

# Clustering Preliminaries

Applications

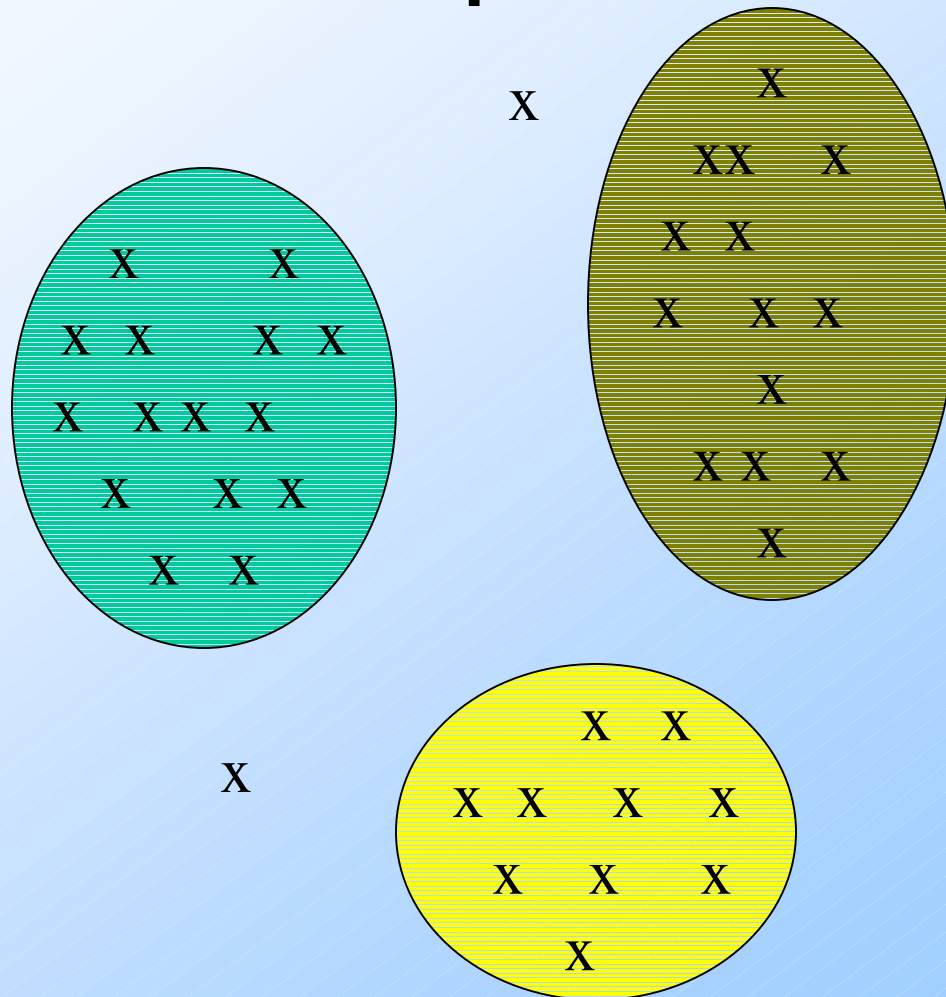
Euclidean/Non-Euclidean Spaces

Distance Measures

# The Problem of Clustering

- ◆ Given a set of points, with a notion of distance between points, group the points into some number of *clusters*, so that members of a cluster are in some sense as close to each other as possible.

# Example



# Problems With Clustering

- ◆ Clustering in two dimensions looks easy.
- ◆ Clustering small amounts of data looks easy.
- ◆ And in most cases, looks are *not* deceiving.

# The Curse of Dimensionality

- ◆ Many applications involve not 2, but 10 or 10,000 dimensions.
- ◆ High-dimensional spaces look different: almost all pairs of points are at about the same distance.
  - ◆ **Example:** assume random points within a bounding box, e.g., values between 0 and 1 in each dimension.

