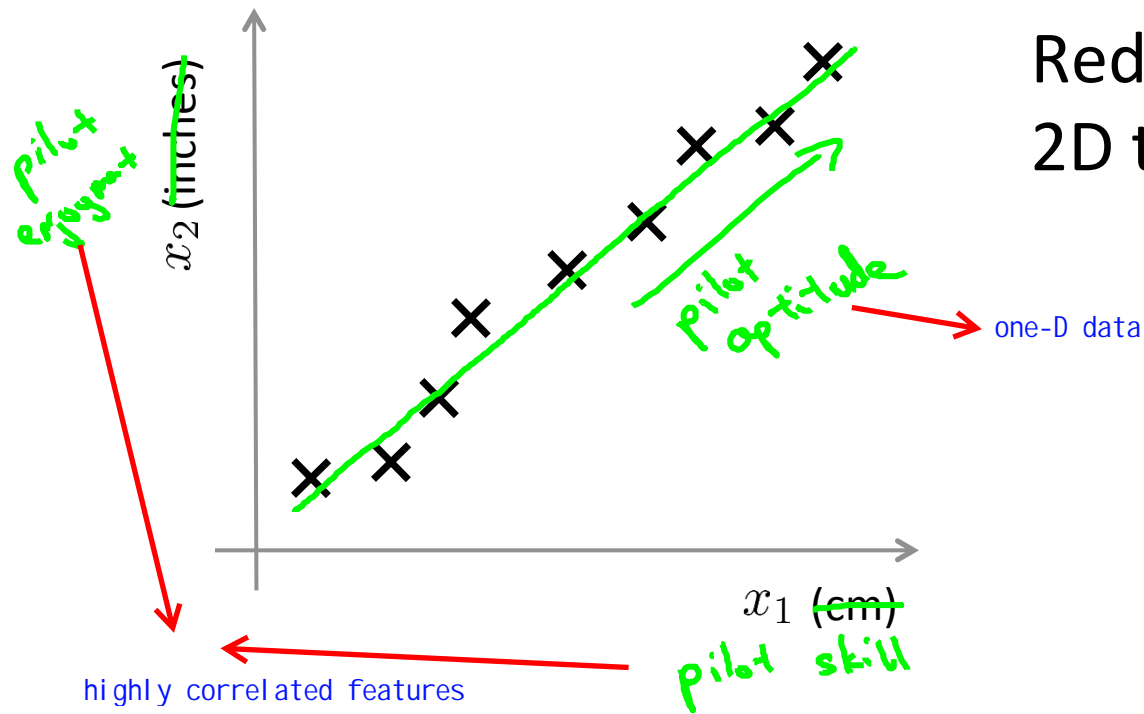# Dimensionality Reduction

## Motivation I: Data Compression

This can help our algorithm to run faster!

Machine Learning

# Data Compression



Reduce data from 2D to 1D

pilot enrolment

highly correlated features

pilot skill

pilot aptitude

one-D data

$x_2$ (inches)

$x_1$ (cm)

Andrew Ng

# Data Compression

project the data onto this line, then we can represent the data using 1 number, therefore reduce the dimension.

Reduce data from 2D to 1D

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

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

$\vdots$

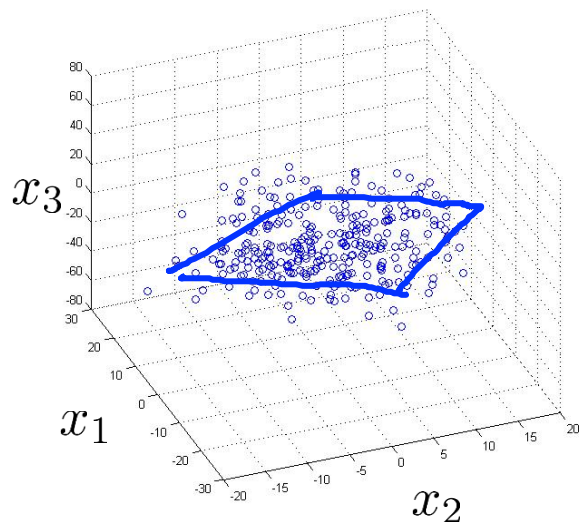$x^{(m)} \in \mathbb{R}^2 \rightarrow z^{(m)} \in \mathbb{R}$

$x_2$ (inches)

$x^{(1)}$

$x^{(2)}$

$z^{(2)}$   $z^{(1)}$

$x_1$ (cm)

