



N

Nebrask	а
Linco	dr

ephen Sco

ntroduction

lustering

Dutline Clustering

Clustering Background

Types of clustering problems:

- Hard (crisp): partition data into non-overlapping clusters; each instance belongs in exactly one cluster
- Fuzzy: Each instance could be a member of multiple clusters, with a real-valued function indicating the degree of membership
- Hierarchical: partition instances into numerous small clusters, then group the clusters into larger ones, and so on (applicable to phylogeny)
 - End up with a tree with instances at leaves

lebiaska Lincoln	Clustering Background (Dis-)similarity Measures: Between Instances
CSCE 478/878 Lecture 8: Clustering ephen Scott	Dissimilarity measure: Weighted L_p norm: $L_p(\mathbf{x}, \mathbf{y}) = \left(\sum_{i=1}^n w_i x_i - y_i ^p\right)^{1/p}$
utline ustering easures: int-Point	Special cases include weighted Euclidian distance $(p = 2)$, weighted Manhattan distance
easures: Point-Set easures: Set-Set Weans ustering	$L_1(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^n w_i x_i - y_i ,$
ustering	and weighted L_{∞} norm
	$L_{\infty}(\mathbf{x}, \mathbf{y}) = \max_{1 \le i \le n} \{w_i x_i - y_i \}$

Similarity measure: Dot product between two vectors (kernel)

・ロト・西ト・ボン・ボン・ボー シック



troductior

ustering

Clustering Background (Dis-)similarity Measures: Between Instances (cont'd)

If attributes come from $\{0, \ldots, k-1\}$, can use measures for real-valued attributes, plus:

- Hamming distance: DM measuring number of places where x and y differ
- Tanimoto measure: SM measuring number of places where x and y are same, divided by total number of places
 - Ignore places *i* where $x_i = y_i = 0$
 - Useful for ordinal features where x_i is degree to which x possesses ith feature

Nebraska Lindin (Dis-)similarity Measures: Between Instance and Set

- Might want to measure proximity of point **x** to existing cluster *C*
- Can measure proximity α by using all points of C or by using a representative of C
- If all points of C used, common choices:

$$\begin{split} \alpha^{ps}_{max}(\mathbf{x}, C) &= \max_{\mathbf{y} \in C} \left\{ \alpha(\mathbf{x}, \mathbf{y}) \right\} \\ \alpha^{ps}_{min}(\mathbf{x}, C) &= \min_{\mathbf{y} \in C} \left\{ \alpha(\mathbf{x}, \mathbf{y}) \right\} \\ \alpha^{ps}_{avg}(\mathbf{x}, C) &= \frac{1}{|C|} \sum_{\mathbf{y} \in C} \alpha(\mathbf{x}, \mathbf{y}) \end{split}$$

where $\alpha(\mathbf{x}, \mathbf{y})$ is any measure between \mathbf{x} and \mathbf{y}



troduction

Clustering Measures: Point-Point

Nebraska Lincoln	k-Means Clustering	Nebraska Lincoln	k-Means Cluster
CSCE 478/878 Lecture 8: Clustering Stephen Scott Introduction Outline Clustering Agenthm Example Hierarchical Clustering	 Very popular clustering algorithm Represents cluster <i>i</i> (out of <i>k</i> total) by specifying its representative m_i (not necessarily part of the original set of instances X) Each instance x ∈ X is assigned to the cluster with nearest representative Goal is to find a set of <i>k</i> representatives such that sum of distances between instances and their representatives is minimized NP-hard (intractable) in general Will use an algorithm that alternates between determining representatives and assigning clusters until convergence (in the style of the EM algorithm) 	CSCE 478/878 Lecture 8: Clustering Stephen Scott Introduction Outline Clustering Agaithm Example Hierarchical Clustering	 Choose value fot Initialize k arbitra E.g., k rando Repeat until rep For all x ∈ X Assign x measure I.e., nea Por each j ∈

EVALUATE: A subset of the set of the set

Nebiaska

k-Means Clustering Example with *k* = 2



Nebraska Hierarchical Clustering

Outline

lustering

Hierarchica

Clustering

• Useful in capturing hierarchical relationships, e.g., evolutionary tree of biological sequences

- End result is a **sequence** (hierarchy) of clusterings
- Two types of algorithms:
 - Agglomerative: Repeatedly merge two clusters into one
 - Divisive: Repeatedly divide one cluster into two

