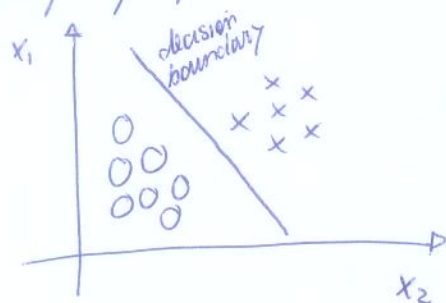# Introduction to Unsupervised learning
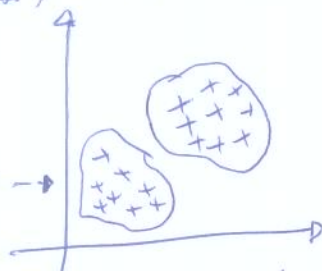
- Learning from unlabelled data.

Recall some typical examples for problems with _supervised learning_



we usually had a training set like this $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \ldots, (x^{(m)}, y^{(m)})\}$.
Our objective was to find a decision boundary separating $O$ and $x$.

In _unsupervised learning_, all we have is a training set of $\{x^{(1)}, x^{(2)}, \ldots, x^{(m)}\}$ with _no_
labels



Then we ~~develop~~ develop algorithms to find _structures_ in the _data_. The _simplest structure_
is a group. Therefore we have many algorithms to find ~~structures~~ _groups_ or _clusters_ in
our data.

Clustering has several applications such as

① Image Retrieval / Info Retrieval.
② Image classification using visual dictionaries.
③ market segmentation (find groups of similar customers).
④ Social networks (find groups of people with similar interests, ~~etc~~ coherent groups of friends, etc.)
⑤ Organization of datacenters (how to connect computers faster to attend demand, etc.)
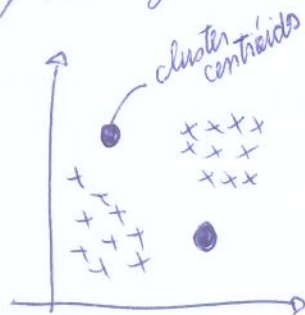⑥ Astronomical data analysis (understanding galaxy formations).

# K-means algorithm

→ most popular and widely used algorithm **to date**.



① ② Randomly initialize clusters.

cluster centroids

Suppose we want to group our data into two **clusters** (groups).
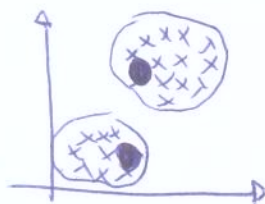


③ cluster assignment
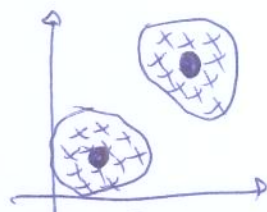
\* K-means has **two main steps**
Iterative algorithm.

① cluster assignment step

② ~~the~~ move centroids step



④ move centroids and new assignments

⑤ move centroids and new assignments

More **formally** : the K-means algorithm takes as input :

- Ⓚ (number of clusters) . ↝ for now, just assume ~~you~~ we knew this by construction.

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

- $x^{(i)} \in \mathbb{R}^n$ ( $\underline{x_0 = 1}$ by convention) used before [is **not**] used here).

# K-means algorithm

① Randomly initialize $K$ cluster centroids $\mu_1, \mu_2, \ldots, \mu_K \in \mathbb{R}^n$.

② Repeat {

cluster assignments
$\left\{\begin{array}{l} \text{for } i \leftarrow ① \text{ to } ⓜ \ \{ \\ \quad c^{(i)} \leftarrow \text{index (from 1 to } k) \text{ of cluster centroid } \underline{closest} \text{ to } x^{(i)}. \\ \quad \text{est.b.k.} \end{array}\right.$

move centroid
$\left\{\begin{array}{l} \text{for } k \leftarrow 1 \text{ to } \underline{K} \ \{ \quad \text{big K} \\ \quad \mu_k \leftarrow \underline{\text{mean of points assigned to cluster } ⓚ}. \\ \} \end{array}\right.$

}

$\quad \underset{k}{\min} \| x^{(i)} - \mu_k \|^2$  ← norm

the ② (square) here is by convention only.

For instance, suppose $\mu_2$ is the currently the center of cluster ② with the points just assigned $x^{(1)}, x^{(2)}, x^{(3)}$. Then the new centroid is going to be $\mu_2 = \sum_{i=1}^{3} \frac{x^{(i)}}{3} = \frac{x^{(1)} + x^{(2)} + x^{(3)}}{3}$;

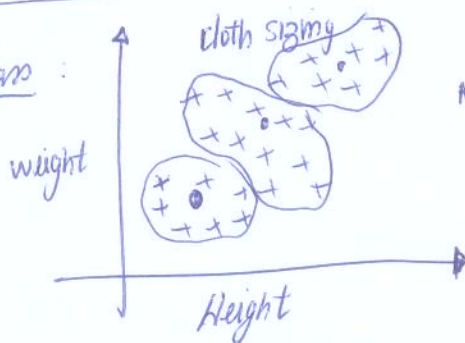$\mu_2 \in \mathbb{R}^m$ just like before.

$\mu_2 = \dfrac{\text{sum over members of } ②}{\text{size of } ②}.$

<u>Question</u>: what if a cluster doesn't have any element?

↳ Eliminate such cluster reducing the number of clusters from ⓚ to ⓚ-1.

↳ Another solution without eliminating the cluster is just reinitialize this cluster's centroid.

What happens when we don't have very well separated clusters?

Example from Andrew Ng's class:



cloth sizing

market segmentation example.

weight / Height

suppose we want ③ sizes ⓢ ⓜ, and ⓛ

# Optimization ~~step~~ objective

~ for all algorithms we studied so far, ~~step goes~~ we had an optimization function or _cost function_ even genetic algorithms have it in order to measure _adaptation_ and _fitness_.