$z_1 \longrightarrow$ new feature

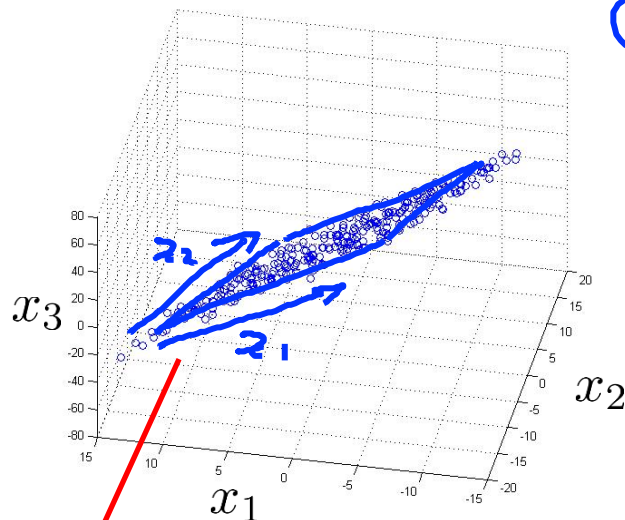Andrew Ng

**Data Compression**

In a typical setting, we may have 10000 dimensions

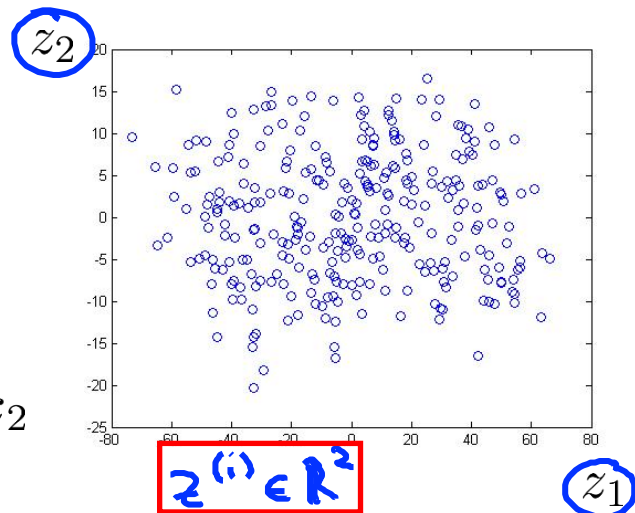$10000 \rightarrow 100D$

Reduce data from <u>3D to 2D</u>

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

$z_2$

$z_1$

Project the data into this plane

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

$z = \begin{bmatrix} z_1 \\ z_2 \end{bmatrix}$

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

Andrew Ng

# Dimensionality Reduction

## Motivation II: Data Visualization

Machine Learning

Fact about some countries in the world

## Data Visualization

$x \in \mathbb{R}^{50}$

50 features

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

| Country | $x_1$ GDP (trillions of US$) | $x_2$ Per capita GDP (thousands of intl. $) | $x_3$ Human Develop- ment 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

Andrew Ng

# Data Visualization



per. person
GDP
(economic activity)