Nebraska

478/878 ectur<u>e 8</u>

Stephen Sco

Introduction

Clustering

ustering

Hierarchical Clustering

- Let $C_t = \{C_1, \ldots, C_{m_t}\}$ be a **level**-*t* clustering of $X = \{\mathbf{x}_1, \ldots, \mathbf{x}_N\}$, where C_t meets definition of hard clustering
- C_t is **nested** in $C_{t'}$ (written $C_t \sqsubset C_{t'}$) if each cluster in C_t is a subset of a cluster in $C_{t'}$ and at least one cluster in C_t is a proper subset of some cluster in $C_{t'}$

$$\begin{aligned} \mathcal{C}_1 &= \{\{\mathbf{x}_1, \mathbf{x}_3\}, \{\mathbf{x}_4\}, \{\mathbf{x}_2, \mathbf{x}_5\}\} \sqsubset \{\{\mathbf{x}_1, \mathbf{x}_3, \mathbf{x}_4\}, \{\mathbf{x}_2, \mathbf{x}_5\}\}\\ \mathcal{C}_1 \not\sqsubset \{\{\mathbf{x}_1, \mathbf{x}_4\}, \{\mathbf{x}_3\}, \{\mathbf{x}_2, \mathbf{x}_5\}\}\end{aligned}$$

Hierarchical Clustering Nebraska Definitions (cont'd) 478/878 Lecture 8 Clustering Agglomerative algorithms start with $C_0 = \{\{\mathbf{x}_1\}, \dots, \{\mathbf{x}_N\}\}$ and at each step *t* merge two tephen Sc clusters into one, yielding $|C_{t+1}| = |C_t| - 1$ and $C_t \sqsubset C_{t+1}$ ntroduction • At final step (step N - 1) have hierarchy: Dutline Clustering $\mathcal{C}_0 = \{\{\mathbf{x}_1\}, \ldots, \{\mathbf{x}_N\}\} \sqsubset \mathcal{C}_1 \sqsubset \cdots \sqsubset \mathcal{C}_{N-1} = \{\{\mathbf{x}_1, \ldots, \mathbf{x}_N\}\}$ -Means ustering • Divisive algorithms start with $C_0 = \{\{x_1, \dots, x_N\}\}$ and at each step t split one cluster into two, yielding $|\mathcal{C}_{t+1}| = |\mathcal{C}_t| + 1$ and $\mathcal{C}_{t+1} \sqsubset \mathcal{C}_t$ • At step N - 1 have hierarchy: $\mathcal{C}_{N-1} = \{\{\mathbf{x}_1\}, \ldots, \{\mathbf{x}_N\}\} \sqsubset \cdots \sqsubset \mathcal{C}_0 = \{\{\mathbf{x}_1, \ldots, \mathbf{x}_N\}\}$

<ロ> <畳> < 三> < 三> < 三> < 三> < 三> < 三> のへの

・ロト・西ト・モー・モー ひゃう

Nebraska Lincoln	Hierarchical Clustering
CSCE 478/878 Lecture 8: Clustering Stephen Scott Introduction	 Initialize C₀ = {{x₁},, {x_N}}, t = 0 For t = 1 to N - 1
Outline Clustering <i>k</i> -Means Clustering Hierarchical Clustering Definitions	 Find closest pair of clusters: (C_i, C_j) = argmin_{C_i,C_r∈C_{r-1}, r≠s} {d (C_s, C_r)} (C_i = (C_{t-1} - {C_i, C_j}) ∪ {{C_i ∪ C_j}} and update representatives if necessary
Pseudocode Example	If SM used, replace argmin with argmax Number of calls to $d(C_k, C_r)$ is $\Theta(N^3)$

Nebraska	Hierarchical Clustering
CSCE 478/878 Lecture 8: Clustering Stephen Scott Introduction Outline	$ \begin{aligned} \mathbf{x}_1 &= [1, 1]^T, \mathbf{x}_2 = [2, 1]^T, \mathbf{x}_3 = [5, 4]^T, \mathbf{x}_4 = [6, 5]^T, \\ \mathbf{x}_5 &= [6.5, 6]^T, DM = Euclidian/\alpha_{\min}^{ss} \\ An \left(N - t \right) \times \left(N - t \right) proximity matrix P_t \text{ gives the proximity} \\ between all pairs of clusters at level (iteration) t \end{aligned} $
Clustering <i>k</i> -Means Clustering Hierarchical Clustering Definitions Pasudocode Example	$P_0 = \begin{bmatrix} 0 & 1 & 5 & 6.4 & 7.4 \\ 1 & 0 & 4.2 & 5.7 & 6.7 \\ 5 & 4.2 & 0 & 1.4 & 2.5 \\ 6.4 & 5.7 & 1.4 & 0 & 1.1 \\ 7.4 & 6.7 & 2.5 & 1.1 & 0 \end{bmatrix}$
	Each iteration, find minimum off-diagonal element (i, j) in P_{t-1} , merge clusters <i>i</i> and <i>j</i> , remove rows/columns <i>i</i> and <i>j</i> from P_{t-1} , and add new row/column for new cluster to get P_t

Nebraska Linon Hierarchical Clustering Pseudocode (cont'd)

CSCE 478/878 Lecture 8: Clustering Stephen Scott Introduction Outline Clustering A-Means Clustering Hierarchical Clustering Definitions Peeudocote Example

A **proximity dendogram** is a tree that indicates hierarchy of clusterings, including the proximity between two clusters when they are merged $x_1 \quad x_2 \quad x_3 \quad x_4 \quad x_5$



Cutting the dendogram at any level yields a single clustering