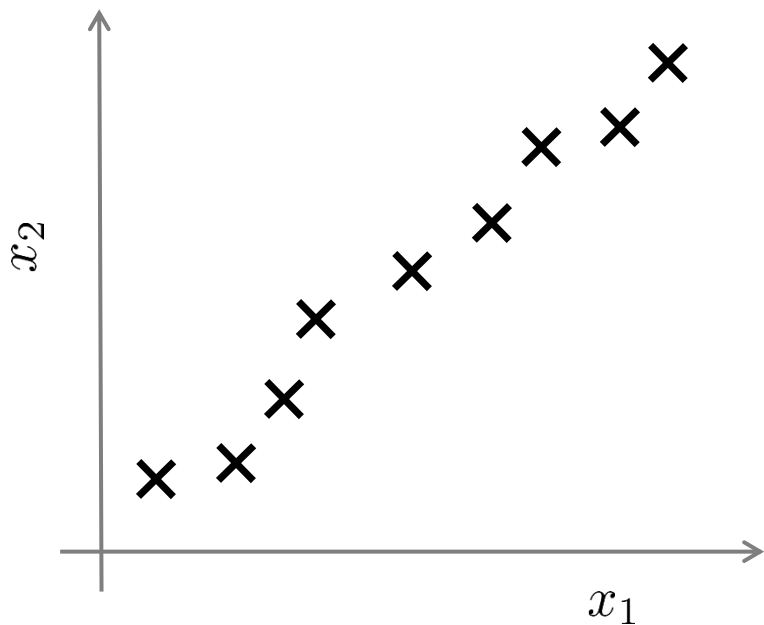# Dimensionality Reduction

Machine Learning (AIM 5002-41)

Joon Hee Choi
Sungkyunkwan University
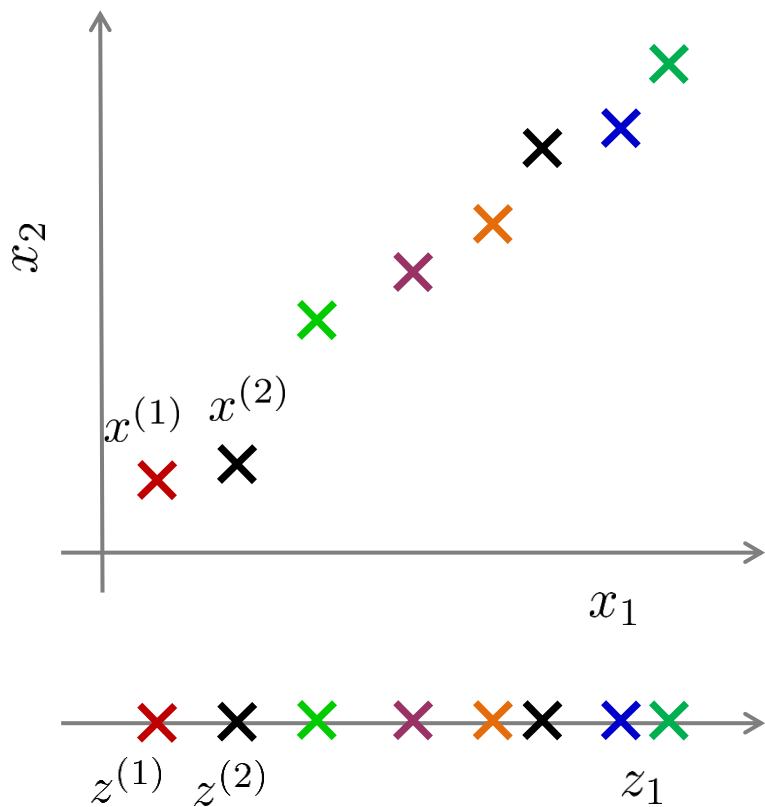
# Motivation I:
# Data Compression