$z_2$

$z^{(i)} \in \mathbb{R}$
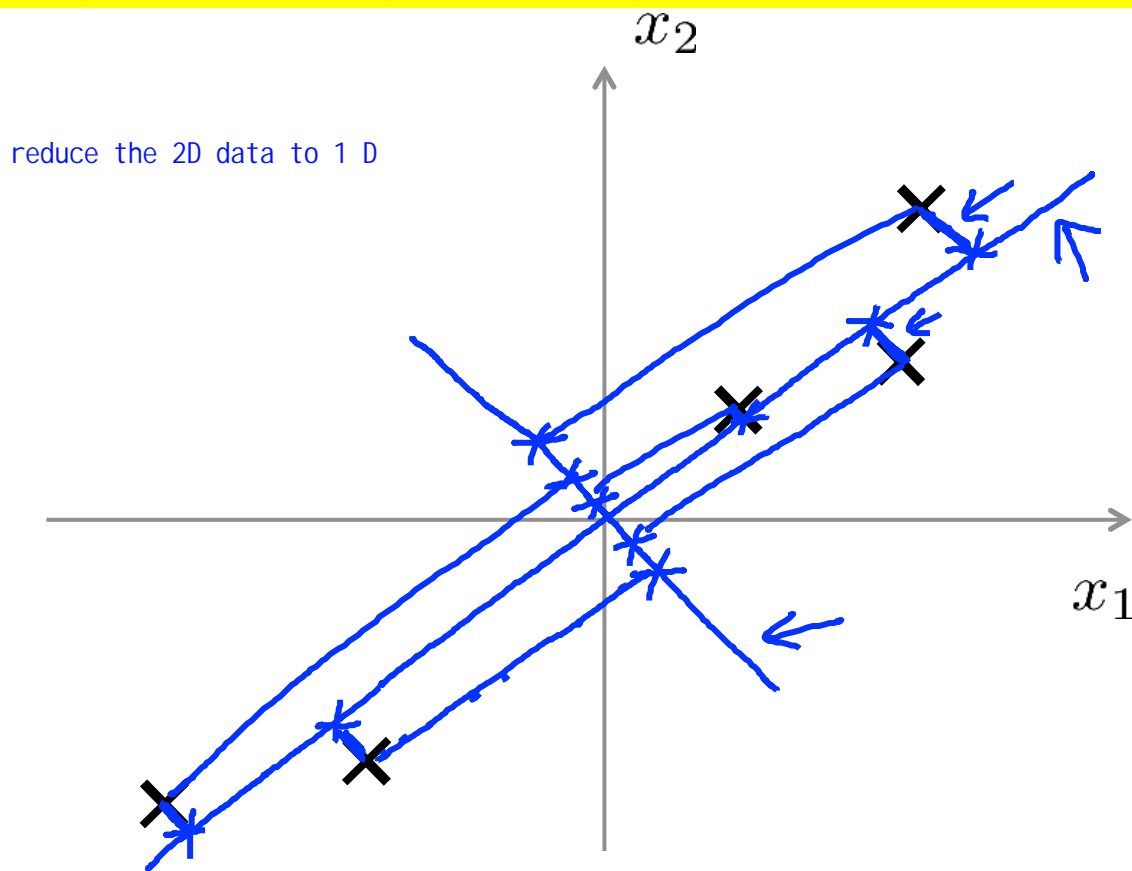
$z_1$

Singapore

USA

Country size /
GDP

1

# Dimensionality Reduction

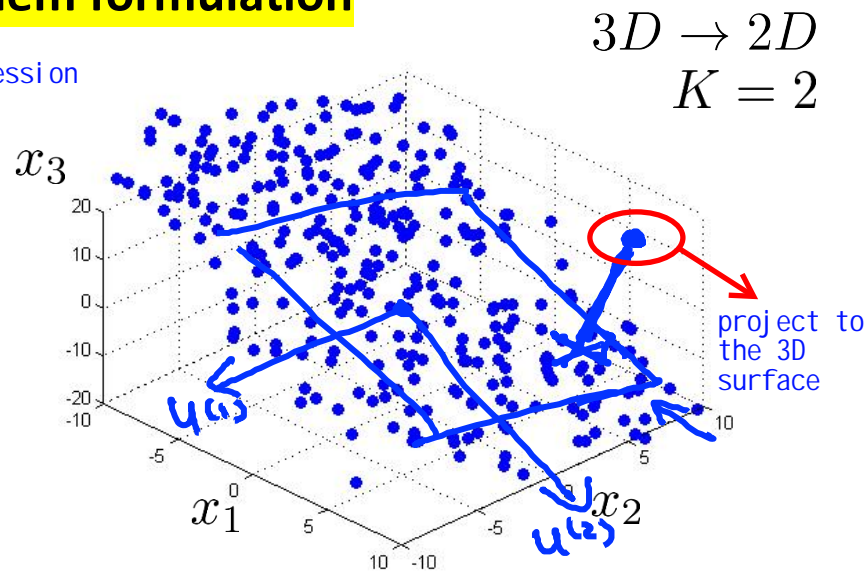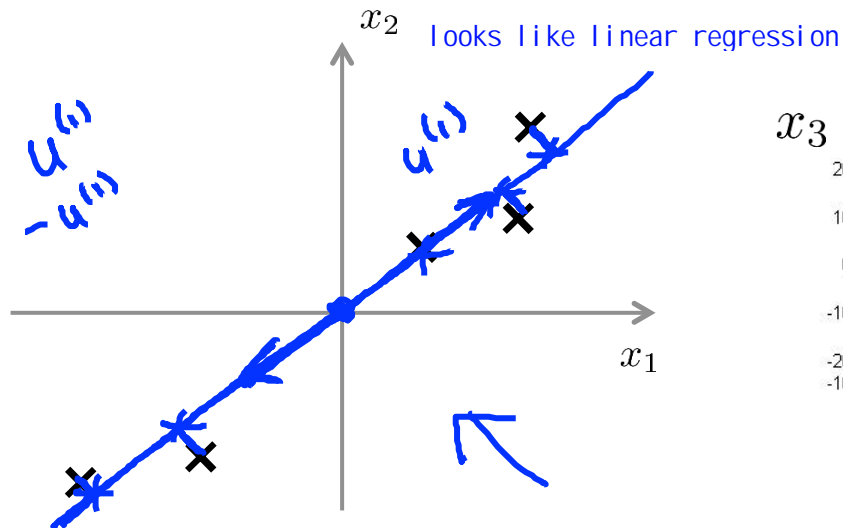## Principal Component Analysis problem formulation

