# CS145: INTRODUCTION TO DATA MINING

## 7: Vector Data: K Nearest Neighbor

**Instructor: Yizhou Sun**

yzsun@cs.ucla.edu

October 23, 2018

# Methods to Learn: Last Lecture

| | Vector Data | Set Data | Sequence Data | Text Data |
|---|---|---|---|---|
| **Classification** | **Logistic Regression; Decision Tree**; KNN **SVM**; **NN** | | | Naïve Bayes for Text |
| **Clustering** | K-means; hierarchical clustering; DBSCAN; Mixture Models | | | PLSA |
| **Prediction** | **Linear Regression** GLM* | | | |
| **Frequent Pattern Mining** | | Apriori; FP growth | GSP; PrefixSpan | |
| **Similarity Search** | | | DTW | |

# Methods to Learn

| | Vector Data | Set Data | Sequence Data | Text Data |
|---|---|---|---|---|
| **Classification** | **Logistic Regression; Decision Tree; KNN SVM; NN** | | | Naïve Bayes for Text |
| **Clustering** | K-means; hierarchical clustering; DBSCAN; Mixture Models | | | PLSA |
| **Prediction** | **Linear Regression** GLM* | | | |
| **Frequent Pattern Mining** | | Apriori; FP growth | GSP; PrefixSpan | |
| **Similarity Search** | | | DTW | |

# K Nearest Neighbor

- Introduction ⬅
- kNN
- Similarity and Dissimilarity
- Summary

# Lazy vs. Eager Learning

- Lazy vs. eager learning
  - **Lazy learning** (e.g., instance-based learning): Simply stores training data (or only minor processing) and waits until it is given a test tuple
  - **Eager learning** (the above discussed methods): Given a set of training tuples, constructs a classification model before receiving new (e.g., test) data to classify
- Lazy: less time in training but more time in predicting
- Accuracy
  - Lazy method effectively uses a richer hypothesis space since it uses many local linear functions to form an implicit global approximation to the target function
  - Eager: must commit to a single hypothesis that covers the entire instance space

# Lazy Learner: Instance-Based Methods

- Instance-based learning:
  - Store training examples and delay the processing ("lazy evaluation") until a new instance must be classified
- Typical approaches
  - *k*-nearest neighbor approach
    - Instances represented as points in, e.g., a Euclidean space.
  - Locally weighted regression
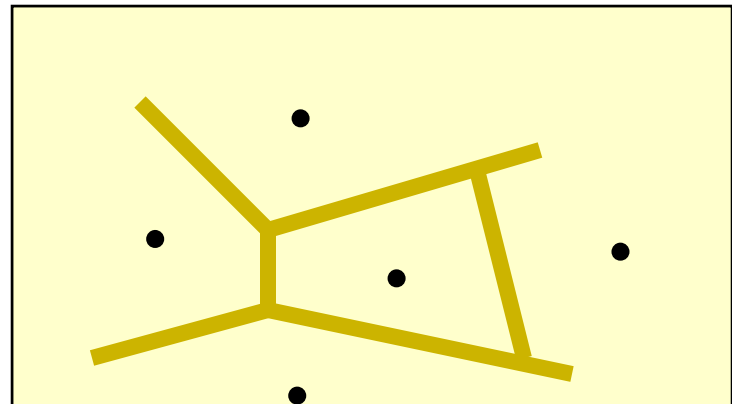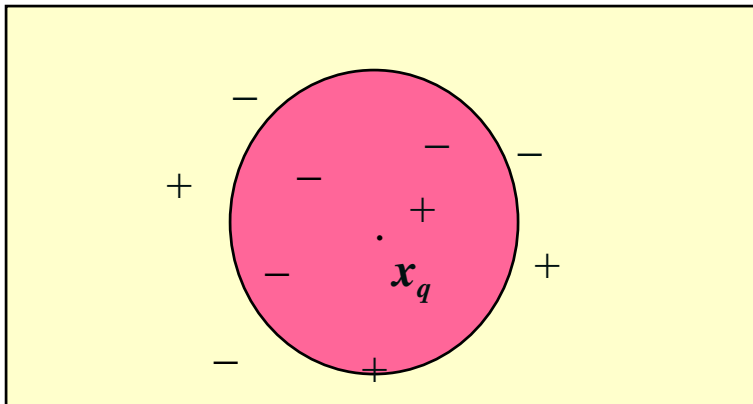    - Constructs local approximation