~ K-means also _has_ a _cost function_.

## K-means optimization objective

$c^{(i)}$ = index of clusters $(1, 2, \cdots, k)$ to which example $x^{(i)}$ is currently assigned

$M_K$ = cluster centroid $k$ $(\mu_K \in \mathbb{R}^n)$.   $\}$ $\circledR$ big number of clusters an $\circledR$ little ~~is not~~ an index.

$M_{C^{(i)}}$ ~~scribbled~~ = cluster centroid of cluster to which example $x^{(i)}$ has been assigned.

for instance, suppose $x^{(i)}$ was assigned to cluster $\circled{5}$. It means then ~~that~~ that $c^{(i)} = 5$. finally, $\boxed{M_{c^{(i)}} = \mu_5}$

## Optimization objective:

$$J(c^{(1)}, \cdots, c^{(m)}, M_1, \cdots, M_K) = \frac{1}{m} \sum_{i=1}^{m} \left\| x^{(i)} - M_{c^{(i)}} \right\|^2$$

$\underbrace{\quad}$ square distance between $x^{(i)}$ and its cluster center/centroid.

therefore we want

$$\min_{\substack{c^{(1)}, \cdots, c^{(m)} \\ \mu^{(1)}, \cdots, \mu^{(m)}}} J\left(c^{(1)}, \cdots, c^{(m)}, M_1, \cdots M_K\right)$$

~ This cost function sometimes is called $\circled{Distortion}$

If we recall our K-means algorithm, we have that:
- ① step of cluster assignment minimizes $J(\cdots)$ keeping $M_1, \cdots, M_K$ fixed. Finds $c^{(1)}, \cdots, c^{(k)}$
- ② step of move centroid minimizes $J(\cdots)$ keeping $c^{(i)}$ fixed. Finds $M_1, \cdots, M_K$.

we can $\boxed{use}$ this cost function to find better clusters and avoid some _local minima_.
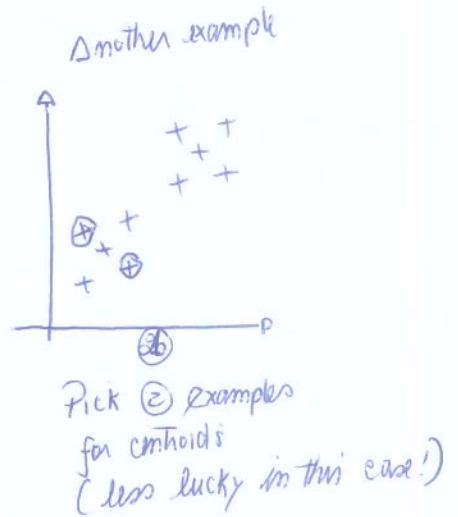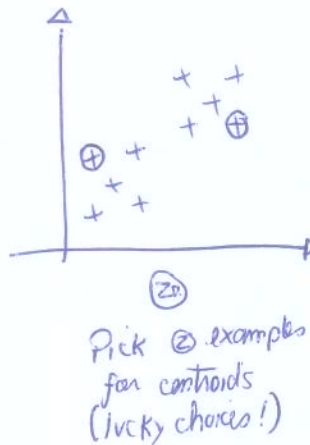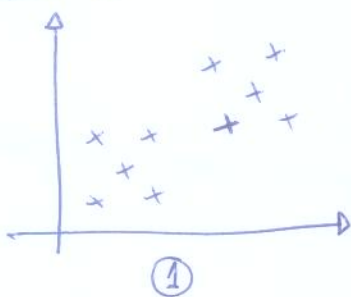
# Random Initialization

Recalling the algorithm:

① Randomly initialize $K$ cluster centroids $\mu_1, \ldots, \mu_K \in \mathbb{R}^n$.

② Repeat {

    for $i \leftarrow 1$ to $m$

      $c^{(i)} \leftarrow$ index (from $1$ to $K$) of cluster centroids closest to $x^{(i)}$

    for $K \leftarrow 1$ to $K$

      $\mu_K \leftarrow$ mean of points assigned to cluster $K$.
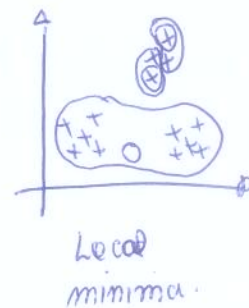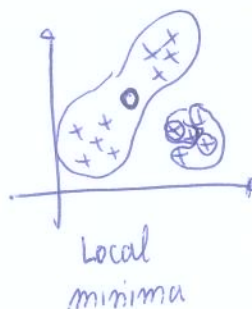
}

How to do this **Random** ?

① Should have $K < m$.

② Randomly pick $K$ training examples

③ Set $\mu_1, \ldots, \mu_K$ equal to these $K$ examples

Example $K = 2$.

Another example



①

Pick ② examples
for centroids
(lucky choices!)

Pick ② examples
for centroids
( less lucky in this case!)

$$\begin{cases} \mu_1 = x^{(i)} \\ \mu_2 = x^{(j)} \\ \quad \vdots \\ \mu_K = x^{(K)}. \end{cases}$$

One thing is obvious : depending on the initialization K-means will converge _differently_. In ~~particular~~ particular, it can converge to _local optima_.



Good
Convergence.

Local
minima

Local
minima.

↳ To avoid such cases or to diminish the chances of that happening, we can run k-means different times with different initializations.

Random initialization{

    for i=1 to 1000{

        Randomly initialize k-means.

        Run k-means.

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

        Compute cost function (distortion) $J(c^{(1)}, \ldots, c^{(m)}, \mu_1, \ldots, \mu_K)$.

    }

    Pick clustering that gave lowest cost $J(\ldots)$.

