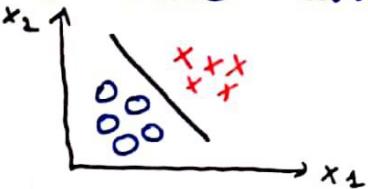


## W8 - Unsupervised Learning Intro

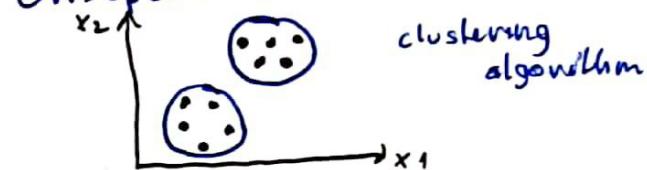
1

- o İlk unsupervised learning algoritmimiz olan Clustering'den bahsedeceğiz.
- o Bildiğimiz gibi, we learn from unlabeled data for unsupervised learning.

### Supervised Learning



### Unsupervised Learning



- o Training Set:  $\{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$
- o Goal is to find a decision boundary that separates positive and negative labeled examples.

- o Training Set:  $\{x^{(1)}, x^{(2)}, x^{(3)}, \dots, x^{(m)}\}$
- o No labels ( $y$ 's).
- o Our goal is to make algorithm find some structure in the data.
- o Algorithm that finds classes like above is called a clustering algorithm

- o Our first unsupervised learning algorithm is Clustering. Although there will be other types of unsupervised learning algorithms that finds other types of structure or other types of patterns in the data other than clusters.

### What is clustering good for?

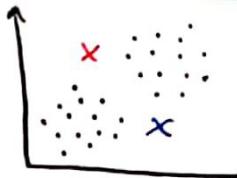
- o Market Segmentation (ticiliyor, gruplara ve farklı davranışları)
- o Social Network Analysis (etickeciyor, grupları)
- o Organize Data Centers (bilgisayarlar ne tenehde birlikte çalışır)
- o Astronomical Data Analysis (galaksi formasyonunu理解)

- o In the next video we'll start to talk about specific clustering algorithm.

## W8 - Clustering - (K-means algorithm)

- In the clustering problem we are given an unlabeled data-set and we would like to have an algorithm automatically group the data into coherent subsets or into coherent clusters for us.
- The K-means algorithm is by far the most popular and widely used clustering algorithm.

### K-MEANS ALGORITHMS

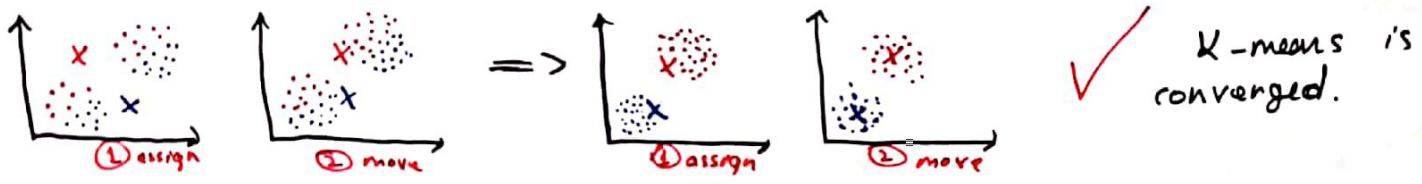


- Diyelim ki soldaki gibi bir data set olsun. Bunu iki cluster'a ayırmak istiyoruz.
- K-means algorithm, önce iki adet noktası randomly initialize eder. Bu nedenle cluster centroids denir.

- K-means is an iterative algorithm and it does 2 things.

Step 1: cluster assignment step (her training ex'i yakınığını gide bir gruba dahil eder)

Step 2: move centroid step (centroid'i yeni elemanların mean'ine laşır)



Input:

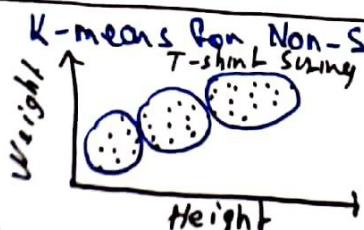
- $K$  (number of clusters)
- Training set  $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$

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

what K-means algorithm does:

- Randomly initialize  $K$  cluster centroids  $M_1, M_2, \dots, M_K \in \mathbb{R}^n$
- Repeat {
  - cluster assignment
    - for  $i = 1$  to  $m$   $\rightarrow \|x^{(i)} - M_k\|^2$  nin hangi k için min olduğu bilenir ve  $c^{(i)}$ 'ye bu değer atanır!
    - $c^{(i)}$  : index (from 1 to  $K$ ) of cluster centroid closest to  $x^{(i)}$
  - move centroid
    - for  $k = 1$  to  $K$ 
      - $M_{k+}$  = average(mean) of points assigned to cluster  $k$ 
        - $\rightarrow$  Diyelim ki  $M_2$  centroidi ile  $x^{(1)}, x^{(5)}, x^{(6)}, x^{(10)}$  ait olsalar  $c^{(1)}=2, c^{(5)}=2, c^{(6)}=2, c^{(10)}=2$
        - $\rightarrow$  Then  $M_2 = \frac{1}{4} [x^{(1)} + x^{(5)} + x^{(6)} + x^{(10)}] \in \mathbb{R}^n$

- ?) what if there is a cluster with "0" points assigned to it? In that case we can eliminate that centroid (common thing to do.). Second option is you can just randomly reinitialize that cluster centroid.



