# Pencast

Scan to open on Studocu

```
▶ Run
1  import numpy as np
2
3  x = np.array([2, 3, 1, 3, 3, 4, 5, 2])
4
5  print(x < 3)
6  print(x == 3)
```

↳ This "<" will tell the code to check the value of each elements of x

⟶ Output of this code set would b

[ True  False  True  False  False ,  False  False  True ]

[ False  True  False  True  True  False  false  false ]

To print it as number  "Make output number"

```
▶ Run
1  import numpy as np
2
3  x = np.array([2, 3, 1, 3, 3, 4, 5, 2])
4  y = np.array([0, 1, 2, 3, 4, 5, 6, 7])
5
6  print(y[x < 3])
7  print(y[x == 3])
```

use        line 6  :   print (y[x<3])
           line 7  :   print (y[x==3])


           . mean ()     :  mean function
           '{:.2f}'.format ( ... )     :  2 decimal place format

# Argmin and Argmax

                        minimum        maximum

numpy  has  methods  .min() and  .max()  values in array
      we can also use  .argmin()    .argmax()  ⟶  to know the location of
                                         min  &  max  value

```
1  import numpy as np
2
3  x = np.array([2, -1, 3, 8])
4
5  print(x.min())     I
6  print(x.max())
```

```
▶ Run
1  import numpy as np
2
3 #      index:  0
4  x = np.array([2, -1, 3, 8])
5
   print(x.argmin())
```
from index 1
from index 3

-1
8

Instead of calculating distance with the inner product → np.sqrt(np.inner(a-b, a-b))

use :   np.linalg.norm(a-b)   which calculate distance between vectors a and b

```python
1  import numpy as np
2
3  a = np.array([1, 3, 4])
4
5  x1 = np.array([4, 3, 5])
6  x2 = np.array([0.4, 10, 50])
7  x3 = np.array([1, 4, 10])
8  x4 = np.array([30, 40, 50])
9
10 x = [x1, x2, x3, x4]
11 dist = np.zeros(4)
12
13 for i in range(4):
14     dist[i] = np.linalg.norm(a - x[i])
15     print('Distance between a and x{}: {:.2f}'.format(i+1, dist[i]))
16
17 print('The nearest neighbour is x{}'.format(dist.argmin()+1))
```

use loop distances
dist[i] = np.linalg.norm(a - x[i])
and .argmin()

# Loading data

## using pandas library

" pd.read_csv(file.name) "

→ • Name of data frame



```
► Run                                          PYTHON
1  import pandas as pd
2
3  marketing = pd.read_csv('/course/data/DirectMarketing.csv')
4
5  print(marketing)
```

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| 995 | Young | Female | Rent | Single | ... | 1 | NaN | 18 | 304 |
| 996 | Middle | Male | Rent | Single | ... | 1 | NaN | 18 | 1673 |
| 997 | Old | Male | Own | Single | ... | 0 | Medium | 24 | 1417 |
| 998 | Middle | Male | Own | Married | ... | 2 | Medium | 18 | 671 |
| 999 | Young | Male | Rent | Married | ... | 1 | Medium | 24 | 973 |

[1000 rows x 10 columns]

✓ Program exited with code 0

head() : give us the top of the file the frame

⎧ Data frame is a type of object that pandas made when it loads in our data.
⎩ → can think of it as data table or excel table

```
► Run                                          PYTHON
1  import pandas as pd
2
3  marketing = pd.read_csv('/course/data/DirectMarketing.csv')
4
5  print(marketing.head(10))
```
Note: the first row is row 0. This is because Python indexing starts at 0. This means we start at row 0 and go up to but not including row 10.

→ 10 heads

# Selecting Colums

**DataFrame[column]**

multiple    column

**Data Frame [[ column1 , ..... ] ]** → list

```
▶ Run                                      python [ ]
  1  import pandas as pd
  2
  3  marketing = pd.read_csv('/course/data/DirectMarketing.csv')
  4
  5  print(type(marketing['Age']))
<class 'pandas.core.series.Series'>
```

Here is another example showing you how to extract out the Age, Married and AmountSpent columns.

```
▶ Run                                      python [ ]
  1  import pandas as pd
  2
  3  marketing = pd.read_csv('/course/data/DirectMarketing.csv')
  4
  5  print(marketing[['Age', 'Married', 'AmountSpent']])
  6  print(type(marketing))

995   Young    Single     384
996   Middle   Single    1973
997   Old      Single    1417
998   Middle   Married    671
999   Young    Married    973

[1000 rows x 3 columns]
<class 'pandas.core.frame.DataFrame'>
```

it's a Data frame
because it has multiple column

# From Pandas to numpy

**DataFrame . to_numpy ()**

Ex :

```
▶ Run                                      python [ ]
  1  import pandas as pd
  2
  3  marketing = pd.read_csv('/course/data/DirectMarketing.csv')
  4  salary = marketing['Salary'].to_numpy()
  5
  6  print(type(salary))
  7  print(salary)

<class 'numpy.ndarray'>
[ 47500  36600  13500  85500  68400  30400  43100  58100  61900  83700
  41700  11500  46100 111400 110900  53100  52800  70100  39900  14000
 14300  49600  68200  43100  42700  52300  52800  72200  65200  21700
 30700  72200  87700 104100  67700  27600  54100  47500  41300  25100
 54800  66800  76200  81500  43600 111500  77700  58900 126500  81300
 96500  37500 105500  66700  43600  65800  66300  67800  47500  35700
 62700  29900  15500  14600  62600  66700  72500  13200  91300  72200
 13600  96400  71500  16400  19400  57600  96600  64300  24400  22100 ]
```