}

$\overbrace{\qquad}^{>10 \text{ for instance}}$

↳ with many clusters this doesn't mean that much but with a few clusters only, it may $\overbrace{\qquad}^{<10}$ help substantially.

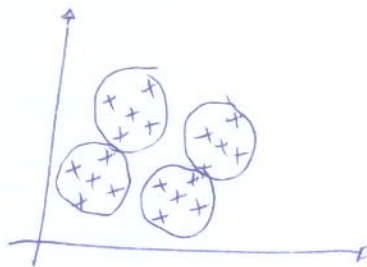Alright, but how to choose Ⓚ the number of clusters?

① The most common answer is manually through experience, using visualization techniques, analyzing the outputs of clustering, etc.

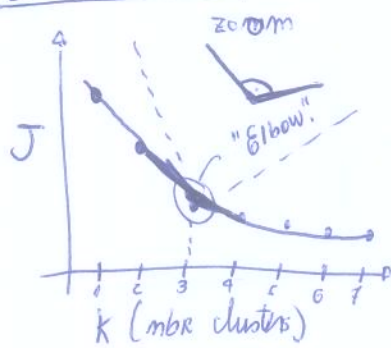The main problem here is ambiguity



VS

And this spirit of ambiguity is common in ⓈⓁ since we don't have labels to rely upon.
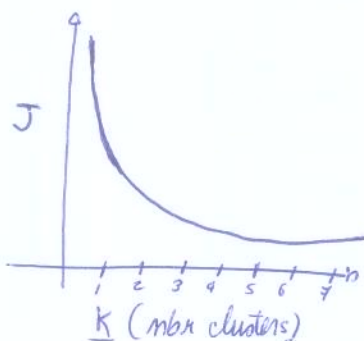
In spite of this natural ambiguity, we can use a method known as the

## Elbow method



Distortion goes rapidly down
until 3 and start stabilizing.

however, sometimes
it is ambiguous

Do not have **high expectations** with this method.

## Another way of choosing $K$.

Sometimes we use k-means for a later **purpose**. Evaluate k-means based on a metric for how well it performs for that later purpose. E.g.: image classification using visual dictionaries (choose a dictionary size and check classif. accuracy on a set of examples).

Other examples: 
- image segmentation: choose $k$, cluster, segment image, check results.
- image compression: choose $k$, compress image, measure quality.
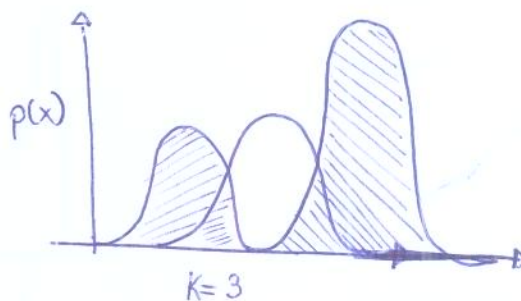
---

## Observations.

① There are many other clustering methods such as 
- k-medoids clustering
- hierarchical clustering
- agglomerative clustering
- etc.

② In k-means we had what we call a **hard assignment policy**: each data point $x^{(i)}$ can belong to only one cluster. what about relaxing this constraint? sometimes, for vectors near the boundary the hard assignment policy may be a poor choice.

③ Let's now consider a **soft-assignment policy** where the strength of the assignment depends on the distance.
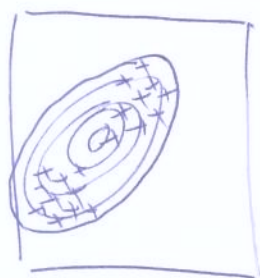
# Gaussian Mixture Model (GMM)

Combine simple models into a complex model

$$p(x) = \sum_{k=1}^{K} \underbrace{\pi_k \, \mathcal{N}(x \mid \mu_k, \Sigma_k)}_{}$$

component

↳ mixing coefficient
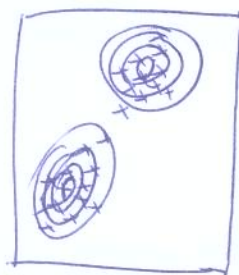


$$\forall k: \pi_k \geqslant 0, \quad \sum_{k=1}^{k} \pi_k = 1.$$



VS

single Gaussian          Mixture of two
                         Gaussians

## Cost function for fitting a GMM

for a point $x^{(i)}$,     $p(x^{(i)}) = \sum_{k=1}^{K} \pi_k \, \mathcal{N}\left(x^{(i)} \mid \overset{\rightarrow}{\mu_k}, \Sigma_k\right)$

centroids    covariances ↝ as we model as Gaussians
                           we need the Gs.

The boxed likelihood of the GMM for (M) points (assuming independence) is

$$\prod_{i=1}^{N} p(x^{(i)}) = \prod_{i=1}^{M} \sum_{k=1}^{K} \pi_k \, \mathcal{N}(x^{(i)} \mid \overset{\rightarrow}{\mu_k}, \Sigma_k).$$

→ we can simplify this in order to optimize by taking the (negative) log-likelihood

$$\mathcal{L}(\theta) = -\sum_{i=1}^{m} \ln \sum_{k=1}^{K} \pi_k \, \mathcal{N}(x^{(i)} \mid \overset{\rightarrow}{\mu_k}, \Sigma_k).$$

where $\theta$ are the parameters we wish to estimate ( $\overset{\rightarrow}{\mu_k}$ and $\Sigma_k$ ).

$\mathcal{L}(\theta)$ is our $J(\theta)$ here written as $\mathcal{L}(\theta)$ to remember its association with
the term log-likelihood.