# Data Compression



Reduce data from 2D to 1D

# Data Compression

Reduce data from 2D to 1D

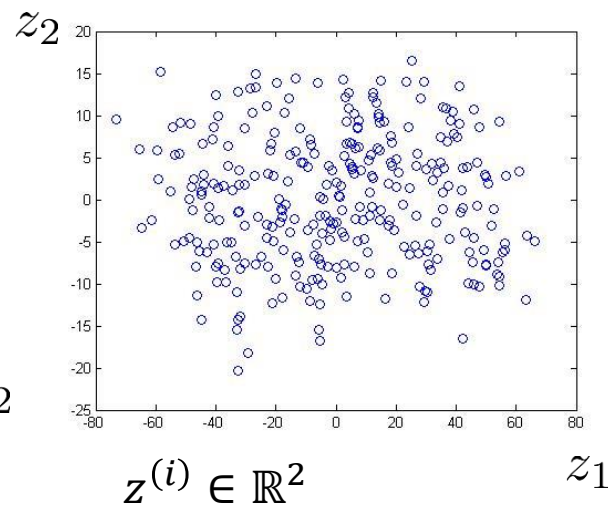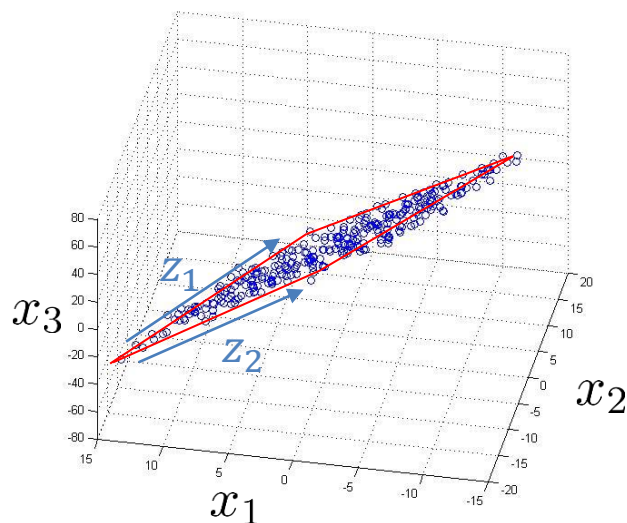$$x^{(1)} \in \mathbb{R}^2 \quad \rightarrow z^{(1)} \in \mathbb{R}$$

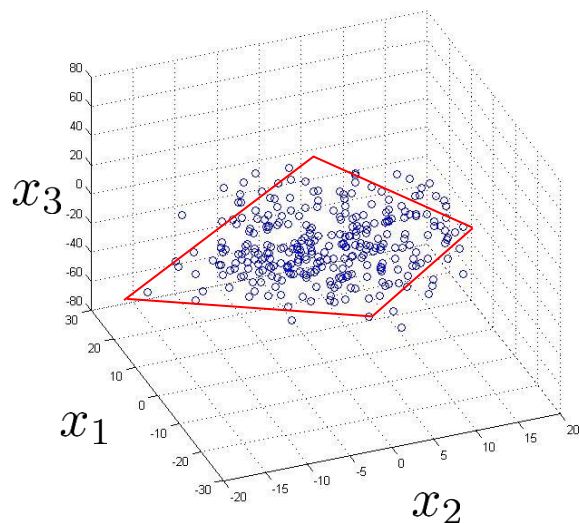$$x^{(2)} \in \mathbb{R}^2 \quad \rightarrow z^{(2)} \in \mathbb{R}$$

$$\vdots$$

$$x^{(m)} \in \mathbb{R}^2 \quad \rightarrow z^{(m)} \in \mathbb{R}$$

# Data Compression

## Reduce data from 3D to 2D



$x^{(i)} \in \mathbb{R}^3$

$z^{(i)} \in \mathbb{R}^2$

$$z = \begin{bmatrix} z_1 \\ z_2 \end{bmatrix} \qquad z^{(i)} = \begin{bmatrix} z_1^{(i)} \\ z_2^{(i)} \end{bmatrix}$$

