

2022.2.7

## Lecture 4

## Jaccard Similarity

	$w_1$	$w_2$	$w_d$
$x$	1	0	
$y$	1	1	

$w_i$  为 word,  $x, y$  为 documents. 若  $w_i$  出现在  $x$  中,  $w_i$  为 1.  
算 distance?

$$L_1(x, y) = \sum_{i=1}^d |x_i - y_i| \text{ 有多少词不同}$$

↓  
称为 Manhattan distance.

not useful! 例:  $x_1, y_1$  极长, 相似, 只有 1 个词不同

$x_2, y_2$  极短, 也只有 1 个词不同, 但相似程度不同

## Jaccard Similarity

$$J_{sim}(x, y) = \frac{|x \cap y|}{|x \cup y|} \quad J_{Dist}(x, y) = 1 - \frac{|x \cap y|}{|x \cup y|}$$

Distance 表示 dissimilarity. dis 与 sim 相反. 因此  $1 - \text{similarity} = \text{distance}$

clustering

将 data 分组: similar data points 在同组中

不相似的 data points 在不同组中



# Clustering 的方式

• partitional

每个 object 只属于 1 个 cluster (将  $n$  个 objects 分成  $k$  个 clusters)

计算每组中每对 points 之间的 distance 之和, 使整体最小.

$$\sum_k \sum_{x_i, x_j \in C_k} d(x_i, x_j)$$

每个 cluster 中每对 possible pair 间 distance 之和.

↓ 表示 cluster  $k$

太难算. instead, 找  $k$  个 points, 使  $\sum_i \sum_{x \in C_i} d(x, \mu_i)$  最小.

Lloyd's Algorithm:

1. randomly 选  $k$  个 centers 做  $\mu_i$

2. 将 每个 data points 归入 距 center 最小的那一组.

3. 计算各 cluster 的 mean, 成为 每组 cluster 新的 center.

4. 重复 2 & 3 直到 每个 data point 不再更改组别

每一次 iteration 中  $\sum_i \sum_{x \in C_i} d(x, \mu_i)$  会变小.

properties:

① Always converge!

② 不一定为最优 cluster

进阶: K-means++

选 1 个 random [data point] 做 1st center, 每次选下 1 个 center 时, 算当前 center 和 其它 data points 距离. 对每个下 1 个 potential centers, 其被选中概率 is proportional 到 其与当前 center 的距离.