# K Nearest Neighbor

- Introduction
- kNN
- Similarity and Dissimilarity
- Summary

# The *k*-Nearest Neighbor Algorithm

- All instances correspond to points in the n-D space
- The nearest neighbor are defined in terms of a distance measure, dist($\mathbf{X_1}$, $\mathbf{X_2}$)
- Target function could be discrete- or real- valued
- For discrete-valued, *k*-NN returns the <span style="color:red">most common value</span> among the *k* training examples nearest to $x_q$
- Vonoroi diagram: the decision surface induced by 1-NN for a typical set of training examples
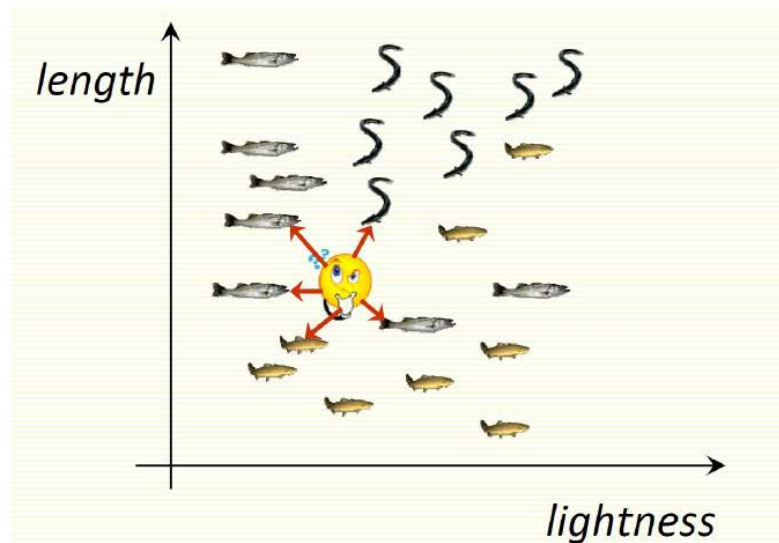
# kNN Example

$X = (length, lightness)$

Classes $= \{$salmon, sea bass, eel$\}$

Task: Identify fish given its (length, lightness)



$K = 5 : 3$ sea bass, $1$ eel, $1$ salmon $\Rightarrow$ sea bass

# kNN Algorithm Summary

- Choose K

- For a given new instance $X_{new}$ , find K closest training points w.r.t. a distance measure
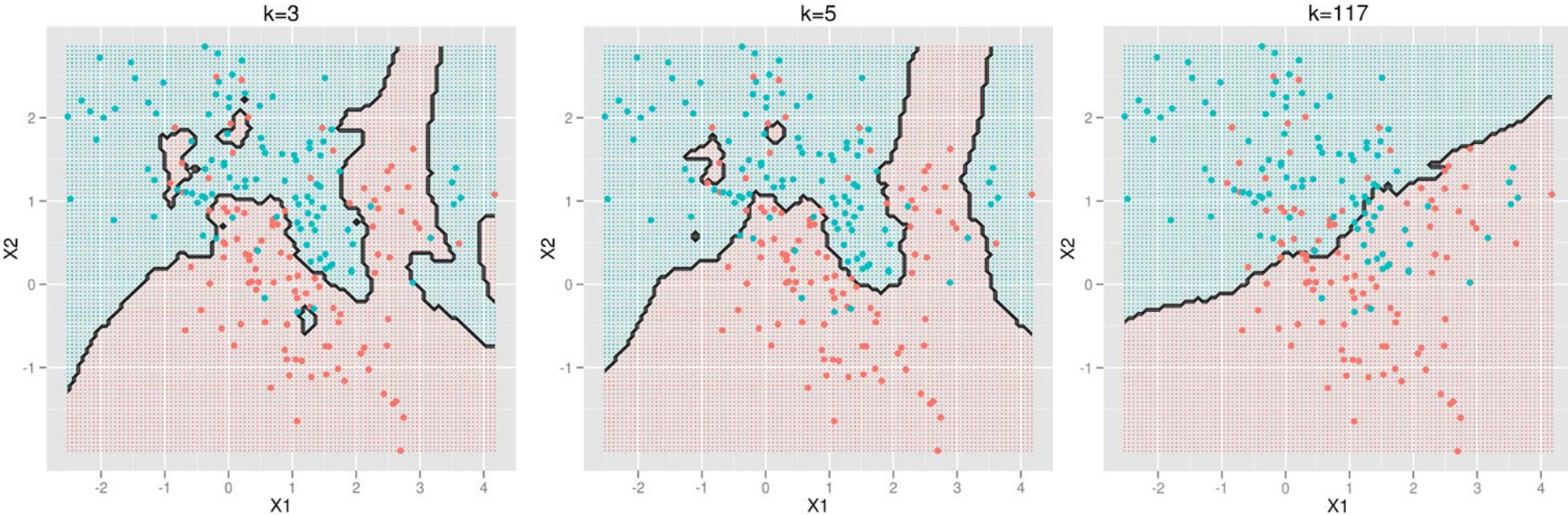
- Classify $X_{new}$ = majority vote among the K points

