

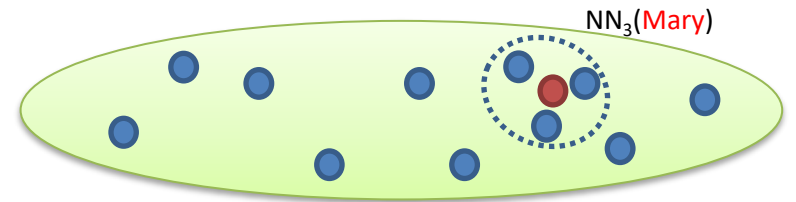
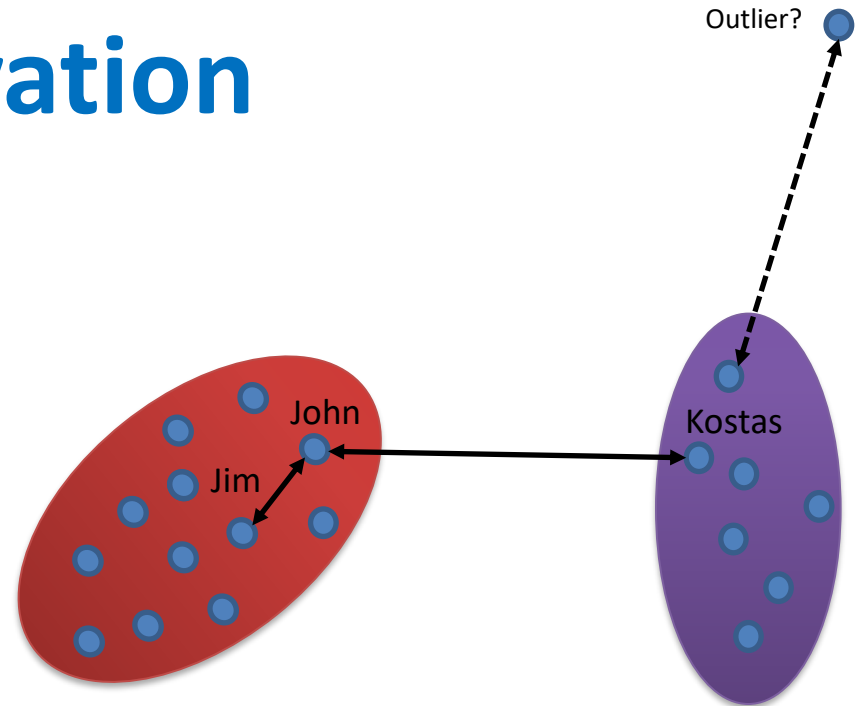
Working with Data

Yannis Kotidis
Department of
Informatics
Athens University of
Economics and Business



Motivation

- We often need ways to assess how **similar** or **dissimilar** objects are in comparison to one another
- Examples: **clustering**, **outlier** analysis, **nearest-neighbor** search, **recommendation**, **visualization**, **classification**



Simple Running Example

- Car dealership
 - a customer inquired about car #1 that was sold
 - which of the other cars is she most likely to buy?

Car	Color	Condition	Mileage (*1000)
1	Blue	Excellent	45
2	Green	Fair	22
3	Red	Good	64
4	Blue	Excellent	28

Roadmap

- We will first discuss simple similarity metrics for common data types: **nominal**, **ordinal** and **numerical** attributes
- We will then extend our techniques to address more complex scenarios
 - Hierarchical domains (e.g. product \Leftrightarrow categories)
 - Sets (e.g. basket data), Bags
 - Vectors (multidimensional data)
 - Strings
 - Time Series
 - Graphs

Preliminaries

- Let **sim(a,b)** denote the similarity of two values a, b of a data attribute
- In order to have a common basis, we will **normalize** values: $\text{sim}(a,b)$ in range $[0..1]$
 - $\text{sim}(a,b) = 1$, iff a and b are identical
 - $\text{sim}(a,b) = 0$, iff a and b are unlike
- **Dissimilarity**: $d(a,b) = 1 - \text{sim}(a,b)$
 - Formula assumes $\text{sim}(a,b) \in [0..1]$
 - Notice that, depending on the internals of the $\text{sim}()$ function, dissimilarity is not necessarily a distance function (i.e. triangle inequality may not hold)

Running Example

- Dataset describing used cars
- 3 known attributes : *color, condition, mileage* (in 1000 Km)

Car	Color	Condition	Mileage (*1000 Km)
1	Blue	Excellent	45
2	Green	Fair	22
3	Red	Good	64
4	Blue	Excellent	28

Nominal Attributes

- The values are **symbols** (e.g., names of things) that often represent some **category** or **state**
 - Are a type of categorical attributes when there is **no ordering**/importance implied by the values
 - Cats are not better than dogs or vice-versa
 - Can be used to differentiate records (e.g., pet="cat" vs pet="dog")
- Numerical attributes (e.g., product-ids) may also be treated as nominal

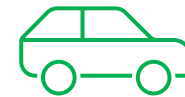
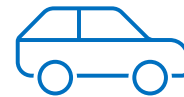
Dissimilarity of nominal attributes

- Let us define

$$d(a,b) = 1 \text{ if } a \neq b, 0 \text{ otherwise}$$

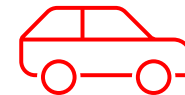
- Examples

– $d(\text{Blue}, \text{Green}) = 1$



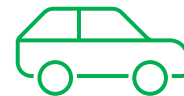
X

– $d(\text{Green}, \text{Red}) = 1$



X

– $d(\text{Green}, \text{Green}) = 0$



✓

Dissimilarity of nominal attributes

- We can form a *dissimilarity matrix* for **Color**:

Car 2 →

0	1	1	0
1	0	1	1
1	1	0	1
0	1	1	0

↑
Car 1

Car	Color	Condition	Mileage (*1000)
1	Blue	Excellent	45
2	Green	Fair	22
3	Red	Good	64
4	Blue	Excellent	28

Note

- Even for medium-sized datasets, the dissimilarity matrix may not fit in memory as it requires $O(n^2)$ space
- We use it in our examples in order to visually inspect the computed pair-wise dissimilarities

Dissimilarity of *ordinal* attributes

- In this case there is an ordering of the symbols
- Consider **rank** of different values
 - Fair (**1**) < Good (**2**) < Excellent (**3**)

- Let

$$d(a, b) = \frac{|rank(a) - rank(b)|}{maxrank - minrank}$$

- Examples

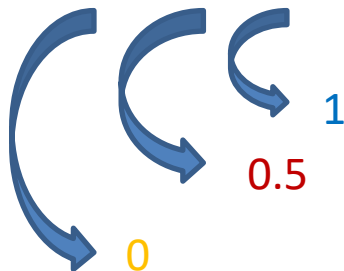
- $d(\text{Fair}, \text{Good}) = \frac{|1-2|}{3-1} = 0.5$

- $d(\text{Excellent}, \text{Fair}) = \frac{|3-1|}{3-1} = 1$

Dissimilarity of ordinal attributes

- Dissimilarity matrix for **Condition**:

0			
1	0		
0.5	0.5	0	
0	1	0.5	0



Car	Color	Condition	Mileage (*1000)
1	Blue	Excellent	45
2	Green	Fair	22
3	Red	Good	64
4	Blue	Excellent	28

Dissimilarity of numerical attributes

- Compute:

$$d(a, b) = \frac{|a - b|}{\text{maxvalue} - \text{minvalue}}$$

- Example

$$- d(45, 22) = \frac{|45 - 22|}{64 - 22} = 0.55$$

Car	Color	Condition	Mileage (*1000)
1	Blue	Excellent	45
2	Green	Fair	22
3	Red	Good	64
4	Blue	Excellent	28

Dissimilarity of numerical attributes

- Dissimilarity matrix for **Mileage**:

0			
0.55	0		
0.45	1	0	
0.40	0.14	0.86	0

Car	Color	Condition	Mileage (*1000)
1	Blue	Excellent	45
2	Green	Fair	22
3	Red	Good	64
4	Blue	Excellent	28

Combining scores

- Simplest approach: take average of computed dissimilarities

$$d(r_a, r_b) = \frac{1}{3} * (d_{\text{color}}(r_a, r_b) + d_{\text{condition}}(r_a, r_b) + d_{\text{mileage}}(r_a, r_b))$$

- Use weights to **prioritize** certain attributes
 - e.g. user prioritizes mileage over color ($w_3 > w_1$)

$$d(r_a, r_b) = \frac{w_1 * d_{\text{color}}(r_a, r_b) + w_2 * d_{\text{condition}}(r_a, r_b) + w_3 * d_{\text{mileage}}(r_a, r_b)}{w_1 + w_2 + w_3}$$

Example (weighted average)

- Assume **mileage** is more (**2x**) important than color, condition
 - $w_{\text{color}}=1$, $w_{\text{condition}}=1$, $w_{\text{mileage}}=2$

$$d(r_a, r_b) = \frac{d_{\text{color}}(r_a, r_b) + d_{\text{condition}}(r_a, r_b) + 2 * d_{\text{mileage}}(r_a, r_b)}{1 + 1 + 2}$$

Take avg of dissimilarity matrices

$$\begin{bmatrix} 0 & & & \\ 1 & 0 & & \\ 1 & 1 & 0 & \\ 0 & 1 & 1 & 0 \end{bmatrix} + \begin{bmatrix} 0 & & & \\ 1 & 0 & & \\ 0.5 & 0.5 & 0 & \\ 0 & 1 & 0.5 & 0 \end{bmatrix} + \begin{bmatrix} 0 & & & \\ 0.55 & 0 & & \\ 0.45 & 1 & 0 & \\ 0.40 & 0.14 & 0.86 & 0 \end{bmatrix}$$

3



$$\begin{bmatrix} 0 & & & \\ 0.85 & 0 & & \\ 0.65 & 0.83 & 0 & \\ 0.13 & 0.71 & 0.79 & 0 \end{bmatrix}$$

Outcome

- Data Matrix:

Car	Color	Condition	Millage (*1000)
1	Blue	Excellent	45
2	Green	Fair	22
3	Red	Good	64
4	Blue	Excellent	28

- Dissimilarity Matrix:

0			
0.85	0		
0.65	0.83	0	
0.13	0.71	0.79	0

Most similar pair of cars?

Most similar pair of cars

- Data Matrix:

Car	Color	Condition	Millage (*1000)
1	Blue	Excellent	45
2	Green	Fair	22
3	Red	Good	64
4	Blue	Excellent	28

- Dissimilarity Matrix:

0			
0.85	0		
0.65	0.83	0	
0.13	0.71	0.79	0

John: I like Car #3

- Data Matrix:

Car	Color	Condition	Millage (*1000)
1	Blue	Excellent	45
2	Green	Fair	22
3	Red	Good	64
4	Blue	Excellent	28

- Dissimilarity Matrix:

0			
0.85	0		
0.65	0.83	0	
0.13	0.71	0.79	0

Can we rank the cars on our lot based on their dissimilarities to car #3?

Nearest Neighbors of Car #3

- Data Matrix:

Car	Color	Condition	Millage (*1000)
1	Blue	Excellent	45
2	Green	Fair	22
3	Red	Good	64
4	Blue	Excellent	28

- Dissimilarity Matrix:

0			
0.85	0		
0.65	0.83	0	
0.13	0.71	0.79	0

NN(Car #3):

Car#1 ($d(1,3) = 0.65$)

Car#4 ($d(4,3) = 0.79$)

Car#2 ($d(3,2) = 0.83$)

Most similar



Less similar

Extended Data Matrix

Car	Color	Condition	Millage (*1000)	Type
1	Blue	Excellent	45	Supermini
2	Green	Fair	22	Crossover
3	Red	Good	64	SUV
4	Blue	Excellent	28	Small family

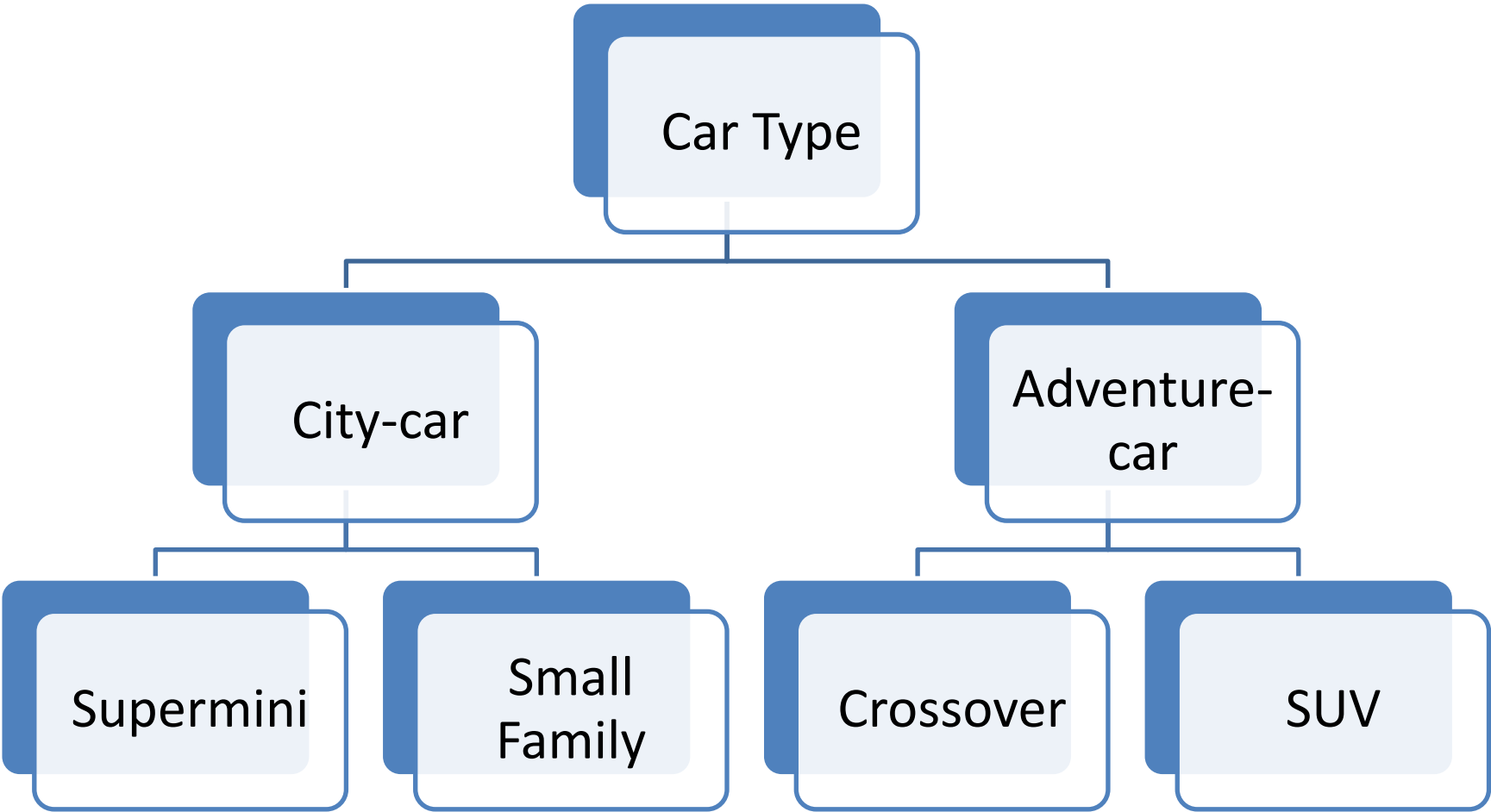
- How do we treat the newly added **Type** attribute?

Extended Data Matrix

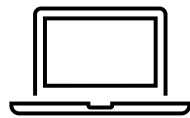
Car	Color	Condition	Millage (*1000)	Type
1	Blue	Excellent	45	Supermini
2	Green	Fair	22	Crossover
3	Red	Good	64	SUV
4	Blue	Excellent	28	Small family

- How do we treat the newly added **Type** attribute?
 - Small family cars more similar to superminis?
 - SUVs more similar to Crossover?

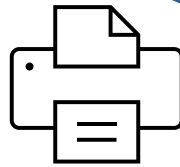
Car-Type Hierarchy?



