# Unsupervised learning -introduction : Clustering

## Unsupervised learning
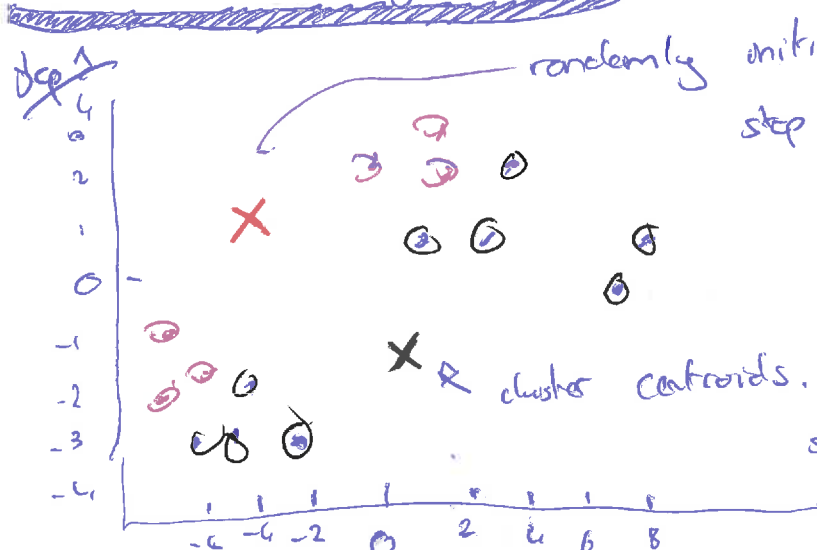


2 different clusters

Training set: $\{x^{(1)}, x^{(2)}, x^{(3)} \ldots x^{(m)}\}$

## Applications of clustering

- Market segmentation
- Social network Analysis
- Organise computing clusters
- astronomical Analysis

# K-means Algorithm

step 1



randomly initialize two cluster center points

X cluster centroids.

step 1) Cluster assignment step

So if a data point is closer to a cluster centroid → it is assigned to that cluster centroid.

step 2) move centroid step

→ compute average of all the assigned cluster datapoints of a cluster and move the centroid to this new centrepoint.

→ after awhile the centroid locations no longer change.

step 2



now closer together - centroids "moved"

## K-means algorithm

**Input:**
- $k$ (number of clusters) ✓
- training set $\{x^{(1)}, x^{(2)} \dots\}$ ✓ (no $y \dots$)

$x^{(i)} \in \mathbb{R}^n$ (drop $x_0 = 1$

---

## K-means algorithm

total number of centroids

Randomly initialize $k$ cluster centroids $\mu_1, \mu_2, \dots, \mu_k \in \mathbb{R}^n$

Repeat {

cluster assignment step

for $i = 1$ to $m$

$c^{(i)} :=$ index (from 1 to $k$) of cluster centroid closest to $x^{(i)}$

$\min_k \| x^{(i)} - \mu_k \|^2$ length.

move centroid step.

for $k = 1$ to $k$

$\mu_k :=$ average (mean) of points assigned to cluster $k$

↳ let assume this cluster ($\mu_2$) has the values $x^{(1)}, x^{(5)}, x^{(6)}, x^{(10)}$ assigned to it.
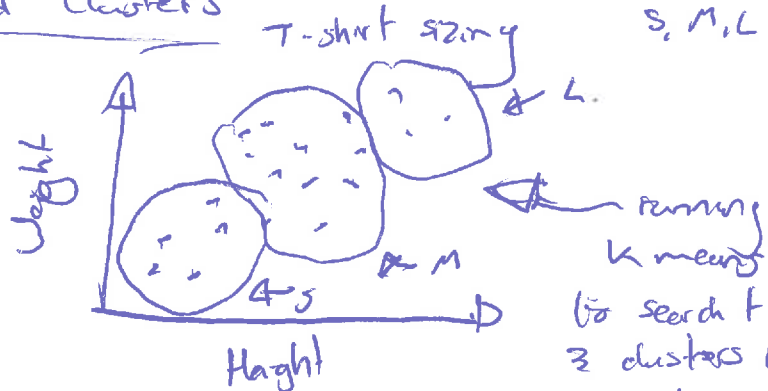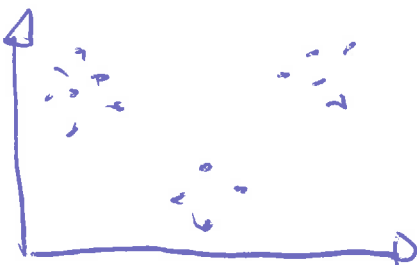
$\therefore c^{(1)} = 2, c^{(5)} = 2, c^{(6)} = 2,$

