

K-Mean Clustering

M T W T F S

DATE: _____

M D W

- Section Overview
 - Understanding clustering
 - History and Intuition of K-mean
 - Mathematical Theory of K-Means
 - Example of K-mean
 - K-mean project Exercise -

Clustering General overview

Clustering uses unlabelled data and looks for similarities between groups (clusters) in order to segment the data into separate clusters -

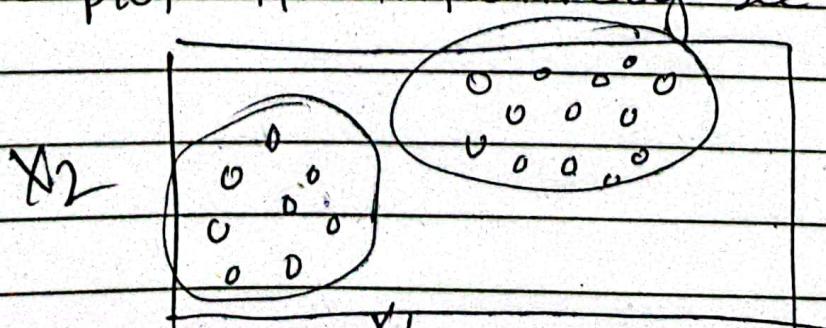
- Keep in mind that we don't actually know the true correct label for this data -

Imagine an example dataset

x_1	x_2
2	4
6	3
5	9
2	1

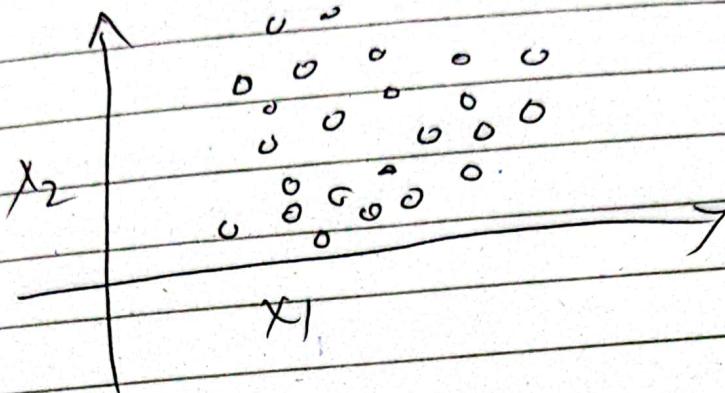
• distance is the intuitive metric

If we plot it intuitively see 2 grouping



METHODS

- 2 or 3 clusters could be reasonable.



Different methods can be used to decide number of clusters -

- Some methods themselves decide and for some we need to decide the number of clusters -

(will name ~~methods~~ group 0, group 1, group 2 methods)

Main clustering Ideas:-

- Use features to decide which points are most similar to other points -
- Realize that there are no final correct label to compare cluster results to.
- We can think of clustering as an unsupervised learning process that discovers potential labels -

M T W T F S

DATE: _____

Unsupervised paradigm shift

How we assign a new datapoint to a cluster?

- Different approaches depending on the unsupervised learning algorithm used -
- use features to assign most appropriate cluster.
- just as before, ^{there is} no way to measure if this was the "correct" assignment
- If we have discovered these new clusters labels, could we use that as a y for supervised training?
 - Yes we can use unsupervised learning to discover possible labels, then apply supervised learning on new data points.

What is trade off?

- Clustering doesn't tell you what these new cluster labels represent, no real way of knowing if these clusters are truly significant.
- Clustering ideas still to come:
 - How do we decide which number of cluster is based best?
 - Do we decide or let the algorithm decide-
 - How can we measure "goodness of fit" for clustering without a y label for comparison -

- Machine learning as an art:

- What is ground truth

- What tradeoffs are we making by using unsupervised learning as a substitute for ground truth of they label that was not given?