To minimize $L(\Theta)$ or $J(\Theta)$, we differentiate w.r.t. $\mu_k$

$$\frac{\partial J(\Theta)}{\partial \mu_k} = \sum_{i=1}^{m} \underbrace{\frac{\pi_k \cdot \mathcal{N}(x^{(i)} | \vec{\mu}_k, \Sigma_k)}{\sum_{j=1}^{k} \pi_j \, \mathcal{N}(x^{(i)} | \vec{\mu}_j, \Sigma_j)}}_{\gamma_{ik}} \overset{\text{covariance}}{\Sigma_k^{-1}} (x^{(i)} - \vec{\mu}_k)$$

### Rearranging

$$\sum_{j=1}^{m} \gamma_{ik} \vec{\mu}_k = \sum_{j=1}^{m} \gamma_{ik} x^{(i)}.$$

therefore $\quad \vec{\mu}_k = \frac{1}{M_k} \sum_{i=1}^{M} \gamma_{ik} \cdot x^{(i)} \qquad$ // weighted mean

$\underset{\text{examples in group k}}{M_k}$

where

$$\gamma_{ik} = \frac{\pi_k \, \mathcal{N}(\vec{x}^{(i)} | \vec{\mu}_k, \Sigma_k)}{\sum_{j=1}^{k} \pi_j \, \mathcal{N}(x^{(i)} | \vec{\mu}_j, \Sigma_j)} \quad , \quad M_k = \sum_{i=1}^{m} \gamma_{ik}.$$

~ and $\gamma_{ik}$ are called _responsibilities_ of mixture component $k$ for vector $\vec{x}^{(i)}$.

$M_k$ is the effective number of vectors assigned to group $k$.

~ $\gamma_{ik}$ play a similar role to the assignment variables 1 or 0 multiplying a group in hard K-means. There we could think of $\gamma_{ik}$ as 1 if a datapoint is in group $k$ and zero to the other groups.

Differentiating $\Sigma_k$ gives

$$\Sigma_k = \frac{1}{M_k} \cdot \sum_{i=1}^{m} \gamma_{ik} \left( \vec{x}^{(i)} - \vec{\mu}_k \right) \left( \vec{x}^{(i)} - \vec{\mu}_k \right)^T \qquad \text{weighted covariance}$$

and w.r.t. $\pi_k$ enforcing the constraint $\sum_{k} \pi_k = 1$ gives $\pi_k = \frac{M_k}{M}$.

sum

which is the average responsibility for the component.

# Expectation maximization (EM) algorithm

Step ① : _Expectation_
 Compute responsibilities using current parameters $\vec{\mu}_K, \Sigma_K$ (assignment)

$$\gamma_{ik} = \frac{\pi_k \cdot N\left(\vec{x}^{(i)} \mid \vec{\mu}_K, \Sigma_K\right)}{\overset{K}{\underset{j=1}{\sum}} \pi_j N\left(\vec{x}^{(i)} \mid \vec{\mu}_j, \Sigma_j\right)}$$

n.b.k
groups

Step ② : _Maximization_ : Re-estimate parameters using computed responsibilities

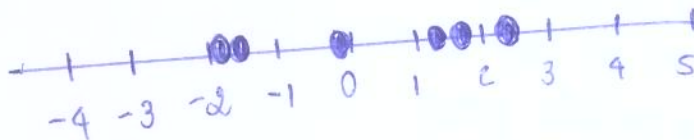$$\vec{\mu}_K = \frac{1}{M_K} \sum_{i=1}^{M} \gamma_{ik} \vec{x}^{(i)}$$

$$\Sigma_K = \frac{1}{M_K} \sum_{i=1}^{M} \gamma_{ik} \left(\vec{x}^{(i)} - \vec{\mu}_K\right) \cdot \left(\vec{x}^{(i)} - \vec{\mu}_K\right)^T$$

$$\pi_K = \frac{M_K}{M} \quad \text{where} \quad M_K = \sum_{i=1}^{M} \gamma_{ik}.$$
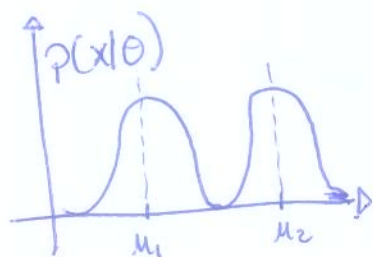
Repeat until convergence.

# Example in ①D :

Data : $x = (x_1, x_2, \ldots, x_M)$.



Objective : fit GMM with $k=2$ components

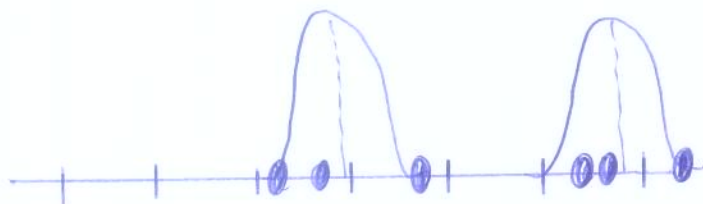Model : $p(x_i \mid \Theta) = \sum_K \pi_k \cdot N(x_i \mid \vec{\mu}_K, \sigma_K)$ where $\sum_{k=1}^{K} \pi_k = 1$
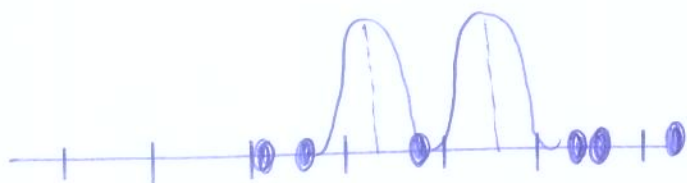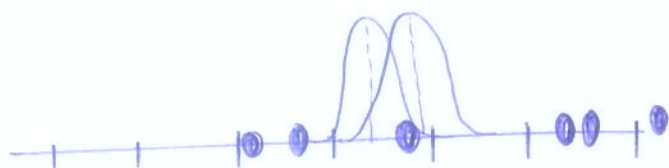
Parameters : $\Theta = \{\pi, \mu, \sigma\}$. keep $\pi, \sigma$ fixed and estimate $\mu$.
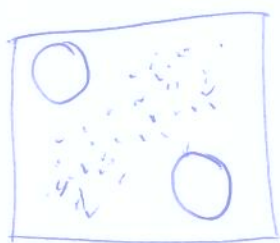
$p(x|\theta)$

$\mu_1$  $\mu_2$

Intuition of EM :

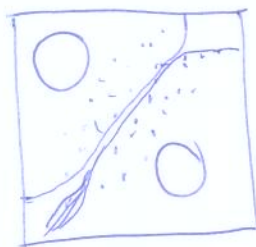E-Step : compute soft assignment of the points using current $\theta$.

