



decture 19

Clustering - Kmeans

Lg Clustering HWG SOLS 5,6

Lecture Notes IX - Clustering

Marina Meilă mmp@stat.washington.edu

> Department of Statistics University of Washington

> > May, 2023

Paradigms for clustering

Parametric clustering algorithms (K given)

Cost based / hard clustering **Figure** K-means clustering and the quadratic distortion **Content** Model based / soft clustering **Content**

Issues in parametric clustering • • Selecting *K*

Reading: Ch. 18

What is clustering? Problem and Notation

- Informal definition Clustering = Finding groups in data
- ▶ Notation $\mathcal{D} = \{x_1, x_2, \dots, x_n\}$ a data set n = number of data points K = number of clusters (K << n) $\Delta = \{C_1, C_2, \dots, C_K\}$ a partition of \mathcal{D} into disjoint subsets k(i) = the label of point i $\mathcal{L}(\Delta) =$ cost (loss) of Δ (to be minimized)
- Second informal definition Clustering = given n data points, separate them into K clusters
- Hard vs. soft clusterings
 - Hard clustering Δ: an item belongs to only 1 cluster
 - **Soft** clustering $\gamma = {\gamma_{ki}}_{k=1:K}^{i=1:n}$

 γ_{ki} = the degree of membership of point *i* to cluster *k*

$$\sum_k \gamma_{ki} = 1 \quad \text{for all } i$$

(usually associated with a probabilistic model)



from Carreira-Perpinan, 2006



image segmentation n = #pixals K = 5

xⁱ = pixel i = ?) grey value = [0,25] or 2) other <u>feature</u>

Paradigms

Depend on type of data, type of clustering, type of cost (probabilistic or not), and constraints (about K, shape of clusters)



Parametric clustering algorithms

- Cost based
 - Single linkage (min spanning tree)
 - Min diameter
 - Fastest first traversal (HS initialization)
 - K-medians
 - K-means
- Model based (cost is derived from likelihood)
 - EM algorithm
 - "Computer science" /" Probably correct" algorithms

[Supplement: Single Linkage Clustering]

Algorithm Single-Linkage

Input Data $\mathcal{D} = \{x_i\}_{i=1:n}$, number clusters K

- 1. Construct the Minimum Spanning Tree (MST) of ${\cal D}$
- 2. Delete the largest K 1 edges
- ► **Cost** $\mathcal{L}(\Delta) = -\min_{k,k'} \operatorname{distance}(C_k, C_{k'})$ where $\operatorname{distance}(A, B) = \operatorname*{argmin}_{x \in A, y \in B} ||x - y||$
- ▶ Running time $O(n^2)$ one of the very few costs \mathcal{L} that can be optimized in polynomial time
- Sensitive to outliers!

[Supplement: Single Linkage Clustering]



Observations

[Supplement: Minimum diameter clustering]

• Cost
$$\mathcal{L}(\Delta) = \max_k \max_{\substack{i,j \in C_k}} ||x_i - x_j||$$

diameter

Mimimize the diameter of the clusters

Optimizing this cost is NP-hard

Algorithms

Fastest First Traversal – a factor 2 approximation for the min cost For every \mathcal{D} , FFT produces a Δ so that

$$\mathcal{L}^{opt} \leq \mathcal{L}(\Delta) \leq 2\mathcal{L}^{opt}$$

rediscovered many times

[Supplement: Minimum diameter clustering]

Algorithm Fastest First Traversal Input Data $\mathcal{D} = \{x_i\}_{i=1:n}$, number clusters Kdefines centers $\mu_{1:K} \in \mathcal{D}$ (many other clustering algorithms use centers) 1. pick μ_1 at random from \mathcal{D} 2. for k = 2: K $\mu_k \leftarrow \operatorname{argmax} \operatorname{distance}(x_i, \{\mu_{1:k-1}\})$ 3. for i = 1: n (assign points to centers) k(i) = k if μ_k is the nearest center to x_i

[Supplement: K-medians clustering]

• Cost $\mathcal{L}(\Delta) = \sum_k \sum_i i \in C_k ||x_i - \mu_k||$ with $\mu_k \in \mathcal{D}$