# Discussion on the *k*-NN Algorithm

- *k*-NN for <u>real-valued prediction</u> for a given unknown tuple
  - Returns the mean values of the $k$ nearest neighbors
- <u>Distance-weighted</u> nearest neighbor algorithm
  - Weight the contribution of each of the $k$ neighbors according to their distance to the query $x_q$
    - Give greater weight to closer neighbors $e.g., w_i = \dfrac{1}{d(x_q, x_i)^2}$
    - $y_q = \dfrac{\sum w_i y_i}{\sum w_i}$, where $x_i$'s are $x_q$'s nearest neighbors $\qquad w_i = \exp(-d(x_q, x_i)^2 / 2\sigma^2)$
- <u>Robust</u> to noisy data by averaging *k*-nearest neighbors
- <u>Curse of dimensionality</u>: distance between neighbors could be dominated by irrelevant attributes
  - To overcome it, axes stretch or elimination of the least relevant attributes

# Selection of k for kNN

- The number of neighbors k
  - Small k: overfitting (high var., low bias)
  - Big k: bringing too many irrelevant points (high bias, low var.)



  - More discussions:
http://scott.fortmann-roe.com/docs/BiasVariance.html

# K Nearest Neighbor

- Introduction

- kNN

- Similarity and Dissimilarity ⬅

- Summary

# Similarity and Dissimilarity

- **Similarity**
  - Numerical measure of how alike two data objects are
  - Value is higher when objects are more alike
  - Often falls in the range [0,1]
- **Dissimilarity** (e.g., distance)
  - Numerical measure of how different two data objects are
  - Lower when objects are more alike
  - Minimum dissimilarity is often 0
  - Upper limit varies
- **Proximity** refers to a similarity or dissimilarity

# Data Matrix and Dissimilarity Matrix

- Data matrix
  - n data points with p dimensions
  - Two modes

$$\begin{bmatrix} x_{11} & \cdots & x_{1f} & \cdots & x_{1p} \\ \cdots & \cdots & \cdots & \cdots & \cdots \\ x_{i1} & \cdots & x_{if} & \cdots & x_{ip} \\ \cdots & \cdots & \cdots & \cdots & \cdots \\ x_{n1} & \cdots & x_{nf} & \cdots & x_{np} \end{bmatrix}$$
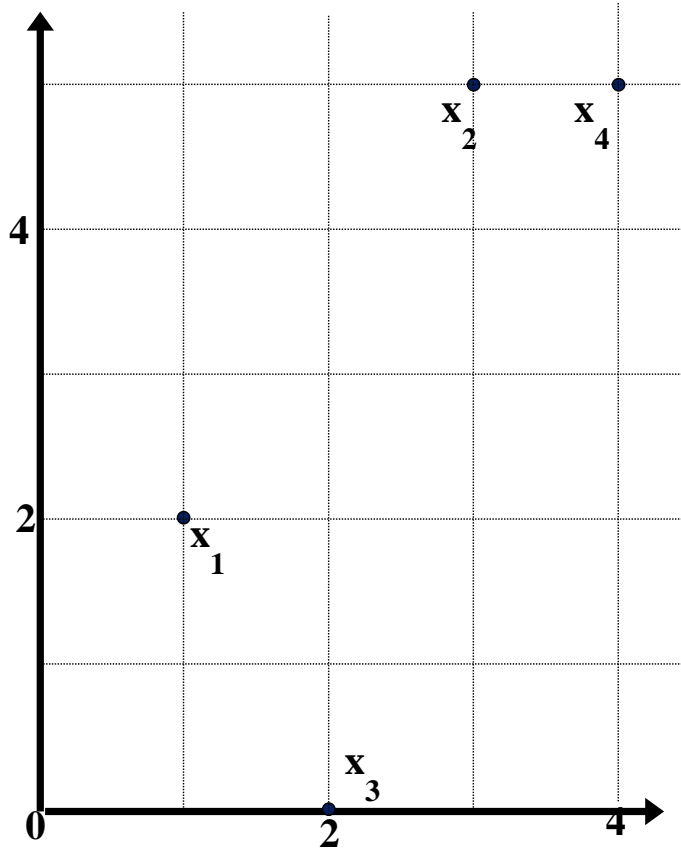
- Dissimilarity matrix
  - n data points, but registers only the distance
  - A triangular matrix
  - Single mode