# Example: SkyCat

- ◆ A catalog of 2 billion “sky objects” represents objects by their radiation in 9 dimensions (frequency bands).
- ◆ **Problem:** cluster into similar objects, e.g., galaxies, nearby stars, quasars, etc.
- ◆ Sloan Sky Survey is a newer, better version.

# Example: Clustering CD's (Collaborative Filtering)

- ◆ Intuitively: music divides into categories, and customers prefer a few categories.
  - ◆ But what are categories really?
- ◆ Represent a CD by the customers who bought it.
- ◆ Similar CD's have similar sets of customers, and vice-versa.

# The Space of CD's

- ◆ Think of a space with one dimension for each customer.
  - ◆ Values in a dimension may be 0 or 1 only.
- ◆ A CD's point in this space is  $(x_1, x_2, \dots, x_k)$ , where  $x_i = 1$  iff the  $i^{\text{th}}$  customer bought the CD.
  - ◆ Compare with the "shingle/signature" matrix: rows = customers; cols. = CD's.



## Space of CD's --- (2)

- ◆ For Amazon, the dimension count is tens of millions.
- ◆ An option: use minhashing/LSH to get Jaccard similarity between “close” CD's.
- ◆ 1 minus Jaccard similarity can serve as a (non-Euclidean) distance.

# Example: Clustering Documents

- ◆ Represent a document by a vector  $(x_1, x_2, \dots, x_k)$ , where  $x_i = 1$  iff the  $i^{\text{th}}$  word (in some order) appears in the document.
  - ◆ It actually doesn't matter if  $k$  is infinite; i.e., we don't limit the set of words.
- ◆ Documents with similar sets of words may be about the same topic.

# Example: Gene Sequences

- ◆ Objects are sequences of  $\{C,A,T,G\}$ .
- ◆ Distance between sequences is *edit distance*, the minimum number of inserts and deletes needed to turn one into the other.
- ◆ Note there is a “distance,” but no convenient space in which points “live.”

# Distance Measures

- ◆ Each clustering problem is based on some kind of “distance” between points.
- ◆ Two major classes of distance measure:
  1. *Euclidean*
  2. *Non-Euclidean*

# Euclidean Vs. Non-Euclidean

- ◆ A *Euclidean space* has some number of real-valued dimensions and “dense” points.
  - ◆ There is a notion of “average” of two points.
  - ◆ A *Euclidean distance* is based on the locations of points in such a space.
- ◆ A *Non-Euclidean distance* is based on properties of points, but not their “location” in a space.

