

Unsupervised Learning

Given: Unlabeled data $D = \{x_1, \dots, x_n\} \subseteq \mathbb{R}^d$ ← No labels y_i
 Unsupervised learning attempts to find (hidden) structure in your data.

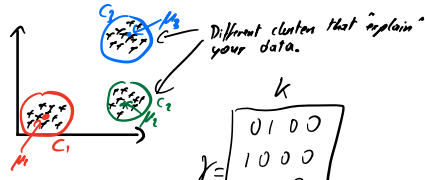
Clustering: K-Means

(MacQueen 1967)

Assumptions: - We know there are k clusters
 - Euclidean distance measures similarity
 - clusters are spherical
 - clusters have similar variance

Let data indices be $\{1, \dots, n\}$.

Cluster assignment: $y_{ie} = \begin{cases} 1 & \text{if } i \text{ falls in cluster } l \\ 0 & \text{otherwise} \end{cases}$



$$Y = \begin{matrix} & \begin{matrix} 1 & 2 & 3 \end{matrix} \\ \begin{matrix} 1 \\ 2 \\ 3 \\ 4 \\ 5 \\ 6 \\ 7 \end{matrix} & \begin{bmatrix} 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 \end{bmatrix} \end{matrix}$$

Binary assignment. Each x_i is in exactly one cluster. $y_{ie} \in \{0, 1\}$

Cluster center: $\vec{\mu}_l = \frac{1}{m_l} \sum_{i=1}^n y_{ie} x_i$

$m_l = \sum_{i=1}^n y_{ie}$
 m_l denotes the number of points assigned to cluster l .

Loss function: $\mathcal{L}(y, \vec{\mu}_1, \dots, \vec{\mu}_k) = \frac{1}{n} \sum_{l=1}^k \sum_{i=1}^n y_{ie} \|x_i - \vec{\mu}_l\|_2^2$

← Average distance from each point to its cluster center.
 ↳ if you want to represent each x_i with its corresponding $\vec{\mu}_l$ instead, this would be the avg. squared error.

① Minimize \mathcal{L} if y is fixed: (Assign $\vec{\mu}_1, \dots, \vec{\mu}_k$)

$$\vec{\mu}_l = \underset{\vec{\mu}_l}{\text{argmin}} \mathcal{L}(y) \Rightarrow \vec{\mu}_l = \frac{1}{m_l} \sum_{i=1}^n y_{ie} x_i$$

mean of all x_i in l

For each $\vec{\mu}_l$, the contribution to the loss is $\sum_{i=1}^n y_{ie} (\vec{x}_i - \vec{\mu}_l)^2$ ← parabola

② Minimize \mathcal{L} if $\vec{\mu}_1, \dots, \vec{\mu}_k$ are fixed: (Assign y)

$$y_{ie} = \begin{cases} 1 & \text{if } l = \underset{l}{\text{argmin}} (\vec{x}_i - \vec{\mu}_l)^2 \leq \dots \\ 0 & \text{otherwise} \end{cases}$$

All points that are closest to $\vec{\mu}_l$.
 For each x_i you can contribute $(x_i - \vec{\mu}_1)^2$ or $(x_i - \vec{\mu}_2)^2$ or ... or $(x_i - \vec{\mu}_k)^2$.
 For each x_i pick the smallest one to minimize \mathcal{L} .

K-means algorithm: 1) Initialize $\vec{\mu}_1, \dots, \vec{\mu}_k$ somehow (e.g. random points, or heuristic)

- 2) Assign y with ②
 - 3) Assign $\vec{\mu}_1, \dots, \vec{\mu}_k$ using ①
- Repeat until convergence

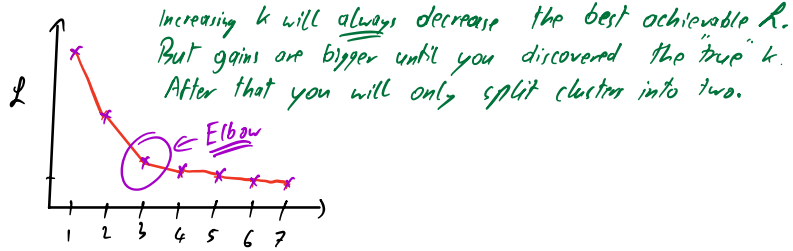
Easy to show: \mathcal{L} can never increase.