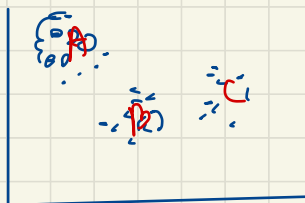Mutiple column

**DataFram . to_numpy ()        : save**

## Multiple columns

We are also able to extract out more than one column at a time.

**Question!** What are the dimensions of the subset array generate by the code below?

```
▶ Run                                      python [ ]
  1  import pandas as pd
  2
  3  marketing = pd.read_csv('/course/data/DirectMarketing.csv')
  4  subset = marketing[['Salary', 'Children']].to_numpy()
  5
  6  print(subset)
```

2 columns

# Clustering

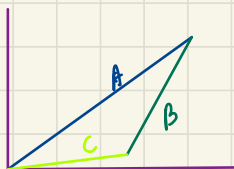idea of clustering is take some datas and form groups automatly

A cluster is a group of sample which we say are ==similar==.



# Calculating Similarity with distance

### Mathematically

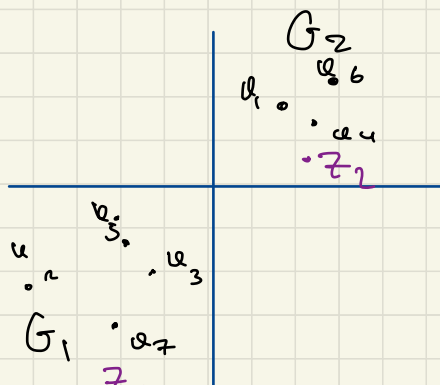$$\|x-y\| = \sqrt{(x_1-y_1)^2 + \ldots + (x_n-y_n)^2}$$

$$= \sqrt{(x-y)^T(x-y)}$$



the points that are closer
we expect them to be
similar
further ⟹ diff

# Notation

$$==J^{clust} = \frac{1}{N}\sum_{i=1}^{N} \|x_i - z_{c_i}\|^2==$$

$i$ is group : $1, \ldots, K$

⟹ $N = 7$    $G_j \in \{1, \ldots, N\}$

$K = 2$

$G_1 = \{2, 3, 5, 7\}$    $G_2 \doteq \{1, 4, 6\}$

" Our goal is to assign each data point in a group such that we minimize the distance between each point and its cluster respresentitive "

It mean that if $x_q$ in group 2 $\Rightarrow$ we try to minimize the distance between $x_q$ & $z_2$ which is $\|x_q - z_2\|$

We done when wth change

# Visual walkthrough

k-mean is a very simple clustering method.
It performs surprisingly well in practice and its simplicity means that it can be applied to massive datasets with ease.

— Clustering algorithm:
1. Set the number k
2. initialise centroid of each cluster (guess or use a heuristic)
3. Assign each data point to its closet centroid (cluster assigment)
4. Set new centroids as the mean of each other
5. Repeat 3-4 until convergene    no change in cluster assigment

— Clustering algorithm (mathematic notation):

Given $x_1, \ldots, x_N \in R^n$ and $z_1, \ldots, z_k \in R^n$
repeat
  update partition: assign i to $G_j$, $j = \arg\min \; j' = 1, \ldots, k \|x_i - z_i\|^2$
  update centroids: $z_j = \frac{1}{|G_j|} \sum_{i \in G_j} x_i$

until    $z_1, \ldots, z_k$ stop changing

# Break

is a way to exit the loop

# Pencast W6:

## Matrices:
        is  a  grid  of  number
        a  list  of  vector
        a  list  of  column vector

## To print matrices

```
▶ Run
1 import numpy as np
2
3 matrix = np.array([
4     [1, 4, 5, 7],
5     [2, 8, 6, 3],
6     [3, 1, 4, 9]
7 ])
8
9 print(matrix)
10 print(matrix.shape)  → tell us the demension

[[1 4 5 7]
 [2 8 6 3]
 [3 1 4 9]]
(3, 4)  → columns
  ↓
 rows
```

USE: print (matrix .

## matrix- scalar operation

```
▶ Run
1 import numpy as np
2
3 matrix = np.array([
4     [1, 4, 5, 7],
5     [2, 8, 6, 3],
6     [3, 1, 4, 9]
7 ])
8
9 print(matrix + 100)    → scalar

[101 104 105 107]
[102 108 106 103]
[103 101 104 109]]
```

## Matrix - matrix operations

```
1 import numpy as np
2
3 A = np.array([
4     [1, 4],
5     [2, 8]
6 ])
7
8 B = np.array([
9     [0, 1],
10     [0, 1]
11 ])
12
13 print(A + B)

[[1 5]
 [2
```