T-shirt sınıflarının farklı meşalelerin boyutunu baki 3 şebeke olarak 3 cluster.

- Data işi ayıratmaya acık gibi görünse de K-means ayırmak istiklal kümeye göre small, diğerine göre large kümeye göre large t-shirt lasırlayabiliriz.

## W8 - Clustering Optimization Example

(3)

Most of the supervised learning algorithms we've seen (lin. reg., log. reg. and so on) have an optimization objective or some cost function that the algorithm was trying to minimize. K-means also has a cost function that it's trying to minimize.

It helps us in 2 ways: 1st: It will help us debug the learning algorithm. 2nd: To help k-means find better costs for this and avoid the local optima. We'll see 2nd one later.

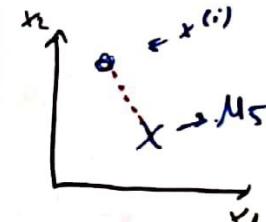
While K-means is running we gonna keep track of 2 sets of variables:

- $c^{(i)}$ : index of cluster ( $1, 2, \dots, K$ ) to which example  $x^{(i)}$  is currently assigned.
- $M_k$ : cluster centroid  $k$  ( $M_k \in \mathbb{R}^n$ )  $\rightarrow k: \# \text{ of cluster (total)}$   
 $\rightarrow k: 1, 2, \dots, K$

- $M_{c(i)}$ : to denote cluster centroid of cluster to which example  $x^{(i)}$  has been assigned.  
If  $x^{(i)}$  is assigned to cluster 5. Then  $c^{(i)} = 5$ . Thus  $M_{c(i)} = M_5$

### Optimization Objective

$$\boxed{\begin{aligned} J(c^{(1)}, \dots, c^{(m)}, M_1, \dots, M_K) &= \frac{1}{m} \sum_{i=1}^m \|x^{(i)} - M_{c(i)}\|^2 \\ \min_{\substack{c^{(1)}, \dots, c^{(m)} \\ M_1, \dots, M_K}} J(c^{(1)}, \dots, c^{(m)}, M_1, \dots, M_K) \end{aligned}}$$



Nasıl inceden best  $\Theta$  buluyorsak burada da aynı best  $c^{(1)}, \dots, c^{(m)}, M_1, \dots, M_K$  buluyoruz  
This cost function is also called Distortion

K-means Algorithm was:  
Randomly initialize  $K$  cluster centroids  $M_1, M_2, \dots, M_K \in \mathbb{R}^n$   
Repeat {  $\uparrow$  Cluster assignment: minimize  $J(\dots)$  w.r.t.  $c^{(1)}, c^{(2)}, \dots, c^{(m)}$   
holding  $M_1, M_2, \dots, M_K$  fixed }

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

for  $k = 1$  to  $K$   
 $M_k :=$  mean of points assigned to cluster  $k$   
 $\downarrow$  Move centroid: minimize  $J(\dots)$  w.r.t.  $M_1, \dots, M_K$

(\*) Burada x falan yok dolayısıyla Cost Function asla on tane. 660 adında osalır.

Yani cost function aslında her training ex.'in kendi centroidine olan uzaklığının karesinin toplamı. Buju minimize etmek işliyor.

W8 - Clustering Random Initialization

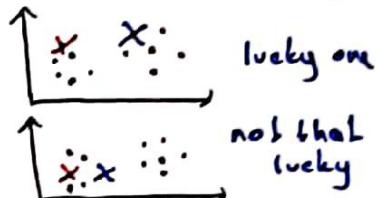
- We'll talk about how to initialize K-means and more importantly this will lead us into a discussion how to make K-means avoid local optima as well.

K-means algorithm'ın ilk adımı guydu:

- Randomly initialize K cluster centers:  $M_1, M_2, \dots, M_K \in \mathbb{R}^n$
- Bu sistem için birebir doğru değil var, fakat one method is much more recommended than the others. Let's talk about that option:

Random Initialization

Let's say  $K=2$



- Should have  $K < m$ . Randomly pick K training examples
- Set  $M_1, \dots, M_K$  equal to these K examples.
- $M_1 = x^{(i)}$ ,  $M_2 = x^{(j)}$  randomly chosen examples.

