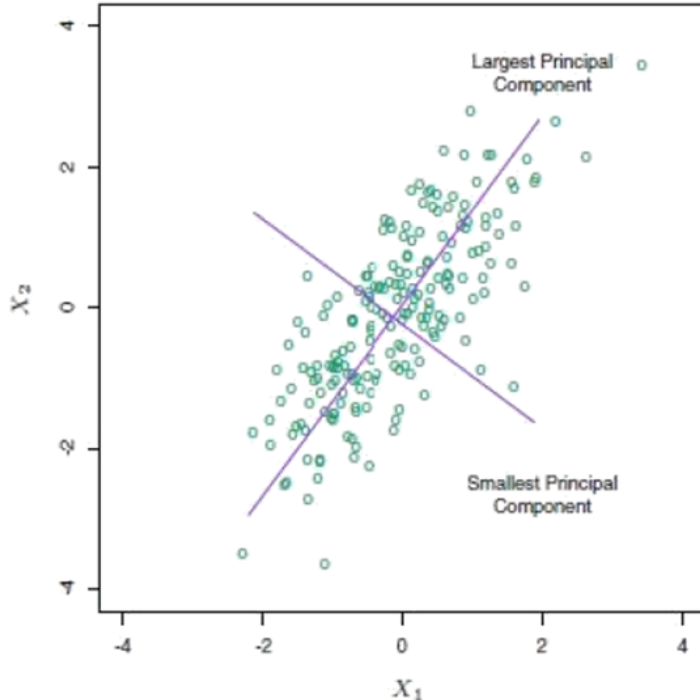
- Principal Component Analysis (PCA)
- Linear Discriminant Analysis (LDA)
- Canonical Correlation Analysis (CCA)

- **Non-linear**

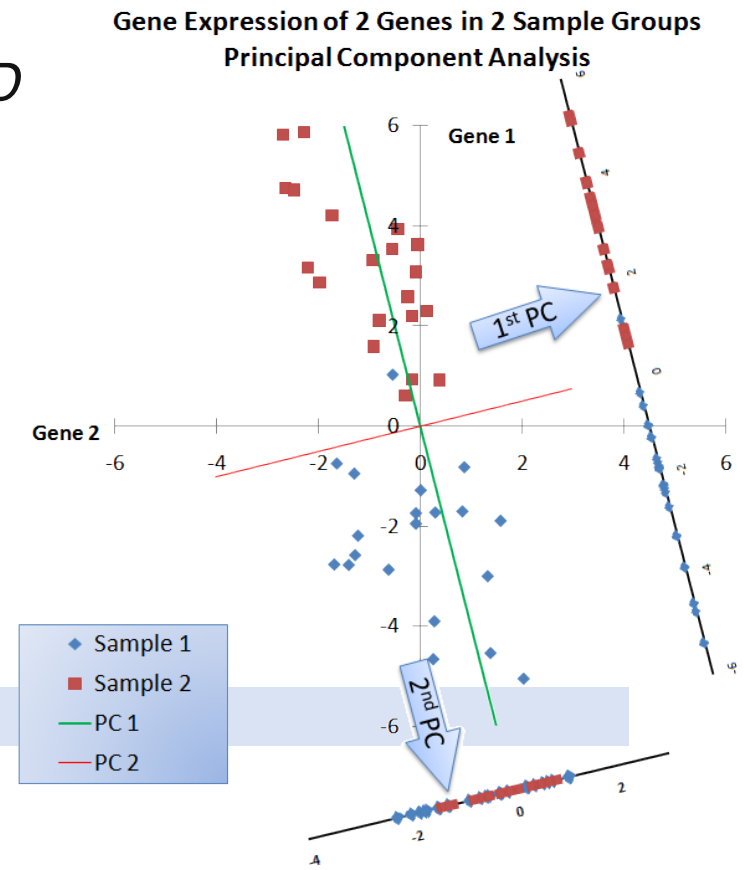
- Nonlinear feature reduction using kernels
- Manifold learning

# Principal Component Analysis (PCA)

- 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,
- Our goal is to project the data onto a space having dimensionality  $M < D$



ariance



# Principal Component Analysis (PCA)

- 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

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

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



# PCA

- Maximize  $\text{Var}(z)$  subject to  $\|\mathbf{w}\|=1$

Using  
Lagrange  
multipliers

$\Sigma \mathbf{w}_1 = \alpha \mathbf{w}_1$  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  $\mathbf{w}_1$

How?

$\Sigma \mathbf{w}_2 = \alpha \mathbf{w}_2$  that is,  $\mathbf{w}_2$  is another eigenvector of  $\Sigma$

and so on.