K-mean Clustering History, Intuition and Theory

- Hugo Steinhaus was the mathematician who began laying out groundwork for clustering methods
- 1940s - 1950s : Hugo Steinhaus polish mathematician and professor -
 - During WWII he was in hiding from nazi occupation, yet still taught clandestine math classes from memory, since higher education for poles was forbidden -
- He even established statistical method for estimating German casualty rates based ~~on~~ simply based on relative frequencies of obituaries in local newspapers.
- After WWII he helped established the Mathematics department at university

M T W T F S

DATE: _____

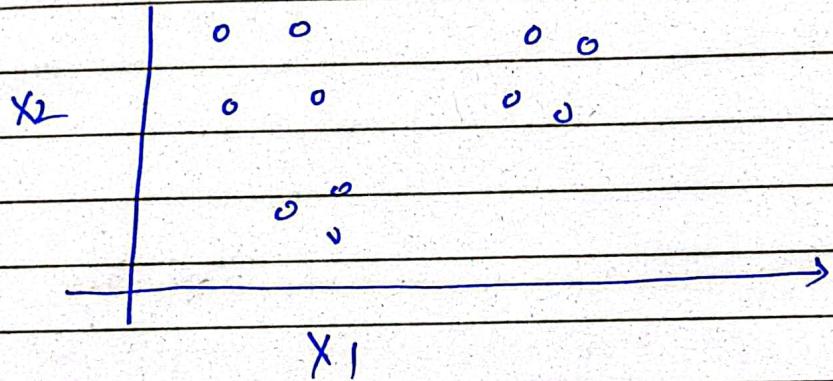
Wroclaw -

- It was at the University of Wroclaw where he began to develop the idea that would become the foundation of k-mean clustering!
- Check out his wikipedia page for more interesting stories from his incredible life-

- o 1957 : Stuart Lloyd of Bell Labs developed the standard algorithmic technique for clustering related to k-Means.
- o However Lloyd did not publish his technique in a journal until 1982.
- o 1965 Edward W. Forgy published the clustering techniques, thus both Lloyd and Forgy are credited with the development.
- o 1967 : James MacQueen of UCLA developed and published "some methods for classification and analysis of multivariate observation".
 - o This is the first publication to introduce the term "k-Mean" as a sequential clustering algorithm

- First a set of properties point must satisfy:
 - Each point must belong to a cluster.
 - Each point can only belong to one cluster (no single point can belong to multiple clusters)

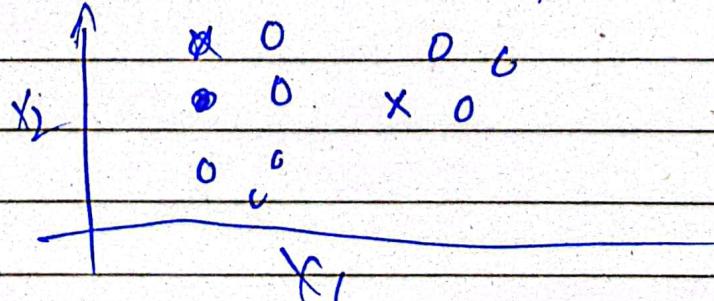
How to kmean cluster works



- step 1: choose the number of cluster to create (this is K value)

: we have choose $K=3$. Note in most situation you wont visualize the data.

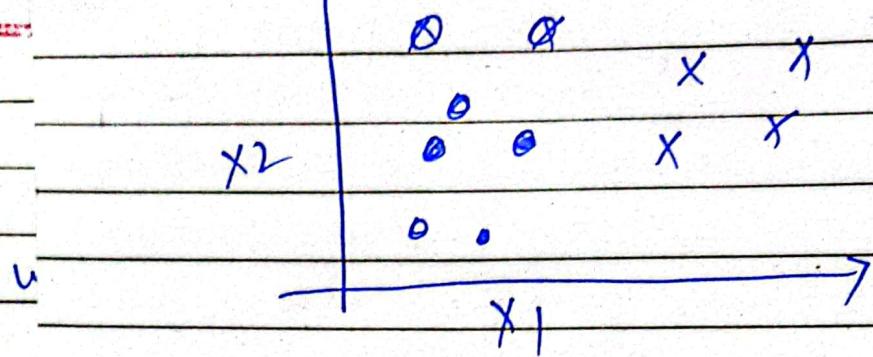
Step 2: Random Select K distinct data points - Our $K=3$:



Step 2: Assign each point
we meet these new K points
our "cluster points" -

DATE: _____

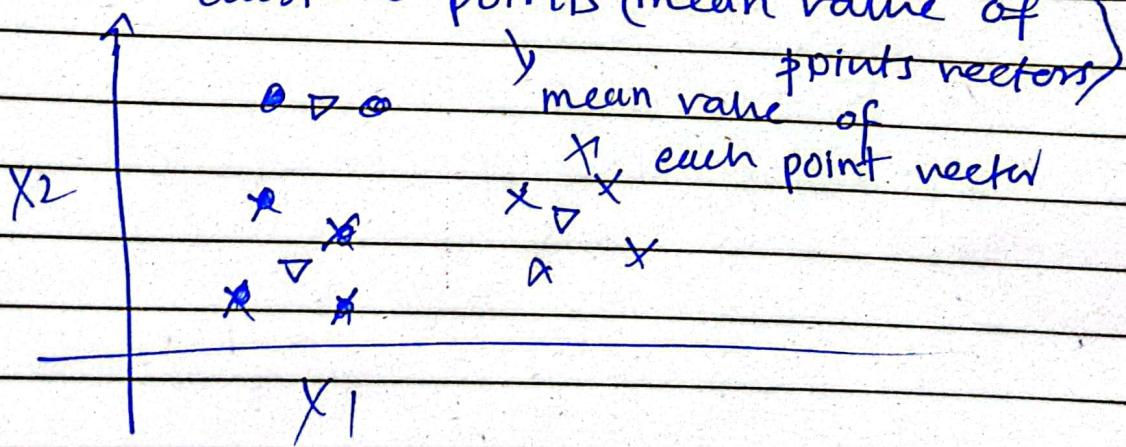
M T W T F S



Step 3: Assign each remaining points
to the 'nearest' "cluster" point)

• Note this is using a distance metric
to judge the nearest points

Step 4: Calculate the center of the
clusters points (mean value of)



Step 5: Now assign each point to the nearest
cluster center -

6. - we repeat step 4 and 5 until there
are no more cluster reassignments

Upcoming considerations:

- How do we choose a reasonable value for K number of clusters?
- Is there any way we can evaluate how good our current K value is at determining clusters?

K-mean clustering coding # part 1

```
df = pd.read_csv(".../Data/Band-full.csv")  
df.info()  
df.describe()
```

two things - domain experience

```
sns.histplot(data=df, x="age", kde=True)  
plt.figure(figsize=(20,10), dpi=200)
```

```
sns.histplot(data=df, x="age", hue="label")
```

```
sns.histplot(data=df, x="pdays", bins=40)
```

```
sns.histplot(data=df[df['pdays'] == 999],  
             x='pdays', bins=40)
```

```
df['contact'].unique()
```

```
df['duration']
```

`sns.histplot(data=df, x='duration',
hue='contact')
plt.xlim(0, 1000)`

`sns.countplot(data=df, x='contact')`

~~df~~ `sns.countplot(data=df, x='job',
plt.xticks(rotation=90): order=5
plt.figure(''))`

same
for

education

↳ check with hue='default'

`df['default'].value_counts()`

`sns.pairplot(df)`

Coding part #2

`X = pd.get_dummies(df)`

`X.head()`

`st = StandardScaler()`

`= st.fit_transform(X)`

`scaled_X`

from sklearn.cluster import KMeans
 help(KMeans)