(usually) assumes centers chosen from the data points (analogy to median) Exercise Show that in 1D $\operatorname{argmin} \sum_{i} |x_i - \mu|$ is the median of $\{x_i\}$

optimizing this cost is NP-hard

has attracted a lot of interest in theoretical CS (general from called "Facility location"

K-means clustering

Algorithm K-Means



K-means clustering

Algorithm K-Means

Input Data $\mathcal{D} = \{x_i\}_{i=1:n}$, number clusters K tialize centers $\mu_1, \mu_2, \dots \mu_K \in \mathbb{R}^d$ at random terate until convergence 1. for i = 1 : n (assign points to clusters \Rightarrow new clustering) · Convergence If arignment $k(i) = \operatorname{argmin}_{i} ||x_i - \mu_k||$ 2. for k = 1: K (recalculate centers) don't change (1) $\mu_k = \frac{1}{|C_k|} \sum_{i \in C_k} x_i$ 12 11 C1 10 3L Hims -10 3

The K-means
$$\cot z = \log (z^{(n)} \cdot is^{(n)}) \cdot how \text{ gook is } \underline{K}^{?}$$

 $\cdot \mathcal{L}(\Delta) = \sum_{k=1}^{K} \sum_{i \in C_{k}} ||x_{i} - \mu_{k}||^{2}$
 $\wedge \mathcal{L}(\Delta) = \sum_{k=1}^{K} \sum_{i \in C_{k}} ||x_{i} - \mu_{k}||^{2}$
 $\wedge \mathcal{L}(\Delta) = \sum_{k=1}^{K} \sum_{i \in C_{k}} ||x_{i} - \mu_{k}||^{2}$
 $\wedge \mathcal{L}(\Delta) = \int C_{1,j} \cdot C_{k}^{?} (2)$
 $\wedge \mathcal{L}(\Delta) = \int C_{1,j} \cdot C_{k} \cdot C_{k}$

Sketch of proof

- step 1: reassigning the labels can only decrease \mathcal{L}
- step 2: reassigning the centers µ_k can only decrease L because µ_k as given by (1) is the solution to

$$\mu_k = \min_{\mu \in \mathbb{R}^d} \sum_{i \in C_k} ||\mathbf{x}_i - \mu||^2 \tag{3}$$

T.

10

 $t=2 \text{ assign } X_{1:n} \to \text{ obsers } (\mu_{1:K}^{2}) \Rightarrow \text{ obsers } \Lambda^{2} \leftarrow \text{from } \mu_{1:K}^{4}$ $for X_{i}: \text{ same } C_{K} \qquad \|X_{i} - \mu_{K}\|^{2} = \min_{k'=k} \|X_{i} - \mu_{K}\|^{2}$ $\int_{i=k}^{k} |X_{i} - \mu_{K}|^{2} \leq \|X_{i} - \mu_{K}\|^{2} \qquad \int_{i=k}^{k} |X_{i} - \mu_{K}|^{2} \leq \|X_{i} - \mu_{K}\|^{2}$

$$\Delta : x_i \in C_k \quad \text{another } C_k' \iff \|x_i - \mu_k\| \quad \text{decreases} \ (k_i = k \quad k_i = k \quad k_i = k \quad \text{decreases} \ (k_i = k \quad k_i = k \quad k_i = k \quad k_i = \mu_{k_i} \quad$$

$$\Rightarrow \lambda(\Delta^{2}, \mu_{1:k}^{1}) \leq \lambda(\Delta^{1}, \mu_{1:k}^{1})$$

$$< \mathcal{H} \Delta^{1} \neq \Delta^{2}$$

$$t=3 \quad \text{recalculate} \quad \mu's : \mu_{1:k}^{2} \in \cdots$$

$$\qquad \lambda(\Delta^{2}, \mu_{1:k}^{2}) \leq \lambda(\Delta^{2}, \mu_{1:k}^{1})$$

$$\text{recalculate} \quad \Delta$$

$$\qquad \lambda(\Delta^{3}, \mu_{1:k}^{3}) \leq \lambda(\Delta^{2}, \mu_{1:k}^{2})$$

⇒ Convergence
$$(=) \land (\land, \mu)$$
 doesn't champe
~ This is Locar Min \land is best only w.r.t small champ
at conver no

· How to choose & after several initializations? for b = 1:B Initialilize $\mu_{1:k}$ Run K-means algorithm $\Rightarrow \Delta^{(b)}, \mu_{1:k}^{(b)}$ Output arguin $\mathcal{L}(\mathcal{L}^{(b)}, \mu_{q:k}^{(b)})$ How to guess us initially?