Initialization: - Can be unlucky. Perform multiple runs, pick solution with lowest \mathcal{L} .
 - k-means++ (Arthur et al. 2009): - pick $\vec{\mu}_1$ randomly from D

- For $l=2$ to k :
 - $\forall x_i \in D$ define $d_i = \min_{j \neq l} (x_i - \mu_j)^2$ ← squared Euclidean distance to closest existing center
 - pick μ_l randomly from D proportional to d_i^2 ← Avoids cluster centers that are very close

How to find k if it is unknown?

Run k-means for $k=1,2,\dots$ monitor the loss L
 (For each k you have to average over multiple runs, due to the randomness of the initialization.)



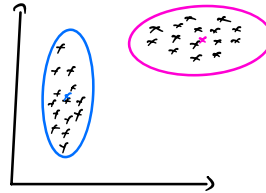
- Strengths:
- Easy to implement
 - fast for small k
 - easily parallelizable

- Limitations:
- sensitive to outliers
 - slow for large k
 - assumes all clusters are spherical
 - makes only hard cluster assignments

Generalization: Gaussian Mixture Model (GMM)

Each cluster becomes a Gaussian $N(\mu_c, \Sigma_c)$

Almost the same as k-means, but each point has a distribution over clusters $\gamma_i = P(C_c | x_i)$.



0.5	0.2	0.1	0.1
0.1	0.3	0.9	0.2
0.9	0.0	0.0	0.1
0.8	0.1	0.05	0.05
0.7	0.1	0.1	0.1
0.0	0.0	0.9	0.1
1.0	0.0	0.0	0.0
0.3	0.2	0.5	0.0

$\gamma =$
 soft assignment:
 $\gamma_i \in [0,1] \quad \sum \gamma_i = 1$

$$\vec{\mu}_c \leftarrow \frac{1}{\sum_i \gamma_{ci}} \sum_{i=1}^n \gamma_{ci} \vec{x}_i \leftarrow \text{weighted average / expected center of cluster } c$$

$$\vec{\Sigma}_c \leftarrow \frac{1}{\sum_i \gamma_{ci}} \sum_{i=1}^n \gamma_{ci} (\vec{x}_i - \vec{\mu}_c)(\vec{x}_i - \vec{\mu}_c)^T$$

$$\pi_c \leftarrow \frac{\sum_i \gamma_{ci}}{n} \leftarrow \pi_c = P(C_c)$$

$$\gamma_{ci} \leftarrow \frac{P(x_i; \mu_c, \Sigma_c) \pi_c}{\sum_{c'} P(x_i; \mu_{c'}, \Sigma_{c'}) \pi_{c'}}$$

M-Step
 Maximization

E-Step
 Estimation

- Strengths:
- more flexible than k-means
 - probabilities tell you how well a sample fits into clusters
- Weaknesses:
- slower than k-means
 - can suffer from singularities (single point clusters)

(if $\Sigma_c = I\sigma^2$ GMM becomes k-means as $\sigma^2 \rightarrow 0$.)

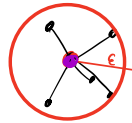
spherical

hard assignments

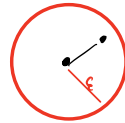
DBSCAN: 1. connect each point with all neighbors within ϵ -radius.
 (Density based spatial clustering with noise)
 Ester et. al 1996

Assumption:

- Clusters have similar density
- Similar points share clusters



core point
 $\geq M$ neighbors



non-core point
 $< M$ neighbors

Repeat:

- Start with random unassigned core-point.
- assign it to a new cluster ID l .
- Until no longer possible:
 - Assign all core points with a core-point neighbor in cluster l also to l .
- Assign all non-core points with a core point neighbor in cluster l to l .
- Assign all non-core points without core-point neighbor as noise

Apply DBSCAN on this ϵ -nearest neighbor graph:

