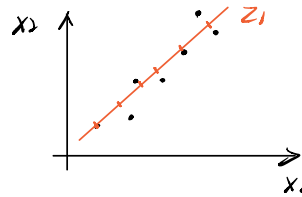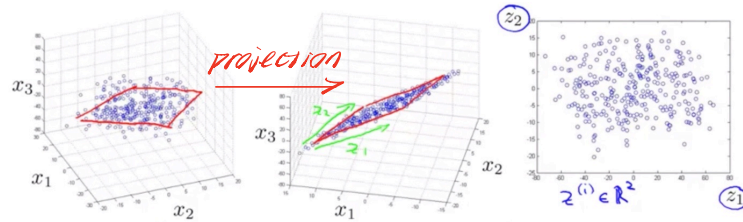② Dimension reduction

application: ① Data compression     Choose $k$ by variation %



can reduce data from 2D to 1D.

$(x_1, x_2) \rightarrow z_1$



② Data visualization.
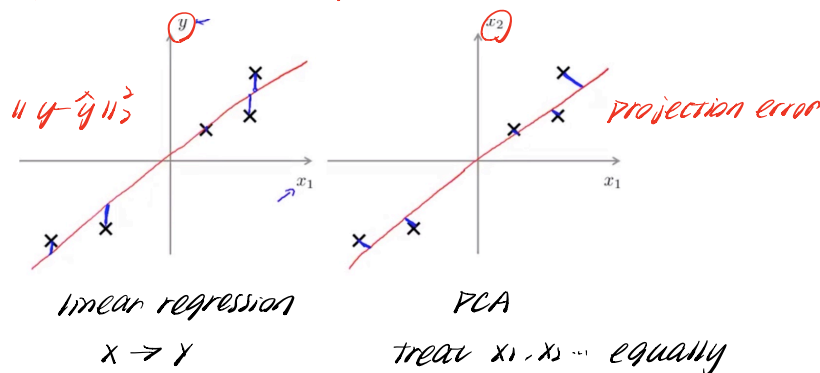
$$R^m \rightarrow R^3 \text{ or } R^2.$$

algorithm: Principle Component Analysis (PCA).

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

PCA is NOT linear regression



$\|y - \hat{y}\|_2^2$

projection error

linear regression
$X \rightarrow Y$

PCA
treat $x_1, x_2 \dots$ equally

### Data preprocessing

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

Preprocessing (feature scaling/mean normalization):

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

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

② 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 summary

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

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

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

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

$$Sigma = (1/m) * X' * X;$$

First k columns

Singular Value Decomposition

### 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

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

$\alpha$

"99% of variance is retained"

Algorithm:

Try PCA with $k=1$

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

Check if

$$\frac{\frac{1}{m} \sum_{i=1}^{m} \|x^{(i)} - x_{approx}^{(i)}\|^2}{\frac{1}{m} \sum_{i=1}^{m} \|x^{(i)}\|^2} \leq 0.01?$$
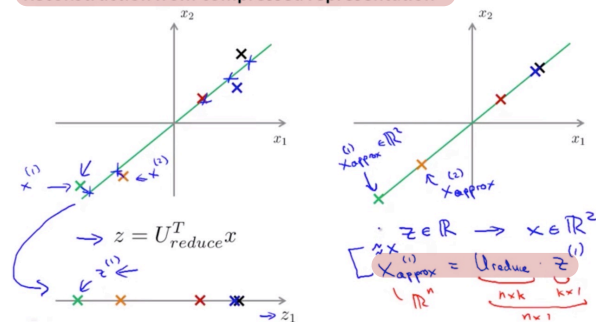
k = 17

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

$$S = \begin{bmatrix} S_{11} & & & \\ & S_{22} & & \\ & & S_{33} & \\ & & & S_{nn} \end{bmatrix}$$

For given k         k=3

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

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

### Reconstruction from compressed representation



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

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

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

$\mathbb{R}^n$

### Supervised learning speedup

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

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

Extract inputs:

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

$\downarrow PCA$     $U_{reduce}$

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

New training set:

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

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

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.

$x \rightarrow z$

### Bad use of PCA: To prevent overfitting

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

Thus, fewer features, less likely to overfit.

Don't know y: Might throw important information.       Bad!

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

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

Misuse of PCA in ML: Use raw data first, try PCA when some issues arise
e.g. memory leak, speed ...