[Supplement: Equivalent and similar cost functions]

The distortion can also be expressed using intracluster distances

$$\mathcal{L}(\Delta) = \sum_{k=1}^{K} \frac{1}{n_k} \sum_{i,j \in C_k} ||x_i - x_j||^2$$
(4)

Correlation clustering is defined as optimizing the related criterion

$$\mathcal{L}(\Delta) = \sum_{k=1}^{K} \sum_{i,j \in C_k} ||x_i - x_j||^2$$

This cost is equivalent to the (negative) sum of (squared) intercluster distances

$$\mathcal{L}(\Delta) = -\sum_{k=1}^{K} \sum_{i \in C_k} \sum_{j \notin C_k} ||x_i - x_j||^2 + \text{constant}$$
(5)

Proof of (6) Replace μ_k as expressed in (1) in the expression of \mathcal{L} , then rearrange the terms **Proof of (5)** $\sum_k \sum_{i,j \in C_k} ||x_i - x_j||^2 = \sum_{i=1}^n \sum_{j=1}^n ||x_i - x_j||^2 - \sum_k \sum_{i \in C_k} \sum_{j \notin C_k} ||x_i - x_j||^2$ independent of Δ

[Supplement: The K-means cost in matrix form – the assignment matrix]

• \mathcal{L} as sum of squared intracluster distances

$$\mathcal{L}(\Delta) = \sum_{k=1}^{K} \frac{1}{|C_k|} \sum_{i,j \in C_k} ||x_i - x_j||^2$$
(6)

• Define the assignment matrix associated with Δ by $Z(\Delta)$ Let $\Delta = \{C_1 = \{1, 2, 3\}, C_2 = \{4, 5\}\}$

$$Z^{unnorm}(\Delta) = \begin{bmatrix} C_1 & C_2 & & & C_1 & C_2 \\ 1 & 0 \\ 1 & 0 \\ 0 & 1 \\ 0 & 1 \end{bmatrix} \text{ point } i \quad Z(\Delta) = \begin{bmatrix} C_1 & C_2 \\ 1/\sqrt{3} & 0 \\ 1/\sqrt{3} & 0 \\ 1/\sqrt{3} & 0 \\ 0 & 1/\sqrt{2} \\ 0 & 1/\sqrt{2} \end{bmatrix}$$

Then Z is an orthogonal matrix (columns are orthornormal) and

$$\mathcal{L}(\Delta) = \operatorname{trace} Z^T D Z \quad \text{with } D_{ij} = ||x_i - x_j||^2$$
(7)

Let $\mathcal{Z} = \{ Z \in \mathbb{R}^{n \times K}, K \text{ orthonormal} \}$

Proof of (7) Start from (2) and note that trace $Z^T A Z = \sum_k \sum_{i,j \in C_k} Z_{ik} Z_{jk} A_{ij} = \sum_k \sum_{i,j \in C_k} \frac{1}{|C_k|} A_{ij}$

[Supplement: The K-means cost in matrix form – the co-ocurrence matrix]

$$n = 5, \ \Delta = (1, 1, 1, 2, 2),$$

$$X(\Delta) = \begin{bmatrix} \frac{1}{2} & \frac{1}{2} & \frac{1}{2} & 0 & 0 \\ \frac{1}{2} & \frac{1}{2} & \frac{1}{2} & 0 & 0 \\ \frac{1}{3} & \frac{1}{3} & \frac{1}{3} & 0 & 0 \\ 0 & 0 & 0 & \frac{1}{2} & \frac{1}{2} \\ 0 & 0 & 0 & 0 & \frac{1}{2} & \frac{1}{2} \end{bmatrix}$$