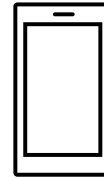
Six products



laptop



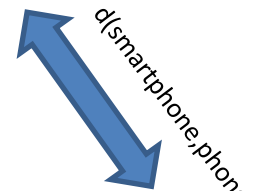
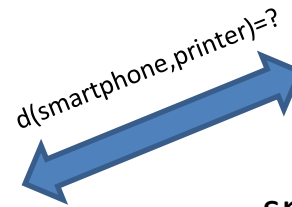
printer



smartphone



Phone



sneakers



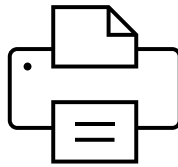
high heels

Groups/Categories

Computers



laptop



printer

Communication Devices



smartphone



Phone



sneakers



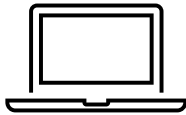
high heels

Shoes

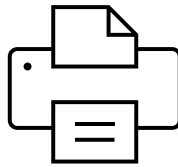
Higher-level Categories

Electronics

Computers

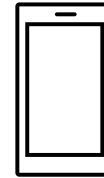


laptop



printer

Communication Devices



smartphone



Phone



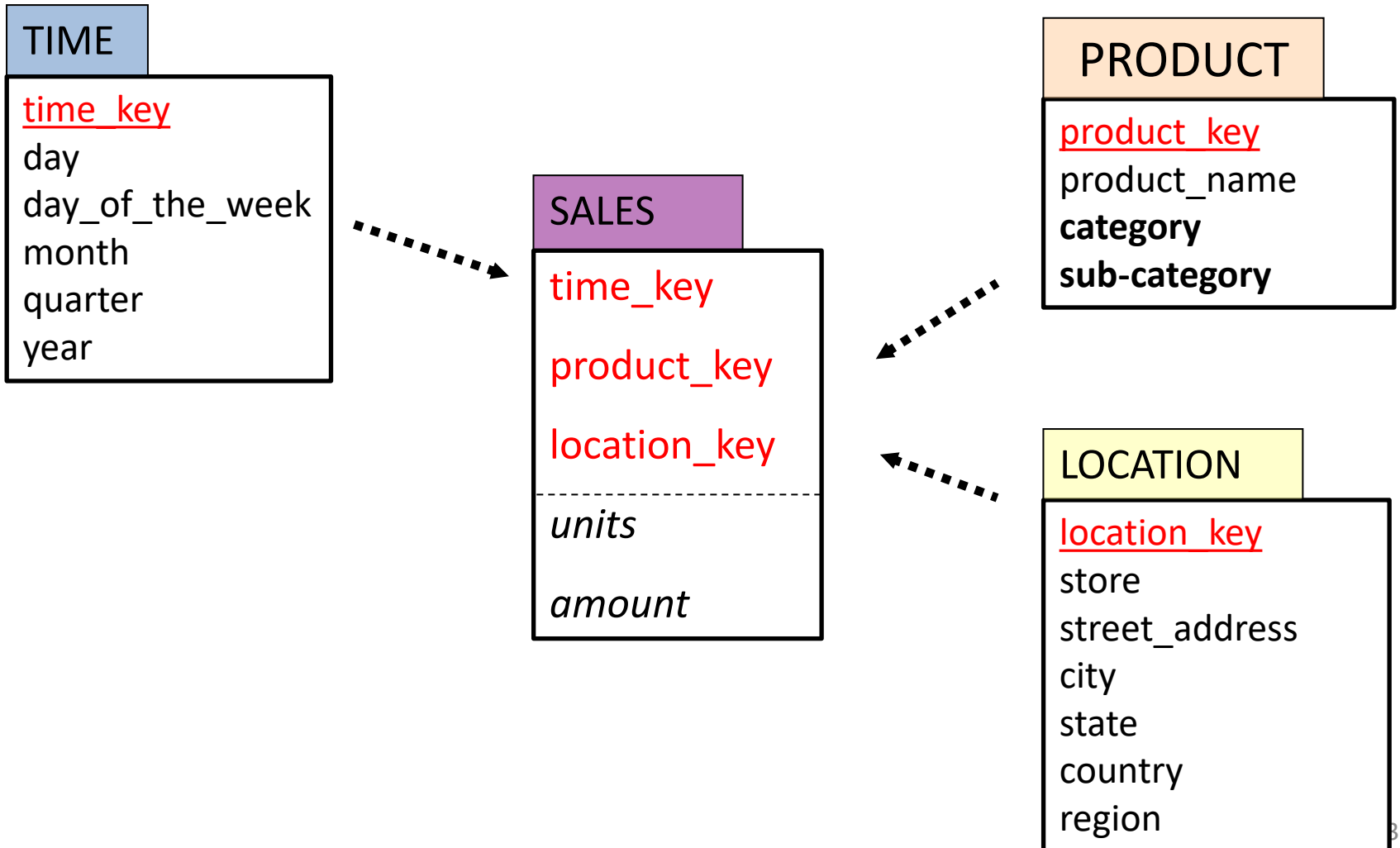
sneakers



high heels

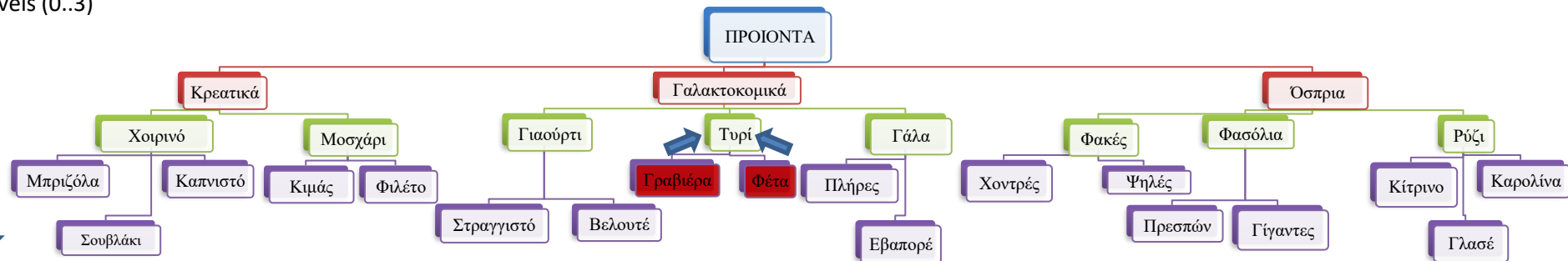
Shoes

Utilize the Star Schema



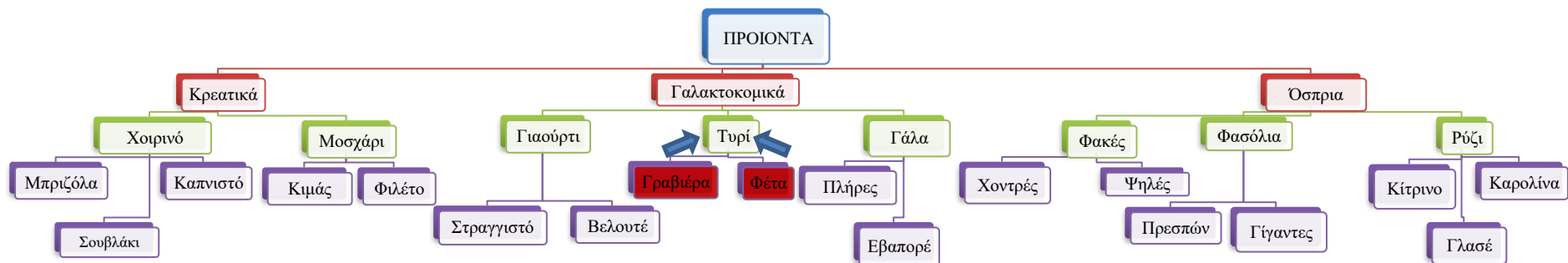
Using Hierarchies (simplified)

levels (0..3)



- Assume a,b are leaves and lca is their *lowest common ancestor* in the hierarchy
- Examples
 - $lca(\text{"Γραβιέρα"}, \text{"Φέτα"}) = \text{"Τυρί"}$
 - $lca(\text{"Γραβιέρα"}, \text{"Εβαπορέ"}) = \text{"Γαλακτοκομικά"}$
 - $lca(\text{"Γραβιέρα"}, \text{"Γίγαντες"}) = \text{"ΠΡΟΙΟΝΤΑ"}$

Using Hierarchies (simplified)



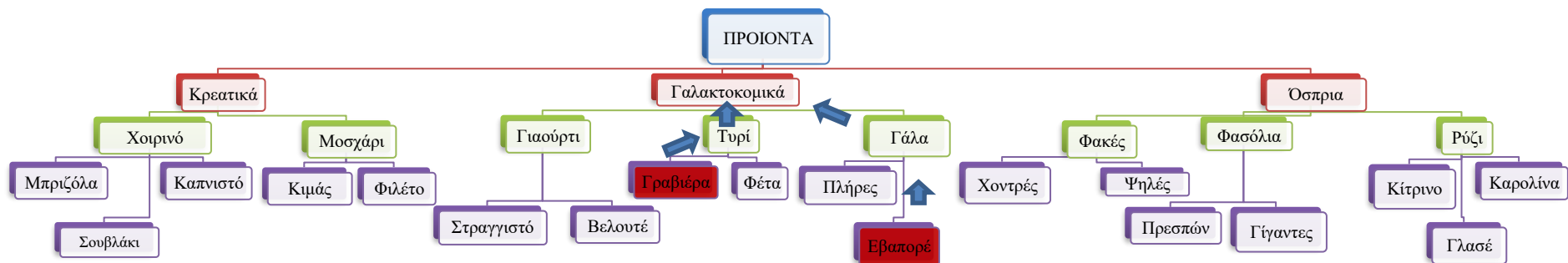
- Assume a, b are leaves and $|p_{lca}|$ = the length of the path towards their *lowest common ancestor* (lca)

– Define $d(a, b) = \frac{|p_{lca}|}{tree_height}$ if $a \neq b$, 0 otherwise

- Example

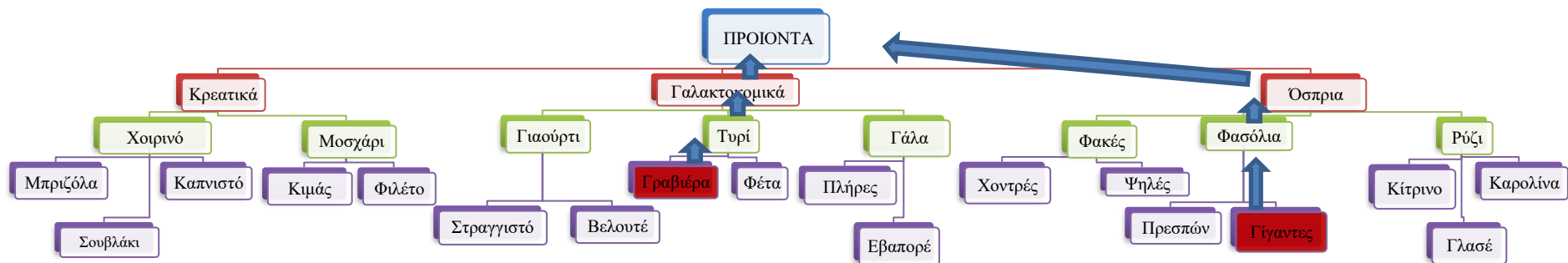
– $d(\text{"Γραβιέρα"}, \text{"Φέτα"}) = \frac{1}{3}$ ($lca = \text{"Τυρί"}$)

Using Hierarchies (simplified)



- Assume a, b are leaves and $|p_{lca}|$ = the length of the path towards their *lowest common ancestor* (lca)
 - Define $d(a, b) = \frac{|p_{lca}|}{tree_height}$ if $a \neq b$, 0 otherwise
- Example
 - $d(\text{"Γραβιέρα"}, \text{"Εβαπορέ"}) = \frac{2}{3}$ ($lca = \text{"Γαλακτοκομικά"}$)

Using Hierarchies (simplified)



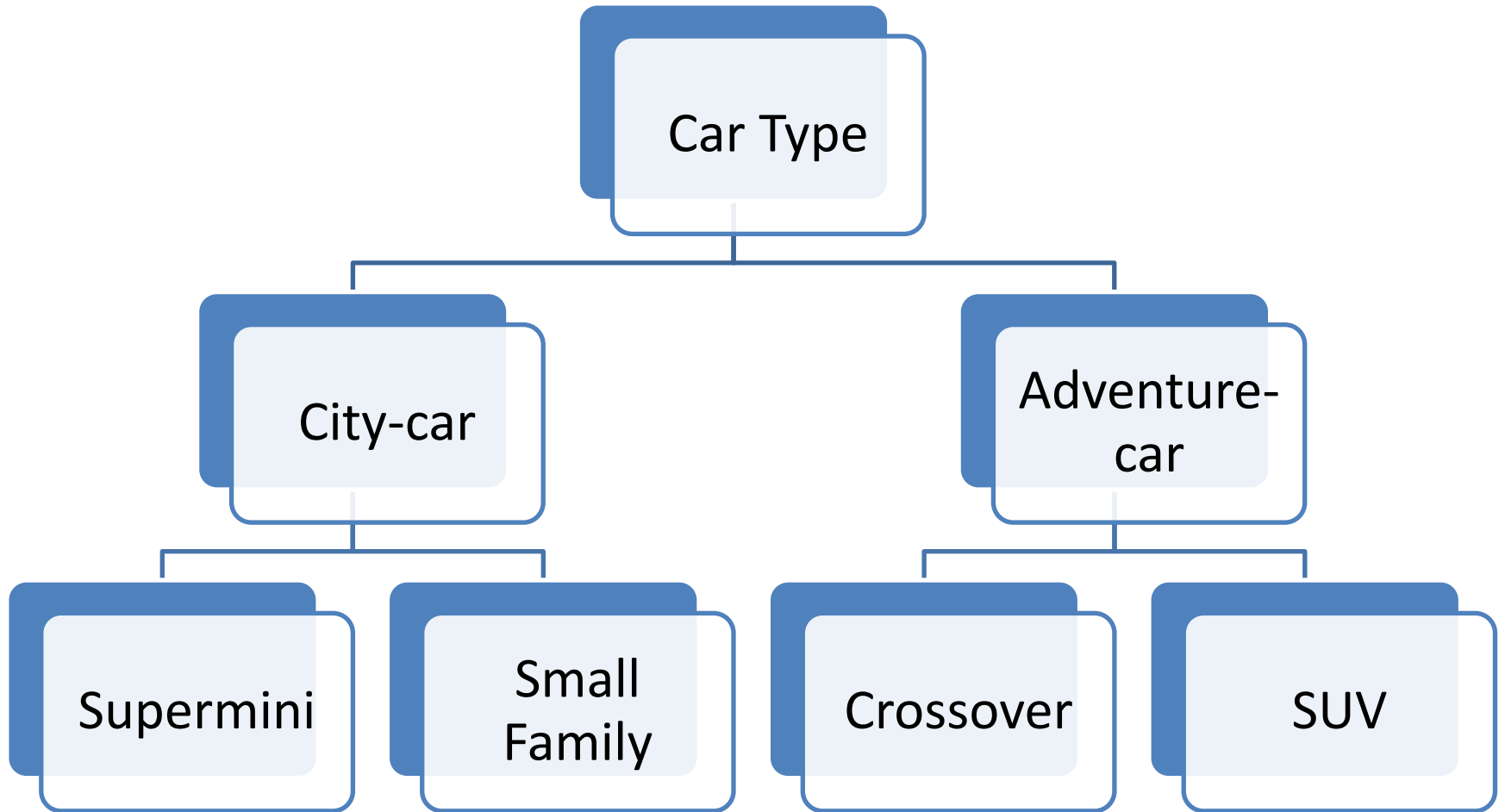
- Assume a, b are leaves and $|p_{lca}|$ = the length of the path towards their *lowest common ancestor* (lca)

– Define $d(a, b) = \frac{|p_{lca}|}{tree_height}$ if $a \neq b$, 0 otherwise

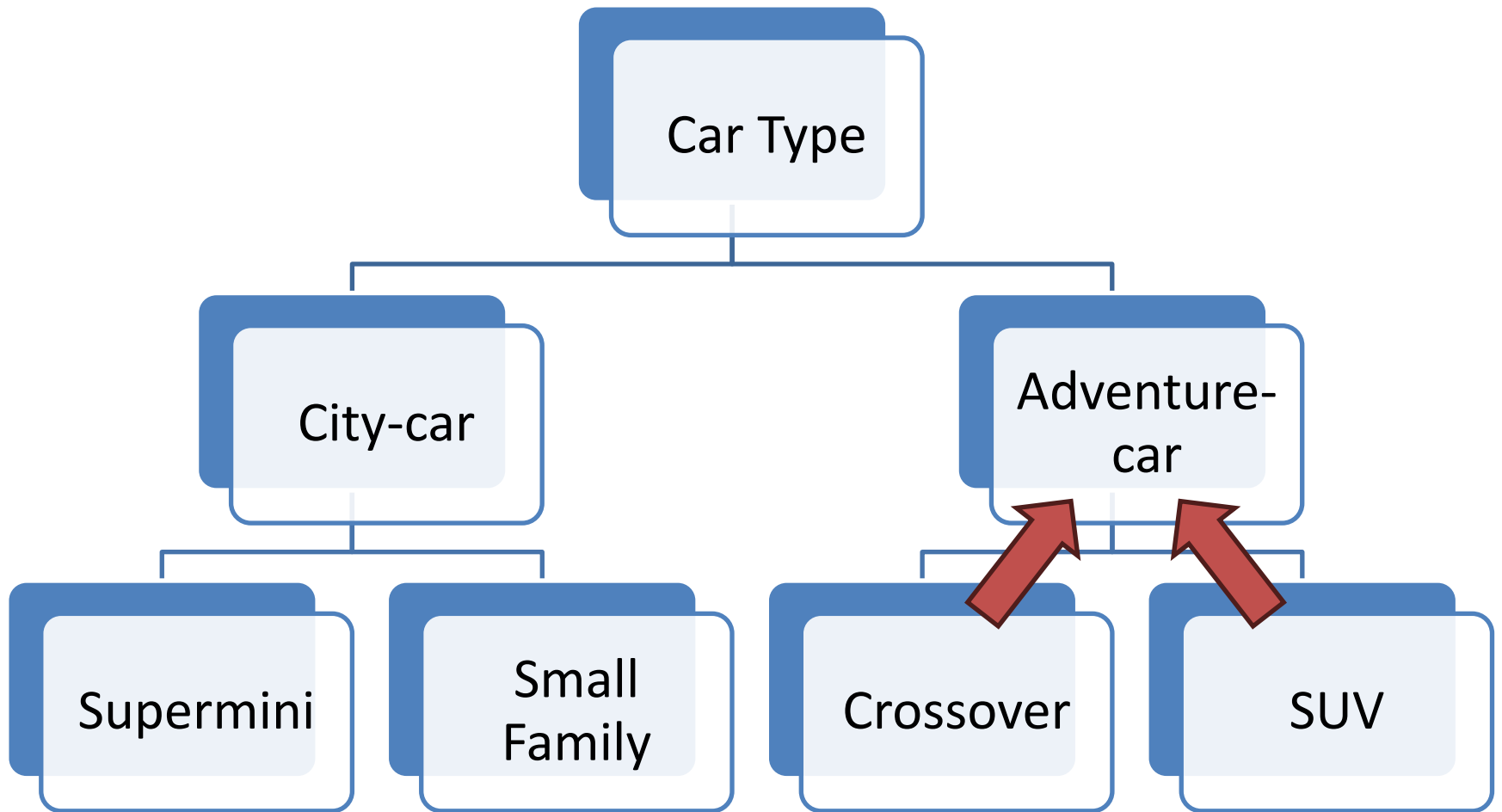
- Example

– $d(\text{"Γραβιέρα"}, \text{"Γίγαντες"}) = \frac{3}{3}$ ($lca = \text{"ΠΡΟΙΟΝΤΑ"}$)

d(Crossover, SUV)=?