- Thus K-means can end up different solutions depending on the initialization.
- So depending on the initialization K-means can end up at Local Optima or Global Optima (desired):



- Local Optima dedigimiz  $J(c^{(1)}, \dots, c^{(m)}, M_1, \dots, M_K)$ 'nın local optima'sı.

- Thus, to avoid getting stuck in bad local optima's we can try multiple random initializations.

Way of applying multiple random Initialization:

- for  $i = 1$  to 100
  - Randomly initialize K-means
  - Run K-means. Get  $c^{(1)}, \dots, c^{(m)}, M_1, \dots, M_K$
  - Compute cost function (distortion) :  $J(\dots)$  (100 farklı  $J(\dots)$  buluyor)

}

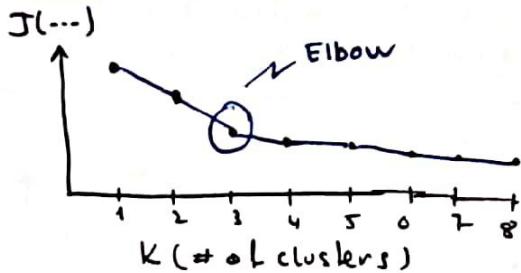
- Pick clustering that gave lowest cost  $J(\dots)$

- If  $K=2-10$  doing multiple random initialization can sometimes make you find a better local optima. But if  $K$  is large having multiple random initialization is less likely to make a huge difference, muhtemelen your data işin sonuc verir ama bir çok farklı olmasın sor.

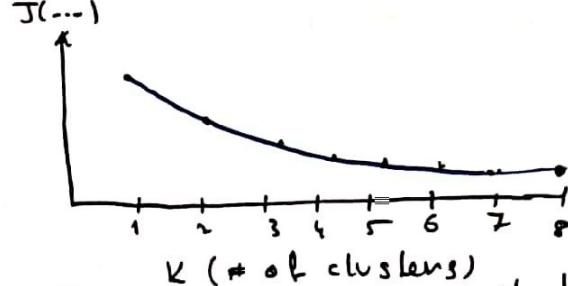
① To be honest, there isn't a great way of automatically choosing  $K$ . Most common way is still choosing it manually by looking at visualizations or by looking at the output of the clustering algorithm or something else.

② Bir data'ya baksak bile koc cluster olmasa gerekliginin mutlak bir cevabi yoktur.

**Elbow Method to Choose  $K$ :** People sometimes use it to determine  $K$ .



③ Elbow'a kadar distortion goes down rapidly then goes down slowly. Maybe  $K=3$  is the right choice.



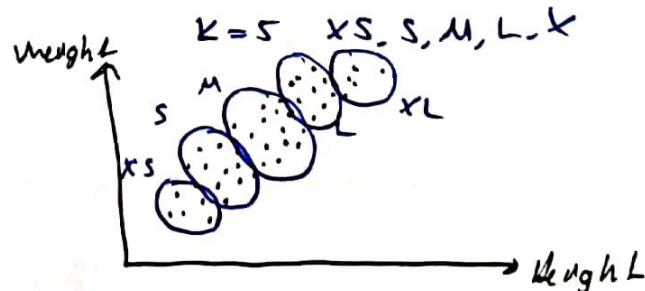
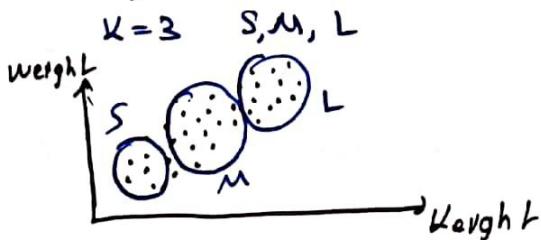
④ Genelde böyle bir graph ike karşılaşırsınız burada kesişen bir elbow yok. Söldənə plotta karışıklığınlığı bir sey kolayca seegebilirsin.

~~⑤~~ **Summary of Elbow Method:** Graph'i dekənə soldəki gibi çəkəsa nə ələ fakat it is a long shot.

~~⑥~~ **Graph'ın şübhə olsalması MUST deyil belki  $K=5$  rən  $K=3$ 'tən daha yüksək bir cost- vənətblir. Bu demədə ki kətə bir local minimum'a təkildi, random initiallərlən multiple dəneyip  $K$ -means'i təkərən fəaliyyətdən.**

**Another Method to Choose  $K$ :** Sometimes, you're running  $K$ -means ~~to~~ to get ~~clusters~~ clusters to use for some later/downstream purpose. Evaluate  $K$ -means based on a metric for how well it performs for that later purpose.

⑤ Let's say you wanna use  $K$ -means for market segmentation like in the T-shirt sizing example before.



~~⑥~~ Burada onluk kənd vəməh van hangisinin dəha əyr oldugu na bəsiləmə cluster olun məsələkən dəha vəyi mi cysun yoxsa bəzi məsələlərənən aradığın, bulamasisın ama dəha ucuz mi maledəlim?