$c^{(10)} = 2.$

$\mu_2 \leftarrow \frac{1}{4} \left[ x^{(1)} + x^{(5)} + x^{(6)} + x^{(10)} \right) \in \mathbb{R}^n$

- cluster assignment step ← distance between $i$th training example and cluster centroid
↳ so ∴ picking the cluster centroid $c^{(i)}$ that is closest to our training example $x^{(i)}$

notes: delete delete a cluster if it has no points assigned to it
↳ or reinitialize it through randomization

---

## K-means for non-separated clusters

T-shirt sizing    S, M, L



Weight

Height

← L.

← running k means
↳ to search for 3 clusters for S, M, L.

# Clustering — optimization objective

## K-means optimization objective

$c^{(i)}$ = index of cluster $(1, 2, \ldots, K)$ to which example $x^{(i)}$ is currently assigned.

$\mu_k$ = cluster centroid $K$ $(\mu_k \in \mathbb{R}^n)$

> lower case $k$ is $k = \{1, 2, \ldots k\}$
> an index into

$\mu_{c(i)}$ = cluster centroid of cluster to which $x^{(i)}$ has been assigned

$$x^{(i)} \to 5 \quad \therefore \quad c^{(i)} = 5$$

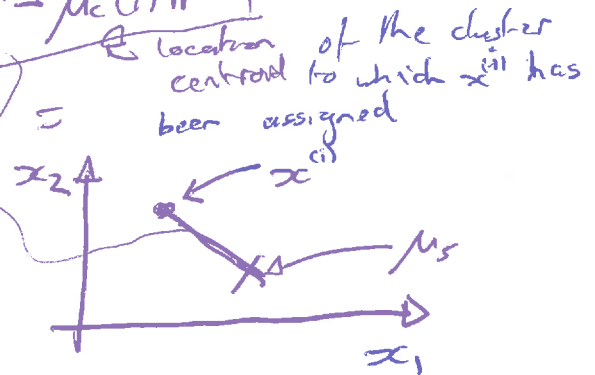$\therefore x$ has been assigned to cluster 5

$$\mu_{c^{(i)}} = \mu_5 = ?$$

optimization objective:

$$J(c^{(1)}, \ldots, c^{(m)}, \mu_1, \ldots, \mu_K) = \frac{1}{m} \sum_{i=1}^{m} \| x^{(i)} - \mu_{c(i)} \|^2$$

$\underset{\substack{c^{(1)}, \ldots, c^{(m)} \\ \mu_1, \ldots, \mu_K}}{\min} J(c^{(1)}, \ldots, c^{(m)}, \mu_1, \ldots, \mu_K)$

← location of the cluster centroid to which $x^{(i)}$ has been assigned

⟶ Sometimes called the "distortion cost function"



⟶ step 1 | The cluster assignment step in minimizing $J(\ldots)$ with respect to $c^{(1)}, c^{(2)}$ etc. ... $c^m$ (holding $\mu_1, \ldots, \mu_K$ fixed)

⟶ step 2 | minimises $J(\ldots)$ w.r.t. $\mu_1 \ldots \mu_K$.