Computational cost :  $O(MD^2)$  for  $M$  principal components

# Minimum-error formulation

- introduce a complete orthonormal set of  $D$ -dimensional basis vectors  $\{\mathbf{u}_i\}$

$$\mathbf{u}_i^T \mathbf{u}_j = \delta_{ij}.$$

$$\mathbf{x}_n = \sum_{i=1}^D \alpha_{ni} \mathbf{u}_i$$

$$\mathbf{x}_n = \sum_{i=1}^D (\mathbf{x}_n^T \mathbf{u}_i) \mathbf{u}_i.$$

# Minimum-error formulation

- introduce a complete orthonormal set of  $D$ -dimensional basis vectors  $\{\mathbf{u}_i\}$

$$\mathbf{x}_n = \sum_{i=1}^D (\mathbf{x}_n^T \mathbf{u}_i) \mathbf{u}_i.$$

- Approximate this data point using a representation involving a restricted number  $M < D$  of variables corresponding to a projection onto a lower-dimensional subspace.

$$\tilde{\mathbf{x}}_n = \sum_{i=1}^M z_{ni} \mathbf{u}_i + \sum_{i=M+1}^D b_i \mathbf{u}_i$$

- Minimize the distortion introduced by the reduction in dimensionality.

$$J = \frac{1}{N} \sum_{n=1}^N \|\mathbf{x}_n - \tilde{\mathbf{x}}_n\|^2.$$

# Minimum-error formulation

- Minimize the distortion introduced by the reduction in dimensionality.

$$J = \frac{1}{N} \sum_{n=1}^N \|\mathbf{x}_n - \tilde{\mathbf{x}}_n\|^2.$$

$$\mathbf{x}_n = \sum_{i=1}^D (\mathbf{x}_n^T \mathbf{u}_i) \mathbf{u}_i.$$

- derivative with respect to  $z_{nj}$   $z_{nj} = \mathbf{x}_n^T \mathbf{u}_j$
- Derivative wrt  $b_j$   $b_j = \bar{\mathbf{x}}^T \mathbf{u}_j$

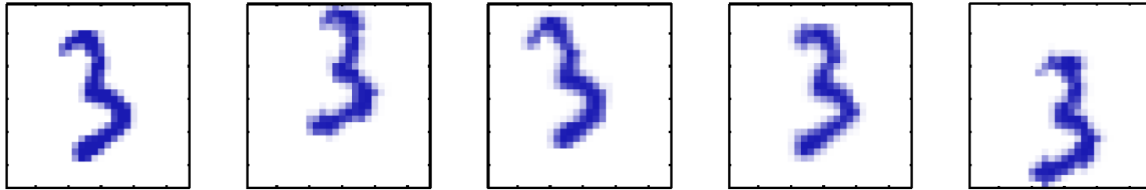
- Displacement vector lies in the space orthogonal to the principal subspace

$$\mathbf{x}_n - \tilde{\mathbf{x}}_n = \sum_{i=M+1}^D \{(\mathbf{x}_n - \bar{\mathbf{x}})^T \mathbf{u}_i\} \mathbf{u}_i$$

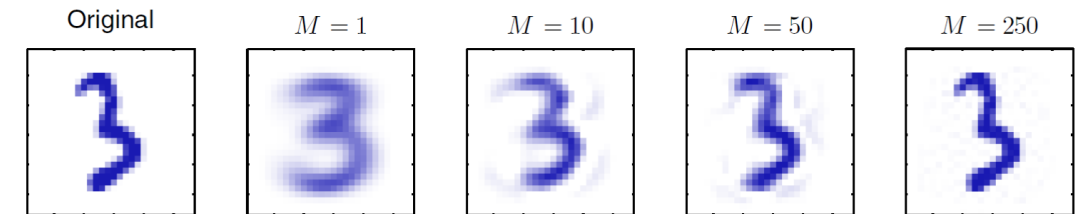
- distortion measure  $J$

$$J = \frac{1}{N} \sum_{n=1}^N \sum_{i=M+1}^D (\mathbf{x}_n^T \mathbf{u}_i - \bar{\mathbf{x}}^T \mathbf{u}_i)^2 = \sum_{i=M+1}^D \mathbf{u}_i^T \mathbf{S} \mathbf{u}_i.$$

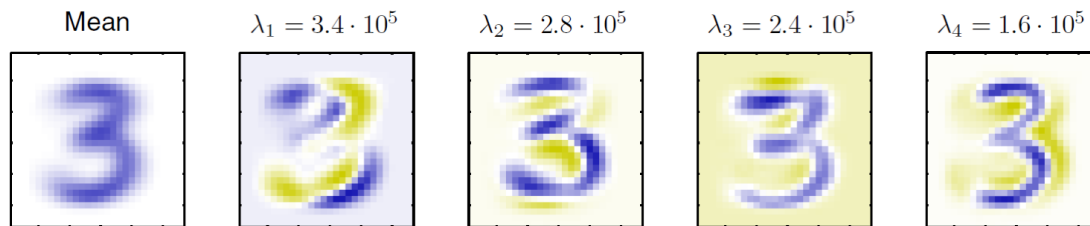
# PCA



A synthetic data set obtained by taking one of the off-line digit images and creating multiple copies in each of which the digit has undergone a random displacement and rotation within some larger image field. The resulting images each have  $100 \times 100 = 10,000$



An original example from the off-line digits data set together with its PCA reconstructions obtained by retaining  $M$  principal components for various values of  $M$ . As  $M$  increases the reconstruction becomes more accurate and would become perfect when  $M = D = 28 \times 28 = 784$ .