We'll start talking about a second type of unsupervised learning problem called dimensionality reduction.

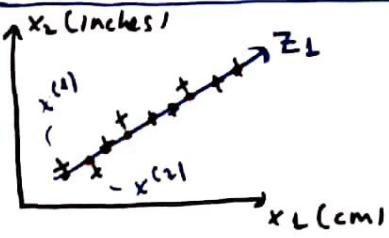
There are couple of different reasons why one might want to do dimensionality reduction:

### Motivation I: Data Compression

- Data compression allow us to use less computer memory and also speed up our learning algorithms.

Let's talk about what is dimensionality reduction?

2D to 1D



- Bir bilmiyor ama bir çok feature içinden  $x_1$  ve  $x_2$  aynı sadexe birimi farklı. This gives us a highly redundant representation.
- Aynı feature'ı ölçmeli olduğum hem  $x_1$  hem  $x_2$  kullanmak yerine tek feature ile işe halledebilirim.

cm ve inch verileri rounded olduğu için tam bir çizgi taripli olmuyorlar.

Dimensionality reduction iken datanın taripliği line belirlenir ve datalar tek tek bu line'a project edilir. Sonra  $x_1$ - $x_2$  yerine  $z_1$  adında bir feature elde edilir.

~~x<sub>1</sub> ile x<sub>2</sub> arasında bu kadar bir直 line straight line (inches to cm) olmasa bu straight featureler highly correlated olabilir, straight line'a yakin (matematiksel olarak olmasa sən deqil). Yine dim. red. iş yapar.~~

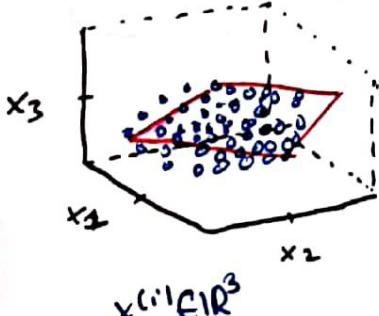
$$x^{(1)} \in \mathbb{R}^2 \rightarrow z^{(1)} \in \mathbb{R} \quad x^{(2)} \in \mathbb{R}^2 \rightarrow z^{(2)} \in \mathbb{R} \quad \dots \quad x^{(m)} \in \mathbb{R}^2 \rightarrow z^{(m)} \in \mathbb{R}$$

~~Summary: If we allow ourselves to approximate the original data set by projecting all of my original examples onto the blue line then I need only one number to specify the position of a point on the line and therefore I can represent each training example by 1 feature instead of 2. This is a approximation.~~

### 3D to 2D Reduction

10000 to 1000 gibi olsa burada orası mass 2D ve 3D

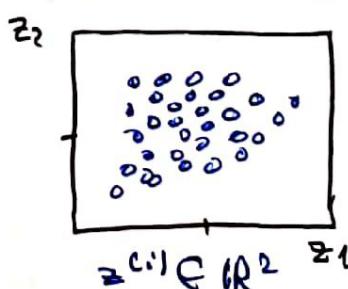
Gensehle bu reduction kullanılıyor.



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

Maybe all of these data roughly lies on the red plane.

- What we can do with dimensionality reduction is project the data to a 2D plane



- Now I can represent any training ex. by  $z_1$  and  $z_2$  instead of  $x_1, x_2, x_3$ .

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

○ In the last video we talked about dimensionality reduction for the purpose of compressing the data.

○ In this video we are gonna talk about the second application of dimensional reduction and that is to visualize the data. Data visualization bin got problems! & we'll solve them later.

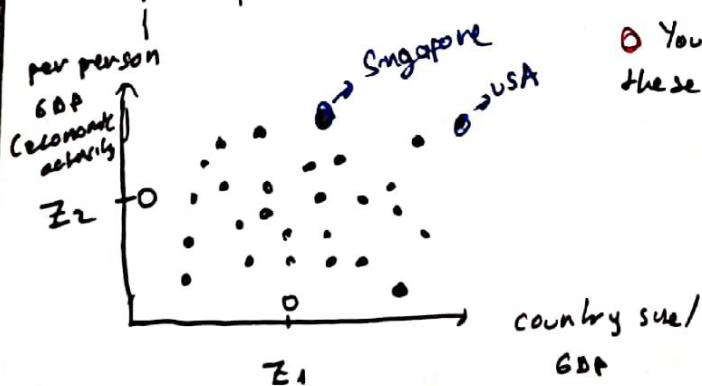
Country	$x_1$ GDP	$x_2$ Per capita GDP	$x_3$ Human Development Index	$x_4$ Life expectancy	$x_5$ Poverty Index	$x_6$ Mean household income	...	$\times \mathbb{R}^{50}$
Canada	---	---	-	-	-	-	-	
China	:	:	:	:	:	:	:	
India	:	:	:	:	:	:	:	
Russia	:	:	:	:	:	:	:	
USA	:	:	:	:	:	:	:	
;								

○ How do we visualize this data if we have 50 features for each country. It is very difficult to plot 50 dimensional data.