$$\begin{bmatrix} 0 & & & & \\ d(2,1) & 0 & & & \\ d(3,1) & d(3,2) & 0 & & \\ \vdots & \vdots & \vdots & & \\ d(n,1) & d(n,2) & \cdots & \cdots & 0 \end{bmatrix}$$

# Example:
# Data Matrix and Dissimilarity Matrix



## Data Matrix

| point | attribute1 | attribute2 |
|-------|-----------|-----------|
| x1 | 1 | 2 |
| x2 | 3 | 5 |
| x3 | 2 | 0 |
| x4 | 4 | 5 |

## Dissimilarity Matrix

### (with Euclidean Distance)

|    | x1 | x2 | x3 | x4 |
|----|------|-----|------|---|
| x1 | 0 |  |  |  |
| x2 | 3.61 | 0 |  |  |
| x3 | 2.24 | 5.1 | 0 |  |
| x4 | 4.24 | 1 | 5.39 | 0 |

16

# Distance on Numeric Data: Minkowski Distance

- *Minkowski distance*: A popular distance measure

$$d(i, j) = \sqrt[h]{|x_{i1} - x_{j1}|^h + |x_{i2} - x_{j2}|^h + \cdots + |x_{ip} - x_{jp}|^h}$$

where $i = (x_{i1}, x_{i2}, ..., x_{ip})$ and $j = (x_{j1}, x_{j2}, ..., x_{jp})$ are two $p$-dimensional data objects, and $h$ is the order (the distance so defined is also called L-$h$ norm)

- Properties

  - $d(i, j) > 0$ if $i \neq j$, and $d(i, i) = 0$ (Positive definiteness)

  - $d(i, j) = d(j, i)$ (Symmetry)

  - $d(i, j) \leq d(i, k) + d(k, j)$ (Triangle Inequality)

- A distance that satisfies these properties is a metric

# Special Cases of Minkowski Distance

- $h = 1$: Manhattan (city block, $L_1$ norm) distance
    - E.g., the Hamming distance: the number of bits that are different between two binary vectors

$$d(i,j) = |x_{i1} - x_{j1}| + |x_{i2} - x_{j2}| + ... + |x_{ip} - x_{jp}|$$

- $h = 2$: ($L_2$ norm) Euclidean distance

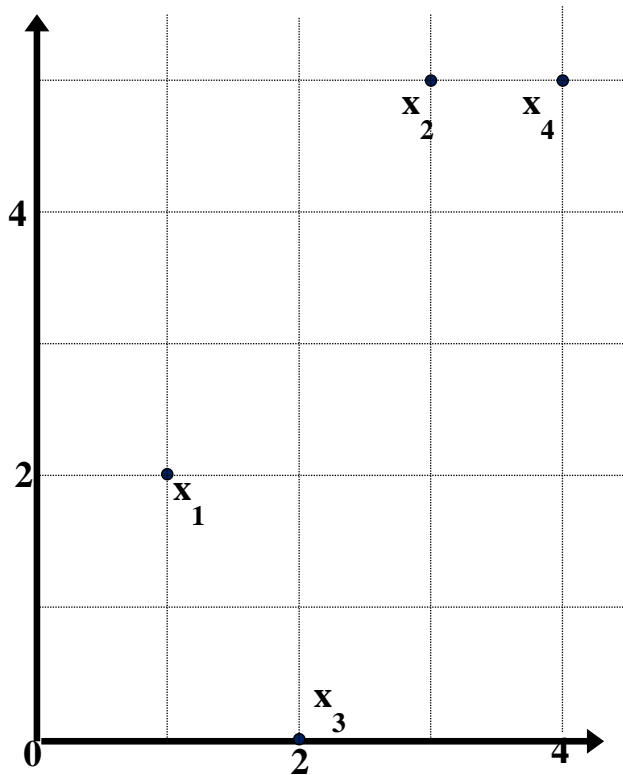$$d(i,j) = \sqrt{(|x_{i1} - x_{j1}|^2 + |x_{i2} - x_{j2}|^2 + ... + |x_{ip} - x_{jp}|^2)}$$