Andrew Ng

# Motivation II:
# Data Visualization

# Data Visualization

$$x \in \mathbb{R}^{50} \qquad x^{(i)} \in \mathbb{R}^{50}$$

| Country | $x_1$ GDP (trillions of US$) | $x_2$ Per capita GDP (thousands of intl. $) | $x_3$ Human Development Index | $x_4$ Life expectancy | $x_5$ Poverty Index (Gini as percentage) | $x_6$ Mean household income (thousands of US$) | ... |
|---------|------|-------|-------|------|------|--------|-----|
| Canada | 1.577 | 39.17 | 0.908 | 80.7 | 32.6 | 67.293 | ... |
| China | 5.878 | 7.54 | 0.687 | 73 | 46.9 | 10.22 | ... |
| India | 1.632 | 3.41 | 0.547 | 64.7 | 36.8 | 0.735 | ... |
| Russia | 1.48 | 19.84 | 0.755 | 65.5 | 39.9 | 0.72 | ... |
| Singapore | 0.223 | 56.69 | 0.866 | 80 | 42.5 | 67.1 | ... |
| USA | 14.527 | 46.86 | 0.91 | 78.3 | 40.8 | 84.3 | ... |
| ... | ... | ... | ... | ... | ... | ... | ... |

[resources from en.wikipedia.org]

Andrew Ng

# Data Visualization

$z^{(i)} \in \mathbb{R}^2$

| Country | $z_1$ | $z_2$ |
|---|---|---|
| Canada | 1.6 | 1.2 |
| China | 1.7 | 0.3 |
| India | 1.6 | 0.2 |
| Russia | 1.4 | 0.5 |
| Singapore | 0.5 | 1.7 |
| USA | 2 | 1.5 |
| ... | ... | ... |

Reduce data from 50D to 2D

# Data Visualization



Singapore

USA

per-person GDP (economic activity)
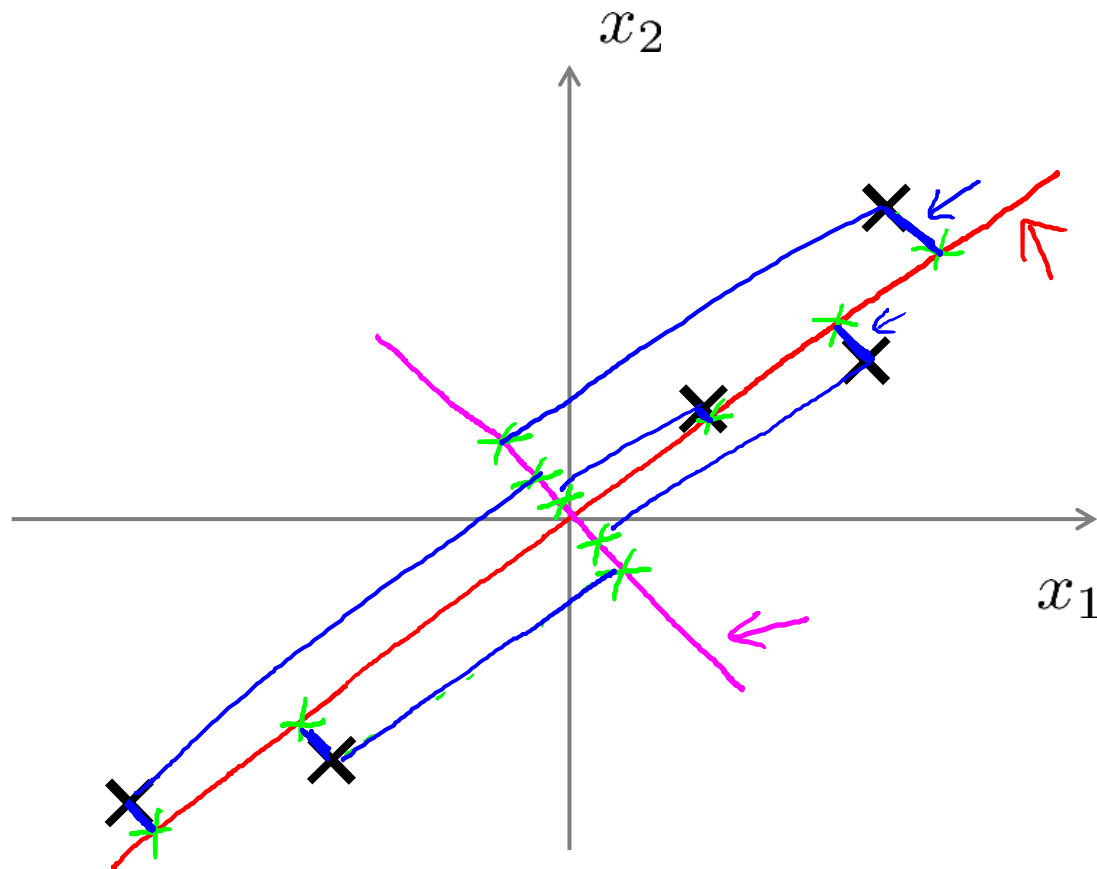
$z_2$