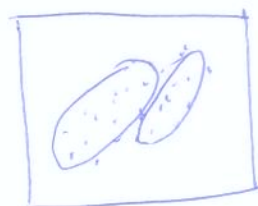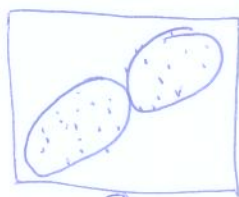m-step : update parameters $\theta$ using current responsibilities.

2-D example

① ② ③ ④

⑤

## Practical Aspects:

1. Normally we use K-means to initialize $EM$ algorithms

2. choice of $K$

3. Also can converge to local min.

————————— x —————————

# Dimensionality Reduction

Two main reasons:
- Data compression { - for speeding up learning algorithms
  - compress data for memory efficiency purposes.
- Visualization
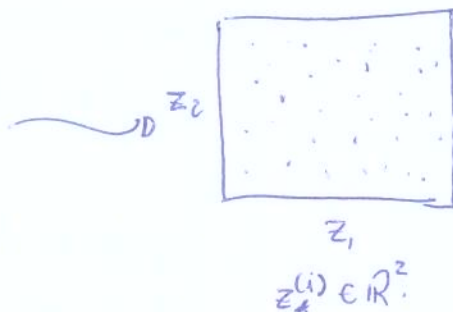
Suppose we have a dataset with <u>many</u> features (here are two)



- Reduce data from ②D to ①D.
- Data may be redundant (cm × meters × inches from different sources).

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

# Data compression 3D ↝ 2D



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

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

# Data visualization:

Suppose we have a dataset such that each $x^{(i)} \in \mathbb{R}^{100}$ and we want to quickly visualize some of its <u>properties</u>:

| país | GDP(tri) | GDP ($) per capita | ... | GINI |
|------|--------|----------|-----|------|
| canada | 1.6 | 39.2 | | 32 |
| china | 5.9 | 7.5 | | 46 |
| : | | | : | |
| singapore | 0.22 | 56.7 | | 42 |
| USD | 14.5 | 46.8. | | 40 |

$x^{(i)} \in \mathbb{R}^{100} \leadsto z^{(i)} \in \mathbb{R}^2.$

※ In $\mathbb{R}^2$ it is up to use to understand what the data is telling us (what each $z^{(i)}$ means).

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

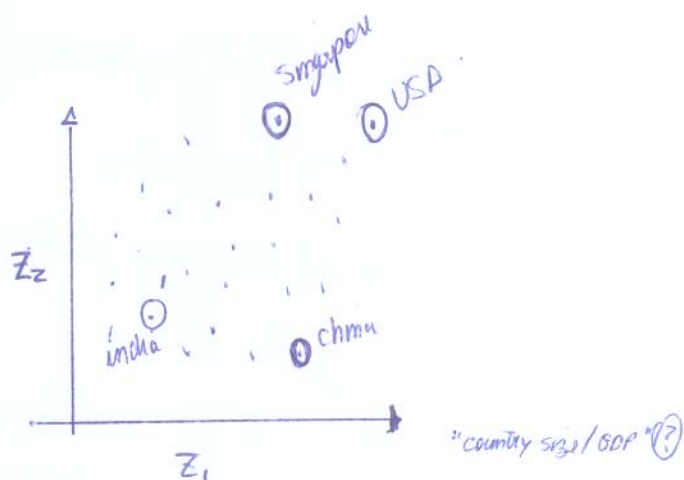| Country | $z_1$ | $z_2$ |
|---|---|---|
| Canada | 1.6 | 1.2 |
| China | 1.7 | 0.3 |
| ⋮ | | |
| Singapore | 0.5 | 1.7 |
| USA | 2 | 1.5 |



Per person GPP / economic activity (?) — $z_2$ axis

"country size / GDP" (?) — $z_1$ axis

Singapore ⊙   ⊙ USA

india ⊙   ⊙ china

---

Before going into details about Dimensionality Reduction, let's see why it is __important__.

## The curse of dimensionality

① Term coined by Bellman in 1961.

② Refers to the problem associated with _multivariate_ data analysis as the dimensionality increases.

Consider a simple example: a 3-class pattern recognition problem.

    ① A simple approach would be divide the feature space into uniform __bins__.

    ② Compute the ratio of examples for each class at each bin and,

    ③ for a new input, find its bin and choose the predominant class in that bin

In our toy problem, we decide to start with one single feature and divide the __real line__ into ③ segments



o o □ , □ □ △ , □ △ △      $x_1 \in \mathbb{R}$.

⟿ After doing this, we notice there exists __overlap__ among the __classes__. The solution? Increase complexity / add more features. Let's start in $\mathbb{R}^2$.

