

Unsupervised Learning

- ↳ How to use ML when no labels are provided?
- ↳ Labels can be expensive to acquire, UL allows for insights despite no labels

Clustering

- k-means
- hierarchical

Gaussian Mix. Models

- Density est.
 - Comparison w/
k-means
- ↳ the same under certain conditions

What is UL used for?

↳ Bayesianism?

- density estimation (est of PDF)
- latent variables
- Dimension reduction
 - ↳ understand data
 - ↳ see if models work well

Why use clustering?

- ↳ can be used on any dataset.
- ↳ Also allows for speed of algorithms if you need to search for similar items
 - ↳ recommendation systems
- ↳ Density estimation
 - ↳ generate new things based on old things (Hidden Markov Models)

↑
known

k-means Cluster-Me

No labels to data

↳ there might still be patterns to find

→ Challenges:

- High dimensions can't be visualised
- Not all interesting info is shown in clustering
 - ↳ How do we evaluate this?

k-means

Based on 2 facts:

- clusters have a center of mass

$$\text{com.} = \frac{1}{c} \sum \vec{z}_i$$

- new points get assigned to closest cluster

More on k-means

- ↳ Start w/ random centroids
 - ↳ map points to closest pt.
 - ↳ recenter to CCM
 - ↳ repeat until converge
- } n centroids
} = k

Soft k-means

"fuzzy" membership

- ↳ membership is not binary but a continuous value

$$\gamma_k^{(n)} = \frac{e^{-\beta d(m_k, x^n)}}{\sum_{j=1}^K e^{-\beta d(m_j, x^n)}}$$

dist between centroid & point

↳ similar to gaussian

$\beta \leftarrow$ "sensitivity term"

Objective Function

$$J = \sum_n \sum_k \|w_k - x^{(n)}\|^2 \quad] \text{ coordinate descent}$$

C will always decrease (but not convex)

↳ multiple local minima

When can k-means fail

- donut problem
- two clusters can also be split awkwardly depending on initialization
- less dense functions that are next to very dense ones

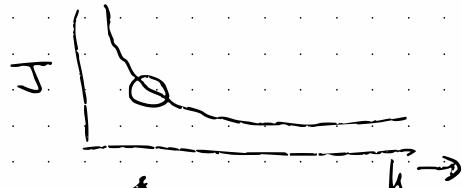
→ Disadv. of k-means:

- k needs to be determined] experimentation
- local min.] needs to be done.
- can only work for spherical clusters

→ How to choose k

→ Plot scree plot

→ Choose k that reduces
the cost the most.



↑ Elbow has
most info per
cluster

Cost function Alternatives

Current inter-cluster dist $\bigcirc \leftrightarrow \bigcirc$

intra cluster dist \bigcirc

↳ Decreases every round

(low \bigcirc , high $\bigcirc - \bigcirc$ is
desirable)

Discuss: problems w/ large dataset

cost is sensitive to scale of
data

cost also sensitive to k

↳ high k overfitting?

Not really - unsupervised
ML - remember?

Davies-Bouldin Index

↳ "internal validation"

$$DBI = \frac{1}{K} \sum_{k=1}^K \max_{j \neq k} \left[\frac{\sigma_j + \sigma_k}{d(c_k, c_j)} \right]$$

Annotations:

- σ_j : dist from each pt to cluster centre
- $d(c_k, c_j)$: dist between cluster's (centre)
- $j \neq k$: different clusters

Hierarchical Clustering

"Agglomerative clustering"

- ↳ grow clusters from bottom-up
- ↳ closest together get clustered first

→ Create dendrogram

Where to decide clusters?

Options for clustering

Distance measurements:

- Euclidean dist
- Sq. Euclid. dist.
- Manhattan
- Max dist
- Mahalanobis dist.

} experiment to
find best results

↳ valid based on certain properties

→ What are some properties distances must have?

- 1) $d(x, y) > 0$
- 2) if 1 then $x = y$
- 3) $d(x, y) = d(y, x)$
- 4) $d(x, z) \leq d(x, y) + d(y, z)$

→ How to join clusters together?