## Clustering - Random initialization
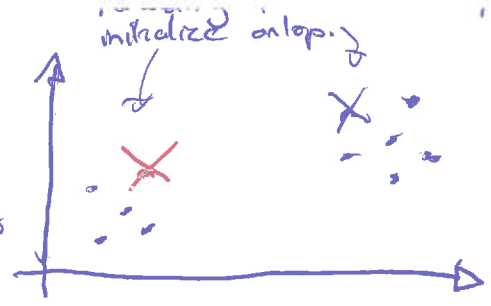
should have $K < m$

Randomly pick $K$ training examples

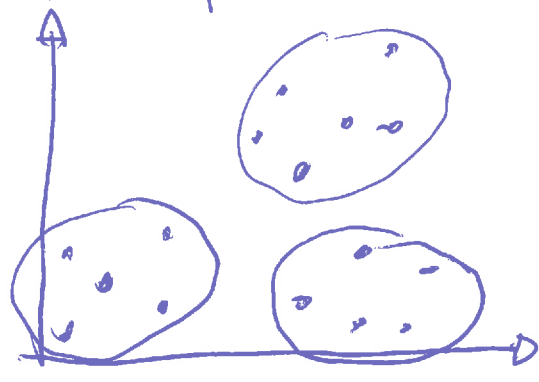Set $\mu_1, \ldots, \mu_K$ equal to these $K$

examples $\mu_1 = x^{(i)}$

$\mu_2 = x^{(j)}$
$\vdots$

lets say:

$\boxed{K = 2}$

two examples



K means can converge to different solutions depending how the clusters were initialized

---

Local optima. $\rightarrow$ can get stuck at different local optima.



initialize $K$ means 100s of times (randomly) to increase the chances of finding the global minimum.

---

Random initialization                typical $50 \rightarrow 1000$

for $i = 1$ to $100$ {

    Randomly initialize K-means

    Run K means to get $c^{(1)}, \ldots, c^{(m)}, \mu_1 \cdots \mu_K$

    Compute cost function (distortion)

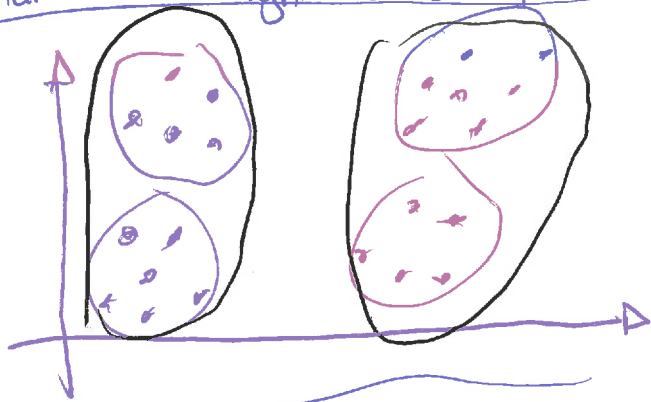    $J(c^{(1)}, \ldots, c^{(m)}, \mu_1, \ldots \mu_K)$

Pick clustering that gives lowest cost $J(c^{(1)}, \ldots, c^{(m)}, \mu_1, \ldots \mu_K)$