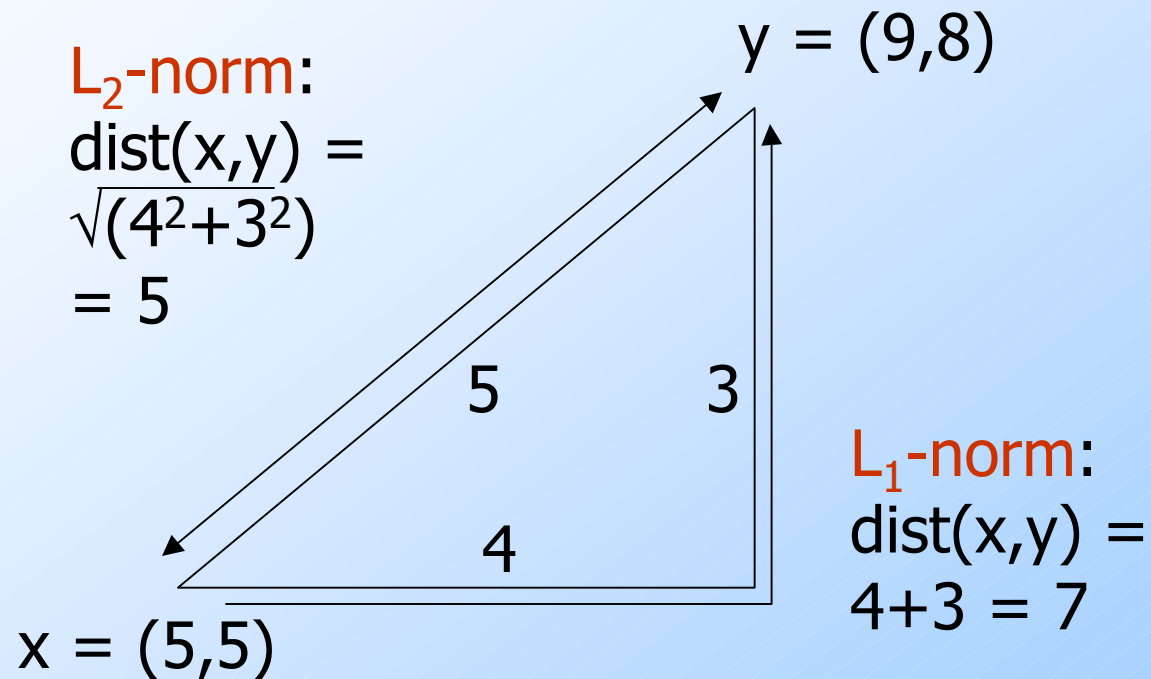
# Axioms of a Distance Measure

- ◆  $d$  is a *distance measure* if it is a function from pairs of points to real numbers such that:
  1.  $d(x,y) \geq 0$ .
  2.  $d(x,y) = 0$  iff  $x = y$ .
  3.  $d(x,y) = d(y,x)$ .
  4.  $d(x,y) \leq d(x,z) + d(z,y)$  (*triangle inequality*).

# Some Euclidean Distances

- ◆  $L_2$  norm :  $d(x,y)$  = square root of the sum of the squares of the differences between  $x$  and  $y$  in each dimension.
  - ◆ The most common notion of “distance.”
- ◆  $L_1$  norm : sum of the differences in each dimension.
  - ◆ *Manhattan distance* = distance if you had to travel along coordinates only.

# Examples of Euclidean Distances





# Another Euclidean Distance

- ◆  $L_\infty$  norm :  $d(x,y)$  = the maximum of the differences between  $x$  and  $y$  in any dimension.
- ◆ Note: the maximum is the limit as  $n$  goes to  $\infty$  of what you get by taking the  $n^{\text{th}}$  power of the differences, summing and taking the  $n^{\text{th}}$  root.

# Non-Euclidean Distances

- ◆ *Jaccard distance* for sets = 1 minus ratio of sizes of intersection and union.
- ◆ *Cosine distance* = angle between vectors from the origin to the points in question.
- ◆ *Edit distance* = number of inserts and deletes to change one string into another.

# Jaccard Distance for Bit-Vectors

- ◆ **Example:**  $p_1 = 10111$ ;  $p_2 = 10011$ .
  - ◆ Size of intersection = 3; size of union = 4, Jaccard similarity (not distance) =  $3/4$ .
- ◆ Need to make a distance function satisfying triangle inequality and other laws.
- ◆  $d(x,y) = 1 - (\text{Jaccard similarity})$  works.

# Why J.D. Is a Distance Measure

- ◆  $d(x,x) = 0$  because  $x \cap x = x \cup x$ .
- ◆  $d(x,y) = d(y,x)$  because union and intersection are symmetric.
- ◆  $d(x,y) \geq 0$  because  $|x \cap y| \leq |x \cup y|$ .
- ◆  $d(x,y) \leq d(x,z) + d(z,y)$  trickier --- next slide.

# Triangle Inequality for J.D.

$$1 - \frac{|x \cap z|}{|x \cup z|} + 1 - \frac{|y \cap z|}{|y \cup z|} \geq 1 - \frac{|x \cap y|}{|x \cup y|}$$

- ◆ Remember:  $|a \cap b|/|a \cup b| =$  probability that  $\text{minhash}(a) = \text{minhash}(b)$ .
- ◆ Thus,  $1 - |a \cap b|/|a \cup b| =$  probability that  $\text{minhash}(a) \neq \text{minhash}(b)$ .