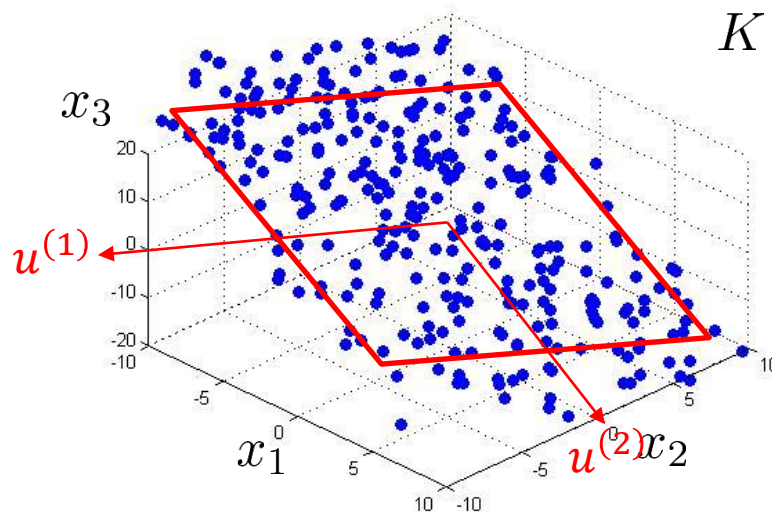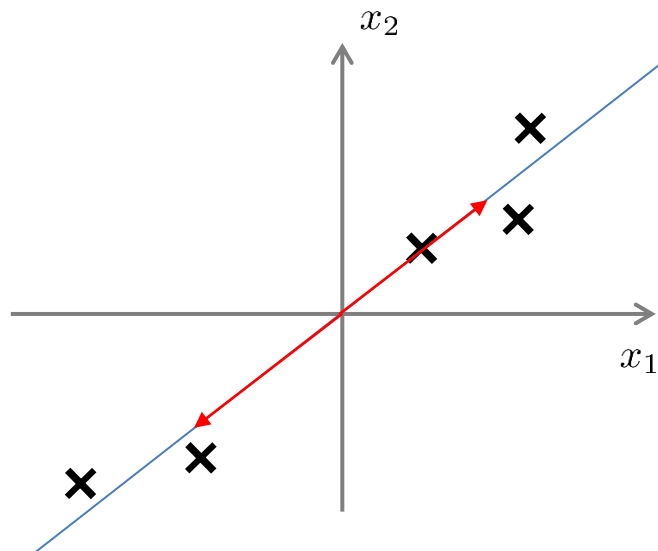Country size / GDP

$z_1$

1

# Principal Component Analysis problem formulation

Andrew Ng

# Principal Component Analysis (PCA) problem formulation



Andrew Ng

# Principal Component Analysis (PCA) problem formulation

$$3D \rightarrow 2D$$
$$K = 2$$



Reduce from 2-dimension to 1-dimension: Find a direction (a vector $u^{(1)} \in \mathbb{R}^n$) onto which to project the data so as to minimize the projection error.

Reduce from n-dimension to k-dimension: Find $k$ vectors $u^{(1)}, u^{(2)}, \ldots, u^{(k)}$ onto which to project the data, so as to minimize the projection error.