- 1. $X(\Delta)$ is symmetric, positive definite, ≥ 0 elements
- 2. $X(\Delta)$ has row sums equal to 1
- 3. trace $X(\Delta) = K$

$$\begin{aligned} \|X(\Delta)\|_{F}^{2} &= \langle X, X \rangle = K\\ X(\Delta) &= Z(\Delta)Z^{T}(\Delta) \end{aligned}$$

$$2\mathcal{L}(\Delta) = \sum_{k=1}^{K} \frac{1}{|C_k|} \sum_{i,j \in C_k} ||x_i - x_j||^2 = \frac{1}{2} \langle D, X(\Delta) \rangle$$

with $D_{ij} = ||x_i - x_j||^2$

[Supplement: Symmetries between costs]

- K-means cost $\mathcal{L}(\Delta) = \min_{\mu_{1:K}} \sum_k \sum_{i \in C_k} ||x_i \mu_k||^2$
- K-medians cost $\mathcal{L}(\Delta) = \min_{\mu_{1:K}} \sum_{k} \sum_{i \in C_k} ||x_i \mu_k||$
- Correlation clustering cost $\mathcal{L}(\Delta) = \sum_k \sum_{i,j \in C_k} ||x_i x_j||^2$
- min Diameter cost $\mathcal{L}^2(\Delta) = \max_k \max_{i,j \in C_k} ||x_i x_j||^2$

▶ Idea 1: start with K points at random

Idea 1: start with K points at random \leftarrow Idea 2: start with K data points at random

_ P5: Ne may be for away from data



- ▶ Idea 1: start with K points at random
- Idea 2: start with K data points at random What's wrong with chosing K data points at random?



The probability of hitting all K clusters with K samples approaches 0 when K > 5

Idea 1: start with K points at random

Idea 2: start with K data points at random What's wrong with chosing K data points at random?



The probability of hitting all K clusters with K samples approaches 0 when K > 5Idea 3: start with K data points using Fastest First Traversal (greedy simple approach to

Idea 3: start with K data points using Fastest First Traversal (greedy simple approach to spread out centers)



More special cases introduce the following description for a covariance matrice in terms of volume, shape, alignment with axes (=determinant, trace, e-vectors). The letters below mean: I=unitary (shape, axes), E=equal (for all k), V=unequal

- EII: equal volume, round shape (spherical covariance)
- VII: varying volume, round shape (spherical covariance)
- EEI: equal volume, equal shape, axis parallel orientation (diagonal covariance)
- VEI: varying volume, equal shape, axis parallel orientation (diagonal covariance)
- EVI: equal volume, varying shape, axis parallel orientation (diagonal covariance)
- VVI: varying volume, varying shape, equal orientation (diagonal covariance)
- EEE: equal volume, equal shape, equal orientation (ellipsoidal covariance)
- EEV: equal volume, equal shape, varying orientation (ellipsoidal covariance)
- VEV: varying volume, equal shape, varying orientation (ellipsoidal covariance)
- VVV: varying volume, varying shape, varying orientation (ellipsoidal covariance)

(from)

EM versus K-means

- Alternates between cluster assignments and parameter estimation
- Cluster assignments γ_{ki} are probabilistic
- Cluster parametrization more flexible



 Converges to local optimum of log-likelihood Initialization recommended by K-logK method

- Modern algorithms with guarantees (for e.g. mixtures of Gaussians)
 - Random projections
 - Projection on principal subspace
 - Two step EM (=K-logK initialization + one more EM iteration)

[Supplement: A two-step EM algorithm]

Similar to K-logK initialization for K-means

Assumes K spherical gaussians, separation $\|\mu_k^{true} - \mu_{k'}^{true} \ge C\sqrt{d}\sigma_k$