$$d(\text{Crossover}, \text{SUV}) = \frac{1}{2}$$



COMBINING EVIDENCE

Combining similarities from difference processes/sources

- Assume we have two separate processes for computing similarities between users
- Process 1: assesses demographic data from the user database (gender, age, marital status, etc.)
 - Reports similarity score s_1 based on demographic data
- Process 2: considers their interaction with our systems (e.g. purchases, logins, etc.)
 - Reports similarity score s_2 based on user activity

Taking weighted averages

- We already saw this computation

$$\text{sim} = \frac{w_1 * s_1 + w_2 * s_2}{w_1 + w_2}$$

- Can be fine-tuned to our preferences or trust on these datasets
 - E.g. if we believe that activity data is more reliable or important, use $w_2 > w_1$

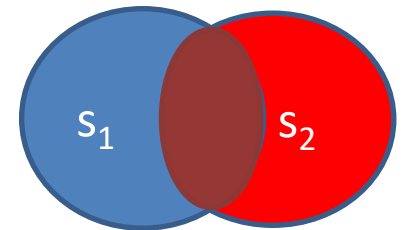
Treating scores as evidence

- One problem with averaging is that low scores from one of the two processes (e.g. due to wrong/missing data) will lower the overall calculation
 - Example: $s_1 = 0.7$, $s_2 = 0.2$, $\text{Average}(s_1, s_2) = 0.45$
- Possible solution: take maximum score
 - Take: $\text{Max}(s_1, s_2) = 0.70$
- Another idea is to treat each score as independent evidence that each boosts our confidence on the similarity between the users

Treating similarities as probabilities

- Assuming independence, combine scores in a probabilistic manner

$$\text{sim}(s_1, s_2) = s_1 + s_2 - s_1 * s_2$$



- In our example
 - $\text{sim}(0.7, 0.2) = 0.7 + 0.2 - 0.14 = 0.76$
- Notice that $\text{sim}(s_1, s_2) \geq \max(s_1, s_2)$

Using additional sources

- This calculation can be extended in case we have more sources suggesting similarity for the customers
 - e.g. based on customer surveys $s_3 = 0.8$
- Combine scores in a probabilistic manner

$$\text{sim}(s_1, s_2, s_3) = \text{sim}(s_1, s_2) + s_3 - \text{sim}(s_1, s_2) * s_3$$

- In our running example
 - $\text{sim}(0.7, 0.2, 0.8) = 0.76 + 0.8 - 0.76 * 0.8 = 0.952$
 - Compare with average $(0.7, 0.2, 0.8) = 0.56$
 - Compare with max $(0.7, 0.2, 0.8) = 0.8$
 - Compare with min $(0.7, 0.2, 0.8) = 0.2$

WORKING WITH SETS

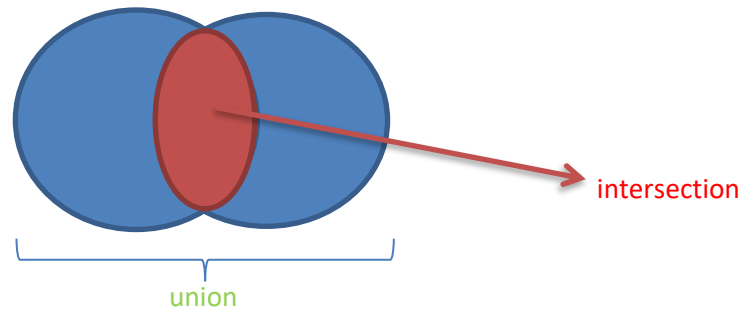
How do we compare sets?

- UserA= {milk, bread, coffee}
- UserB= {milk, bread, donut}
- UserC= {milk, bread, soda, potatoes}

- Straightforward idea: look at their **intersection**
 - Intersection(UserA,UserB) = {milk,bread}
 - Intersection(UserA,UserC) = {milk,bread}
- Intersection not enough!
 - Need to look at their **differences** too

Set similarity: Jaccard Index

- Jaccard(S_1, S_2) = the ratio of the sizes of the **intersection** and **union** of S_1 and S_2
 - Jaccard(S_1, S_2) = $|S_1 \cap S_2| / |S_1 \cup S_2|$



- Note that $|S_1 \cap S_2| \leq |S_1 \cup S_2|$
- Thus, $0 \leq \text{Jaccard}(S_1, S_2) \leq 1$

Jaccard Index Examples

- Recall: $\text{Jaccard}(S1, S2) = \frac{|S1 \cap S2|}{|S1 \cup S2|}$
- $\text{Jaccard}(\{\text{potatoes}, \text{lettuce}\}, \{\text{potatoes}, \text{tomatoes}\}) = \frac{1}{3}$
- $\text{Jaccard}(\{\text{potatoes}, \text{lettuce}, \text{cucumbers}\}, \{\text{potatoes}, \text{tomatoes}, \text{ketchup}\}) = \frac{1}{5}$
- $\text{Jaccard}(\{\text{potatoes}, \text{lettuce}\}, \{\text{potatoes}, \text{lettuce}, \text{tomatoes}\}) = \frac{2}{3}$
- $\text{Jaccard}(\{\text{lettuce}\}, \{\text{milk}, \text{soda}\}) = 0$
- $\text{Jaccard}(\{\text{soda}, \text{milk}\}, \{\text{milk}, \text{soda}\}) = 1$

Toy exercise

(python jupyter notebook in e-class)

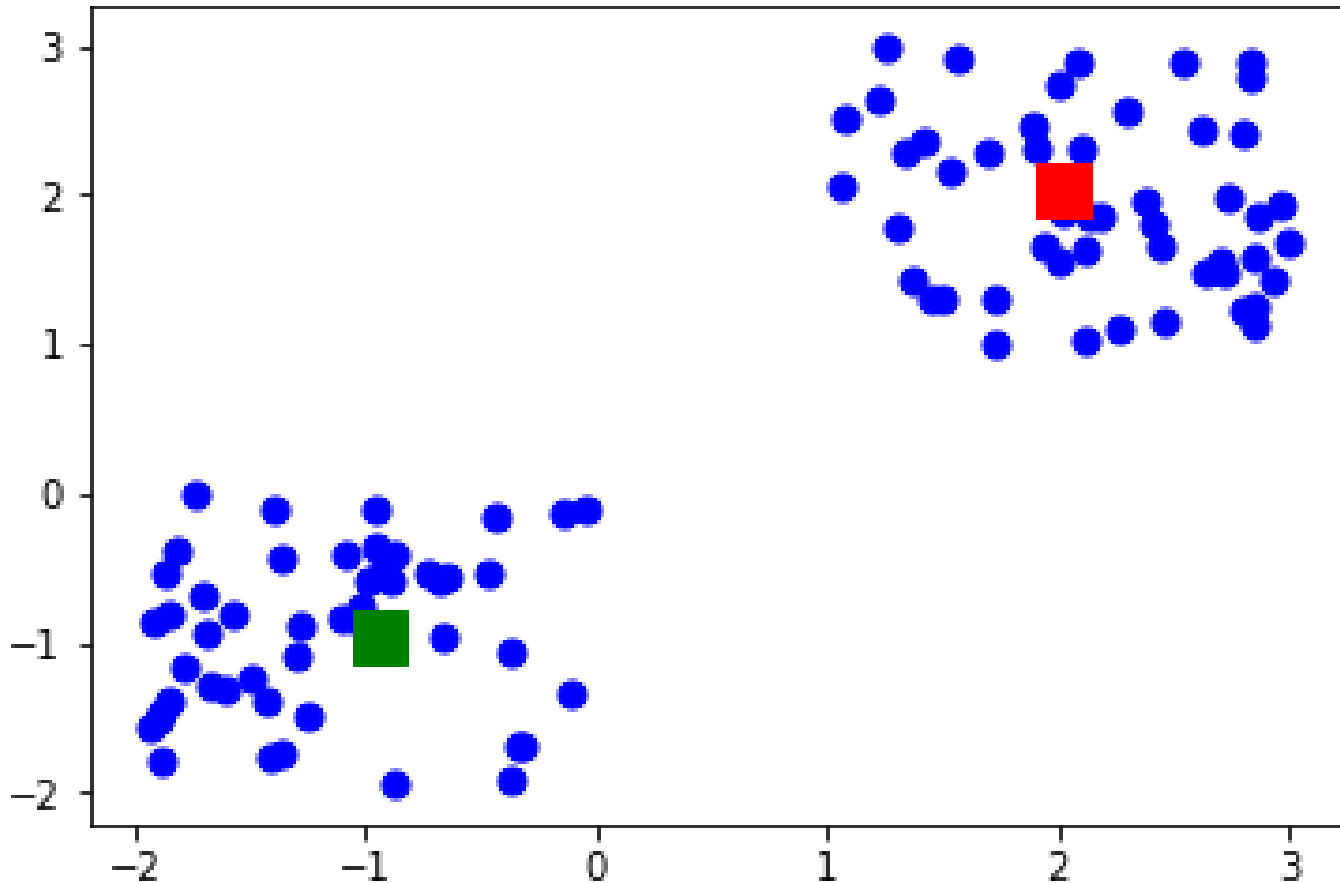
- Assume the following 5 customers with their purchases

```
user1 : ['milk', 'bread', 'coffee']  
user2 : ['milk', 'bread', 'cola']  
user3 : ['cereal', 'milk', 'donut']  
user4 : ['donut', 'cream', 'cola']  
user5 : ['cola', 'milk', 'cereal', 'tea']
```

- Can you group these customers into two clusters?

Clustering

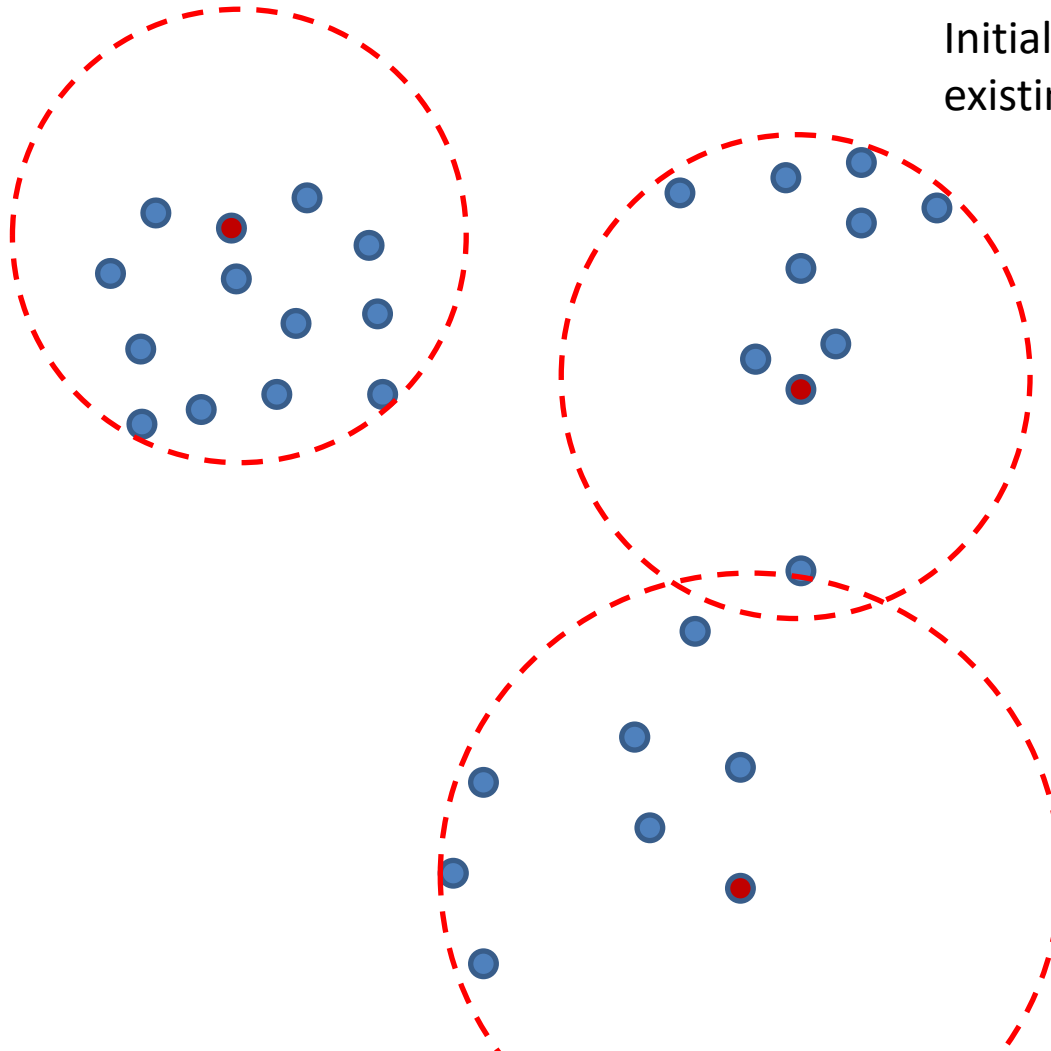
- ❑ Separate data into disjoint groups such that:
 - ❑ Increased similarity among members of the same group (cluster cohesion)
 - ❑ Members of different groups are dissimilar



The famous k-Means algorithm

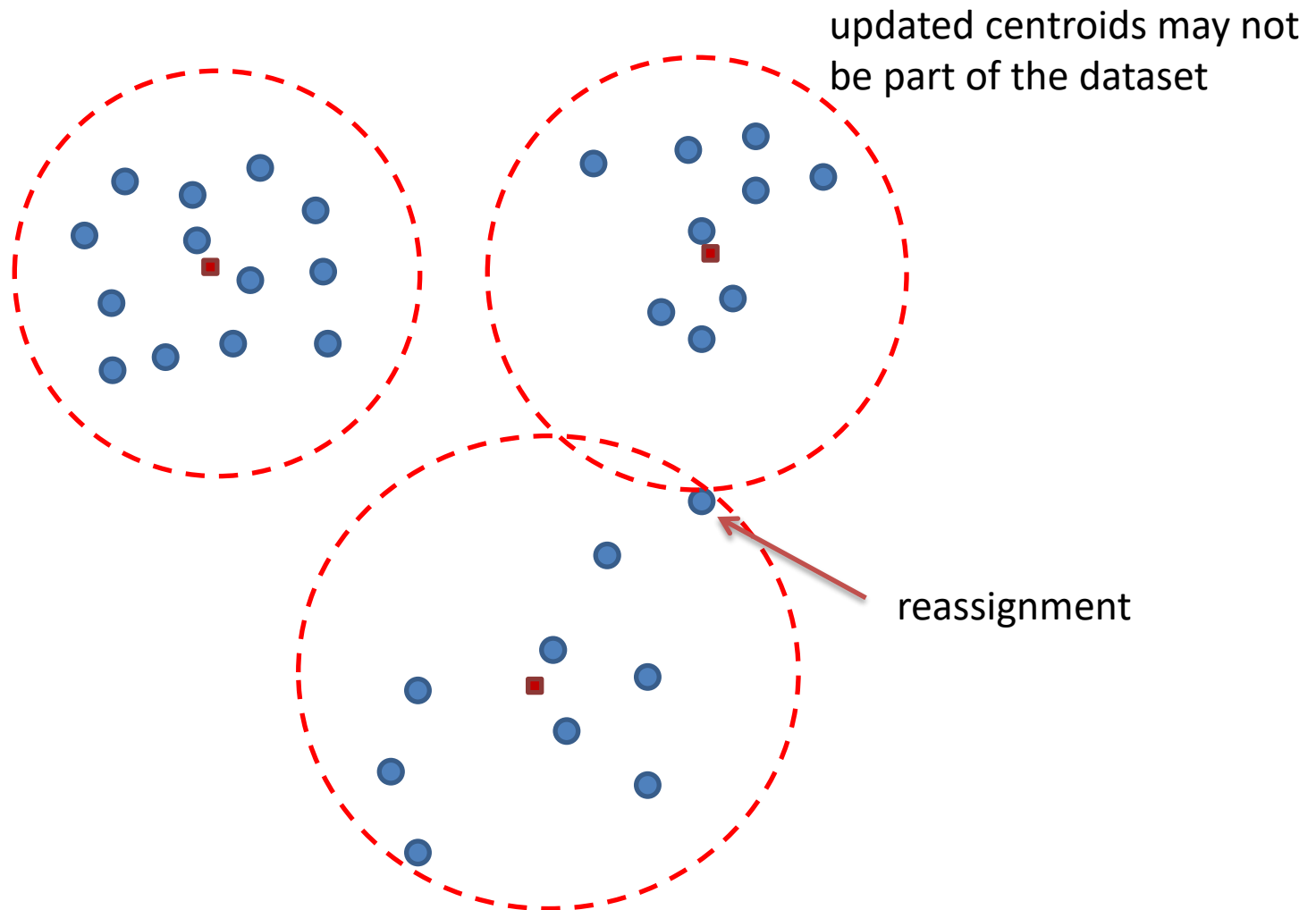
- Assume n points in the Euclidian space and a user-defined value of $k = \# \text{clusters}$
 1. Pick k points (centroids), one per cluster
 2. Assign remaining points to closest centroid
 3. In each cluster update location of its centroid
 4. Reassign points, if necessary
 5. Repeat steps 3-4 until clusters stabilize
- k-Means seeks to minimize the sum of squared distances (thus the variance of the distances) from the centroids
 - the algorithm always converges to some (local) minimum solution