**Data preprocessing**

Training set: $x^{(1)}, x^{(2)}, \ldots, x^{(m)}$

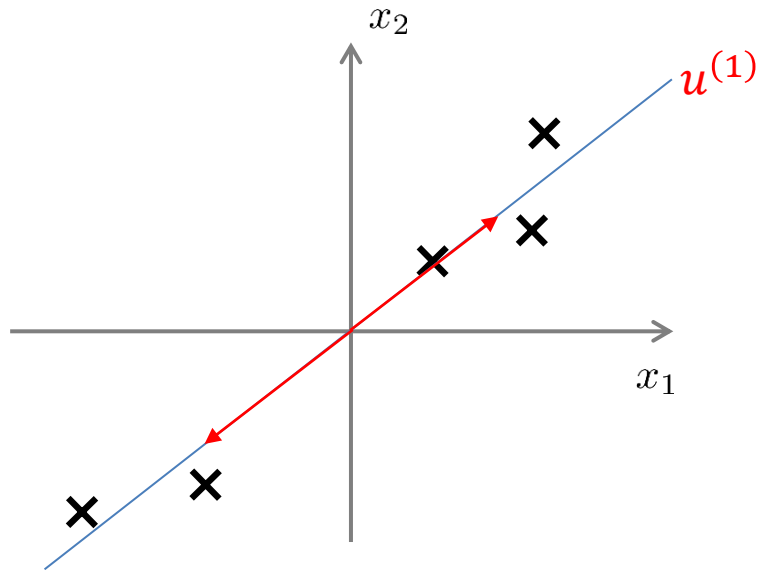Preprocessing (feature scaling/mean normalization):

$$\mu_j = \frac{1}{m} \sum_{i=1}^{m} x_j^{(i)}$$
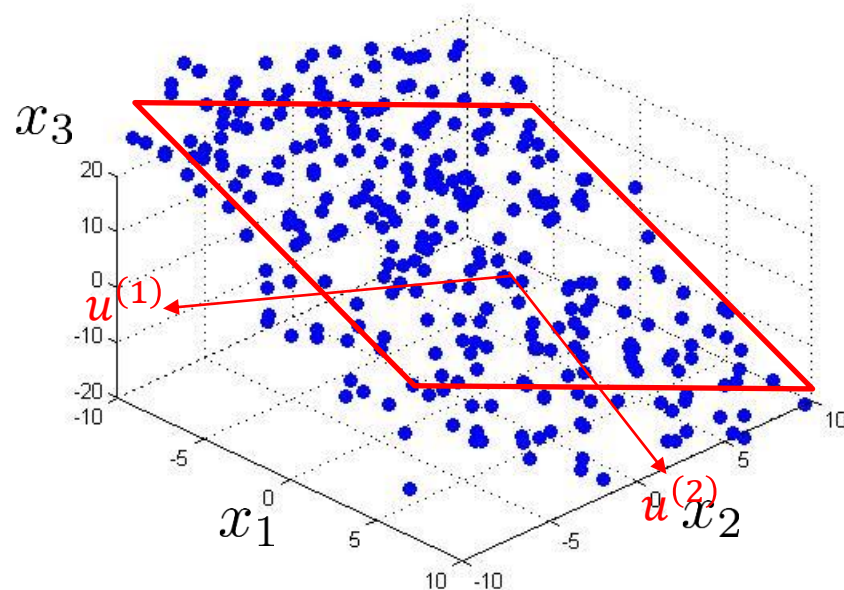
Replace each $x_j^{(i)}$ with $x_j - \mu_j$.

If different features on different scales (e.g., $x_1 =$ size of house, $x_2 =$ number of bedrooms), scale features to have comparable range of values.

$$x_j^{(i)} \leftarrow \frac{x_j^{(i)} - \mu_j}{s_j}$$

# Principal Component Analysis (PCA) algorithm



Reduce data from 2D to 1D

Reduce data from 3D to 2D

# Principal Component Analysis (PCA) algorithm

Reduce data from $n$-dimensions to $k$-dimensions

Compute "covariance matrix":

$$\Sigma = \frac{1}{m} \sum_{i=1}^{n} \boxed{(x^{(i)})(x^{(i)})^T} \quad n \times n$$

$$\underset{n \times 1}{} \quad \underset{1 \times n}{}$$

→ Singular value decomposition

Compute "eigenvectors" of matrix $\Sigma$:

```
[U,S,V] = svd(Sigma);
```

$$U = \begin{bmatrix} | & | & & | \\ u^{(1)} & u^{(2)} & \dots & u^{(n)} \\ | & | & & | \end{bmatrix} \in \mathbb{R}^{n \times n}$$