→ Single linkage (closest points)

→ complete linkage (furthest pts)

→ mean dist (UPGMA)

→ Ward's criterion

↳ Minimise increase in variance

$$\text{cost} = \frac{\sum_{\substack{i \in [a, b] \\ i \neq a, b}} n_i}{n} (\bar{a} - \bar{b})^2$$

↳ Problem w/ single linkage

- always grab nearest obj but miss pattern

Why does this happen?

↳ Always grabbing obj that's very close
so it looks like no new cluster is forming

Gaussian Mixture Models

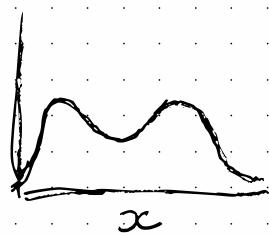
↳ form of density estimation

· approx. of prob. dist. of our data

· used when data is multi-modal

↳ human...

↳ estimate w/ multiple
normal distributions?



yes! $p(X=x) = \sum_{i=0}^n \pi_i N(\mu_i, \Sigma_i)$

↑ ↑ ↑
weighted mean cov matrix
of norm

NB: $\cdot \pi_i$ = prob that X belongs to i th norm.

$$\int p(X=x) dx = 1$$

↳ since each N has int of 1

$$\sum \pi_i = 1$$

$\Rightarrow \pi$ is thus its own PMF

→ Training a Gaussian Mixture Model

Similar to training k-means

2 steps:

1) Calculate responsibilities (γ)

$$\gamma_k^{(n)} = \frac{\pi_k N(x^{(n)} | \mu_k, \Sigma_k)}{\sum_{i=1}^k \pi_i N(x^{(n)} | \mu_i, \Sigma_i)}$$

prop of 1 norm over responsibility
of all norms

2) Given responsibilities, recalculate the rest

$$\mu_k = \frac{1}{N_k} \sum_{n=1}^N \gamma_k^{(n)} x^{(n)} \quad \leftarrow \text{mean of resp. pt}$$

$$\Sigma_k = \frac{1}{N_k} \sum_{n=1}^N \gamma_k^{(n)} (x^{(n)} - \mu_k)(x^{(n)} - \mu_k)^T$$

Norm for sample size resp of that mean sq dist. of each pt. from mean

$$\pi_k = \frac{N_k}{N} \quad \text{w/} \quad N_k = \sum_{n=1}^N \gamma_k^{(n)}$$

\curvearrowleft renorm π or something

\rightarrow Comparison betw. AMM & k-means

\hookrightarrow very similar algorithm

1. initialize points
2. evaluate data
3. update
4. rinse & repeat until convergence

k-means \rightarrow no π var, each cluster has equal weight.

Σ vs β \rightarrow Σ is full cover matrix

\hookrightarrow greater flexibility of shape

k-means is gym where π is uniform
& is perfectly spherical

→ Practical issues w/ gmm's

→ What happens if var = 0

↳ If some var/covar in Σ
is 0 then it becomes singular

(ie squishes space into lower dim than
dim of full rank matrix)

Solution use diagonal covar

by assumption that each dim is indep
then all covars are 0 w/ only vars
on diag.

↳ Also assume all covars are equal

Also resolves local min.

→ kernel density estimation

↳ fitting of prob. dist. w/ kernels



→ Simplest version is a histogram!

GMM

↳ kde using a normal dist as kernel

↳ possible
to use
various kernels

→ Expectation maximisation

THE MLE: Max $P(X|\theta)$ or $\log P(X|\theta) = L(\theta)$

↳ add hidden var Z

$$\therefore P(X|\theta) = P(X, Z|\theta)$$

$$= \sum_z P(X, Z|\theta)$$

$$= \sum_z P(X|\theta, Z) P(Z|\theta) \text{ by LOTP}$$

:

2 steps: Adjusting Z & Θ

1. E-step

find dist of $Z|X, \Theta_n$
current setting of Θ

2. M-step

find Θ that maximises joint dist

$$P(X, Z|\Theta)$$

(\hookrightarrow) Max expected log like. of this
over dist. in Step 1

$$\text{ie } \Theta_{n+1} = \underset{\text{max}}{\arg} E[\log P(X, Z|\Theta)]$$

(Just use a library)