By far, the most popular and commonly used unsupervised learning algorithm

Machine Learning

# Principal Component Analysis (PCA) problem formulation

$x_2$

$x \in \mathbb{R}^2$

reduce the 2D data to 1 D

$x_1$

# Principal Component Analysis (PCA) problem formulation
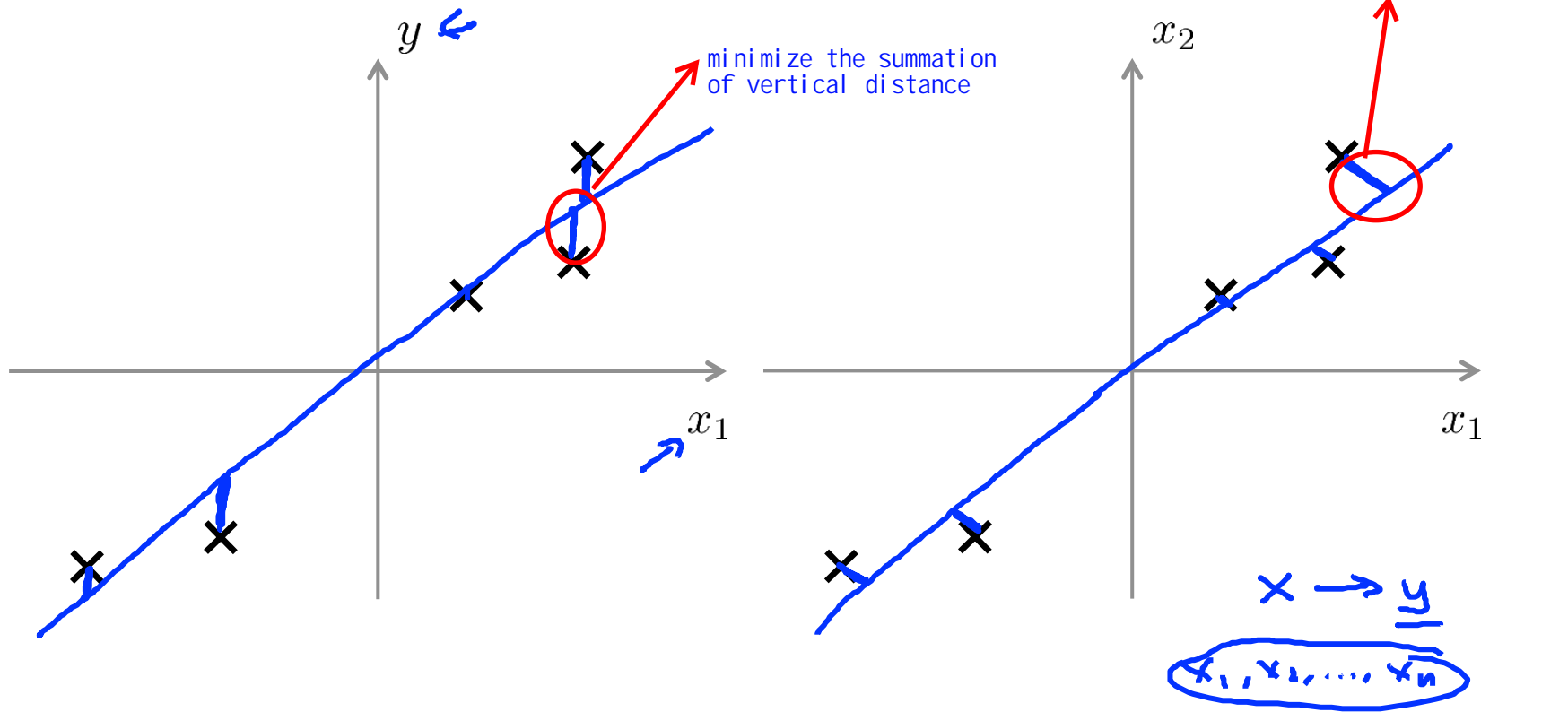


$3D \rightarrow 2D$
$K = 2$

looks like linear regression

project to the 3D surface

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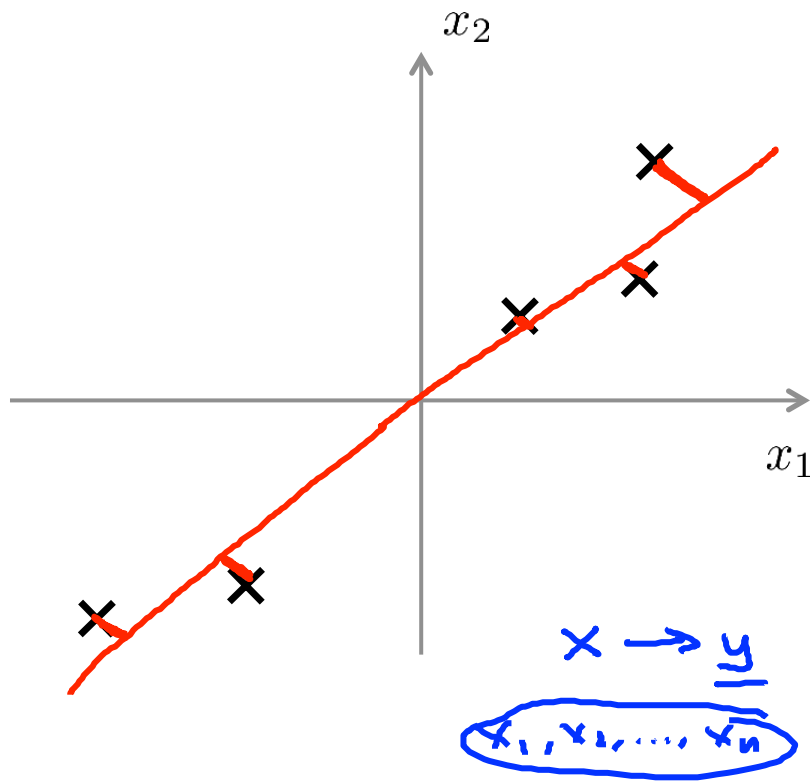
Reduce from n-dimension to k-dimension: Find $k$ vectors $u^{(1)}, u^{(2)}, \ldots, u^{(k)}$ onto which to project the data, so as to minimize the projection error.

Andrew Ng

# PCA is not linear regression



minimize the summation
of vertical distance

minimize this projection
distance summation

$y$

$x_1$

$x_2$

$x_1$

$x \longrightarrow y$

$x_1, x_2, \ldots, x_n$

# PCA is not linear regression



Reduce data from 3D to 2D

# Dimensionality Reduction

## Principal Component Analysis algorithm

Machine Learning

## Data preprocessing

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

Preprocessing (feature scaling/mean normalization):

Before we actually apply PCA

$$\mu_j = \frac{1}{m} \sum_{i=1}^{m} x_j^{(i)}$$

zero mean feature

Replace each $x_j^{(i)}$ with $x_j - \mu_j$.
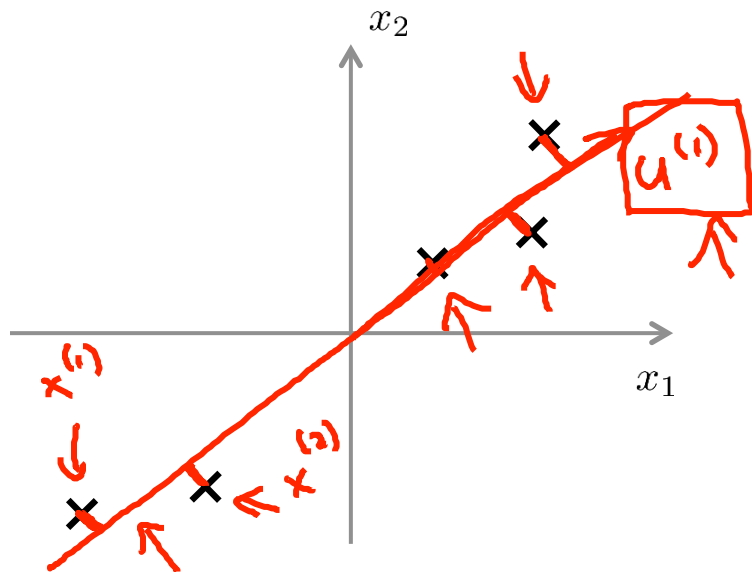
If different features on different scales (e.g., $x_1 =$ size of house, $x_2 =$ number of bedrooms), scale features to have comparable range of values.

feature scaling

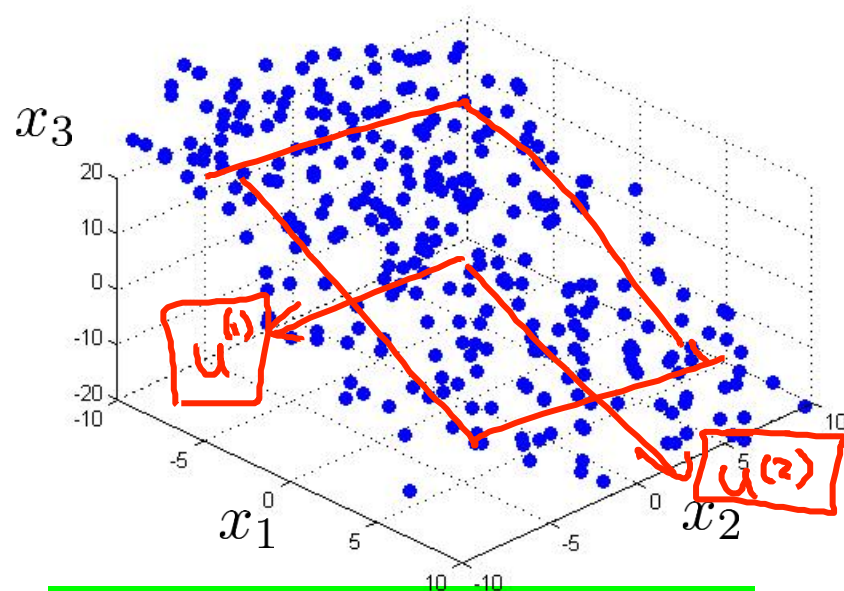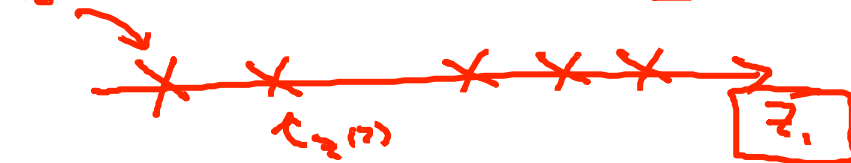$$x_j^{(i)} \leftarrow \frac{x_j^{(i)} - \mu_j}{s_j}$$

some measures of the feature j
e.g. Max value

# 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":

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

$n \times 1$  $1 \times n$  $n \times n$  Sigma

Compute "eigenvectors" of matrix $\Sigma$:

$\rightarrow$ Singular value decomposition

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

eig (Sigma)

$n \times n$ matrix

$U = \begin{bmatrix} | & | & | & & | \\ u^{(1)} & u^{(2)} & u^{(3)} & \cdots & u^{(n)} \\ | & | & | & & | \end{bmatrix}$

$U \in \mathbb{R}^{n \times n}$

$u^{(1)}, \ldots, u^{(k)}$

we mainly require this one

$k$

Andrew Ng

# Principal Component Analysis (PCA) algorithm

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

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

low dimensional representation of x^{i}

K dimensional vector

$$x \in \mathbb{R}^n \longrightarrow z \in \mathbb{R}^k$$

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

$$z \in \mathbb{R}^k$$

$$n \times k$$

$$U_{reduce}$$

$$k \times n$$

$$n \times 1$$

$$k \times 1$$

# Principal Component Analysis (PCA) algorithm summary

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

$$\text{Sigma} = \frac{1}{m} \sum_{i=1}^{m} (x^{(i)})(x^{(i)})^T$$

```
[U,S,V] = svd(Sigma);
Ureduce = U(:,1:k);
z = Ureduce'*x;
```

we grab the 1st k columns

find the z representation

$$X = \begin{bmatrix} \underline{\quad} & x^{(1)T} & \underline{\quad} \\ & \vdots & \\ \underline{\quad} & x^{(m)T} & \underline{\quad} \end{bmatrix}$$

$$\text{Sigma} = (1/m) * X' * X;$$

Vector representation of sigma!!
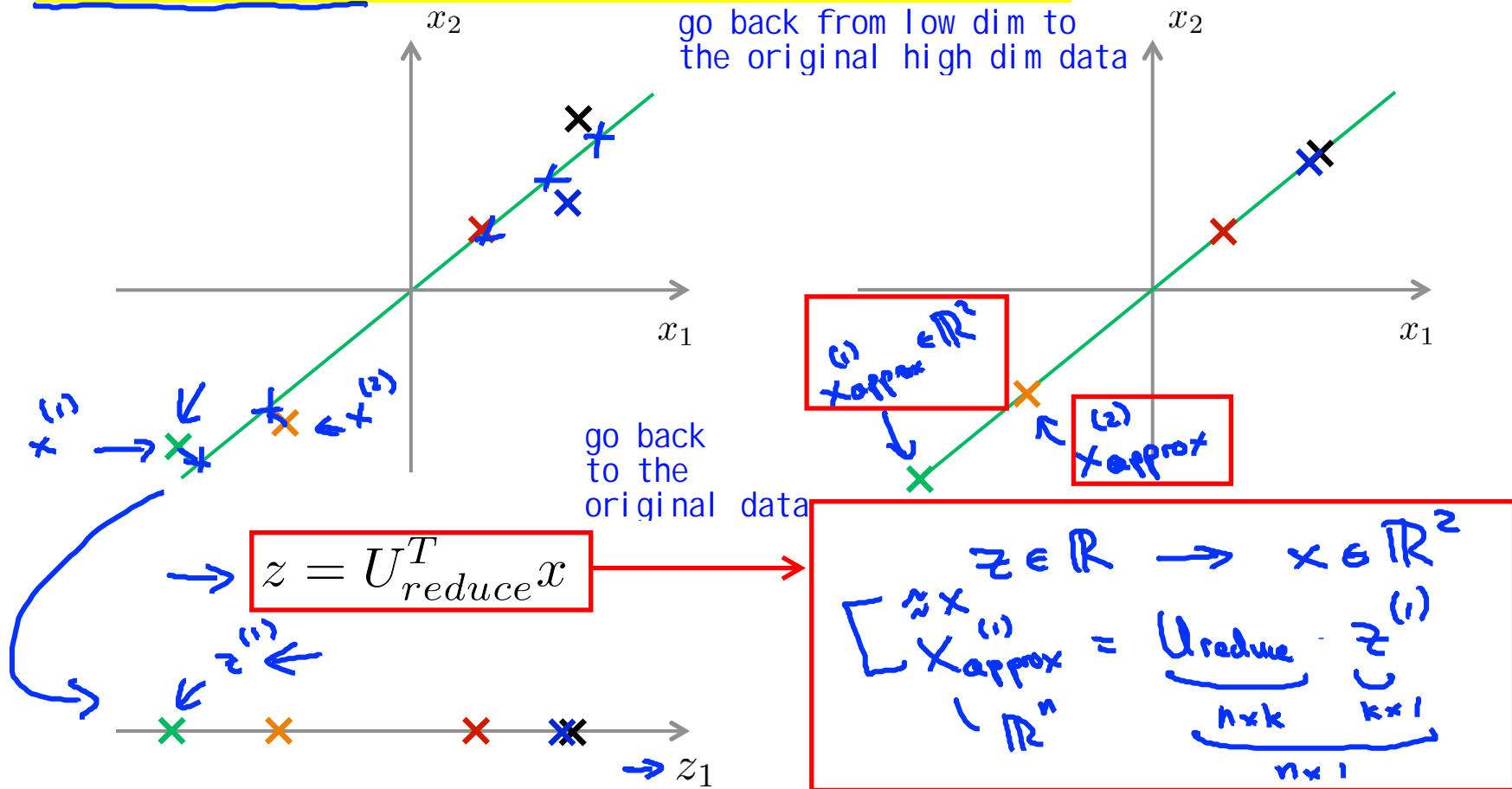
$$x \in \mathbb{R}^n \qquad \cancel{x_0 = 1}$$

# Dimensionality Reduction

Reconstruction from compressed representation

This section shows how do we go back from 100 dimensional data to 1000 dimensional data!