$\underbrace{\qquad\qquad}_{k}$

Andrew Ng

# Principal Component Analysis (PCA) algorithm

From `[U,S,V] = svd(Sigma)`, we get:

$$U = \begin{bmatrix} | & | & & | \\ u^{(1)} & u^{(2)} & \ldots & u^{(n)} \\ | & | & & | \end{bmatrix} \in \mathbb{R}^{n \times n}$$

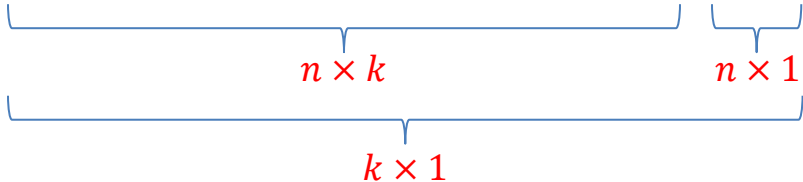$$\underbrace{\phantom{u^{(1)} \quad u^{(2)}}}_{k}$$

$$z^{(i)} = \begin{bmatrix} | & | & & | \\ u^{(1)} & u^{(2)} & \ldots & u^{(k)} \\ | & | & & | \end{bmatrix}^{T} x^{(i)}$$

$$\underbrace{\phantom{u^{(1)} \quad u^{(2)} \quad u^{(k)}}}_{n \times k} \quad \underbrace{\phantom{x^{(i)}}}_{n \times 1}$$

$$\underbrace{\phantom{u^{(1)} \quad u^{(2)} \quad u^{(k)} \quad x}}_{k \times 1}$$

# Principal Component Analysis (PCA) algorithm summary

After mean normalization (ensure every feature has zero mean) and optionally feature scaling:

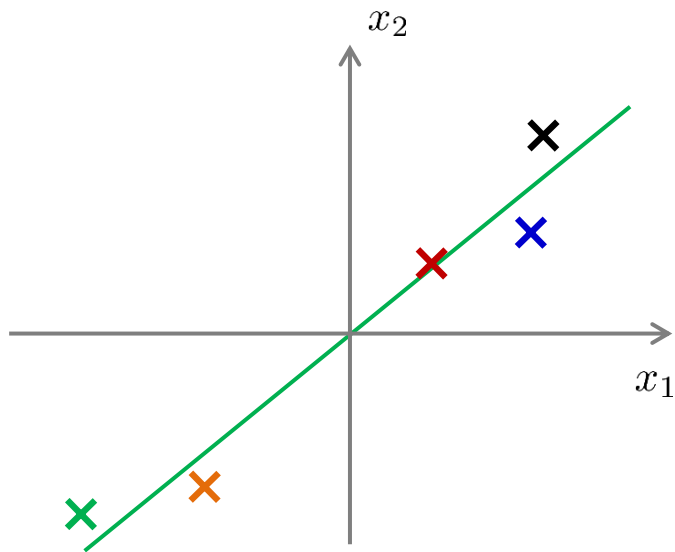**Sigma =** $\dfrac{1}{m}\displaystyle\sum_{i=1}^{m}(x^{(i)})(x^{(i)})^{T}$
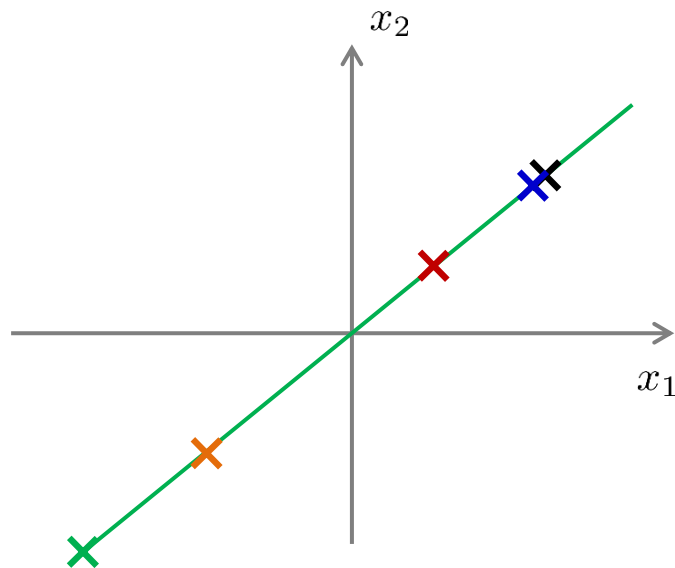
```
[U,S,V] = svd(Sigma);
Ureduce = U(:,1:k);
z = Ureduce'*x;
```