$K = 2 \rightarrow 10$

If you have $K > 10$, multiple random initializations will not make much difference
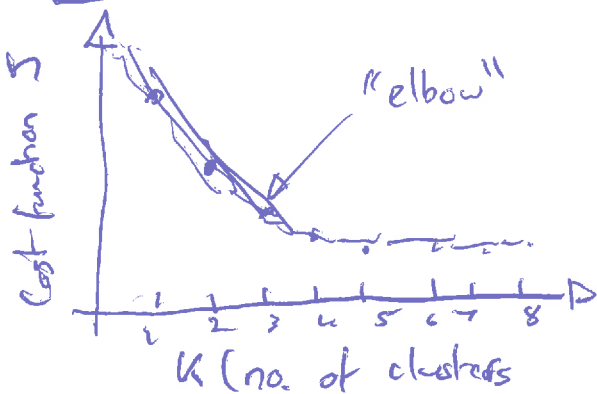
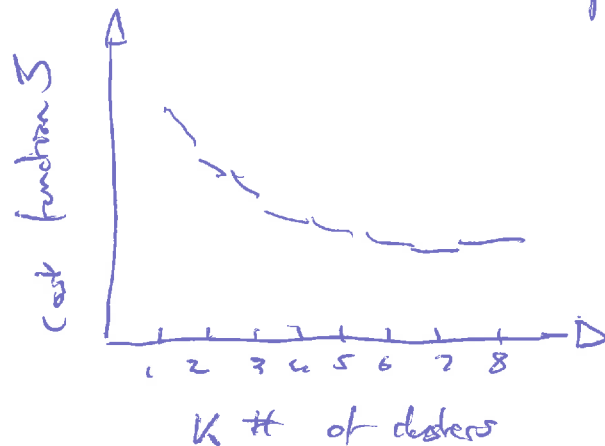# Clustering - choosing the number of clusters

What is the right value of K?



## choosing the value of K
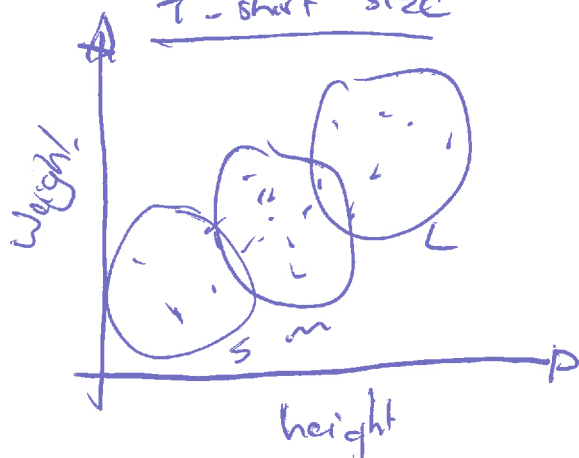
Elbow method:   "worth a shot"



Sometimes this happens:



## Choosing the value of K

Sometimes, you're running k means to get clusters for some later/downstream purpose. Evaluate k-means based on a metric for how well it performs for that later purpose.

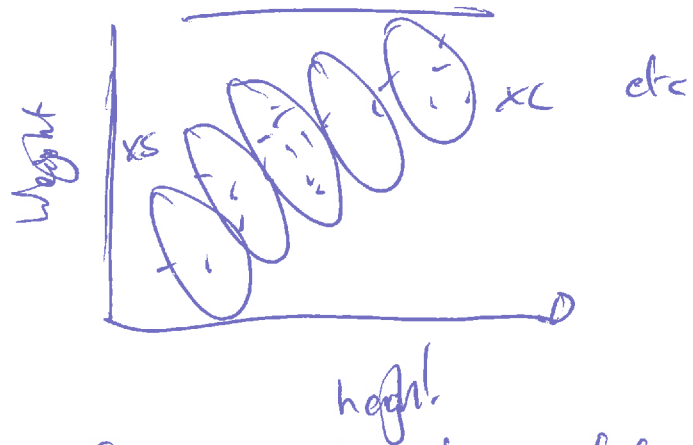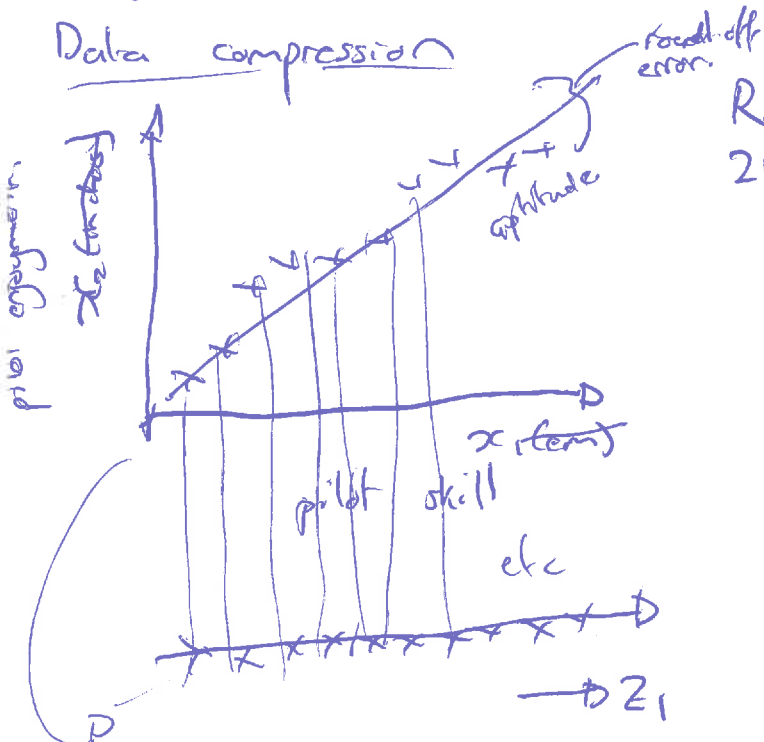$K=3$     S, M, L