Example for $k=3$



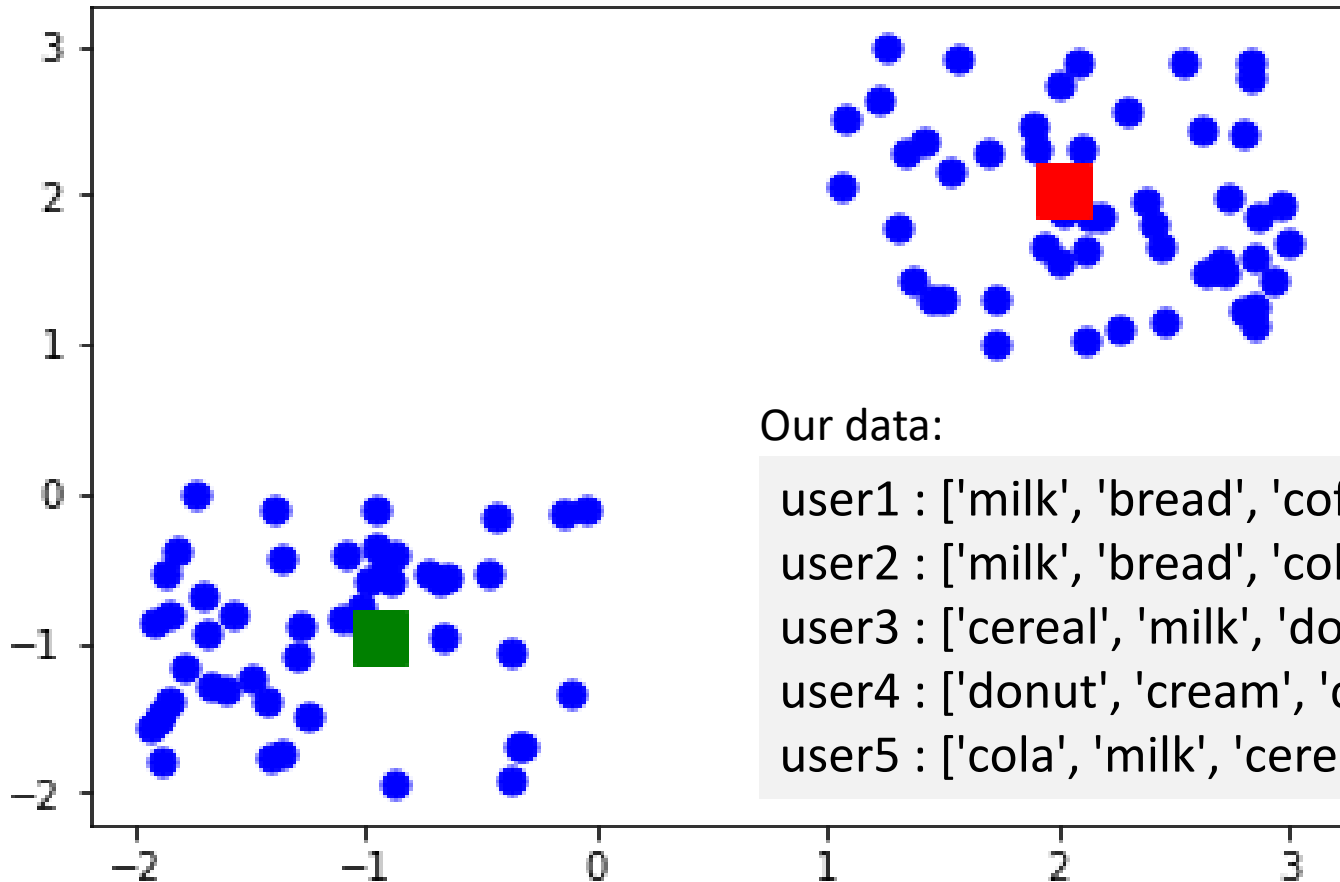
Initial centroids are
existing dataset points

New centroids + reassignment

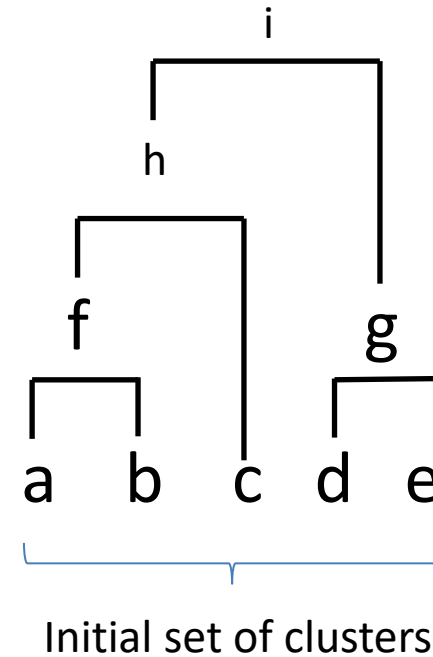
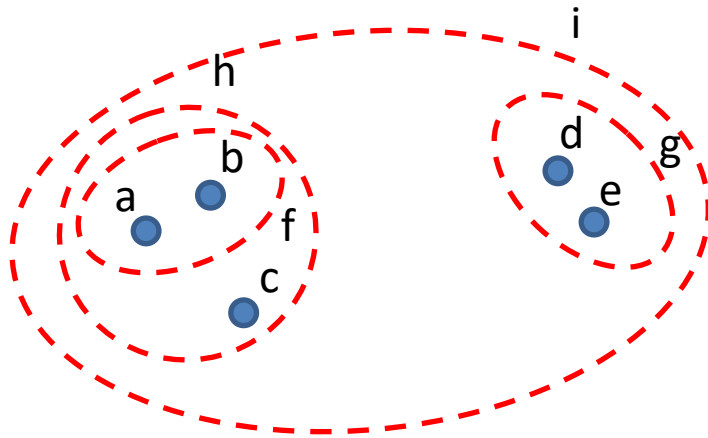


Wait!

- Our dataset is not points in a Euclidian space
 - There is no obvious way to compute a “centroid”



Hierarchical Clustering to the rescue



Executive decision

- Purchases are modelled as **sets of items**
 - Use **Jaccard** for computing customer pair-wise similarity

user1 : ['milk', 'bread', 'coffee']

user2 : ['milk', 'bread', 'cola']

user3 : ['cereal', 'milk', 'donut']

user4 : ['donut', 'cream', 'cola']

user5 : ['cola', 'milk', 'cereal', 'tea']

Jaccard_sim = $2/4 = 50\%$

Jaccard_sim = $1/6 = 16\%$

Jaccard Similarity

- All-pair similarity computation

user1 : ['milk', 'bread', 'coffee']
user2 : ['milk', 'bread', 'cola']
user3 : ['cereal', 'milk', 'donut']
user4 : ['donut', 'cream', 'cola']
user5 : ['cola', 'milk', 'cereal', 'tea']



Jaccard_sim of user1 , user2 is 0.5
Jaccard_sim of user1 , user3 is 0.2
Jaccard_sim of user1 , user4 is 0.0
Jaccard_sim of user1 , user5 is 0.16
Jaccard_sim of user2 , user3 is 0.2
Jaccard_sim of user2 , user4 is 0.2
Jaccard_sim of user2 , user5 is 0.4
Jaccard_sim of user3 , user4 is 0.2
Jaccard_sim of user3 , user5 is 0.4
Jaccard_sim of user4 , user5 is 0.16

Hierarchical Clustering

- Merge most similar pair to form a new cluster

user1 : ['milk', 'bread', 'coffee']
user2 : ['milk', 'bread', 'cola']
user3 : ['cereal', 'milk', 'donut']
user4 : ['donut', 'cream', 'cola']
user5 : ['cola', 'milk', 'cereal', 'tea']



Jaccard_sim of user1 , user2 is 0.5
Jaccard_sim of user1 , user3 is 0.2
Jaccard_sim of user1 , user4 is 0.0
Jaccard_sim of user1 , user5 is 0.16
Jaccard_sim of user2 , user3 is 0.2
Jaccard_sim of user2 , user4 is 0.2
Jaccard_sim of user2 , user5 is 0.4
Jaccard_sim of user3 , user4 is 0.2
Jaccard_sim of user3 , user5 is 0.4
Jaccard_sim of user4 , user5 is 0.16

New state

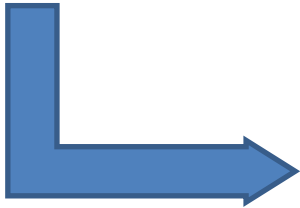
- Merge best pair (user1+user2) to form a new cluster
 - Represent cluster of customers as their **union** (not ideal, other options exist)

```
user1 : ['milk', 'bread', 'coffee']  
user2 : ['milk', 'bread', 'cola']  
user3 : ['cereal', 'milk', 'donut']  
user4 : ['donut', 'cream', 'cola']  
user5 : ['cola', 'milk', 'cereal', 'tea']  
user1+user2 : ['bread', 'cola', 'milk', 'coffee']
```

Next step

(most similar pair: user3, user5)

user3 : ['cereal', 'milk', 'donut']
user4 : ['donut', 'cream', 'cola']
user5 : ['cola', 'milk', 'cereal', 'tea']
user1+user2 : ['bread', 'cola', 'milk', 'coffee']



user4 : ['donut', 'cream', 'cola']
user1+user2 : ['bread', 'cola', 'milk', 'coffee']
user3+user5 : ['cereal', 'donut', 'milk', 'cola', 'tea']

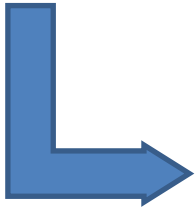
Final step

(most similar pair: user4, user3+user5)

user4 : ['donut', 'cream', 'cola']

user1+user2 : ['bread', 'cola', 'milk', 'coffee']

user3+user5 : ['cereal', 'donut', 'milk', 'cola', 'tea']



user1+user2 : ['bread', 'cola', 'milk', 'coffee']

user4+user3+user5 : {'donut', 'cereal', 'milk', 'cream', 'cola', 'tea'}

Cluster 1

user1 : ['milk', 'bread', 'coffee']

user2 : ['milk', 'bread', 'cola']

Cluster 2

user3 : ['cereal', 'milk', 'donut']

user4 : ['donut', 'cream', 'cola']

user5 : ['cola', 'milk', 'cereal']

Notes

- In this toy example we performed **Hierarchical Clustering** up to 2 clusters without checking the quality of the intermediate clusters
 - Sometimes it is better to stop sooner than later
- To simplify the code, we used as a representative (clustoid) of a cluster the UNION of its members
 - Can you think of examples where this is a bad choice?

Cluster 1

```
user1 : ['milk', 'bread', 'coffee']  
user2 : ['milk', 'bread', 'cola']
```

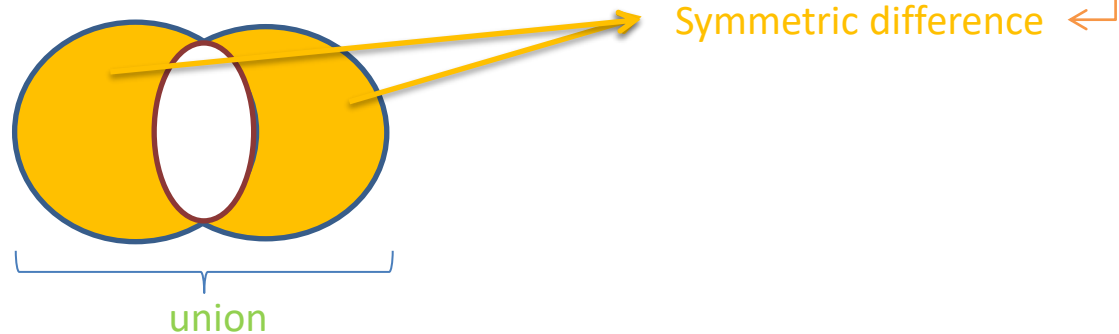
Cluster 2

```
user3 : ['cereal', 'milk', 'donut']  
user4 : ['donut', 'cream', 'cola']  
user5 : ['cola', 'milk', 'cereal']
```

Jaccard Distance between sets

- Can be defined as the complement of their Jaccard similarity

$$- d_{\text{jacc}}(S1, S2) = 1 - \frac{|S1 \cap S2|}{|S1 \cup S2|} = \frac{|S1 \cup S2| - |S1 \cap S2|}{|S1 \cup S2|}$$



How about **bags**?



Bags are “sets” with
repetition of elements
allowed

Jaccard can be extended to
work with bags

Intersection(S1,S2): count an
element n times in the
intersection, where n is the
minimum of the number of times
the element appears in S1 and S2

Union(S1,S2): count the element
the **sum** of the number of times
it appears in S1, S2

Example

- $S1 = \{a,a,a,b\}$, $S2 = \{a,a,b,b,c\}$
- Then, intersection is $\{a,a,b\}$ and union $\{a,a,a,a,a,b,b,b,c\}$
- Bag-similarity is thus, $3/9 = 1/3$
- Note, bag similarity is between 0 and $\frac{1}{2}$ (why?)

Alternative bag similarity

- Count an element n times in the intersection, where n is the **minimum** of the number of times the element appears in $S1$ and $S2$
- In the union, count the element the **max** of the number of times it appears in $S1$, $S2$

Example (alt)

- $S1 = \{a,a,a,b\}$, $S2 = \{a,a,b,b,c\}$
- Then, intersection is $\{a,a,b\}$ and union $\{a,a,a,b,b,c\}$
- Bag-similarity of $S1$, $S2$ is thus, $3/6 = 50\%$
- Note, alternative bag similarity is between 0 and 1 (why?)

Bag Similarity Example

- Movies ratings dataset
 - John: Star_Wars_I:3/5, Avatar: 4/5, Aliens: 2/5
 - Mary: Star_Wars_I: 2/5, Avatar: 5/5, ET: 4/5
 - Nick: Star_Wars_I: 4/5, Aliens: 2/5, ET: 1/5
- Who is the Nearest Neighbor of John?
- Note: if treated as sets
 - $\text{Jaccard}(\text{John}, \text{Mary}) = \text{Jaccard}(\text{John}, \text{Nick}) = 2/4 = 50\%$
 - Let us consider their bag similarity instead!