We decide to preserve the granularity of each axis. Hence the number of bins goes from ③ in $\mathbb{R}$ (1D) to $3^2 = 9$ in $\mathbb{R}^2$ (2D).
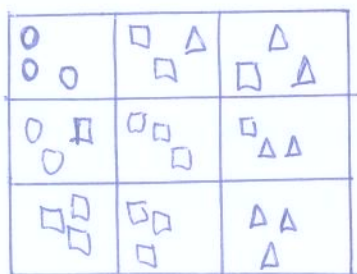
Now we need to decide
- maintain the density of examples per bin — (A)
  or
- keep the number of examples we had for 1D — (B)

<u>In case of (A)</u>

Increases the number of examples from $m = 9$ to $m = 27$ (2D)
(1D)

<u>In case of (B)</u>

Results in a very sparse 2D scatter plot.



Cte density          Const. # examples

what about moving to $\mathbb{R}^3$?
① number of bins $\to 3^3 = 27$.
② cte density $\to 81$ needed examples
③ for the same number of examples (cte # examples), 3D feature space is almost empty.

_____

\* Our approach to divide the sample space is naive but even with more interesting solutions, the sparsity problem still <u>exists</u>.

\* How to deal with the curse of <u>dimensionality</u>?

① By incorporating prior knowledge

② By providing increasing smoothness of the objective function ( e.g. regularization)

③ By reducing dimensionality.

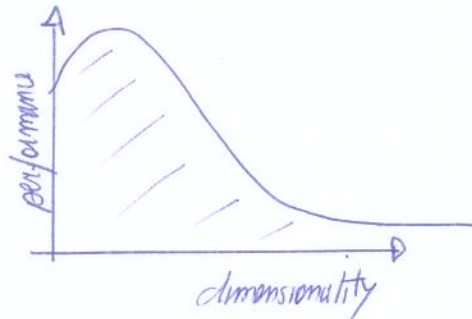In practice, _curse of dim_ means there is, for a given sample size (dataset of a given, size)
a max # of features above which the performance of a classifier will degrade instead of
improve.



## Implications of the curse of dimensionality

① Exponential growth in # of examples required to maintain a given sampling density.
for $M$ examples/bin and $D$ dimensions, # of examples is $M^D$ !

② Exponential growth in the complexity of the _objective function_ ( a density estimate)
with increasing dimensionality.
  ↳ To learn well, a more complex target function requires denser sample points."
  ↳ "a function in high-dimensional space is likely to be more complex than
     one in lower dimensional spaces, and those complications are harder to
     discern" Friedman.

③ what if it's not a _Gaussian_? For ①-D, we have several density functions.
  For $m$-dim spaces $\mathbb{R}^m$, only the multivariate Gaussian is _available_.
    ↳ Good to spot patterns in ①,② and ③ dimensions

④ Humans ⟨ ↳ Good to spot patterns in ①,② and ③ dimensions
            ↳ Pretty bad when $m > 3$.

How to reduce dimensionality then ?

① **feature extraction**: create a subset of new features by combining the existing ones

② **feature selection** : choosing a subset of all the features. (the ones more informative)

$$\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{bmatrix}^T \xrightarrow{\text{selection}} \begin{bmatrix} x^{(s)} \\ x_{i+1}^{(s+1)} \\ \vdots \\ x_d^{(s+k)} \end{bmatrix}^T \qquad \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{bmatrix}^T \xrightarrow{\text{extraction}} \begin{bmatrix} z_1 \\ z_2 \\ \vdots \\ z_d \end{bmatrix}^T = f\left(\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{bmatrix}^T\right)$$

Let's focus on feature extraction:

① Given a feature space $x^{(i)} \in \mathbb{R}^m$, find a mapping $z = f(x): \mathbb{R}^m \to \mathbb{R}^d$ with $d < m$.
Such that the transformed vector $z^{(i)} \in \mathbb{R}^d$ preserves (most of) the information and/or structure in $\mathbb{R}^m$.

② An optimal mapping $z = f(x)$ will be the one that results in low distortion/error.

③ In general, the optimal mapping will be a non-linear function. However, there is no systematic way to generate non-linear functions.

④ For this reason, we often are limited with linear transforms $\boxed{z = W \cdot x}$
   ↳ ② is a linear projection of ⓧ.

⑤ When the mapping is a non-linear function, the reduced space is called a $\boxed{\text{manifold}}$

$$\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{bmatrix} \xrightarrow{\text{linear feature extraction}} \begin{bmatrix} z_1 \\ z_2 \\ \vdots \\ z_d \end{bmatrix} = \begin{bmatrix} w_{11} & w_{12} & \cdots & w_{1N} \\ w_{21} & w_{22} & \cdots & w_{2N} \\ \vdots & \vdots & & \vdots \\ w_{d1} & w_{d2} & \cdots & w_{dm} \end{bmatrix} \begin{bmatrix} x_1 \\ \vdots \\ \vdots \\ x_m \end{bmatrix}$$

⑥ (MLPs) multi-layer perceptions can help with non-linear **mappings**.

① The selection of the feature extraction mapping $z = f(x)$ is guided by an objective function to be <u>minimized</u> or <u>maximized</u>.

② Depending on the criteria; feature extraction techniques can be of one out of two categories:
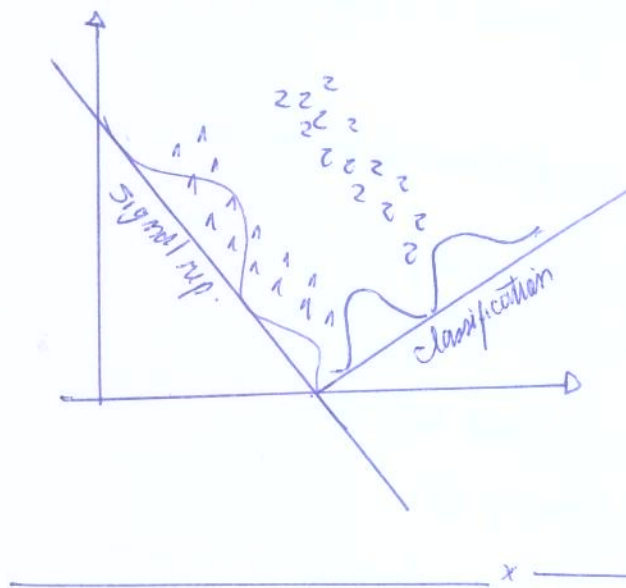
    (a) <u>Signal representation</u> : objective is to represent samples accurately in ↓ lower dimensional space

    (b) <u>Classification</u> : ~~objec~~ objective is to <u>enhance</u> class discrimination in the lower dimensional space.

<u>Two ~~It tt~~ techniques are highlighted</u> :
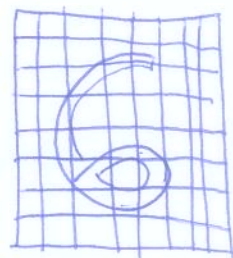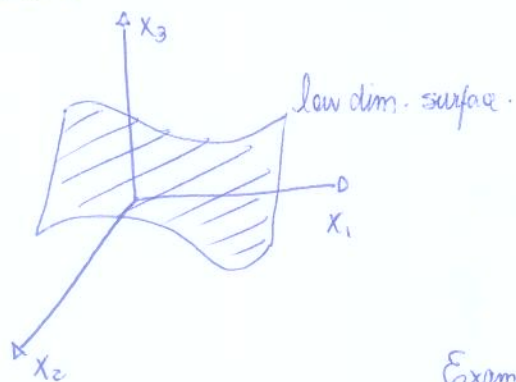
    ① Principal component analysis (example of (a)).

    ② Linear Discriminant Analysis (example of (b)).

# Why reduce dimensionality - Example

① Intrinsic dimension of data: often data is measured in high dimensions but _its_
actual variation lies on a _low dimensional surface_ (plus noise).



low dim. surface.

**Example**     $64 \times 64$ sub-image $\rightarrow \{0, 1\}^{4096}$

① there is irrelevant info/noise.

② feature extraction × selection: what would be
better?

# Projection to lower dimensions

$$\mathbb{R}^m \rightsquigarrow \mathbb{R}^d \quad \text{where} \quad m \gg d \quad \text{and often} \quad d = 2 \ (e.g., \text{when visualizing data}).$$

$$\big[ \ \big] = \big[ d \times m \big] \big[ \ \big].$$

# Principal Component Analysis

① Determine a set of (orthogonal) axes which best represent the data



a linear projection     $\big[ \ \big] = \big[ d \times m \big] \big[ \ \big].$

<u>Steps</u>   determine a set of (orthogonal) axes which best represent the data

<u>Step 1</u> : Compute the vector to the data centroid $\vec{c}$

<u>Step 2</u>: Compute the principal axes $u_\alpha$.

<u>In more details</u>

(1)  Scale and normalize the data.

$$\vec{c} = \frac{1}{N} \sum_{i=1}^{M} x^{(i)}$$

and Transform the data so that $\vec{c}$ becomes the new origin

$$\vec{x}^{(i)} \leftarrow \vec{x}^{(i)} - \vec{c}$$

optionally we can <u>scale</u> (z-norm, max-min, etc)

$$x^{(i)} \leftarrow \frac{x^{(i)} - \vec{c}}{\sigma} \quad \text{or} \quad x^{(i)} \leftarrow \frac{x^{(i)} - \vec{c}}{x_{max} - x_{min}}$$

(2)  Compute the first principal axis: determine the direction that best explains the data. Find a direction (unit vector) $\vec{u}$ such that

$$\sum_i \left( x^{(i)} - (\vec{u}^T \cdot x^{(i)}) \vec{u} \right)^2 \quad \text{is minimized or}$$

$$\sum_i \left( \vec{u}^T x^{(i)} \right)^2 \quad \text{is maximized. This is the direction of } \underline{max} \text{ variation.}$$

To enforce $\|\vec{u}\| = 1$, let's introduce a Lagrange multiplier and find the stationary point of

$$J = \sum_i \left(u^T x^{(i)}\right)^2 + \lambda \left(u^T u\right).$$

with respect to $\vec{u}$.

$$J = \sum_i u^T \left(x^{(i)} x^{(i)T}\right) u + \lambda\left(1 - u^T u\right).$$

$$= u^T \sum_i \left(x^{(i)} x^{(i)T}\right) \cdot u + \lambda\left(1 - u^T u\right)$$

$$= u^T S \, u + \lambda\left(1 - u^T u\right)$$

**Explanation:**

$$\left(A^T B\right)^2 = \left(A^T B\right) \cdot \left(A^T\right.$$
$$= A^T \left(B \cdot B^T\right)\left(A\right.$$
$$= A^T B B^T A$$

where $\circled{S}$ is the $m \times m$ simmetrix matrix $S = \sum_i x^{(i)} x^{(i)T}$

Then

$$\frac{\partial J}{\partial u} = 2 S\vec{u} - 2\lambda\vec{u} \quad \text{and hence} \quad S\vec{u} = \lambda \vec{u}.$$

In other words $\circled{\vec{u}}$ is an eigen vector of $\circled{S}$. Thus the variation

$$\sum_i \left(\vec{u}^T x^{(i)}\right)^2 = \vec{u}^T S \vec{u} = \lambda \underbrace{\vec{u}^T u}_{\text{Identity}} = \circled{\lambda}$$

is maximized by the eigen vector $\circled{u_1}$ corresponding to the largest eigen value $\circled{\lambda_1}$ of $\circled{S}$. $\circled{u_1}$ is then the first principal component.

③ Now compute the next axis, orthogonal to $\circled{u_1}$.

This must be, again, an eigen vector of $\circled{S}$, since $S\vec{u} = \lambda\vec{u}$ gives all the stationary points of variation, and hence is given by $\circled{u_2}$, the eigenvector corresponding to the second largest eigen value $\lambda_2$ of $\circled{S}$. $\circled{u_2}$ is the 2nd principal component.