○ Using dimensionality reduction what we can do is instead of representing each country by  $x^{(i)} \in \mathbb{R}^{50}$  we can come up with a different feature representation  $z^{(i)} \in \mathbb{R}^2$ . If  $z_1$  and  $z_2$  somehow summarizes my numbers we can plot these countries in  $\mathbb{R}^2$  and understand the countries better.

Country	$z_1$	$z_2$
Canada	---	---
China	:	:
India	:	:

○ Reduce data from 50D to 2D  
○ When you do that, the meaning of  $z_1$  and  $z_2$  is not very clear. It is often up to us to figure out roughly what these features  $z_1, z_2$  means



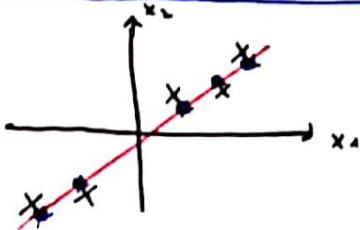
○ You might find that given these 50 features these are the 2 main dimensions of the deviation.

○ That's how we can use dimensionality reduction for visualization purpose.

○ In the next part we will start to develop a specific algorithm called PCA (Principal Component Analysis) which will allow us to do this and also the earlier application we talked about of compressing the data

- For the problem of dimensionality reduction by far the most popular algorithm is principal component analysis (PCA)
- Let's formulate precisely, exactly what we would like PCA to do.

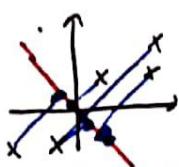
### Principal Component Analysis's Problem Formulation



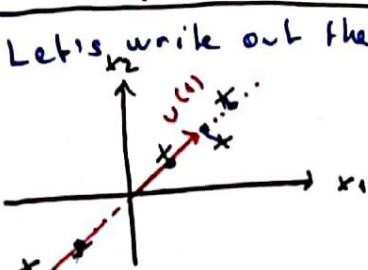
- Let's say we have a data set as here.  $x \in \mathbb{R}^2$
- I want to reduce the dimension from 2D to 1D in other words I would like to find a line onto which to project the data.

Yukarıda çizilen kırıltı çitgi iyi bir linea benzeyen ama neden?  
Çünkü siyah orijinal noktaların project edildiğinde mavi noktaların utaklıklar toplamı çok az. Bu demek oluyor ki kırıltı çitgi dimensionality reduction için iyi bir çitgi

- So what PCA does is, it tries to find lower dimensional surface (a line in this case) onto which to project the data so that the sum of squares of distances b/w original data points and projected ones is minimized.
- Gerekli data  $X$  ile projected data  $\circ$  arasındaki mesafeye basen "projection error."
- Before applying PCA, it's standard practice to first perform mean normalization and feature scaling so that features should have zero mean and should have comparable ranges of values.

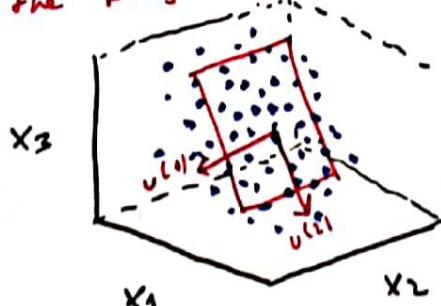


Yukarıdaki yerde yukarıda gibi bir çitgi görüldüğü gibi projection errorsları antarıyor. Bu yüzden bu seviye yarıs olur.



- Let's write out the PCA problem a little more formally.  $\rightarrow \mathbb{R}^2$  in this case
- 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.
  - I'm hoping that PCA will find the vector  $u^{(1)}$  above. So that when I project the data onto the line that I define by extending out the  $u^{(1)}$  vector, and end up with pretty small projection errors.  $\circ u^{(1)}$  veya  $u^{(1)}$  farklınes aynı line!

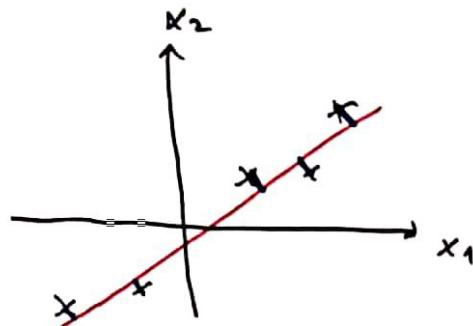
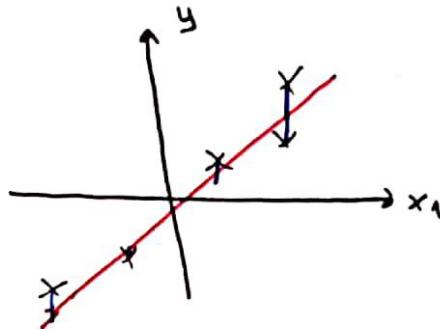
- More General Case Reduce  $n$ -dimension to  $k$ -dimension: Find  $k$  vectors  $u^{(1)}, u^{(2)}, \dots, u^{(k)}$  onto which to project the data so as to minimize the projection error.