$K=5$   T-shirt size   xs, s, m, L, XL.

T-shirt size





deciding what set of sizes for t shirts to manufacture is

# Dimensionality Reduction — Motivation I: data compression

## Data compression



pilot enjoyment

$x_2$ (arbitrary)

round-off error.

amplitude

$x$ (term)

pilot skill

etc

$\rightarrow z_1$

need only $1$ real number to specify a point on a line.

Reduce data from 2D to 1D

$$x^{(1)} \in \mathbb{R}^2 \rightarrow z^{(1)} \in \mathbb{R}$$
$$x^{(2)} \in \mathbb{R}^2 \rightarrow z^{(2)} \subset \mathbb{R}$$
$$\vdots$$
$$x^{(m)} \longrightarrow z^{(m)}$$

## Data Compression

1000 D to 100D

reduce data from 3D to 2D.



project

$\rightarrow$

data onto a 2D plane

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

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

Reduce data from 50D to 2D.



$Z_2$ $0$

$0$

$Z_1$

Singapore

USA.

country
size
GDP-

GDP
per person

# Dimensionality Reduction - principal component analysis algorithm

## 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.

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

$\hookleftarrow$ standard deviation of feature $j$.

## Principal Component Analysis Algorithm

- Reduce data from 2D to 1D
- Reduce data from 3D to 2D.

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

$$z = \begin{bmatrix} z_1 \\ z_2 \end{bmatrix}$$
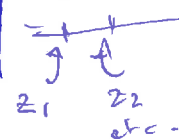
e.g

$z_1$  $z_2$  etc.

## Principal Component Analysis (PCA) Algorithm

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

Compute "covariance matrix"

not to be confused with summation...

$$\Sigma = \frac{1}{m} \sum_{i=1}^{n} \underset{n \times 1}{(x^{(i)})} \underset{1 \times n}{(x^{(i)})^T} \longrightarrow n \times n \text{ matrix}$$

Compute "eigenvectors" of matrix $\Sigma$ $\longrightarrow$ $n \times n$ matrix

$\longrightarrow$ $[U, S, V] = svd(Sigma),$   (or $eig(sigma)$)

$\hookleftarrow$ octave   $\hookrightarrow$ singular value decomposition.

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

$\underbrace{\phantom{aaaa}}_{K}$

$\uparrow$ vectors for plane

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

$u^{(1)}, \ldots, u^{(k)}$ : $\hookleftarrow$ k direction onto which we want to project the data.

# Dimensionality reduction - principal component analysis algorithm

from $[u, s, v] = svd(Sigma)$, we get:

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

$$\underbrace{\phantom{U = \begin{bmatrix} | & | \\ u^{(1)} & u^{(2)} \end{bmatrix}}}_{K.}$$

$$X \in \mathbb{R}^n \longrightarrow Z \in \mathbb{R}^k$$

$$Z = \underbrace{\begin{bmatrix} | & | & & | \\ u^{(1)} & u^{(2)} & \cdots & u^{(k)} \\ | & | & & | \end{bmatrix}^T}_{n \times k.} X^{(i)} \leq \underbrace{\begin{bmatrix} \text{---} & u^{(1)} & \text{---} \\ \text{---} & u & \text{---} \\ \text{---} & (u^{(k)})^T & \text{---} \end{bmatrix}}_{k \times n} \times \underset{n \times 1}{X^{(i)}}$$

$$\underbrace{\phantom{XXXXXXXXXXXXXXXXXXXX}}_{k \times 1}$$

$U_{reduce}$

$Z = \mathbb{R}^k$

---

## PCA Algorithm summary

$$Sigma = \frac{1}{m} \sum_{i=1}^{m} (x^{(i)})(x^{(i)})^T \quad \begin{array}{l}\text{implemented} \\ \text{in octave} \\ \text{looks like}\end{array} \longrightarrow sigma = (\tfrac{1}{m}) * X' * X ;$$

$$\left. \begin{array}{l} [U, S, V] = svd(Sigma); \\ Ureduce = U(:, 1:k); \\ Z = Ureduce' * x ; \end{array} \right\} \begin{array}{l}\text{octave} \\ \text{code}\end{array}$$

$$x \in \mathbb{R}^n \quad \cancel{x_0 = 1}$$

## Dimension Reduction : <u>Reconstruction</u> from compressed representation

in the previous examples we ~~app~~ reduced a 2D (or $N^{th}D$) dimension space, down to 1D using

$(x_2, x_1)$ $\Big\{$ $Z = U^T reduce\ x$
↑
new dimension (1D)

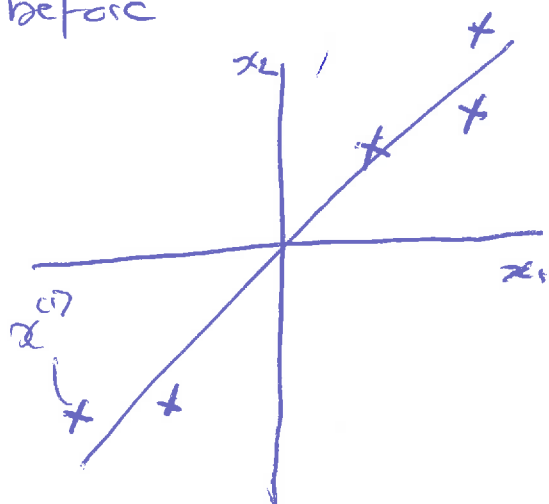$\Big\{$ to reverse it, we need to do the following:

$Z \in \mathbb{R} \longrightarrow x \in \mathbb{R}^2.$

$x_{approx} = \underbrace{U_{reduce}}_{n \times k} \cdot \underbrace{Z}_{k \times 1}$
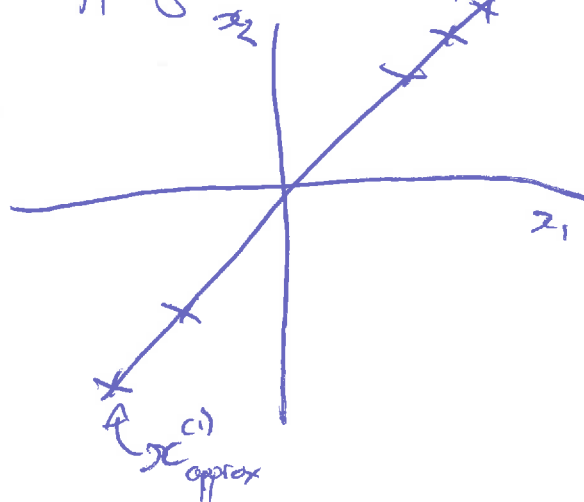
$\approx x$   $\underbrace{\phantom{xxxxxxxxxx}}_{n \times 1}$
↑
original
$x.$