## Example

Data: 3 points $x^{(1)}, x^{(2)},$ and $x^{(3)} \in \mathbb{R}^3$.



$$x^{(1)} = (1, 1, 1)^T.$$
$$x^{(2)} = (2, 2, 1)^T$$
$$x^{(3)} = (3, 3, 1)^T.$$

Centroid $\bar{X} = (2, 2, 1)^T.$

Centered data are now

$$x^{(1)} = (-1, -1, 0)$$
$$x^{(2)} = (0, 0, 0)$$
$$x^{(3)} = (1, 1, 0)$$

writing $X = \left[ x^{(1)}, x^{(2)}, x^{(3)} \right],$ we have $X = \begin{bmatrix} -1 & 0 & 1 \\ -1 & 0 & 1 \\ 0 & 0 & 0 \end{bmatrix}.$

$$S = \sum_i x^{(i)} x^{(i)T} = XX^T = \begin{bmatrix} -1 & 0 & 1 \\ -1 & 0 & 1 \\ 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} -1 & -1 & 0 \\ 0 & 0 & 0 \\ 1 & 1 & 0 \end{bmatrix} = \begin{bmatrix} 2 & 2 & 0 \\ 2 & 2 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

Its eigen-vector decomposition is

$$S = \left[ \vec{u}_1, \vec{u}_2, \vec{u}_3 \right] \begin{bmatrix} \lambda_1 & 0 & 0 \\ 0 & \lambda_2 & 0 \\ 0 & 0 & \lambda_3 \end{bmatrix} \cdot \left[ \vec{u}_1, \vec{u}_2, \vec{u}_3 \right]^T$$

$$= \begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} & 0 \\ \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} & 0 \\ 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} 2 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} & 0 \\ \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

Then $z^{(i)} = U_1^T x^{(i)}$ and $z^{(i)} = \{-\sqrt{2}, 0, \sqrt{2}\}$ for the

points $x^{(i)}$.

Doing for one case: $U_1^T = \begin{bmatrix} 1/\sqrt{2} & 1/\sqrt{2} & 0 \end{bmatrix}$

$$x^{(1)} = \begin{bmatrix} -1 \\ -1 \\ 0 \end{bmatrix}$$

$$z^{(1)} = \begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} & 0 \end{bmatrix} \begin{bmatrix} -1 \\ -1 \\ 0 \end{bmatrix} = -\frac{1}{\sqrt{2}} - \frac{1}{\sqrt{2}} = -\frac{2}{\sqrt{2}} = -\frac{2\sqrt{2}}{2} = \boxed{-\sqrt{2}}$$

## Now we are ready for the algorithm

Given data $\{x^{(1)}, \ldots, x^{(m)}\} \in \mathbb{R}^n$

① Compute the centroide, center the data and optionally __scale__ it

$$\vec{c} = \frac{1}{M} \sum_{i=1}^{M} \vec{x}^{(i)}.$$

and transform the data so that $\vec{c}$ becomes the new origin

$$\vec{x}^{(i)} \leftarrow x^{(i)} - \vec{c}.$$

② write the centered data in matrix form $X = \begin{bmatrix} \vec{x}^{(1)}, & \vec{x}^{(2)}, & \vec{x}^{(3)}, & \ldots, & \vec{x}^{(m)} \end{bmatrix}$ and

compute the covariance matrix ⑤

$$S = \frac{1}{M} \cdot \sum_{i} x^{(i)} x^{(i)T} = \frac{1}{M} \cdot XX^T.$$

③ Compute the eigen decomposition of ⑤

$$S = UDU^T$$

④ Principal components are the columns $\vec{u}_i$ of $U$ ordered by the magnitude of the eigenvalues in ⑩.

⑤ The dimensionality of the data is reduced by the projection: $\vec{z} = U_d^T X$ where $U_d$ are the first $d$ columns of ⑩ and $\vec{z}$ is a $d$-vector.

Observations

① PCA is a linear transformation that rotates the data so that it is maximally decorrelated.

② A limitation of PCA is the linearity — it can't (fit) curved surfaces well.

How much is lost by PCA approximation

The eigenvectors $U$ provide an orthogonal basis for any $\vec{x} \in \mathbb{R}^m$.

$$X = \sum_{j=1}^{m} (\vec{U}_j^T \vec{x}) \vec{U}_j \qquad * \text{ matrices } U, X, U.$$

The PCA approx with $d$ components is

$$\tilde{X} = \sum_{j=1}^{d} (U_j^T X) U_j$$

So, the error is

$$X - \tilde{X} = \sum_{j=d+1}^{m} (U_j^T X) U_j$$

Using $U_j^T U_k = \delta_{jk}$, the squared error is

$$\|X - \tilde{X}\|^2 = \left( \sum_{j=d+1}^{m} (U_j^T \cdot X) U_j \right)^2 = \sum_{d+1}^{m} U_j^T (XX^T) U_j$$

Hence the mean squared error (mse) is

$$\frac{1}{M} \sum_{i}^{M} \|x^{(i)} - \tilde{x}^{(i)}\|^2 = \frac{1}{M} \sum_{d+1}^{N} U_j^T \cdot \left( x^{(i)} x^{(i)T} \right) U_j = \sum_{d+1}^{m} U_j^T \cdot S \cdot U_j$$

and since $S U_j = \lambda U_j$,

$$\boxed{\frac{1}{M} \sum_{i}^{M} \|x^{(i)} - \tilde{x}^{(i)}\| = \sum_{d+1}^{m} U_j^T S U_j = \sum_{d+1}^{m} \lambda_j}$$

(1) what the last equation shows is that the squared reconstruction error is given by the (sum) of _eigenvalues_ of the _unused eigenvectors_.

————————— X —————————