- Burada 3D to 2D dim. red. yapılıyor.  $u^{(1)}$  ve  $u^{(2)}$  vektörleri birebir 2D surface oluşturuyorlar. Data bu surface'e project edilecek. PCA'nın görevi datanın min hatayla project edilecek surface'ı bulmaktır.
- 2D yerine  $k$ D de aynı! 2D surface yerine  $k$ D surface formüllemek oluyor.

PCA is not linear regression:

- Despite some cosmetic similarity, these are actually totally different algorithms.



- Lin. Reg.: we are trying to predict the value of some variable  $y$  given some input features  $x$ . We fit a line and minimize the squared errors b/w points and the line.
- Halaya bakarən məni cıqıllı vətq!
- In contrast in PCA, it tries to minimize the magnitude of different kind of errors (shortest distance b/w point  $x$  and red line)
- As in linear regression there is this distinguished variable  $y$  we are trying to predict whereas in PCA there is no special variable  $y$  that we are trying to predict.

Summary: PCA is trying to find a lower dimensional surface onto which to project the data, so as to minimize the squared projection error. (squared distance b/w each point and its projection)

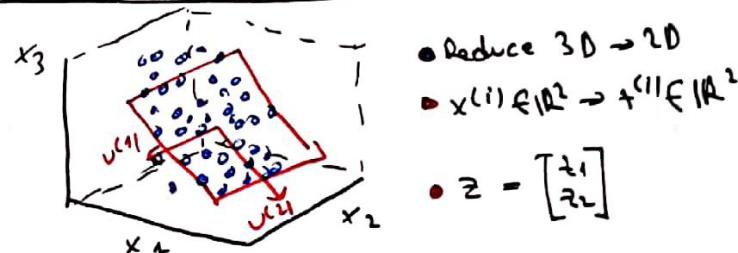
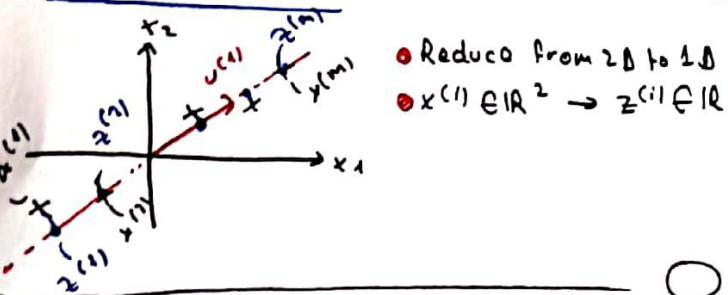
Sıradakı dərəcələr lower dimensional surfaces nasıl bulunur  
onu bəkacığın.

## -W8 - PCA Algorithm -

(5)

- Given a training set of  $m$  unlabeled examples it is important to always perform mean normalization and depending on the data maybe perform feature scaling.
- Training set:  $x^{(1)}, x^{(2)}, \dots, x^{(m)}$ 
  - $\mu_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)}$  Feature Scaling / Mean Normalization:  $\mu_j$ : mean for each feature
  - Replace each  $x_j^{(i)}$  with  $x_j - \mu_j$
  - If different features on different scales perform feature scaling

### DCA Review:



- ~~What PCA has to do~~ → Compute  $u^{(1)}, u^{(2)}, \dots, u^{(k)}$  vectors.  
 → Compute  $z^{(1)}, z^{(2)}, \dots, z^{(m)}$  points.

### PCA Algorithm

Assume we wanna reduce data from  $n$ -D to  $k$ -D

- Compute "covariance matrix" :  $\Sigma = \frac{1}{m} \sum_{i=1}^m (x^{(i)}) \cdot (x^{(i)})^T$
- Compute "eigenvectors" of matrix  $\Sigma$  :  $[U, S, V] = \text{svd}(\Sigma)$
- $U = \begin{bmatrix} | & | & | & \dots & | \\ u^{(1)} & u^{(2)} & u^{(3)} & \dots & u^{(n)} \end{bmatrix}_{n \times n}$  where  $U \in \mathbb{R}^{n \times n}$   $\rightarrow u^{(1)}, \dots, u^{(k)}$  anadığımız surface'ı oluşturur,

$k \rightarrow$  we are gonna use first  $k$  column of  $U$

- Sıradaki adım  $x \in \mathbb{R}^n$  için  $z \in \mathbb{R}^k$ 'yi bulmakız.
- $U_{\text{reduced}} = \begin{bmatrix} | & | & | & \dots & | \\ u^{(1)} & u^{(2)} & u^{(3)} & \dots & u^{(k)} \end{bmatrix}_{n \times k}$
- $z^{(i)} = U_{\text{reduced}}^T \cdot x^{(i)}$