**before**



remapping    buck (after)

# Dimensionality Reduction – choosing the number of principal components

## Choosing k (number of principal components)

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

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

Typically, choose k to be smallest value so that

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

or $0.05$ (5%)

"99% of variance is retained" / "95% of variance retained"

---

## choosing k (number of principal components)

Algorithm:

Try PCA with $k=1$, $k=2$, $k=3$ ... if not less than 0.01, try 2, if not, try... etc

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

check if:

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

$[U,S,V] = svd(sigma);$

$$S = \begin{bmatrix} S_{11} & & & 0 \\ & S_{22} & & \\ & & S_{33} & \\ & & & \ddots \\ 0 & & & S_{nn} \end{bmatrix}$$

For given value k   $k=3$

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

---

## Summary

$\Rightarrow$ Pick smallest value of K for which

$$\frac{\sum_{i=1}^{k}S_{ii}}{\sum_{i=1}^{m}S_{ii}} \geq 0.99 \quad (99\% \text{ variance retained})$$

"choose what dimension to reduce the data too"

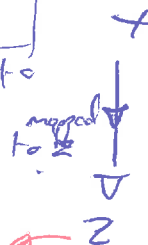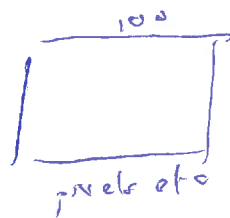# Dimensionality reduction – advice for applying PCA

## Supervised learning speedup.

$$(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)}) \qquad x^{(i)} \in \mathbb{R}^{10000}$$

pixels of o

mapped to $z$

### extract inputs:

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

$\downarrow$ PCA

new representation

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

$z$

### New training set

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

if using logistic regression for example:

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

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.

$U_{reduce}$

parameter learned by PCA $\rightarrow$ fit parameters only to training set. not cv or test set.

apply mapping too train cv & test sets

## Application of PCA

- Compression
  - Reduce memory/disk needed to store data.
  - Speed up learning algorithm
    $\rightarrow$ choose $k$ by % of variance retained.
- Visualization. ✓
  $\rightarrow k=2$ or $k=3$ ✓

# Dimensionality reduction - advice for applying PCA

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

Use $\underline{z^{(i)}}$ instead of $x^{(i)}$ to reduce the number of features $\underline{k < n}$ $\overbrace{\quad}^{1000} \underbrace{\quad}_{1000}$

Thus, fewer features, less likely to overfit.  **Bad !**

This might work ok, but it is not a good way to address overfitting. Use regularization instead.

$\boxed{\lambda}$ ftw.  ($\lambda$ knows what the values of $y$ are ... )
   $\leftarrow$ minimization term

(PCA)
$\rightarrow$ it throws away data, without knowing what the values of $y$ are.

PCA is sometimes used where it should not be.

Design of ML system:
  - ~~Get~~ training set ( $(x^{(1)}, y^{(1)}) \ldots (x^{(m)}, y^{(m)})$ )
  - ~~Run PCA to reduce $x^{(i)}$ in dimension to get $z^{(i)}$~~
  - $\rightarrow$ Train logistic regression on ( $(z^{(1)}, y^{(1)}), \ldots (z^{(m)}, y^{(m)})$ )
  - $\rightarrow$ test on test set: Map $x_{test}^{(i)}$ to $z_{test}^{(i)}$. Run $h_G(z)$ on
     ( $z_{test}^{(1)}, y_{test}^{(1)} \ldots z_{test}^{(m)}, y_{test}^{(m)}$ )

$\boxed{\rightarrow \text{how about the whole thing without PCA}}$

  - $\rightarrow$ try running the normal thing with original / raw data $x^{(i)}$. If the implementation is too slow etc try implementing $z^{(i)}$ etc