- $h \to \infty$. "supremum" ($L_{max}$ norm, $L_\infty$ norm) distance.
    - This is the maximum difference between any component (attribute) of the vectors

$$d(i,j) = \lim_{h \to \infty} \left( \sum_{f=1}^{p} |x_{if} - x_{jf}|^h \right)^{\frac{1}{h}} = \max_{f}^{p} |x_{if} - x_{jf}|$$

# Example: Minkowski Distance

**Dissimilarity Matrices**

| point | attribute 1 | attribute 2 |
|-------|-------------|-------------|
| **x1** | 1 | 2 |
| **x2** | 3 | 5 |
| **x3** | 2 | 0 |
| **x4** | 4 | 5 |

**Manhattan (L$_1$)**

| L | x1 | x2 | x3 | x4 |
|---|----|----|----|----|
| **x1** | 0 | | | |
| **x2** | 5 | 0 | | |
| **x3** | 3 | 6 | 0 | |
| **x4** | 6 | 1 | 7 | 0 |

**Euclidean (L$_2$)**

| L2 | x1 | x2 | x3 | x4 |
|----|----|----|----|----|
| **x1** | 0 | | | |
| **x2** | 3.61 | 0 | | |
| **x3** | 2.24 | 5.1 | 0 | |
| **x4** | 4.24 | 1 | 5.39 | 0 |

**Supremum**

| L$_\infty$ | x1 | x2 | x3 | x4 |
|------------|----|----|----|----|
| **x1** | 0 | | | |
| **x2** | 3 | 0 | | |
| **x3** | 2 | 5 | 0 | |
| **x4** | 3 | 1 | 5 | 0 |

# Standardizing Numeric Data

- Z-score: $z = \dfrac{x - \mu}{\sigma}$

  - X: raw score to be standardized, μ: mean of the population, σ: standard deviation

  - the distance between the raw score and the population mean in units of the standard deviation

  - negative when the raw score is below the mean, "+" when above

- An alternative way: Calculate the mean absolute deviation

$$s_f = \frac{1}{n}(|x_{1f} - m_f| + |x_{2f} - m_f| + ... + |x_{nf} - m_f|)$$

where $m_f = \dfrac{1}{n}(x_{1f} + x_{2f} + ... + x_{nf}).$

  - standardized measure (*z-score*): $z_{if} = \dfrac{x_{if} - m_f}{s_f}$

- Using mean absolute deviation is more robust than using standard deviation

# Proximity Measure for Nominal Attributes

- Can take 2 or more states, e.g., red, yellow, blue, green (generalization of a binary attribute)

- <u>Method 1</u>: Simple matching

  - $m$: # of matches, $p$: total # of variables

$$d(i, j) = \frac{p - m}{p}$$

- <u>Method 2</u>: Use a large number of binary attributes

  - creating a new binary attribute for each of the $M$ nominal states

# Proximity Measure for Binary Attributes

- A contingency table for binary data

- Distance measure for symmetric binary variables:

- Distance measure for asymmetric binary variables:

- Jaccard coefficient (*similarity* measure for *asymmetric* binary variables):

Object $j$

|  | 1 | 0 | sum |
|---|---|---|---|
| Object $i$  1 | $q$ | $r$ | $q+r$ |
| 0 | $s$ | $t$ | $s+t$ |
| sum | $q+s$ | $r+t$ | $p$ |

$$d(i, j) = \frac{r + s}{q + r + s + t}$$

$$d(i, j) = \frac{r + s}{q + r + s}$$

$$sim_{Jaccard}(i, j) = \frac{q}{q + r + s}$$

# Dissimilarity between Binary Variables

- Example

| Name | Gender | Fever | Cough | Test-1 | Test-2 | Test-3 | Test-4 |
|------|--------|-------|-------|--------|--------|--------|--------|
| Jack | M | Y | N | P | N | N | N |
| Mary | F | Y | N | P | N | P | N |
| Jim | M | Y | P | N | N | N | N |