- 1. Pick $K' = \mathcal{O}(K \ln K)$ centers μ_k^0 at random from the data
- 2. Set $\sigma_k^0 = \frac{d}{2} \min_{k \neq k'} ||\mu_k^0 \mu_{k'}^0||^2$, $\pi_k^0 = 1/K'$
- 3. Run one E step and one M step $\implies \{\pi_k^1, \mu_k^1, \sigma_k^1\}_{k=1:K'}$
- 4. Compute "distances" $d(\mu_k^1, \mu_{k'}^1) = \frac{||\mu_k^1 \mu_{k'}^1||}{\sigma_k^1 \sigma_{k'}^1}$
- 5. Prune all clusters with $\pi_k^1 \leq 1/4K'$
- Run Fastest First Traversal with distances d(μ¹_k, μ¹_{k'}) to select K of the remaining centers. Set π¹_k = 1/K.
- 7. Run one E step and one M step $\implies \{\pi_k^2, \mu_k^2, \sigma_k^2\}_{k=1:K}$
- eorem For any $\delta, \varepsilon > 0$ if d large, n large enough, separation $C \ge d^{1/4}$ the Two step EM algorithm obtains centers μ_k so that

$$||\mu_k - \mu_k^{true}|| \le ||\text{mean}(C_k^{true}) - \mu_k^{true}|| + \varepsilon \sigma_k \sqrt{d}$$

Selecting K

- Run clustering algorithm for $K = K_{min} : K_{max}$
 - obtain $\Delta_{K_{min}}, \ldots \Delta_{K_{max}}$ or $\gamma_{K_{min}}, \ldots \gamma_{K_{max}}$
 - choose best Δ_K (or γ_K) from among them
- Typically increasing $K \Rightarrow \text{cost } \mathcal{L} \text{ decreases}$
 - (\mathcal{L} cannot be used to select K)
 - Need to "penalize" L with function of number parameters

Selecting K for mixture models

The BIC (Bayesian Information) Criterion

- let θ_K = parameters for γ_K
- ▶ let $\#\theta_K$ =number independent parameters in θ_K
 - e.g for mixture of Gaussians with full Σ_k 's in d dimensions

$$\#\theta_{K} = \underbrace{K-1}_{\pi_{1:K}} + \underbrace{Kd}_{\mu_{1:K}} + \underbrace{Kd(d-1)/2}_{\Sigma_{1:K}}$$

define

$$BIC(\theta_{K}) = I(\theta_{K}) - \frac{\#\theta_{K}}{2} \ln r$$

- Select K that maximizes $BIC(\theta_K)$
- selects true K for $n \to \infty$ and other technical conditions (e.g parameters in compact set)
- but theoretically not justified (and overpenalizing) for finite n

Number of Clusters vs. BIC EII (A), VII (B), EEI (C), VEI (D),

EEV, 8 Cluster Solution

EVI (E), VVI (F), EEE (G), EEV (H), VEV (I), VVV (J)



(from)

Number of Clusters vs. BIC EII (A), VII (B), EEI (C), VEI (D), EVI (E), VVI (F), EEE (G), EEV (H), VEV (I), VVV (J)

EEV, 8 Cluster Solution



(from)

[Supplement: Stability methods for choosing K]

- like bootstrap, or crossvalidation
- Idea (implemented by)

for each K

- 1. perturb data $\mathcal{D} \rightarrow \mathcal{D}'$
- 2. cluster $\mathcal{D}' \to \Delta'_{\mathcal{K}}$
- 3. compare Δ_K , Δ'_K . Are they similar? If yes, we say Δ_K is stable to perturbations

Fundamental assumption If Δ_K is stable to perturbations then K is the correct number of clusters

- these methods are supported by experiments (not extensive)
- not YET supported by theory ... see for a summary of the area

Clustering with outliers

- What are outliers?
- let p = proportion of outliers (e.g 5%-10%)
- Remedies
 - mixture model: introduce a K + 1-th cluster with large (fixed) Σ_{K+1} , bound Σ_k away from 0
 - K-means and EM
 - robust means and variances
 e.g eliminate smallest and largest pnk/2 samples in mean computation (trimmed mean)
 - K-medians
 - replace Gaussian with a heavier-tailed distribution (e.g. Laplace)
 - single-linkage: do not count clusters with < r points</p>
 - Is K meaningful when outliers present?
 - alternative: non-parametric clustering