↳ n-clusters = 8 default

max_iteration = 300 default

model = KMeans(n_clusters=2)
 model.fit_predict(scaled_x)

clusters_labels = model.fit_predict(scaled_x)
 cluster_labels

array([1, 1, 1, ..., 0, 0, 0])

x['cluster'] = cluster_labels

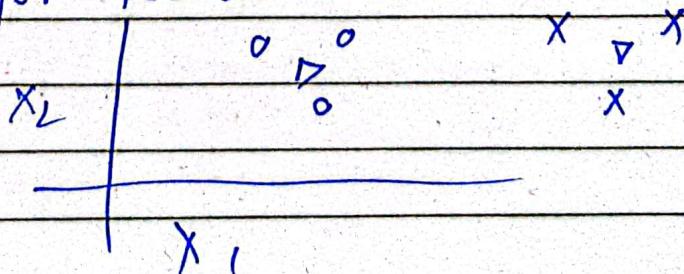
x.corr()['cluster'].iloc[: -1].sort_values()
 .plot(kind='bar')

Q6 # Coding part 3

How to measure 'goodness of fit'

- we could measure the sum of the distances from point to cluster centers.

for K = 2

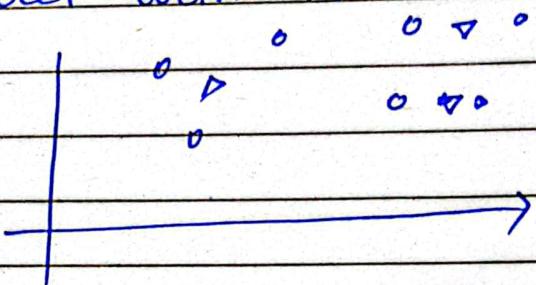


M T W T F S

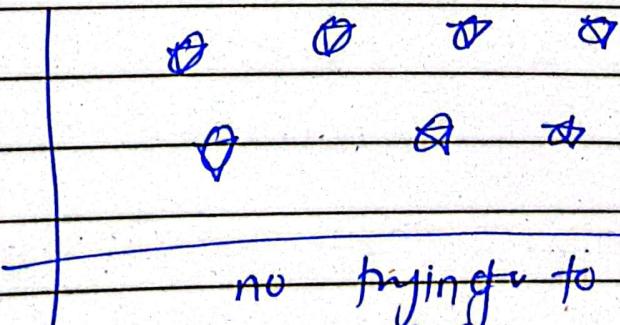
DATE: _____

- for all the points we measure the sum of squared distances from point to the cluster center:

- Then we fit an entirely new k-means model with $K+1$:



- Then measure again the sum of the squared distance (SSD) to center.
- In theory this SSD would go to zero once K is equal to the number of points-
- You would have a cluster for each point! SSD would be perfect at 0!

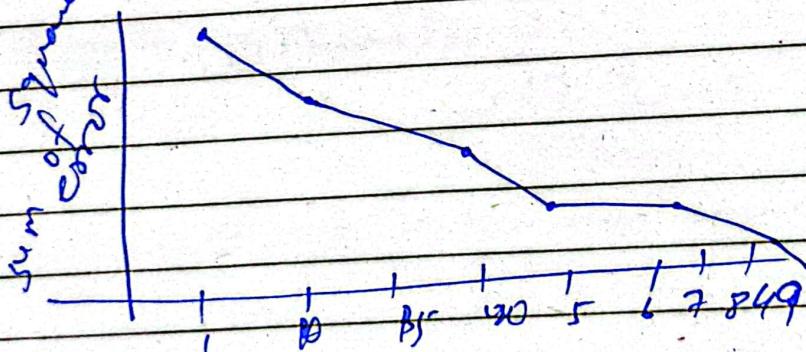


no trying to minimize sum of Squared distance but instead using as a measurement tool

- We keep track of this SSD value for a range of reduction in K until it begins to decline.
- This signifies that adding an extra cluster is not obtaining enough clarity of clusters separation to justify increasing K .

↓ This is known as the "elbow" method since we will track where decrease in SSD begins to flatten out compared to increasing K -values-

- Start with $K=2$
- Increase K and measure SSD



These points are strong indicator that increasing K further is no longer justified as its not revealing more "signal".