Machine Learning

# Reconstruction from compressed representation

$x_2$

go back from low dim to
the original high dim data

$x_2$

$x_1$

$x_1$

$x^{(1)}$  $x^{(2)}$

go back
to the
original data

$x^{(1)}_{approx} \in \mathbb{R}^2$

$x^{(2)}_{approx}$

$$z = U^T_{reduce} x$$

$z^{(1)}$

$z \in \mathbb{R} \longrightarrow x \in \mathbb{R}^2$

$$x^{(1)}_{approx} = \underbrace{U_{reduce}}_{n \times k} \cdot \underbrace{z^{(1)}}_{k \times 1}$$
$$\underbrace{\qquad}_{n \times 1}$$

$\approx x \in \mathbb{R}^n$

$z_1$

# Dimensionality Reduction

## Choosing the number of principal components

Machine Learning

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

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

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

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

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

k is used to estimate x_approx

0.05    5%

0.10    (10%)

"99% of variance is retained"

95% to 90%

Andrew Ng

# Choosing $k$ (number of principal components)

much more efficient

**Algorithm:**

Try PCA with $k = 1$   $k = 2$  $k = 3$  $k = 4$

Compute $U_{reduce}, z^{(1)}, z^{(2)},$

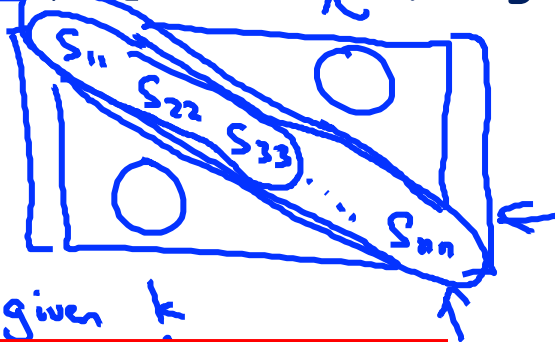$\dots, z^{(m)}, x^{(1)}_{approx}, \dots, x^{(m)}_{approx}$

Check if

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

$k = 17$

Equivalent!!!

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

$S =$

$S_{11}$  $S_{22}$  $S_{33}$  $\dots$  $S_{nn}$

$k = 3$

For given $k$

$$1 - \frac{\sum_{i=1}^{k} S_{ii}}{\sum_{i=1}^{n} S_{ii}} \leq 0.01$$

Simpler Calculation!!!

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

Andrew Ng

# Choosing $k$ (number of principal components)

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

Pick smallest value of $k$ for which

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

$k=100$

(99% of variance retained)

Andrew Ng

# Dimensionality Reduction

## Advice for applying PCA

Machine Learning

One common use of PCA

# Supervised learning speedup

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

$x^{(i)} \in \mathbb{R}^{10,000}$

very large feature

Extract inputs:

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

$\downarrow PCA$

$U_{reduce}$

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

$x$
$\downarrow$
$z$

New training set:

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

$h_\theta(z) = \dfrac{1}{1 + e^{-\theta^T z}}$

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

$x \to z$

Then apply this mapping to cross validation and test sets

# **Application of PCA**

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

  Choose k by % of variance retain

- **Visualization**

  k = 2 or k = 3

# Bad use of PCA: To prevent overfitting

Use $z^{(i)}$ instead of $x^{(i)}$ to reduce the number of features to $k < n$. — 1000 — 10000

Thus, fewer features, less likely to overfit.

*Bad!*

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

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

Andrew Ng

**PCA is sometimes used where it shouldn't be**

Design of ML system:

- - Get training set $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \ldots, (x^{(m)}, y^{(m)})\}$
- - Run PCA to reduce $x^{(i)}$ in dimension to get $z^{(i)}$
- - Train logistic regression on $\{(z^{(1)}, y^{(1)}), \ldots, (z^{(m)}, y^{(m)})\}$
- - Test on test set: Map $x_{test}^{(i)}$ to $z_{test}^{(i)}$. Run $h_\theta(z)$ on
  $\{(z_{test}^{(1)}, y_{test}^{(1)}), \ldots, (z_{test}^{(m)}, y_{test}^{(m)})\}$

How about doing the whole thing without using PCA?

Before implementing PCA, first try running whatever you want to do with the original/raw data $x^{(i)}$. Only if that doesn't do what you want, then implement PCA and consider using $z^{(i)}$.

Andrew Ng