The mean vector  $\bar{x}$  along with the first four PCA eigenvectors  $\mathbf{u}_1, \dots, \mathbf{u}_4$  for the off-line digits data set, together with the corresponding eigenvalues.

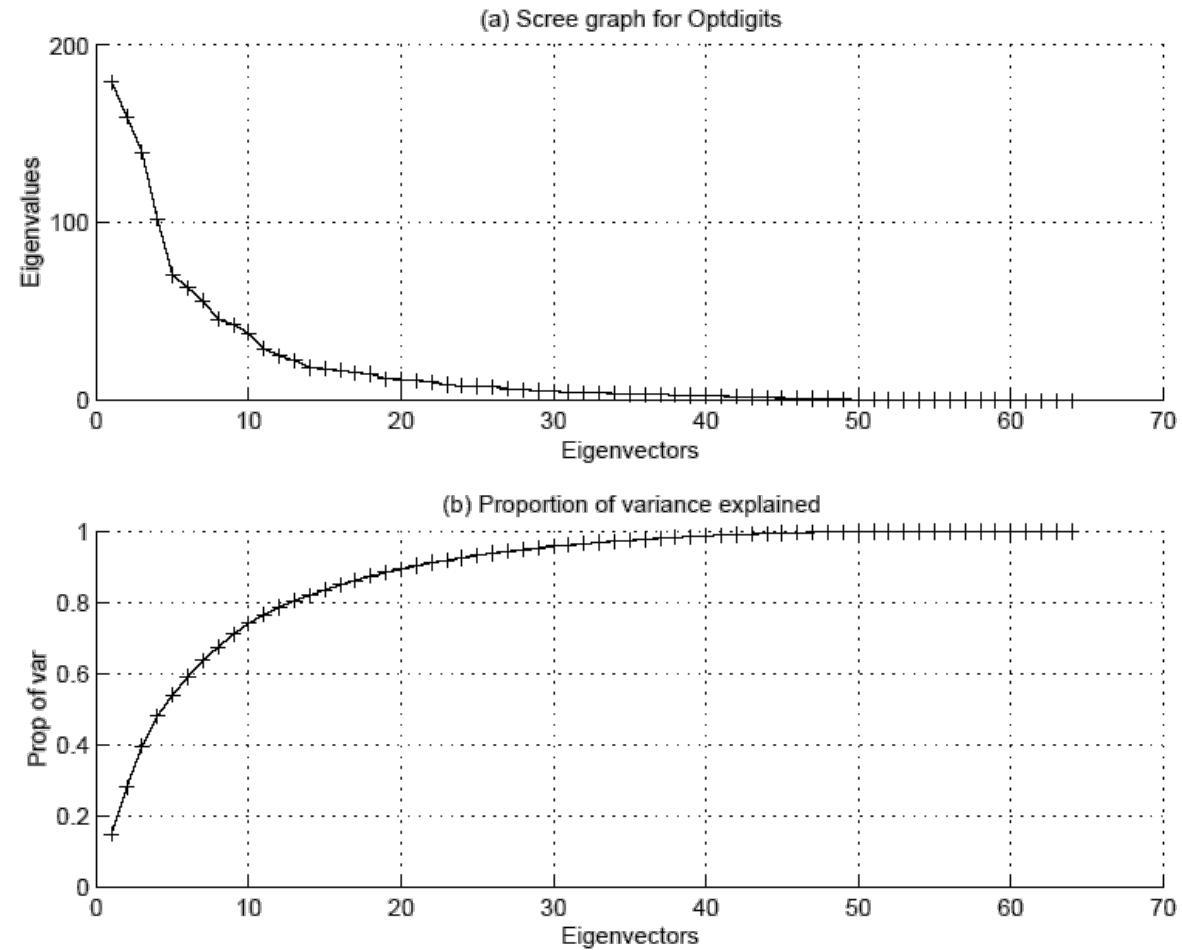
# How to choose $k$

- Proportion of Variance (PoV) explained

when  $\lambda_i$  are sorted in descending order

- Typically, stop at  $\text{PoV} > 0.9$
- See graph plots of PoV vs  $k$ , stop at “elbow”

# Illustration



# Example