Bag Similarity Example

- Convert to bags:
 - John: {Star_Wars_I, Star_Wars_I, Star_Wars_I, Avatar, Avatar, Avatar, Aliens, Aliens}
 - Mary: {Star_Wars_I, Star_Wars_I, Avatar, Avatar, Avatar, Avatar, ET, ET, ET, ET}
 - Nick: {Star_Wars_I, Star_Wars_I, Star_Wars_I, Star_Wars_I, Aliens, Aliens, ET}
- $\text{Bag_similarity_alt}(\text{John}, \text{Mary}) = (2+4)/(3+5+4+2) = 6/14 = 42.9\%$
- $\text{Bag_similarity_alt}(\text{John}, \text{Nick}) = (3+2)/(4+4+2+1) = 5/11 = 45.5\%$

WORKING WITH VECTORS

Basket data example

- Three distinct products:
 - potato (p), lettuce (l),
tomato (t)
- Three users with the following purchases
 - John: 2 potatoes, 1 lettuce
 - Kostas: 1 tomato
 - Mary: 10 potatoes, 6
lettuces

Vector Model

$\langle \#p, \#l, \#t \rangle$

$$\vec{J} = \langle 2, 1, 0 \rangle$$

$$\vec{K} = \langle 0, 0, 1 \rangle$$

$$\vec{M} = \langle 10, 6, 0 \rangle$$



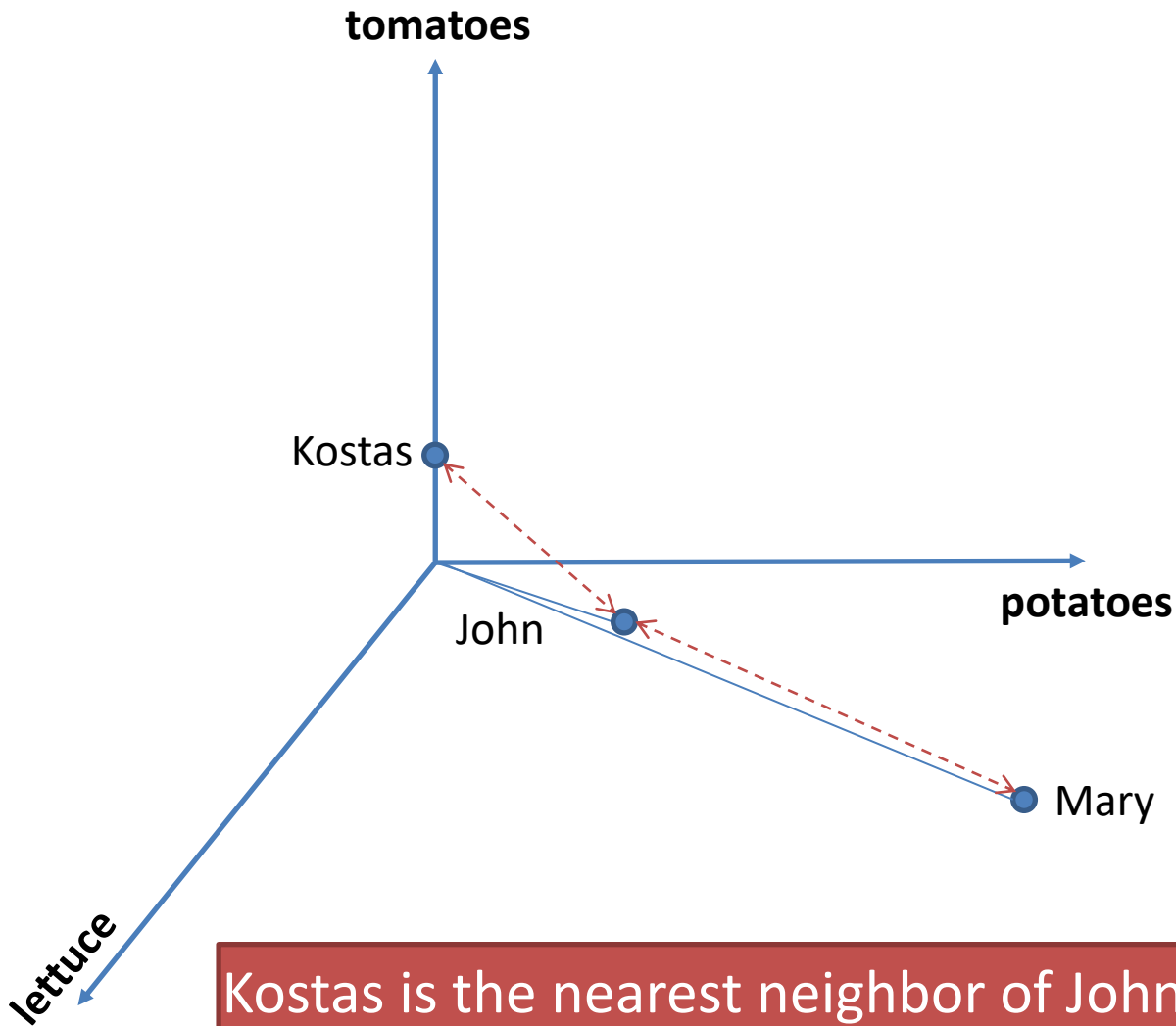
Definition of Euclidean Distance

- $\vec{x} = \langle 2, 1, 0, 5 \rangle$
- $\vec{y} = \langle 5, 6, 1, 10 \rangle$

- Recall that:

$$\begin{aligned}d(\vec{x}, \vec{y}) &= \sqrt{(2 - 5)^2 + (1 - 6)^2 + (0 - 1)^2 + (5 - 10)^2} \\ &= \sqrt{9 + 25 + 1 + 25} = \sqrt{60} = 7.75\end{aligned}$$

Euclidean Distance NN Calculations



Vector Model

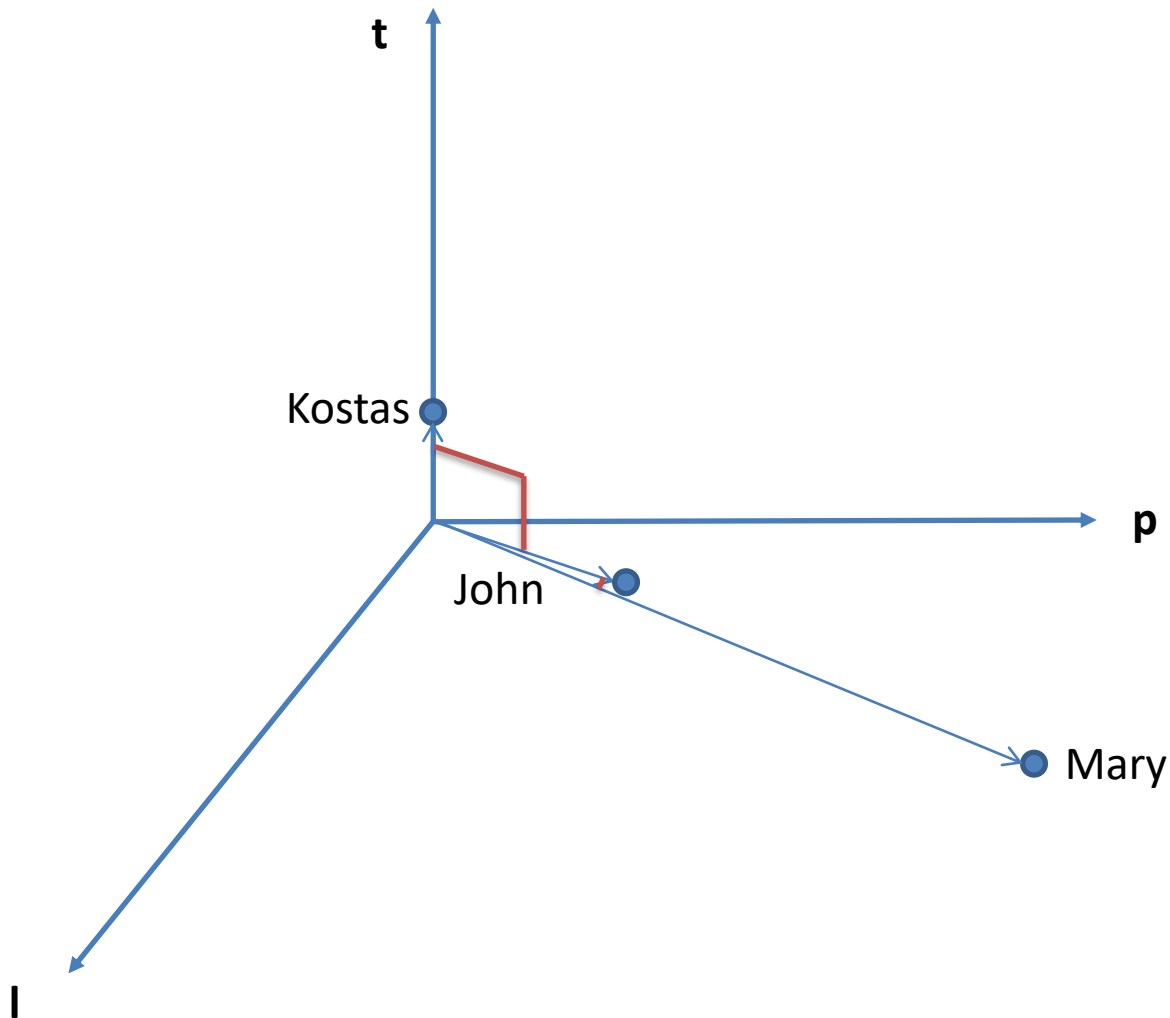
$$\vec{J} = \langle 2, 1, 0 \rangle$$

$$\vec{K} = \langle 0, 0, 1 \rangle$$

$$\vec{M} = \langle 10, 6, 0 \rangle$$

Kostas is the nearest neighbor of John!!!

Angle Calculations: favor direction over length (norm)



Vector Model

$$\vec{J} = \langle 2, 1, 0 \rangle$$

$$\vec{K} = \langle 0, 0, 1 \rangle$$

$$\vec{M} = \langle 10, 6, 0 \rangle$$

$$\theta(\vec{J}, \vec{K}) = 90^\circ$$

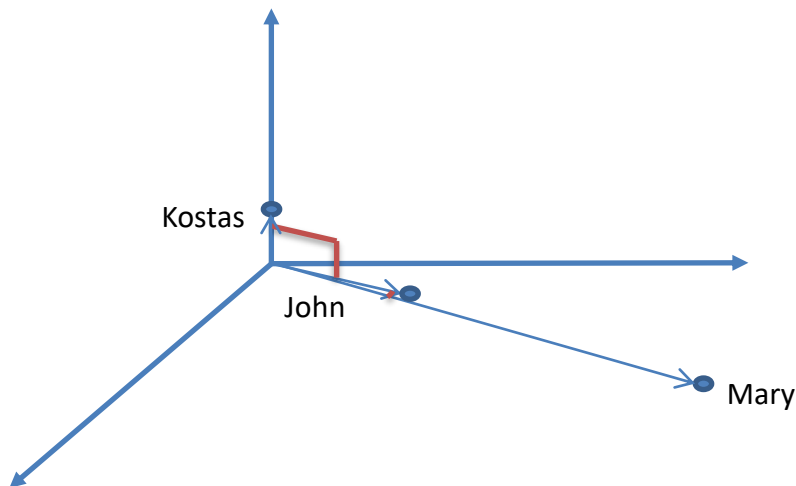
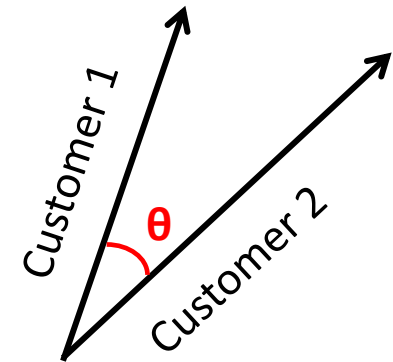
$$\theta(\vec{J}, \vec{M}) = 4.3^\circ$$

When to use Cosine?



Cosine Similarity

- $\text{sim}(\vec{x}, \vec{y}) = \cos(\theta(\vec{x}, \vec{y})) \in [-1..+1]$
 - Used in collaborative filtering
 - Popular in document matching



$$\cos(\theta(\vec{J}, \vec{K})) = \cos(90^\circ) = 0$$

$$\cos(\theta(\vec{J}, \vec{M})) = \cos(4.3^\circ) = 0.997$$

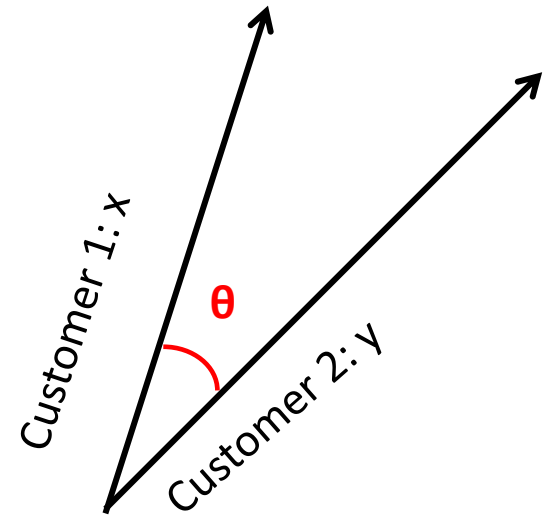
Dot (inner) product between two vectors

- $\vec{x} \cdot \vec{y} = \sum (x_k * y_k)$
- Example:

$$\vec{x} = (1, 3, 0, 5)$$
$$\vec{y} = (1, 0, 1, 6)$$

- Then:

$$\vec{x} \cdot \vec{y} = 1 * 1 + 3 * 0 + 0 * 1 + 5 * 6 = 31$$
$$= |\vec{x}| * |\vec{y}| * \cos(\theta(\vec{x}, \vec{y}))$$



From dot to cosine

- $\cos(\theta(\vec{x}, \vec{y})) = \frac{\vec{x} \cdot \vec{y}}{|\vec{x}| * |\vec{y}|}$ $\vec{x} = (1, 3, 0, 5)$
 $\vec{y} = (1, 0, 1, 6)$
- In this example
- $|\vec{x}| = \sqrt{1^2 + 3^2 + 0^2 + 5^2} = \sqrt{35}$
- $|\vec{y}| = \sqrt{1^2 + 0^2 + 1^2 + 6^2} = \sqrt{38}$
- $\cos(\theta(\vec{x}, \vec{y})) = \frac{31}{\sqrt{35} * \sqrt{38}} = 0.85$

Dot product with **unit** vector

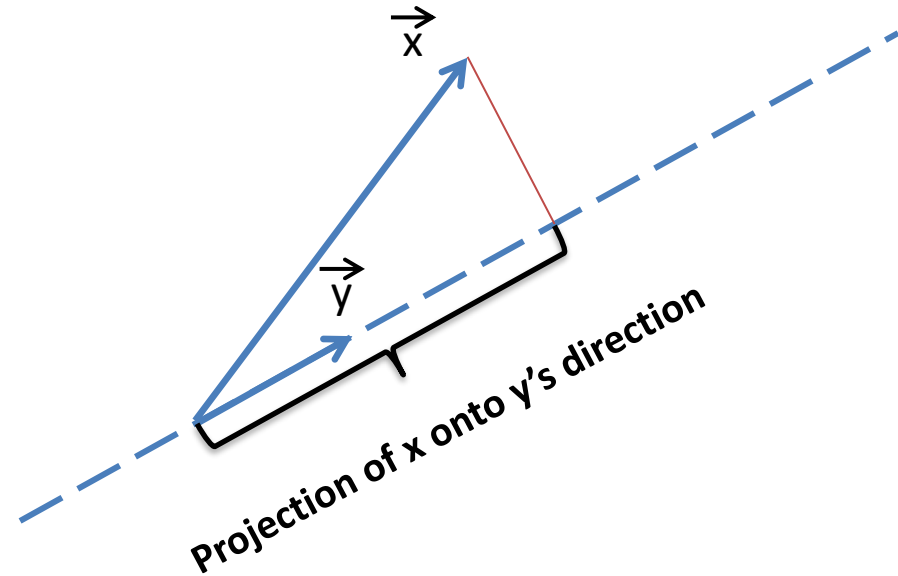
- $\vec{x} \cdot \vec{y} = \sum (x_k * y_k)$
- Example for unit vector \vec{y} :

$$\vec{x} = (1, 3, 0, 5)$$

$$\vec{y} = (1/2, 1/2, 1/2, 1/2)$$

- Notice that $|\vec{y}| = 1$
- Then:

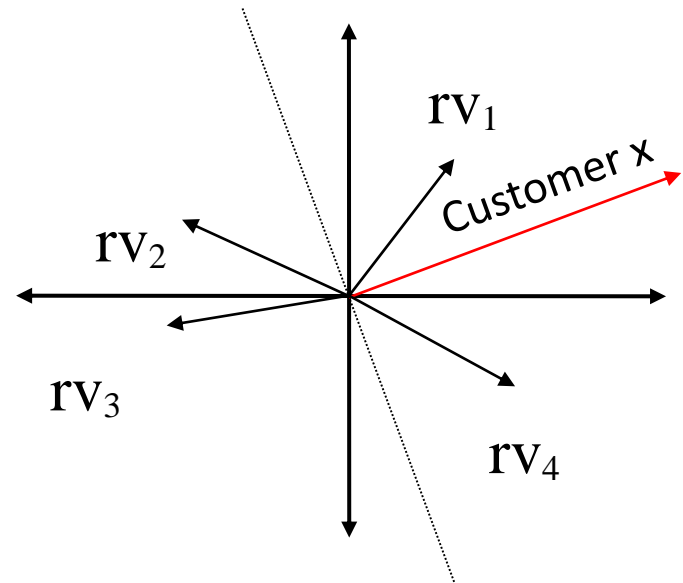
$$\begin{aligned}\vec{x} \cdot \vec{y} &= 1/2 + 3/2 + 5/2 = 9/2 = 4.5 \\ &= |\vec{x}| * 1 * \cos(\theta(\vec{x}, \vec{y}))\end{aligned}$$



Random Hyperplane Projection

(Mining Massive Data Sets, Sec. 3.7.2)

- uses n d-dimensional random vectors (rv_i)
- Generates for each input vector a bitmap of size n as follows:
 - Sets $bit_i=1$ if dot product of input vector with i -th random vector is **positive**
 - Sets $bit_i=0$ if dot product of input vector with i -th random vector is **negative**