### PCA Algorithm Summary

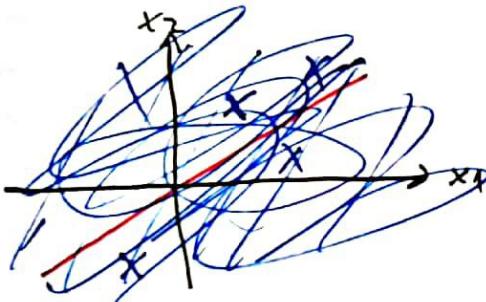
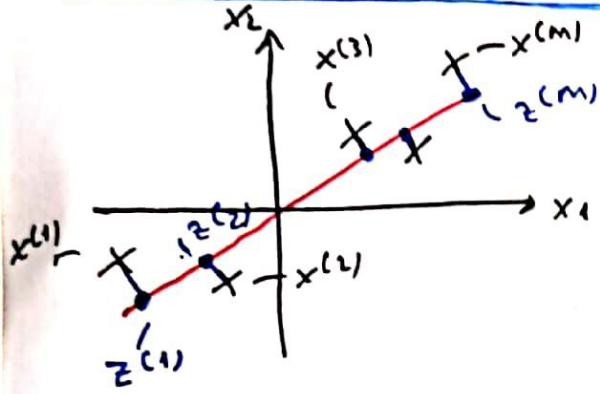
- After mean normalization (ensure every feature has zero mean / std) optionally feature scaling:
- $\Sigma = \frac{1}{m} \sum_{i=1}^m (x^{(i)}) \cdot (x^{(i)})^T$  if  $X$  is given as  $X = \begin{bmatrix} x_1^{(1)} & x_2^{(1)} & x_3^{(1)} & \dots & x_n^{(1)} \\ x_1^{(2)} & x_2^{(2)} & x_3^{(2)} & \dots & x_n^{(2)} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ x_1^{(m)} & x_2^{(m)} & x_3^{(m)} & \dots & x_n^{(m)} \end{bmatrix}_{m \times n}$  then  $\Sigma = \left(\frac{1}{m}\right)^* X^T \cdot X$
- $[U, S, V] = \text{svd}(\Sigma)$ ;
- $U_{\text{reduced}} = U(:, 1:k)$ ; // Grab first  $k$  columns
- $z = U_{\text{reduced}}^T \cdot x$ ; // Similar to K-means  $x \in \mathbb{R}^n$  no  $x_0 = 1$  convention

## W8 - Applying PCA - Reconstruction from compressed representation

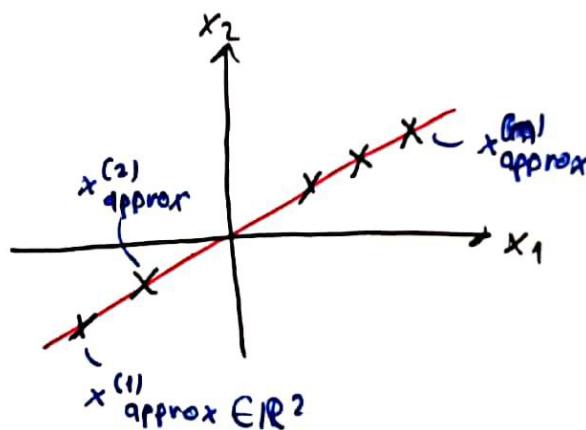
(6)

- I was talking about PCA as a compression algorithm where you may have say, a 1000-D data and compress it to 100-D feature vector.
- There should be a way to go back from this compressed representation back to an approximation of your original high-dimensional data.
- So given  $z^{(i)}$  which may be 100-Dimensional, how do you go back to the original representation  $x^{(i)}$  which maybe a 1000-Dimensional.

### Reconstruction from Compressed Representation



- We know: 
$$z = U^T \cdot x$$
- So if  $z \in \mathbb{R}^k$  is given  $\rightarrow$  what is  $x \in \mathbb{R}^n$
- $$X_{\text{approx}} = U_{\text{reduced}} \cdot z$$



(6)

## WB - Choosing the # of Principal Components

### Choosing k (number of principal components):

- Here are couple of useful concepts to choose  $k$ :

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

◦ Total variation in the data:  $\frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2$  : On average how far are my features from being zero.

- Typically choose  $k$  to be smallest value so that:

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

Baren %5, %10  
da olabilir.

"99% of variance is retained"

- In order to retain 99% of the variance you can often reduce the dimension of the data significantly and retain most of the variance. Because for most real life data-sets many features are highly correlated.