# Reconstruction from compressed representation

# Reconstruction from compressed representation

$$z = U_{reduce}^T x$$

# Choosing the number of principal components

**Choosing $k$ (number of principal components)**

Average squared projection error: $\frac{1}{m}\sum_{i=1}^{m}\left\|x^{(i)} - x_{approx}^{(i)}\right\|^2$

Total variation in the data: $\frac{1}{m}\sum_{i=1}^{m}\left\|x^{(i)}\right\|^2$

Typically, choose $k$ to be smallest value so that

$$\frac{\frac{1}{m}\sum_{i=1}^{m}\left\|x^{(i)} - x_{approx}^{(i)}\right\|^2}{\frac{1}{m}\sum_{i=1}^{m}\left\|x^{(i)}\right\|^2} \leq 0.01 \qquad (1\%)$$

"99% of variance is retained"

# Choosing $k$ (number of principal components)

Algorithm:

Try PCA with $k = 1$

Compute $U_{reduce}, z^{(1)}, z^{(2)},$
$\ldots, z^{(m)}, x^{(1)}_{approx}, \ldots, x^{(m)}_{approx}$

Check if

$$\frac{\frac{1}{m} \sum_{i=1}^{m} \|x^{(i)} - x^{(i)}_{approx}\|^2}{\frac{1}{m} \sum_{i=1}^{m} \|x^{(i)}\|^2} \leq 0.01?$$

`[U,S,V] = svd(Sigma)`

# Choosing $k$ (number of principal components)

`[U,S,V] = svd(Sigma)`

Pick smallest value of $k$ for which

$$\frac{\sum_{i=1}^{k} S_{ii}}{\sum_{i=1}^{m} S_{ii}} \geq 0.99$$

(99% of variance retained)

# Advice for applying PCA

**Supervised learning speedup**

$$(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \ldots, (x^{(m)}, y^{(m)})$$

Extract inputs:

Unlabeled dataset: $x^{(1)}, x^{(2)}, \ldots, x^{(m)} \in \mathbb{R}^{10000}$

$$\downarrow PCA$$

$$z^{(1)}, z^{(2)}, \ldots, z^{(m)} \in \mathbb{R}^{1000}$$

New training set:

$$(z^{(1)}, y^{(1)}), (z^{(2)}, y^{(2)}), \ldots, (z^{(m)}, y^{(m)})$$

Note: Mapping $x^{(i)} \rightarrow z^{(i)}$ should be defined by running PCA only on the training set. This mapping can be applied as well to the examples $x_{cv}^{(i)}$ and $x_{test}^{(i)}$ in the cross validation and test sets.

# Application of PCA

- Compression
    - Reduce memory/disk needed to store data
    - Speed up learning algorithm

- Visualization

**Bad use of PCA: To prevent overfitting**

Use $z^{(i)}$ instead of $x^{(i)}$ to reduce the number of features to $k < n$.

Thus, fewer features, less likely to overfit.

This might work OK, but isn't a good way to address overfitting. Use regularization instead.

$$\min_{\theta} \frac{1}{2m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^{n} \theta_j^2$$