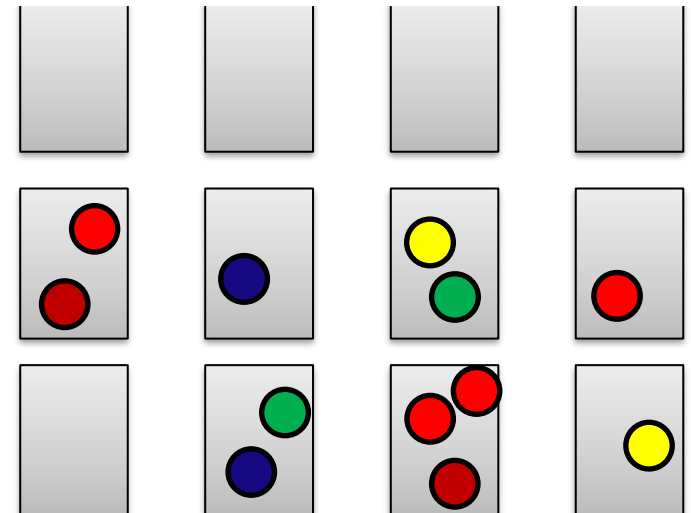
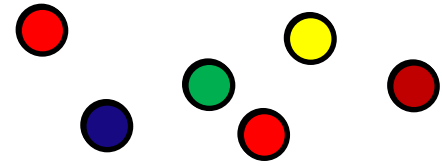


RHP(x)

1	0	0	1
---	---	---	---

Locality Sensitive Hashing (LSH)

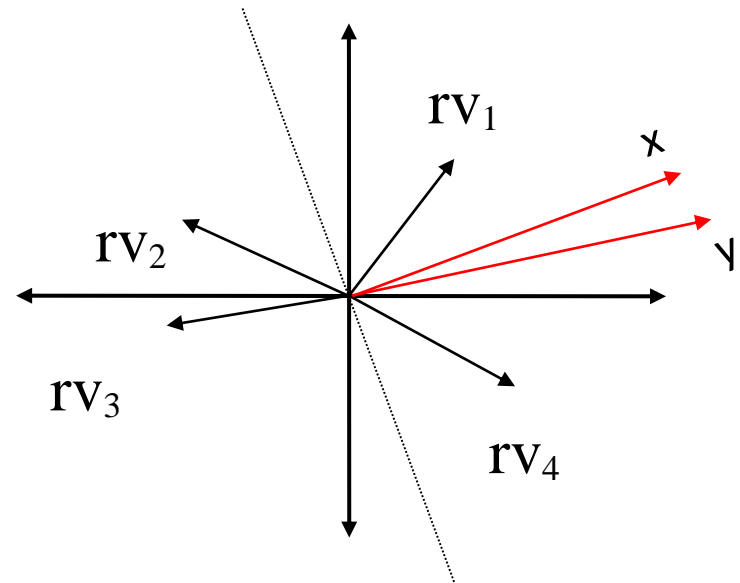
- Assign items to **buckets** using a hash function $h(x)$
 - E.g. $h(\bullet)$ =0110 in binary
 - Details of function $h()$ depend on the preferred similarity metric:
 - Similar objects are hashed to the same bucket with high probability
 - Dissimilar objects are hashed to the same bucket with very small probability
- Repeat several times



Buckets (1-4)

Is RHP a locality-sensitive hashing scheme?

- Assume vectors for customers x and y point (approximately) towards the same direction
 - This means their cosine is close to 1
- We expect that with high probability the RHP values will be identical
- Use RHP encodings as “bucket ids”
 - Similar customers are hashed to the same bucket (with high probability)



$$\text{RHP}(x) = \text{RHP}(y)$$

1	0	0	1
---	---	---	---

Hamming Distance

- The Hamming distance between **two equal-length strings of symbols** is the number of positions at which the corresponding symbols are different (Wikipedia)
 - $D_h('00110101', '10110110') = 3$
 - $D_h('abc', 'acc') = 1$

Approximate Similarity Computation via Hamming Distance of RHP bitmaps

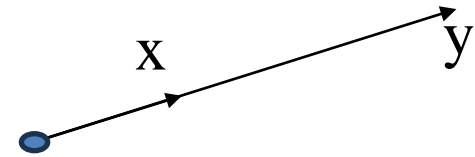
$RHP_{(x)}$:

0	0	0	1	1	1
---	---	---	---	---	---

← n bits →

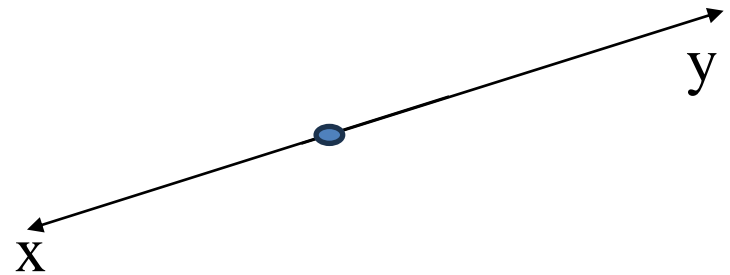
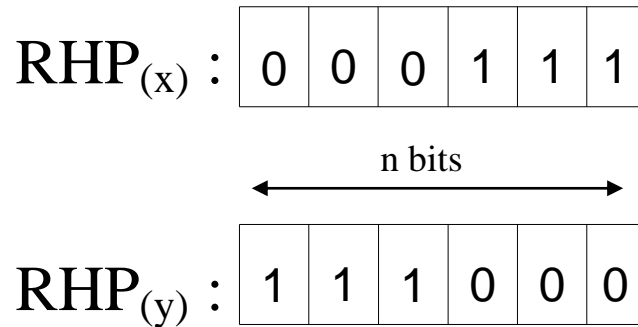
$RHP_{(y)}$:

0	0	0	1	1	1
---	---	---	---	---	---



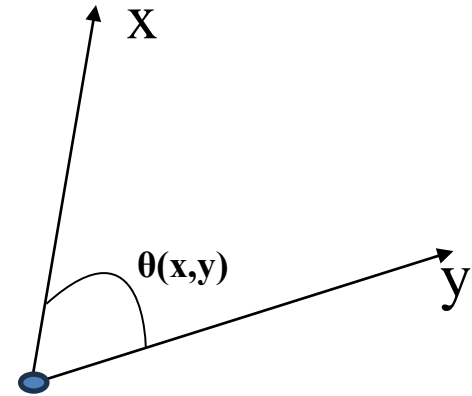
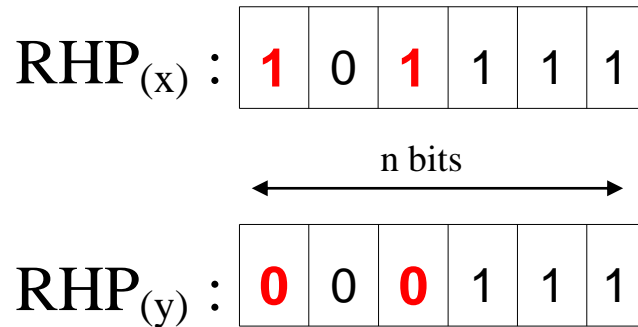
- Vectors are collinear ($\theta(\vec{x}, \vec{y})=0$, cosine similarity = 1)

Approximate Similarity Computation via Hamming Distance of RHP bitmaps



- Vectors are opposite ($\theta(\vec{x}, \vec{y}) = \pi$, cosine similarity = 0)

Approximate Similarity Computation via Hamming Distance of RHP bitmaps

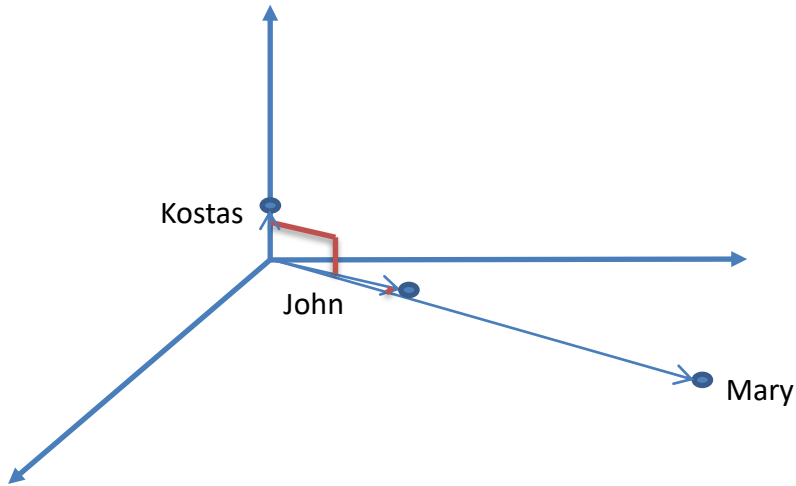


$$\hat{\theta}(\text{RHP}_{(x)}, \text{RHP}_{(y)}) = 2/6 * \pi = \pi/3$$

Estimate $\hat{\theta}(x,y) = D_h(\text{RHP}_{(x)}, \text{RHP}_{(y)}) * \pi/n$

- Also works for the *Pearson correlation*
 - $Cor(x, y) = \text{Cos}(x-\bar{x}, y-\bar{y})$

RHP Example



$$J = \langle 2, 1, 0 \rangle$$

$$K = \langle 0, 0, 1 \rangle$$

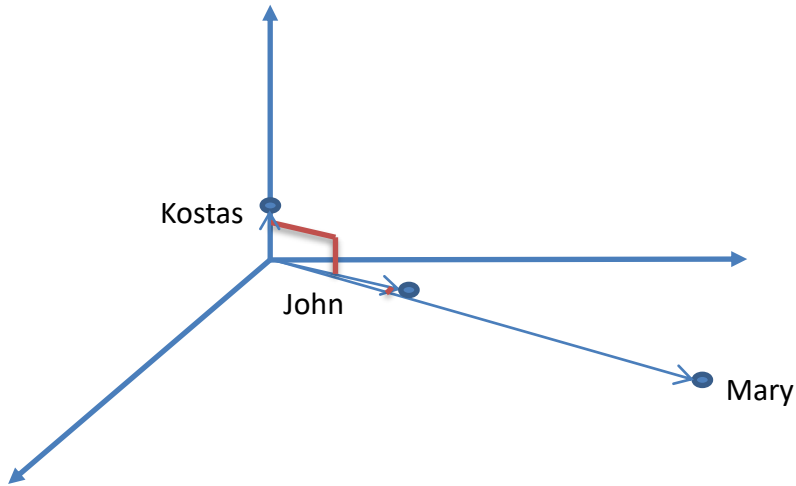
$$M = \langle 10, 6, 0 \rangle$$

$$rv1 = \langle 3, 1, 6 \rangle$$

$$rv2 = \langle -5, 3, 2 \rangle$$

- Calculations for John:
- $\langle 2, 1, 0 \rangle \cdot \langle 3, 1, 6 \rangle = 2 \cdot 3 + 1 \cdot 1 + 0 \cdot 6 = +7 \rightarrow \text{bit} = 1$
- $\langle 2, 1, 0 \rangle \cdot \langle -5, 3, 2 \rangle = -10 + 3 = -7 \rightarrow \text{bit} = 0$
- Thus, **RHP(John) = 10**

RHP Example



$$J = \langle 2, 1, 0 \rangle$$

$$K = \langle 0, 0, 1 \rangle$$

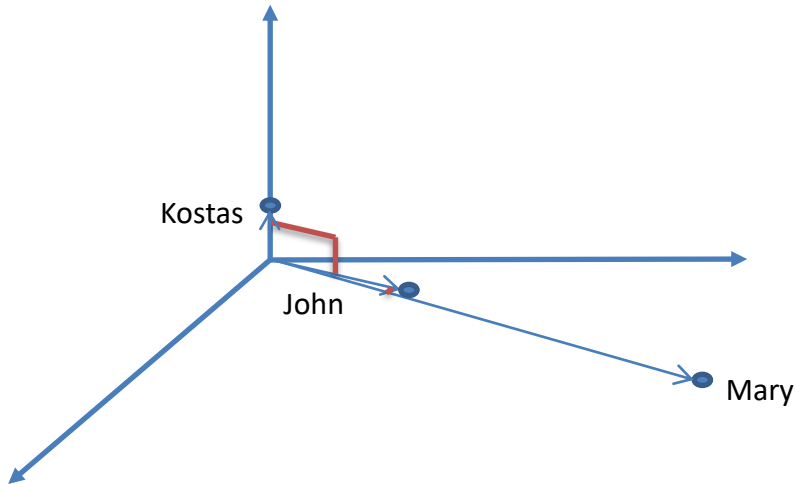
$$M = \langle 10, 6, 0 \rangle$$

$$rv1 = \langle 3, 1, 6 \rangle$$

$$rv2 = \langle -5, 3, 2 \rangle$$

- Calculations for Kostas:
- $\langle 0, 0, 1 \rangle \cdot \langle 3, 1, 6 \rangle = +6 \rightarrow \text{bit} = 1$
- $\langle 0, 0, 1 \rangle \cdot \langle -5, 3, 2 \rangle = +2 \rightarrow \text{bit} = 1$
- Thus, **RHP(Kostas) = 11**

RHP Example



$$J = \langle 2, 1, 0 \rangle$$

$$K = \langle 0, 0, 1 \rangle$$

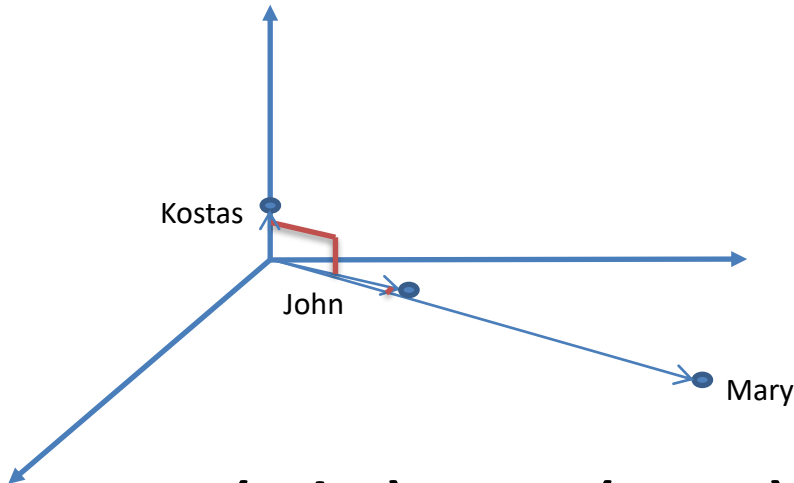
$$M = \langle 10, 6, 0 \rangle$$

$$rv1 = \langle 3, 1, 6 \rangle$$

$$rv2 = \langle -5, 3, 2 \rangle$$

- Calculations for Mary:
- $\langle 10, 6, 0 \rangle \cdot \langle 3, 1, 6 \rangle = +36 \rightarrow \text{bit} = 1$
- $\langle 10, 6, 0 \rangle \cdot \langle -5, 3, 2 \rangle = -32 \rightarrow \text{bit} = 0$
- Thus, **RHP(Mary) = 10**

RHP Example



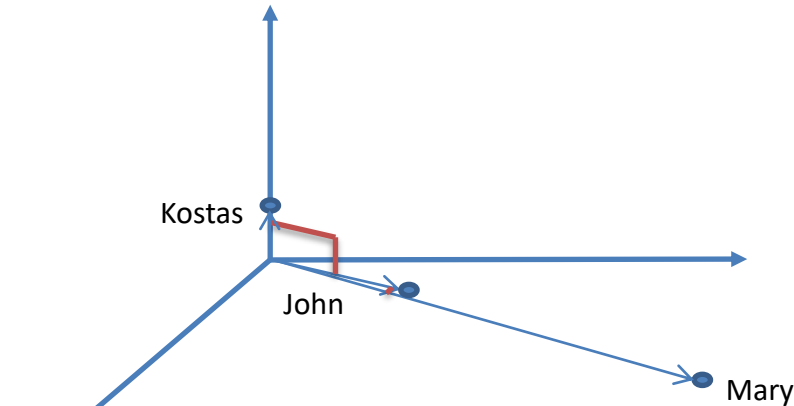
$$J = \langle 2, 1, 0 \rangle$$

$$K = \langle 0, 0, 1 \rangle$$

$$M = \langle 10, 6, 0 \rangle$$

- **$RHP(\text{John}) = RHP(\text{Mary}) = 10$**
- Hamming distance = 0
- Estimated angle is 0
 - Thus, estimated cosine similarity = 1
 - True cosine = 0.997
- Good accuracy by using just two bits!
 - Disclaimer: I am cherry picking favorable examples here

RHP Example



$$J = \langle 2, 1, 0 \rangle$$

$$K = \langle 0, 0, 1 \rangle$$

$$M = \langle 10, 6, 0 \rangle$$

- **RHP(John) = 10, RHP(Kostas) = 11**
- Hamming distance = 1 (out of $n = 2$ bits)
- Estimated angle is $\pi/2 = 90^\circ$
 - Thus, estimated cosine similarity = $\cos(\pi/2) = 0$
 - This is also the true cosine similarity

WORKING WITH STRINGS

String distance computations

- Why it is useful
 - String Matching
 - Spelling Checking
- Examples
 - Fix data entry errors: replace “Yiannis” with “Yannis”
 - Address matching/correction
 - Compare “Patisson” , “Patisation Str”, “Patisation St”
 - Fraud Detection
 - Are “Kotidis123”, “Kotidis554” and “7Kotidis123” the same user?