- You can also measure out this SSD in a barplot

•

Code Section

M T W T F S

DATE: _____

ssd = []

for K in range(2, 10):

model = kmeans(n_clusters=k)

model.fit(scaled_X)

ssd.append(model.inertia_)

↓
#SSD points

cluster center

ssd

plt(range(2, 10), ssd, 'o--')

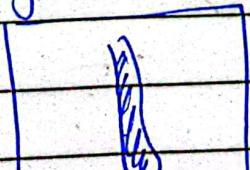
pd.Series(ssd)

pd.Series(ssd).diff()

K-Means Clustering Color Quantization part one

- One image - interesting application of clustering is on image quantization.
- Let discuss images, computers, colors, and quantization to get an idea of how K Mean clustering can be applied to different fields.

- Imagine an image of a single pen stroke!



- This image is in grayscale, meaning the color ranges from black to white.



- A computer will store this information as an array with values between a range.
- Notice 0 is white and 1 is black, with values in between representing gray.

0	0	0	0	0	0	0	0
0	0	0	0.1	0.2	0.3	0.4	0.5
0	0	0	0.9	0.8	0.7	0.6	0.5
0	0	0	0.4	0.3	0.2	0.1	0.0
0	0	0	0.0	0.1	0.2	0.3	0.4

- It is also very common for computer to store value from 0-255 for scales.
- The range 0 to 255 has to do with how computer stores 8-bit numbers.

- Color Images can be represented as array, what about color images.
- Color Images can be represented as a combination of Red, Green, Blue.
- Combining Green, Red, Blue to create colors.
- From a computer perspective, this looks like 3 arrays, each array representing a color channel.

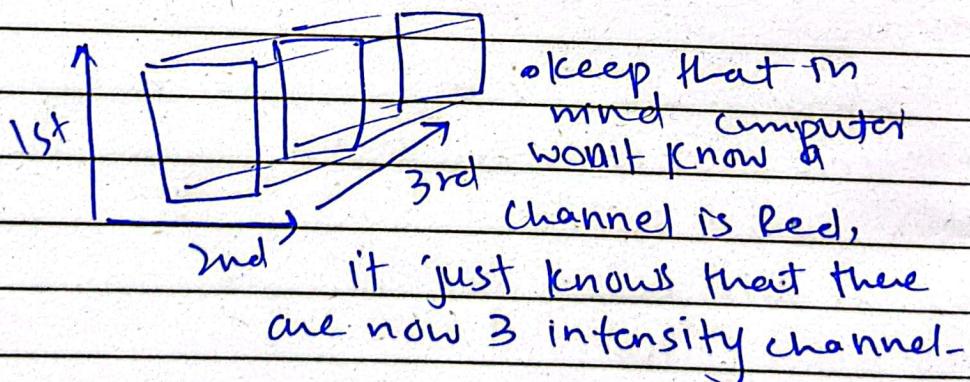
- For example a single pixel (1 by 1 image) here is (213, 111, 56) for (R, G, B)

- The shape of the color arrays there has 3 dimensions -

- Height
- Width
- Color channel

- If you read an image and check its shape - It will look like

- (1280, 720, 3)
- (1280 pixel width)
- (720 pixel height)
- (3 color channels)



- From a computer perspective you simply have an array with 3 dimensions, where a user or display function can attribute each dimension to a color channel (e.g. red "intensity")

- Swapping these array across channels would allow for effect such as color inversion.

- Now how can we apply clustering to RGB color channel and images?

Code part two

```
import numpy as np
```

```
import matplotlib.image as mpimg  
image as array
```

```
mpimg.imread('.../palm-tree.jpg')
```

```
plt.figure(dpi=200)
```

```
plt.imshow(image_as_array)
```

) need
numpy
array

```
(H,W,C) --- 2D
```

```
(H^W, C)
```

```
(h,w,c) = image_as_array.shape
```

```
image_as_array_2d = image_as_array.reshape((h^w, c))
```

```
Image.array
```

```
from sklearn.cluster import KMeans
```

```
model = KMeans(n_clusters=6)
```

```
label = model.fit_predict(image_as_array_2d)
```

```
labels #
```

M T W T F S

DATE: _____

model.cluster_centers_

rgb-codes = model.cluster_centers_.round(0).
as type(int)

rgb-codes[labels]

, np.reshape(rgb-codes[labels], (h, w, c))
quantized
image

plt.imshow(quantized_image)