### Implementation:

#### Algorithm I

- Try PCA with  $k=1$
- Compute Ureduced,  $z^{(1)}, z^{(2)}, \dots, z^{(m)}$ ,  $x_{\text{approx}}^{(1)}, \dots, x_{\text{approx}}^{(m)}$
- Check if 99% of variance is retained:

$$\frac{\frac{1}{m} \sum \dots}{\frac{1}{m} \sum \dots} \leq 0.01 ?$$

- If 99% of var. is. ret. use current  $k$ . Else go back and increase  $k$ . So we are trying to find the minimum  $k$  that provides 99% of variance is retained.

• Very inefficient her adminder PCA uygulayınca çok fazla hesaplama var.

#### Algorithm II (Easen)

- $[U, S, V] = \text{svd}(Sigma)$

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

- For a given  $k$ , retained variance value can be calculated using the  $S$  only.  $S_{ii}$  de yalnızca basta donan bulunanlar yeterlidir:

$$\text{retained variance} = 1 - \frac{\sum_{i=k+1}^n S_{ii}}{\sum_{i=1}^n S_{ii}}$$

$$\text{Mesela } k=2 \text{ iken } 1 - \frac{S_{11} + S_{22}}{S_{11} + S_{22} + \dots + S_{nn}}$$

- Sırayla k yi artır if 99% variance retained sağlanmış dur!

## W8 - Advice for Applying PCA -

In an earlier video, we've said that PCA sometimes can be used to speed up the running time of a learning algorithm.

### Supervised Learning Speedup: Most common use of PCA.

- Let's say we have training set:  $(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots (x^{(m)}, y^{(m)})$
- And let's say examples  $x^{(i)}$  are very high dimensional  $x^{(i)} \in \mathbb{R}^{10000}$ . Mesela computer vision ile oğreniyorum diyelim  $100 \times 100$  px resimlerde uğraşıysak bu 20000 features yapar.
- Feature vectors dimension 10000 olunca running a learning algorithm might be slow. PCA ile dimensionality reduction yapip algoritmayi hızlandırmır.

#### ④ First Extract Inputs:

- Unlabeled dataset:  $x^{(1)}, x^{(2)}, \dots x^{(m)} \in \mathbb{R}^{10000}$
- Apply PCA now I have  $z^{(1)}, z^{(2)}, \dots z^{(m)} \in \mathbb{R}^{1000}$

#### ⑤ Now I have a new training set:

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

- Now I can take this reduced dimension training set and feed it to a learning algorithm maybe a NN or log-reg.  $h_{\theta}(z) = \frac{1}{1 + e^{-\theta^T z}}$
- If we have a new example maybe a new test example  $x$ , map it through the same mapping  $x \rightarrow z \rightarrow h_{\theta}(z) \rightarrow$  prediction

~~Final 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.

- Mapping yaparken bir Ureduced buluruz. Bu Ureduced is like a parameter that's learned by PCA and we should be fitting our parameters only to our training sets and not to our cross validation or test sets.
- Then having found Ureduced or having found the parameters for feature scaling where the mean normalization and the scale that you divide the features by to get them on to comparable scales. Having found all those PARAMETERS, on the training set you can then apply the same mapping to other examples maybe in your cross-validation or test sets.

~~Reducing the data from 10000 to 10000 is not unrealistic.~~  
For most of the problems we can reduce the data by 10x and still retain the most of the variance you're not barely affecting the classification accuracy.

Application of PCA Summary:

- Compression
  - Reduce memory/disk needed to store data
  - Speed up learning algorithm
- Visualization ] choose  $k=2$  or  $k=3$ .

Bad use of PCA : To prevent overfitting:

- People use  $z^{(i)}$  instead of  $x^{(i)}$  to reduce # of features to  $k < n$
  - Thus, fewer features, less likely to overfit.
  - It might actually work OK. But it's not a good way to address overfitting
- Regularization is a much better way for overfitting problem.

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

- PCA doesn't use the values  $y$ . Thus it throws some information without knowing the  $y$  values. It may throw away some valuable information. Regularization will work better because this makes  $\theta$  actually know what the values  $y$  are. Less likely to throw away some valuable information.

PCA is sometimes used where it shouldn't be. Don't use it if you don't need it:

Mesela bir ML System design ediliyor:

- Mesela bir ML System design ediliyor:
  - Get training set  $\{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$
  - Run PCA  $x^{(i)}$  in dimension to get  $z^{(i)}$  ~~X~~ Genetik olabilir!
  - Train logistic regression on  $\{z^{(1)}, y^{(1)}, \dots, (z^{(m)}, y^{(m)})\}$
  - Test on test set: Map  $x_{test}^{(i)}$  to  $z_{test}^{(i)}$ . Run  $h_0(z)$  on  $\{(z_{test}^{(1)}, y_{test}^{(1)}), \dots, (z_{test}^{(m)}, y_{test}^{(m)})\}$

→ 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  $z^{(i)}$ .

PCA is one of the most powerful unsupervised learning algorithms.