String Edit Distance

- The edit distance between strings $x = x_1x_2..x_n$ and $y = y_1y_2..y_m$ is the smallest number of **insertions** and **deletions** of single characters that will convert x to y
- As an example to convert $x = \text{"abcde"}$ to $y = \text{"acfdeg"}$
 1. **delete** b and get "a**bcde**"
 2. **insert** f after c and get "ac**fde**"
 3. **insert** g after e and get "acfde**g**" = y
- Thus, $d_{\text{edit}}(\text{"abcde"}, \text{"acfdeg"}) = 3$

Longest Common Subsequence (LCS)

- The LCS of x and y is the longest common string that is constructed by deleting positions from x and y
 - For $x = \text{“}\underline{a}\underline{b}\underline{c}\underline{d}\underline{e}\text{”}$ to $y = \text{“}\underline{a}\underline{c}\underline{f}\underline{d}\underline{e}\underline{g}\text{”}$
 - $\text{LCS}(x,y) = \text{“}acde\text{”}$
- It holds that
 - $d_{\text{edit}}(x,y) = \text{len}(x) + \text{len}(y) - 2 * \text{len}(\text{LCS}(x,y))$
- In our example $d_{\text{edit}}(x,y) = 5 + 6 - 2 * 4 = 3$

Levenshtein Distance

- In addition to **insertions** and **deletions** of single characters, **Levenshtein distance** also allows **substitutions**
- As an example, for $x = \text{"STALL"}$ and $y = \text{"TABLE"}$,
 $d_{\text{lev}}(x, y) = 3$
 1. (starting with "STALL") delete S and get "TALL",
 2. substitute first L with B and get "TABL",
 3. insert E at the end and get "TABLE"
- In comparison $d_{\text{edit}}(\text{"STALL"}, \text{"TABLE"},) = 4$
 - Notice that 1 substitution \Leftrightarrow 1 deletion + 1 insertion

Note

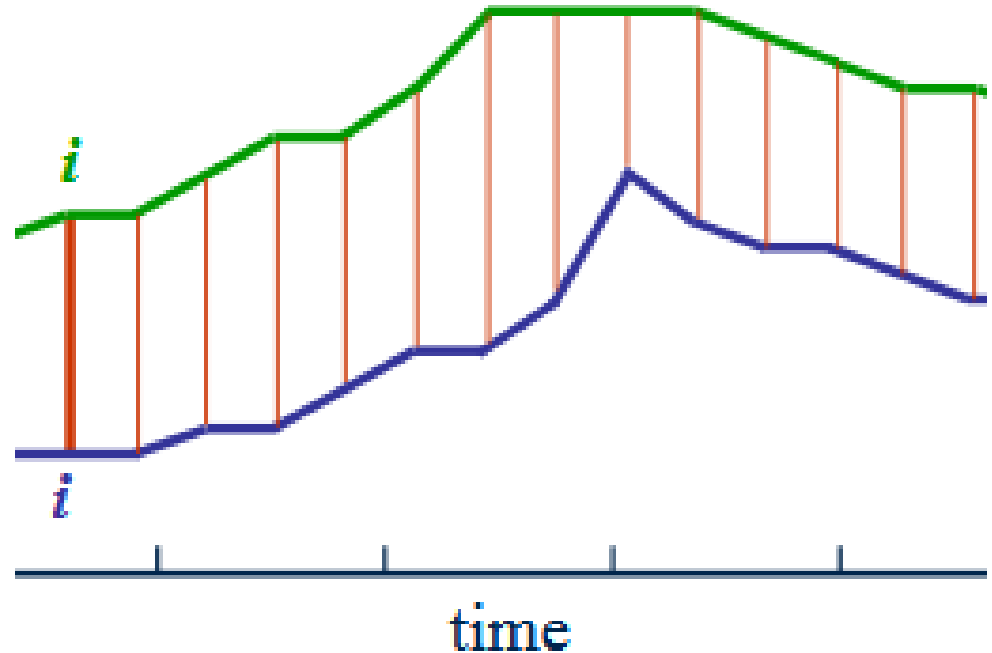
- In the literature sometimes Levenshtein distance is referred as edit distance (e.g. edit distance adjusted to permit insertions, deletions as well as substitutions)

Additional Metrics for strings

- **Damerau–Levenshtein** distance further allows **transportation** between two successive characters
 - Corssroads → Crossroads
- **Jaro** distance only allows transportations

Time Series

- Sequence of data points indexed in time order
 - Examples: financial data, sensor data, speech, etc
 - Univariate (running examples) vs multivariate
- Can be compared with Euclidean distance (given two series of same length)
 - The i^{th} point on one time series is aligned with the i^{th} point on the other
- However, this often gives poor results
- Does not work if series have difference lengths
 - Padding?



Time Series - Euclidean Distance

- $S_x = \langle 2, 1, 0, 1 \rangle$
- $S_y = \langle 2, 0, 2, 3 \rangle$

$$\begin{aligned} \text{Euclidean-distance } d(S_x, S_y) &= \\ &= \sqrt{(2 - 2)^2 + (1 - 0)^2 + (0 - 2)^2 + (1 - 3)^2} \\ &= \sqrt{0^2 + 1^2 + 2^2 + 2^2} \\ &= \sqrt{0 + 1 + 4 + 4} \\ &= \sqrt{9} = 3 \end{aligned}$$

Time Series - Euclidean Distance

- $S_x = \langle 2, 1, 0, 1 \rangle$
- $S_y = \langle 2, 0, 2, 3, 1, 2, 2, 0, 4 \rangle$

What now?

Time Series - Euclidean Distance

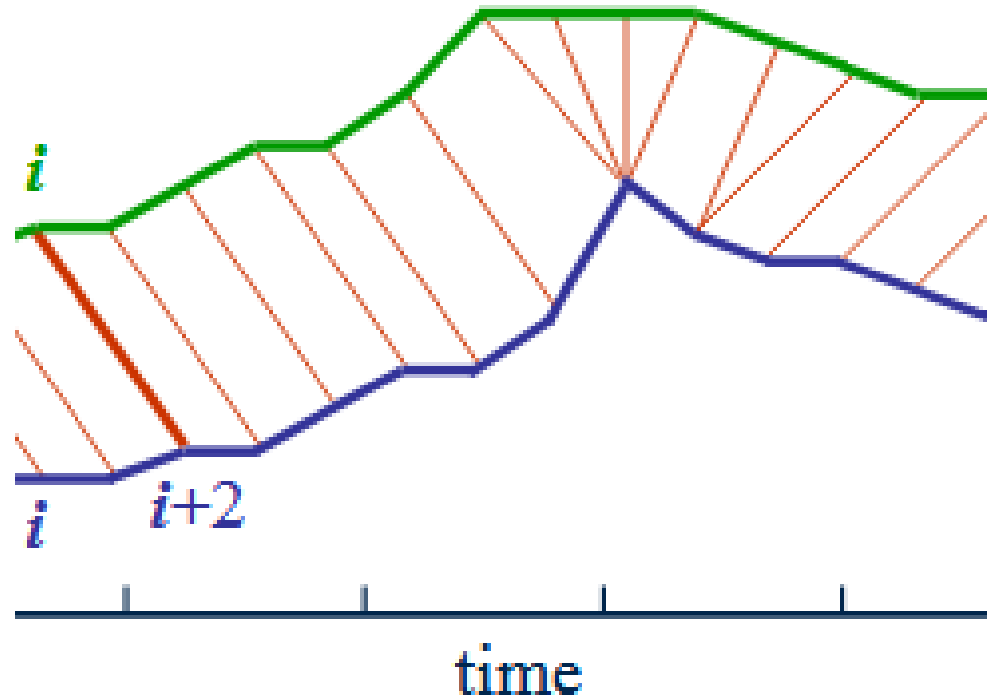
- $S_x = \langle 2, 1, 0, 1 \rangle$
- $S_y = \langle 2, 0, 2, 3, 1, 2, 2, 0, 4 \rangle$

Padding (convert to same length)?

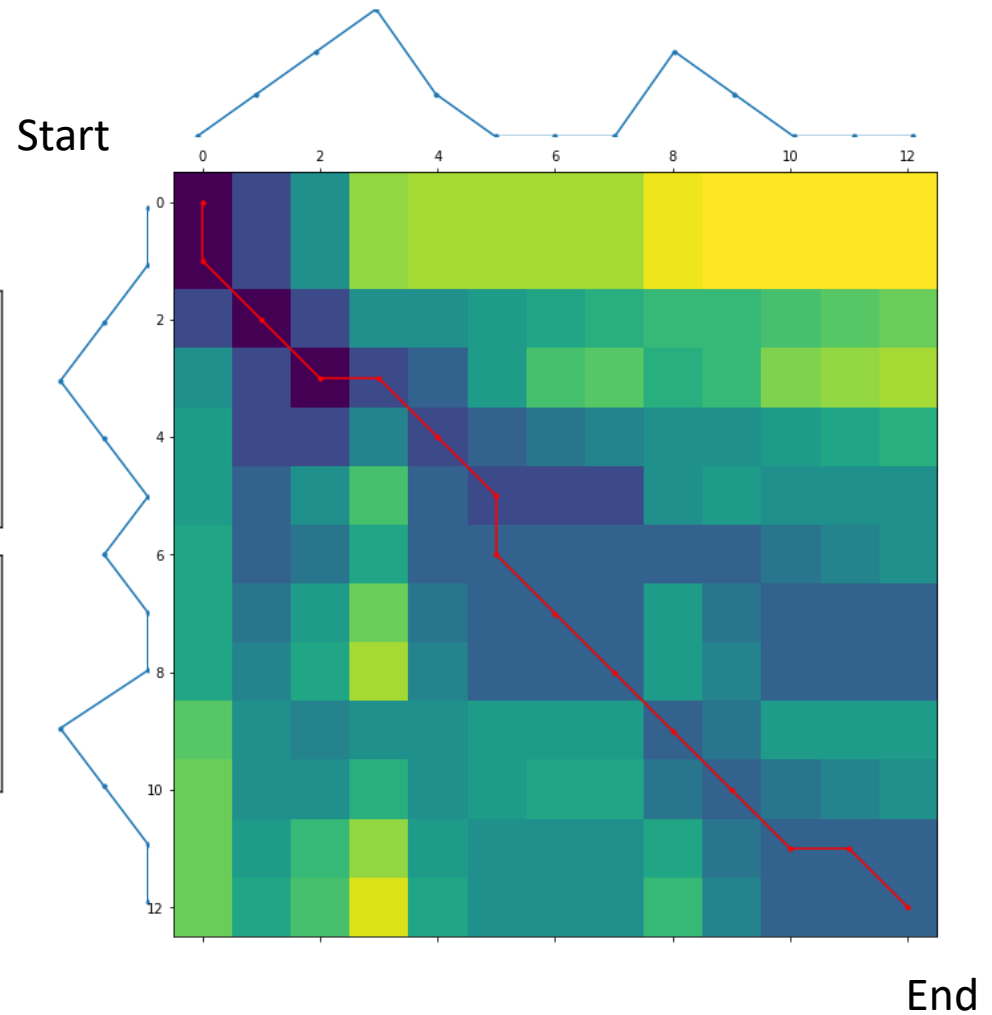
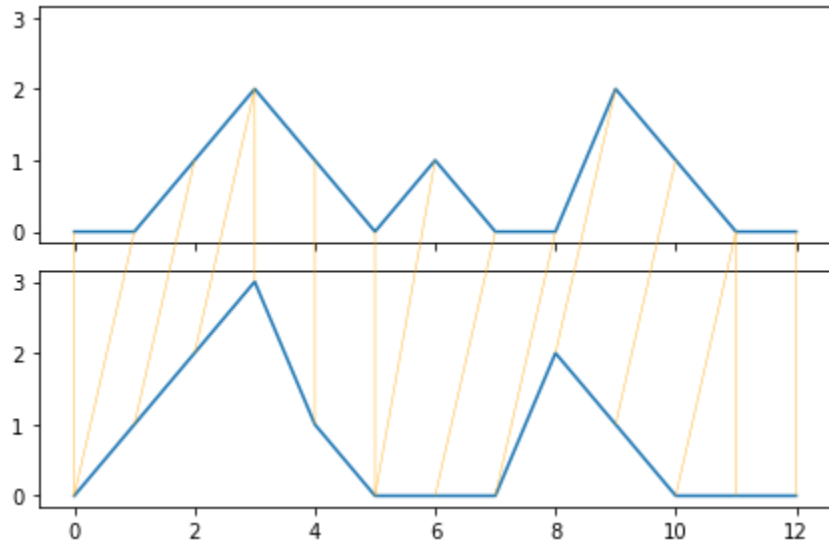
- $S_x = \langle 2, 1, 0, 1, 0, 0, 0, 0, 0 \rangle$
- $S_y = \langle 2, 0, 2, 3, 1, 2, 2, 0, 4 \rangle$

Dynamic Time Wrapping

- **DTW** computes the best **alignment** between the two-time series
 - Works even if the input series have different lengths
 - Useful if series have different frequencies or are out of phase (e.g. lag)
- Has been shown to be superior than Euclidean distance for tasks such as time series classification
- Drawback: quadratic complexity $O(n^2)$



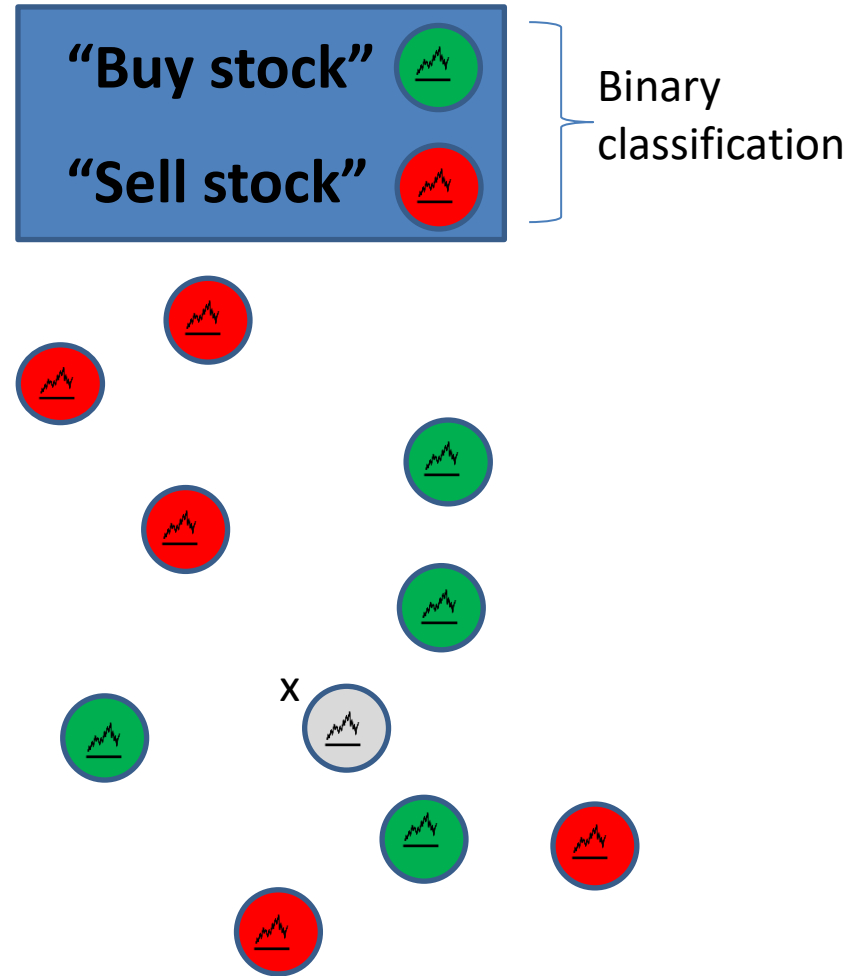
Computation complexity: $O(n*m)$



Example: Time Series Classification

Problem Statement

- Given:
 - n time series x_1, \dots, x_n along with their labels (classes) y_1, \dots, y_n to be used as **training** examples
 - a time series x with an **unknown** label
- Goal:
 - classify x : find the class label of x



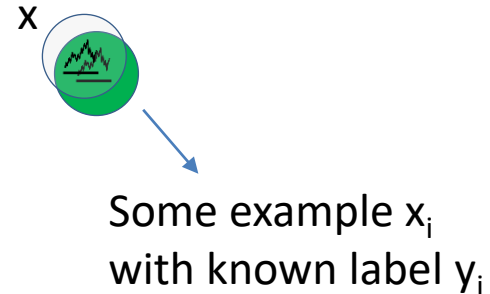
Intuition

In a perfect world:

- Assume there is another data point (time series) x_i that is very similar to the input series x
- I would then pick the label y_i of x_i as my selection
- This decision is **optimal** if $x = x_i$ or, equivalently when $d(x, x_i) \rightarrow 0$