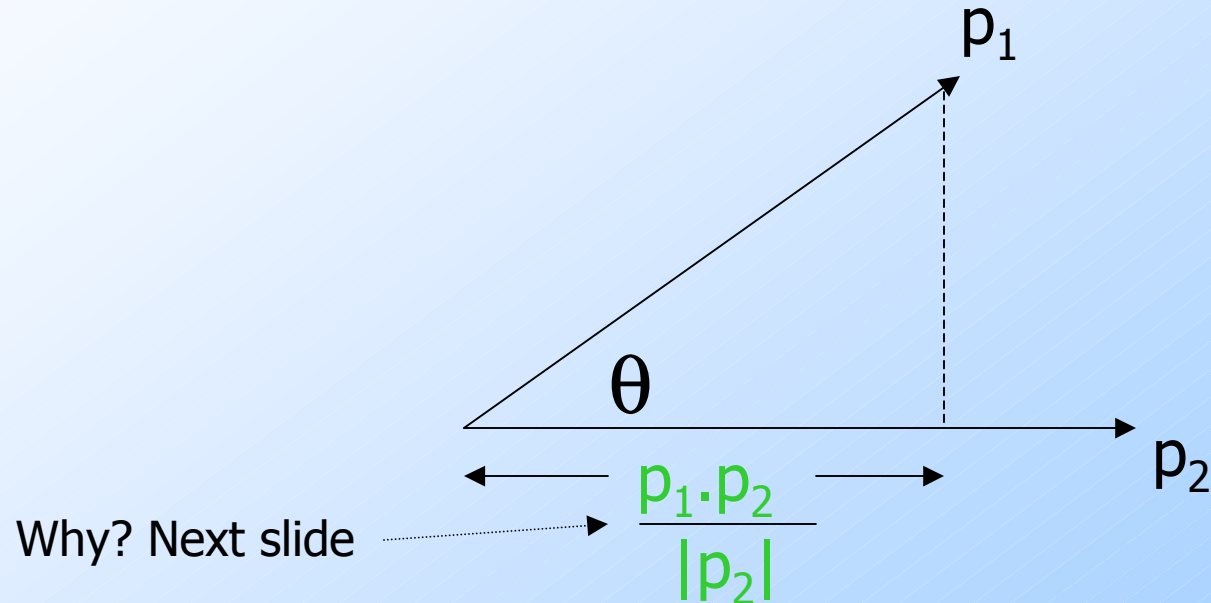
# Triangle Inequality --- (2)

- ◆ Observe that  $\text{prob}[\text{minhash}(x) \neq \text{minhash}(y)] \leq \text{prob}[\text{minhash}(x) \neq \text{minhash}(z)] + \text{prob}[\text{minhash}(z) \neq \text{minhash}(y)]$
- ◆ **Clincher**: whenever  $\text{minhash}(x) \neq \text{minhash}(y)$ , at least one of  $\text{minhash}(x) \neq \text{minhash}(z)$  and  $\text{minhash}(z) \neq \text{minhash}(y)$  must be true.

# Cosine Distance

- ◆ Think of a point as a vector from the origin  $(0,0,\dots,0)$  to its location.
- ◆ Two points' vectors make an angle, whose cosine is the normalized dot-product of the vectors:  $p_1 \cdot p_2 / |p_2| |p_1|$ .
  - ◆ Example  $p_1 = 00111$ ;  $p_2 = 10011$ .
  - ◆  $p_1 \cdot p_2 = 2$ ;  $|p_1| = |p_2| = \sqrt{3}$ .
  - ◆  $\cos(\theta) = 2/3$ ;  $\theta$  is about 48 degrees.

