

# Unsupervised Machine Learning:

## Big Data:

1) too large  $n$  (Long data)

•) Computers are happy and easy to deal with.

But:

•) Large number of observations makes it easy to capture the slightest correlations. (very low  $P$ -value because  $n$  is high)

•) So sometimes some variables are statistically significant without being substantively significant.

2) Large number of variables (wide data)  
( $P$  is large)

•) Modeling is more complex (Computationally and Conceptually)

3)  $n$  and  $P$  are not large  
but the structure of data is very complex.

•) Long Computational process.

~~at~~

✓

## Unsupervised Learning:

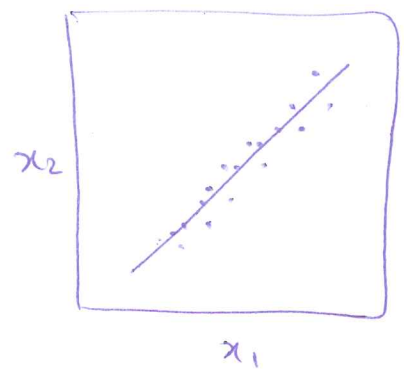
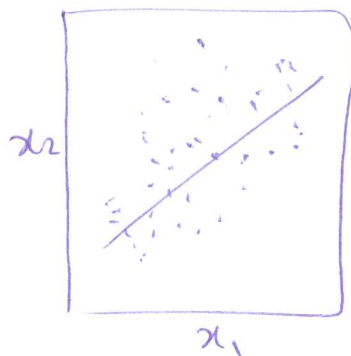
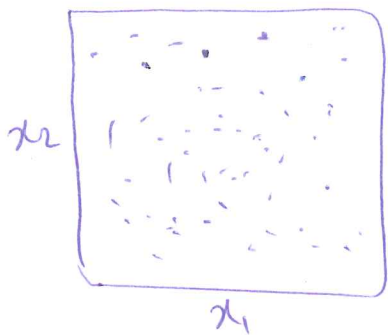
we take a large set of variables and observations and try to detect a smaller set of underlying group. (Factors & clusters) (we do not have  $y$ , all variables are  $x_1, x_2, \dots, x_p$ )

## Supervised Learning:

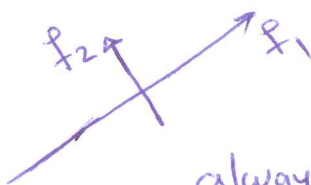
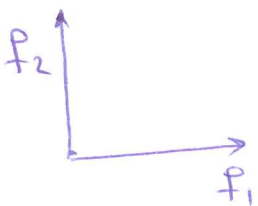
Any technique with the goal to use a large number of variables and observations to model and predict  $Y$ .

## Factor Analysis:

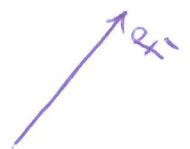
is the search for those underlying dimensions, or factors underneath a large set of variables.



Factors:



always perpendicular



Factors equations  $P^{\text{th}}$  Principal Component =  $C_p$

$$C_p = \phi_{p1} X_1 + \phi_{p2} X_2 + \phi_{p3} X_3 + \dots + \phi_{pp} X_p \quad (P-1 \text{ dimension})$$

↘ a vector containing  $n$  value for variable  $X_1$

\* There are always as many components as there are variables ( $P$ ).

\*  $\{\phi_{p1}, \phi_{p2}, \dots\}$  is the projection of the original variables onto  $C_p$ . So the larger  $\phi_{12}$  is, for instance, the larger the projection of  $X_2$  on  $C_1$ . So the larger  $\phi_{12}$  is the more similar  $X_2$  &  $C_1$  are.

\* Before any analysis, data in all variables must be scaled using normal scaling, or linear, —  
(mean=0, std=1)

\* Two study is being done but they ~~are~~ tend to produce very similar results:

a) Factor Analysis (FA)

b) Principal Component Analysis (PCA)

# Eigen Vectors and Eigen Values:

- ① a  $p \times p$  matrix has exactly  $p$  eigen vectors where each eigen vector is a vector of length  $p$ . ~~Also~~
- ② An eigen vector is defined as that vector such that, when multiplied by its associated matrix gives us a result the exact same vector, except stretched or shrunk by some value known as the eigenvalue.

$$M v_p = \lambda_p v_p$$

real number      vector  $p \times 1$

- ③ Eigen vectors turn up all over mathematics & science.
- ④ Eigen vectors, like Principal Components, give you a reduced version of the data that contains much of the original data.
- ⑤ For temporal Process, they often describe where a temporal structure with lots of feedback end up in <sup>Complex</sup> equilibrium over time.



⑥ Every Matrix  $M_{n \times p}$ , has  $p$  eigenvalues  
eigen vectors and  $p$  eigen values associated with it:  
 $\{v_1, \lambda_1\}, \{v_2, \lambda_2\}, \{v_3, \lambda_3\}, \dots$  (ordered by size of  $\lambda_i$ )

⑦ The eigen vector with the largest  $\lambda_p$  is the first  
eigen vector (Thus the first principal Component) and  
so on for all the rest.

⑧ For our data Analysis: Matrix  $M$  is:

$M = \text{Covariance matrix of all variables.}$

Procedure:

for  
Factor  
Analysis

- ① Standardize the variables: (ex. mean=0  
Sd=1)
- ② Create the Covariance Matrix " $M$ "
- ③ Find the eigenvectors & eigen values " $M$ "
- ④ Choose How many of them we want to  
keep and analyze.

example) \* get data from "psych" package

\* It also has some machine learning function

Library(psych)

Like: Fa  $\rightarrow$  factor analysis.

data(msq)

msq1 <- msq[, 2:72]  $\rightarrow$  all indep. variable  
 $X_1, X_2, \dots, X_{72}$

msq1[msq1 == "9"] = NA

msq1 <- data.frame(msq1)

msq2 <- na.omit(msq1)

names(msq2) # "afraid", "alert", "angry", "anxious"

dim(msq2) # 1747 71 ... "unhappy", "upset", "guilty", "happy", ...

fact <- fa(msq2, nfactors = 2)

fact1 <- fact\$loadings[, 1] # extract the projection  
fact1 <- [order(fact1)] of the variables  
onto the factors.

# so the highest  $\Phi_{ij}$   
are "energetic", "lively" (+)  
and "sluggish", "tired" (-)

$\sim [2] \left\{ \begin{array}{l} \Phi_{11}, \Phi_{12}, \Phi_{13}, \dots \\ \Phi_{21}, \Phi_{22}, \Phi_{23}, \dots \end{array} \right\}$

fact2 <- fact\$loadings[, 2]

fact2 <- [order(fact2)] # "relaxed" "calm" (+)  
"distressed" "frustrated" (-)

② Other Methods:  $\left\{ \begin{array}{l} \text{prcomp} \\ \text{princomp} \end{array} \right.$  built-in R

$\left\{ \begin{array}{l} \text{pcaA} \leftarrow \text{prcomp}(\text{msq2}) \\ \text{pcaA1} \leftarrow \text{pcaA} \$ \text{rotation}[,1] \end{array} \right.$	method 1
$\left\{ \begin{array}{l} \text{pcaB} \leftarrow \text{princomp}(\text{msq2}) \\ \text{pcaB1} \leftarrow \text{pcaB} \$ \text{loadings}[,1] \end{array} \right.$	method 2

③ Manual:

$\left\{ \begin{array}{l} \text{covm} \leftarrow \text{COV}(\text{msq2}) \\ \text{eigenm} \leftarrow \text{eigen}(\text{covm}) \\ \text{eigen1} \leftarrow \text{eigenm} \$ \text{vectors}[,1] \\ \text{eigen-value1} \leftarrow \text{eigenm} \$ \text{values}[,1] \end{array} \right.$
--

\* All these results:  $\{ \text{fact1}, \text{pcaA1}, \text{pcaB1}, \text{eigen1} \}$   
are very similar.

④

Estimation Boundary	$\left\{ \begin{array}{l} \text{tm.msqr2} \leftarrow \text{t(as.matrix(msq2))} \\ \text{m.msqr2} \leftarrow \text{as.matrix(msq2)} \\ \text{va} \leftarrow \text{rnorm(ncol(msqr2))} \\ \text{for } i \text{ in } 1:10 \{ \\ \quad \text{vb} \leftarrow \text{m.msqr2} \%*\% \text{scale(va)} \\ \quad \text{va} \leftarrow \text{tm.msqr2} \%*\% \text{scale(vb)} \end{array} \right.$
------------------------	---

→ first eig vector estimated.

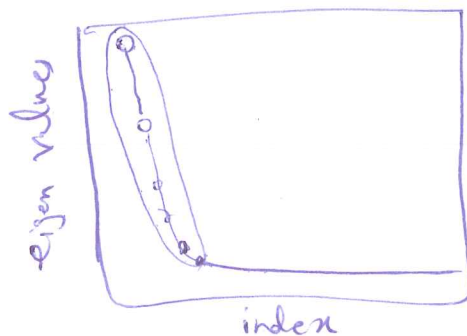
How many of P factors to keep?

①

| `plot(eigenm eigenm$values, type="b")`

look for the "elbow" →

#

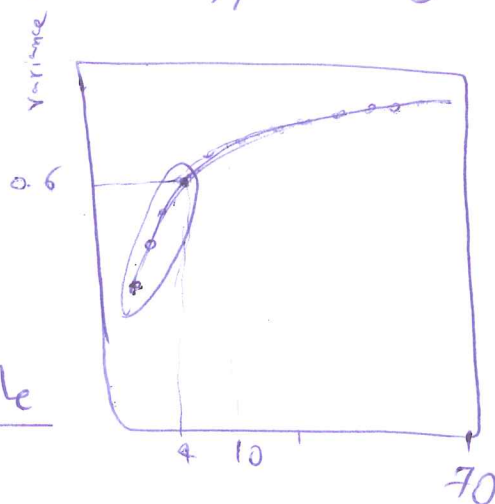


\* The flat area is all the factors that are just noise. So we count from left until the elbow.

② by plotting the proportion of total variance in data explained by the factors as you add them together.

| `plot(cumsum(eigenm$values)/sum(eigenm$values), ylim=c(0,1))`

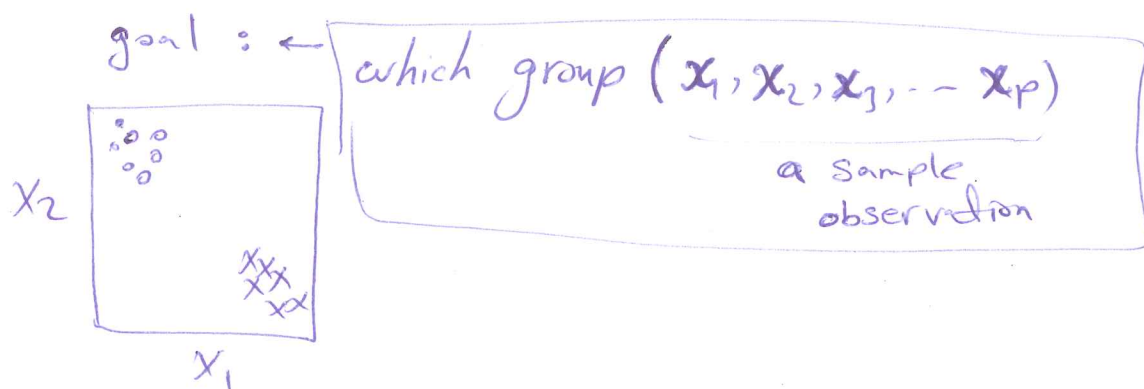
so more than 50% of the total variance in data is explained by just those 4 variables





# Clustering :

- ① unlike factor analysis, cluster analysis is targeted as observation level rather than the variables.
- ② we want to know what natural group the observation fall into:



Clustering Methods :

- ① k-mean
- ② Hierarchical

## ① \* k-means Algorithm :

- ① knowing how many clusters we divide data into
- ② Assign a category at random to every data point.
- ③a for each category calculate the mean of all points in that group. (Centroid of that group)
- ③b after calculating all centroids, reassign each point to the group of the nearest centroid. go back ③a until convergence.

Codes:

```
kout <- kmeans(msg2, center2, nstart = 25) 25 times starts randomly  
centroids <- kout$centers  
tpvars - centroid1 <- centroids[1, order(centroids[1,])]  
tpvars - centroid2 <- centroids[2, order(centroids[2,])]
```

The tutorial code is in the file Computer.

manual

how many clusters?

start with low increase one by one

as total variance ~~will~~ is lowering significantly.

## ② Hierarchical Clustering:

goal: → ① Finding how many  $k$   
(clusters)  
"Dissimilarity":

Complete: the dissimilarity between two cluster A and B is equal to the largest distance between a member of group A and a member of group B

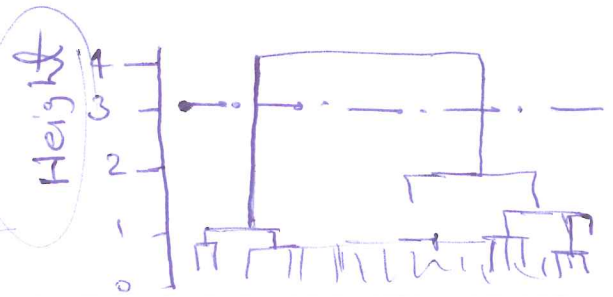
Single: like complete, but sets the dissimilarity between A & B to the smallest distance between a member of A & a member of B

Average: The average dissimilarity between every member of A and every member of B

Centroid: The distance between the centroid of A and the centroid of B

```
hout <- hclust(dist(simdat1[,1:2]), method = "average")
plot(hout)
```

```
abline(a=3, b=0, col="red")
```



every observation is a ~~cluster~~ centroid.

Then we find dissimilarities using "average" method

The two ~~clusters~~ centroids having smallest dissimilarity join together. This continues until they are all united.

▶ The "height" where any two clusters join is equal to the dissimilarity between those two clusters.

▶ draws a line,  $a = \text{intercept}$ ,  $b = \text{inclination}$

```
as.vector(cutree cutree(hout, 2))
```

# This gives a vector

of clusters associated with each observation

# 1 1 1 2 1 1 2 2 2 1 2 1 . .

```
as.vector(cutree(hout, h=3))
```

height

# The same output