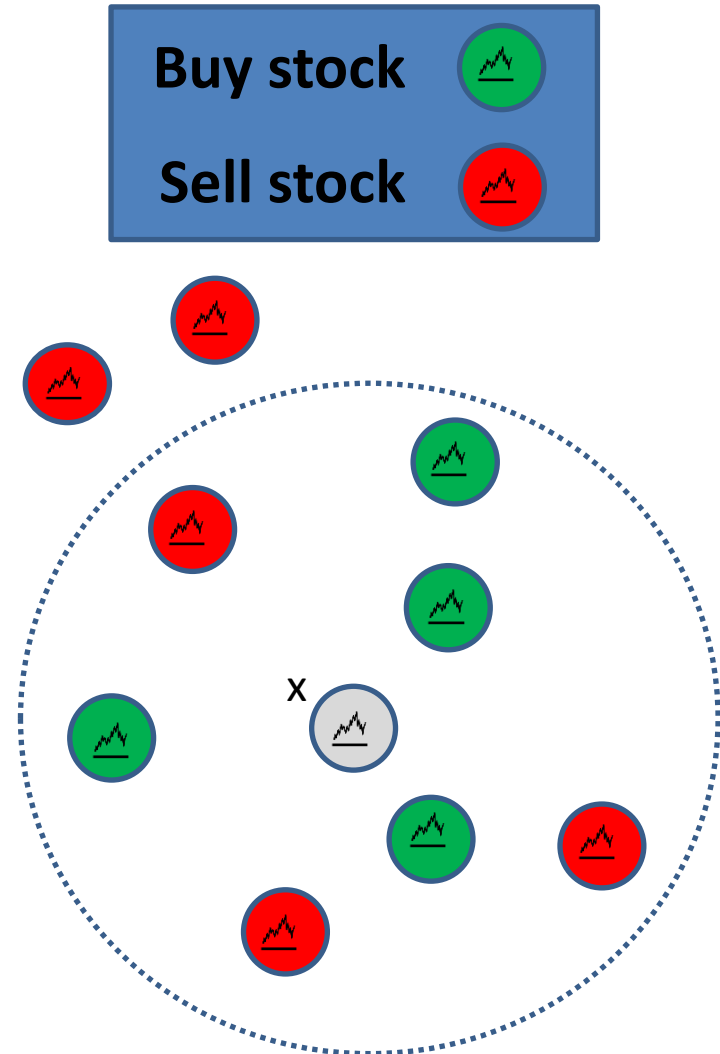
In practice:

- We will look at labeled data from the **neighborhood** of x_i



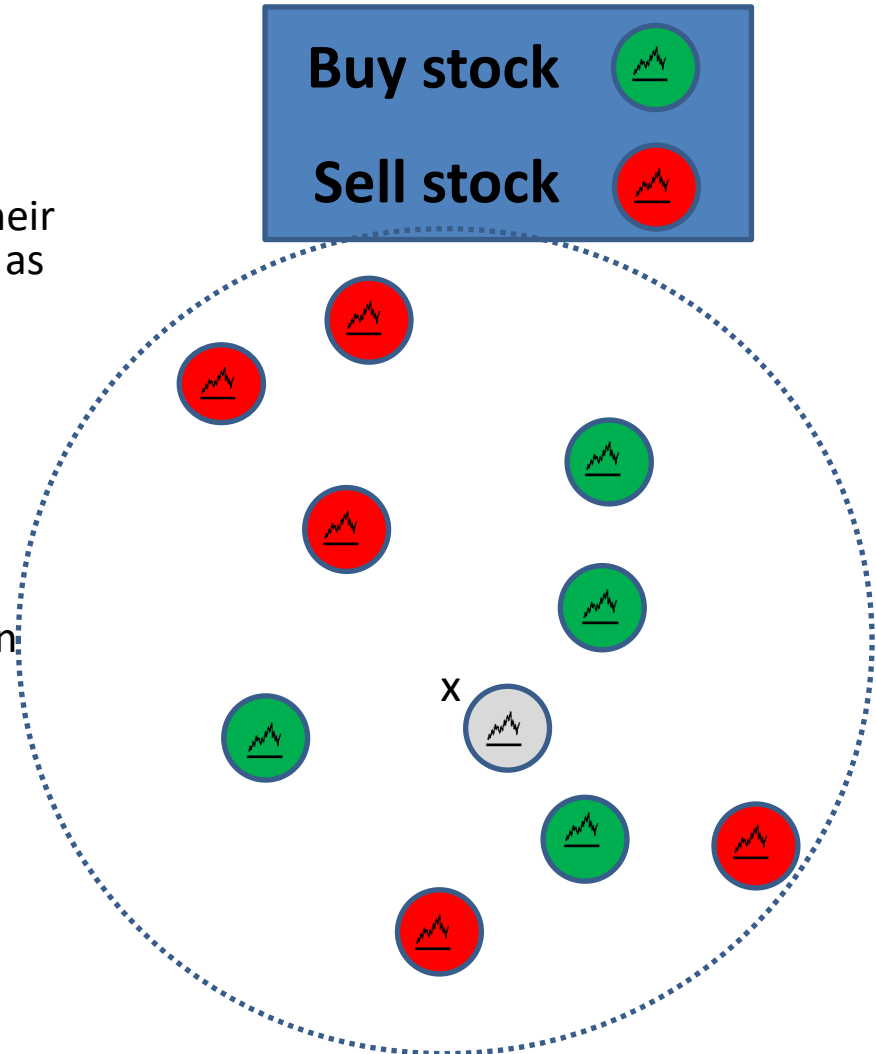
k-NN algorithm

- Given:
 - n time series x_1, \dots, x_n along with their labels (classes) y_1, \dots, y_n to be used as **training** examples
 - a time series x with an unknown label
- Goal:
 - classify x : find the class label of x
- *Intuition:*
 - assign x to the class most common among its k nearest neighbours
- Considerations:
 - selection of k
 - weigh neighbours



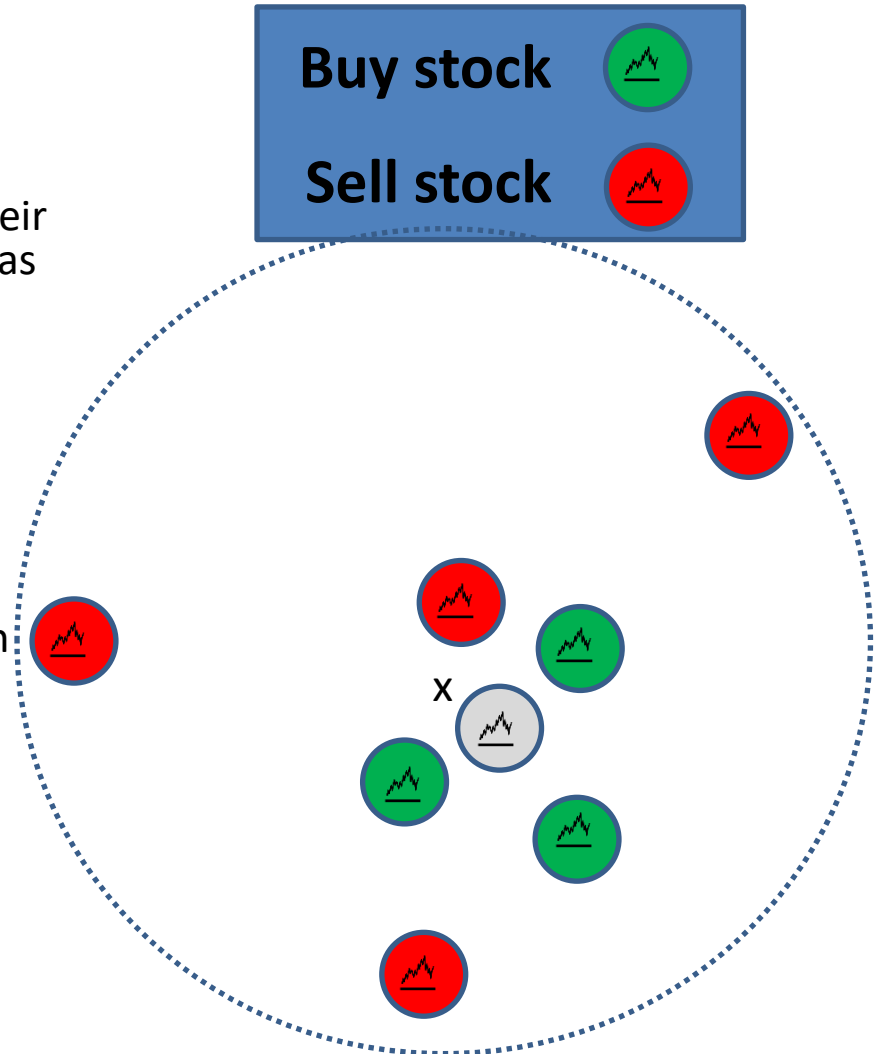
k-NN algorithm

- Given:
 - n time series x_1, \dots, x_n along with their labels (classes) y_1, \dots, y_n to be used as **training** examples
 - a time series x with an unknown label
- Goal:
 - classify x : find the class label of x
- *Intuition:*
 - assign x to the class most common among its k nearest neighbours
- Considerations:
 - selection of k
 - weigh neighbours

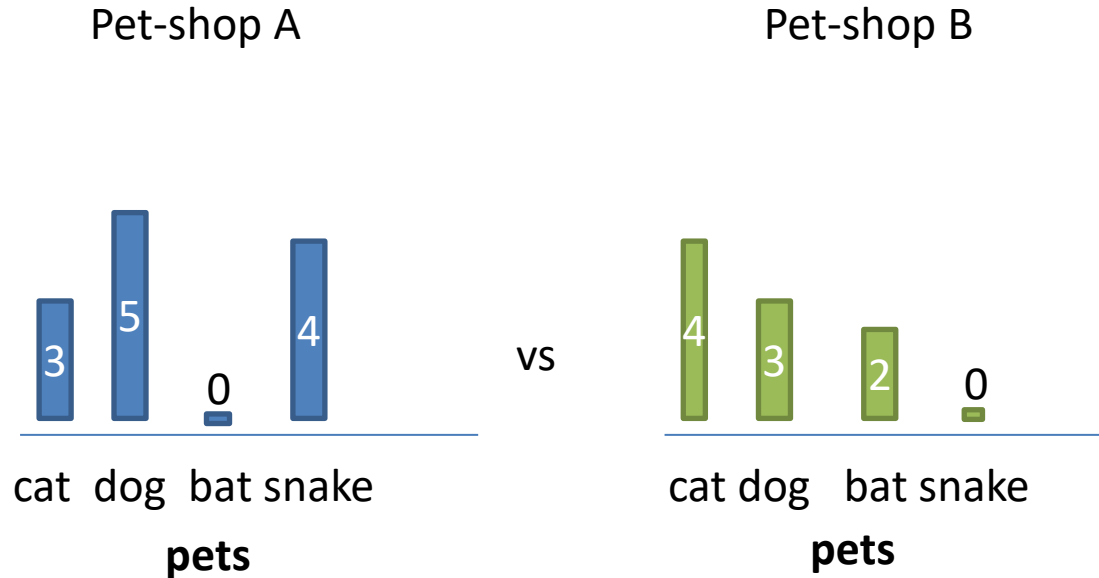


k-NN algorithm

- Given:
 - n time series x_1, \dots, x_n along with their labels (classes) y_1, \dots, y_n to be used as **training** examples
 - a time series x with an unknown label
- Goal:
 - classify x : find the class label of x
- *Intuition:*
 - assign x to the class most common among its k nearest neighbours
- Considerations:
 - selection of k
 - **weigh neighbours**



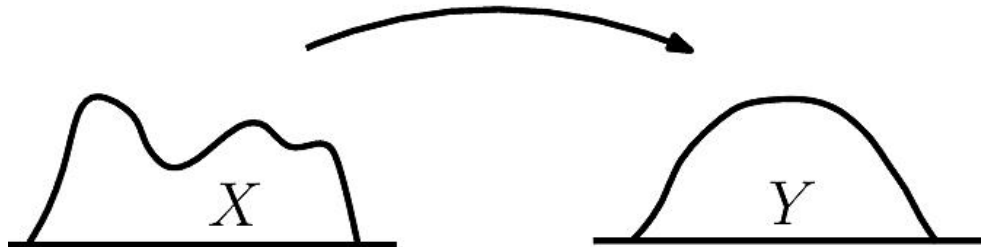
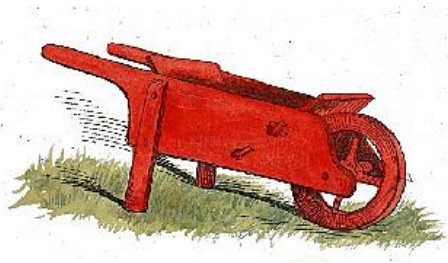
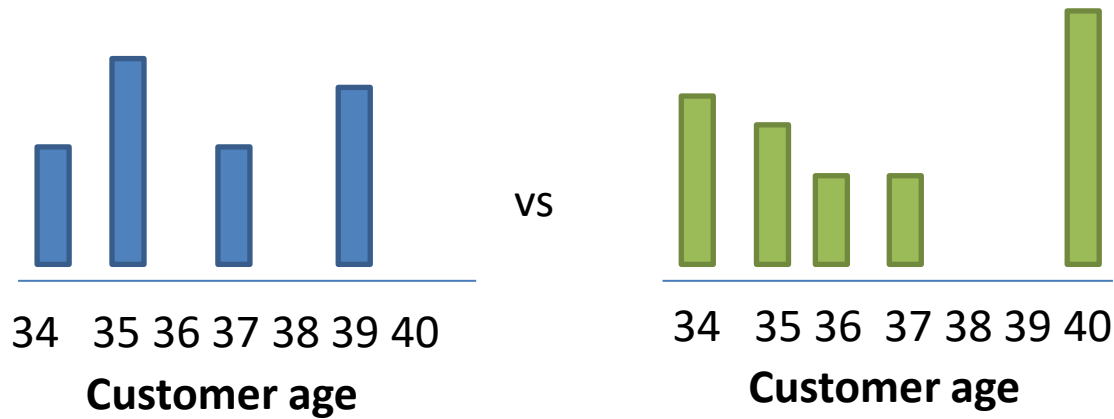
Comparing Distributions (1): Convert to vectors



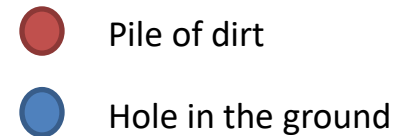
$\langle 3, 5, 0, 4 \rangle$ vs. $\langle 4, 3, 2, 0 \rangle$

- Makes sense for categorical domains

Comparing Distributions (2): Earth Movers Distance

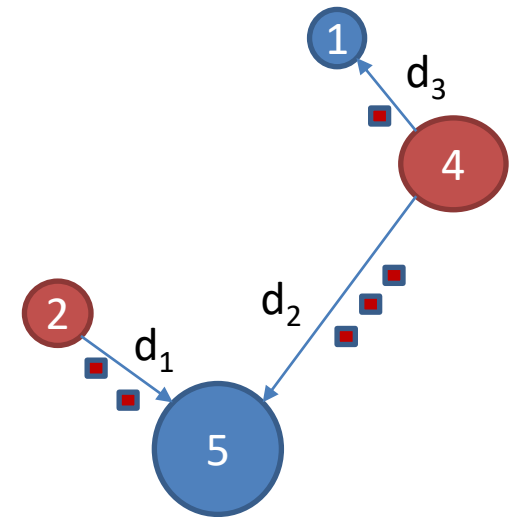


EMD Intuition



- Compute minimum amount of **work** required to change one distribution into the other.
 - **Unit of work**: the amount of work necessary to move one unit of weight by one unit of **ground distance**.
 - Informally: work = amount of dirt moved x distance travelled
 - **Ground distance**: the distance measure between weight locations.

$$EMD(X, Y) = \frac{\min(\text{Work}(X, Y))}{\min(\text{Weight}(X), \text{Weight}(Y))}$$



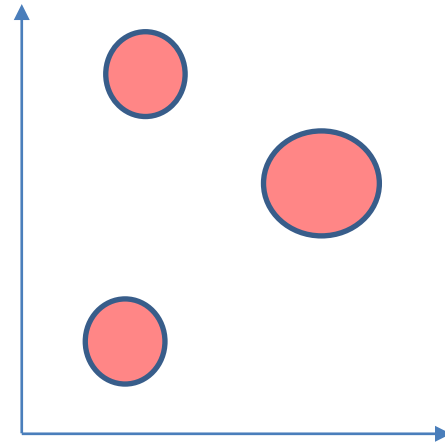
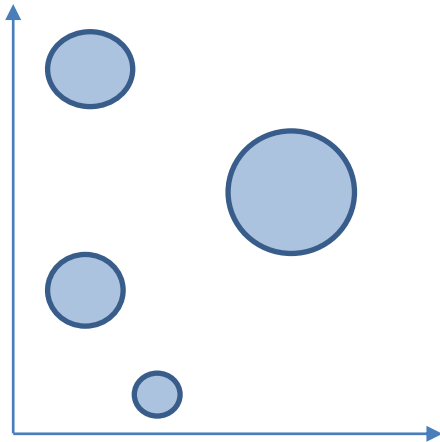
- EMD allows **partial matching** (when cumulative weights don't match): $\text{Weight}(X) \neq \text{Weight}(Y)$
 - all the weight in the **lighter** distribution should be matched to weight in the **heavier** distribution
 - In this case $EMD(x, y)$ is not a distance metric

$$\text{Work} = 2 * d_1 + 3 * d_2 + 1 * d_3$$

$$EMD = \text{WORK} / 6$$

Compare results of Clustering

- Clusters: $\{(x_i, n_i), i=1, \dots, n\}$
 - x_i is the cluster centroid
 - n_i is the size of the cluster



Compare Features Exported from dataset

- Features: $\{(f_i, n_i), i=1, \dots, n\}$
 - f_i : feature i
 - n_i : number of times f_i appears in dataset
 - Ground distance: $\text{dist}(f_i, f_j)$

Neat Application: Word Movers Distance (Kusner et. al.)