# Cosine-Measure Diagram



$$\text{dist}(p_1, p_2) = \theta = \arccos\left(\frac{p_1 \cdot p_2}{|p_2| |p_1|}\right)$$

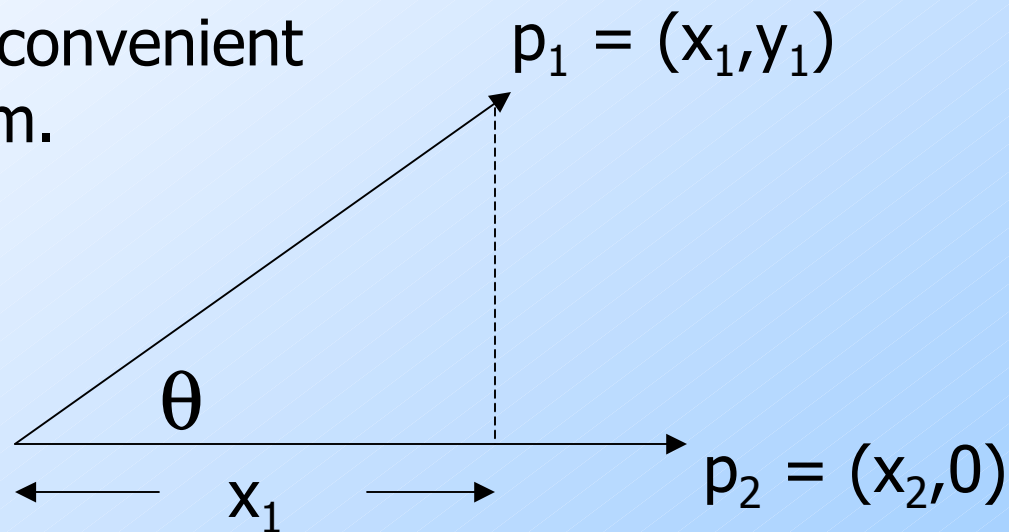


# Why?

Dot product is invariant under rotation, so pick convenient coordinate system.

$$p_1 \cdot p_2 = x_1 x_2.$$

$$|p_2| = x_2.$$



$$x_1 = x_1 x_2 / x_2 = p_1 \cdot p_2 / |p_2|$$

# Why C.D. Is a Distance Measure

- ◆  $d(x,x) = 0$  because  $\arccos(1) = 0$ .
- ◆  $d(x,y) = d(y,x)$  by symmetry.
- ◆  $d(x,y) \geq 0$  because angles are chosen to be in the range 0 to 180 degrees.
- ◆ **Triangle inequality**: physical reasoning.  
If I rotate an angle from  $x$  to  $z$  and then from  $z$  to  $y$ , I can't rotate less than from  $x$  to  $y$ .

# Edit Distance

- ◆ The edit distance of two strings is the number of inserts and deletes of characters needed to turn one into the other.
- ◆ **Equivalently:**  $d(x,y) = |x| + |y| - 2|LCS(x,y)|$ .
  - ◆ LCS = *longest common subsequence* = longest string obtained both by deleting from  $x$  and deleting from  $y$ .

# Example

- ◆  $x = abcde$  ;  $y = bcduve$ .
- ◆ Turn  $x$  into  $y$  by deleting  $a$ , then inserting  $u$  and  $v$  after  $d$ .
  - ◆ Edit-distance = 3.
- ◆ Or,  $\text{LCS}(x,y) = bcde$ .
- ◆  $|x| + |y| - 2|\text{LCS}(x,y)| = 5 + 6 - 2*4 = 3$ .

# Why E.D. Is a Distance Measure

- ◆  $d(x,x) = 0$  because 0 edits suffice.
- ◆  $d(x,y) = d(y,x)$  because insert/delete are inverses of each other.
- ◆  $d(x,y) \geq 0$ : no notion of negative edits.
- ◆ **Triangle inequality**: changing  $x$  to  $z$  and then to  $y$  is one way to change  $x$  to  $y$ .

# Variant Edit Distance

- ◆ Allow insert, delete, and *mutate*.
  - ◆ Change one character into another.
- ◆ Minimum number of inserts, deletes, and mutates also forms a distance measure.