- Gender is a symmetric attribute
- The remaining attributes are asymmetric binary
- Let the values Y and P be 1, and the value N 0

$$d(jack, mary) = \frac{0+1}{2+0+1} = 0.33$$

$$d(jack, jim) = \frac{1+1}{1+1+1} = 0.67$$

$$d(jim, mary) = \frac{1+2}{1+1+2} = 0.75$$

# Ordinal Variables

- Order is important, e.g., rank

- Can be treated like interval-scaled

  - replace $x_{if}$ by their rank $\quad r_{if} \in \{1, \dots, M_f\}$

  - map the range of each variable onto [0, 1] by replacing *i*-th object in the *f*-th variable by

$$z_{if} = \frac{r_{if} - 1}{M_f - 1}$$

  - compute the dissimilarity using methods for interval-scaled variables

# Attributes of Mixed Type

- A database may contain all attribute types
  - Nominal, symmetric binary, asymmetric binary, numeric, ordinal
- One may use a weighted formula to combine their effects

$$d(i,j) = \frac{\Sigma_{f=1}^{p} \delta_{ij}^{(f)} d_{ij}^{(f)}}{\Sigma_{f=1}^{p} \delta_{ij}^{(f)}}$$

- $f$ is binary or nominal:

  $d_{ij}^{(f)} = 0$ if $x_{if} = x_{jf}$, or $d_{ij}^{(f)} = 1$ otherwise
- $f$ is numeric: use the normalized distance
- $f$ is ordinal
  - Compute ranks $r_{if}$ and $\quad z_{if} = \dfrac{r_{if} - 1}{M_f - 1}$
  - Treat $z_{if}$ as interval-scaled

# Cosine Similarity

- A **document** can be represented by thousands of attributes, each recording the *frequency* of a particular word (such as keywords) or phrase in the document.

| Document | team | coach | hockey | baseball | soccer | penalty | score | win | loss | season |
|----------|------|-------|--------|----------|--------|---------|-------|-----|------|--------|
| Document1 | 5 | 0 | 3 | 0 | 2 | 0 | 0 | 2 | 0 | 0 |
| Document2 | 3 | 0 | 2 | 0 | 1 | 1 | 0 | 1 | 0 | 1 |
| Document3 | 0 | 7 | 0 | 2 | 1 | 0 | 0 | 3 | 0 | 0 |
| Document4 | 0 | 1 | 0 | 0 | 1 | 2 | 2 | 0 | 3 | 0 |

- Other vector objects: gene features in micro-arrays, …
- Applications: information retrieval, biologic taxonomy, gene feature mapping, …
- Cosine measure: If $d_1$ and $d_2$ are two vectors (e.g., term-frequency vectors), then
$$\cos(d_1, d_2) = (d_1 \bullet d_2) / ||d_1||\ ||d_2||\ ,$$
   where $\bullet$ indicates vector dot product, $||d||$: the length of vector $d$

# Example: Cosine Similarity

- cos($d_1$, $d_2$) = ($d_1 \bullet d_2$) / ||$d_1$|| ||$d_2$|| ,
  where $\bullet$ indicates vector dot product, $||d|$ : the length of vector $d$

- Ex: Find the **similarity** between documents 1 and 2.

$d_1$ = (5, 0, 3, 0, 2, 0, 0, 2, 0, 0)
$d_2$ = (3, 0, 2, 0, 1, 1, 0, 1, 0, 1)

$d_1 \bullet d_2$ = 5*3+0*0+3*2+0*0+2*1+0*1+0*1+2*1+0*0+0*1 = 25
$||d_1||$ = (5*5+0*0+3*3+0*0+2*2+0*0+0*0+2*2+0*0+0*0)$^{0.5}$=(42)$^{0.5}$ = 6.481
$||d_2||$ = (3*3+0*0+2*2+0*0+1*1+1*1+0*0+1*1+0*0+1*1)$^{0.5}$=(17)$^{0.5}$ = 4.12
cos($d_1$, $d_2$) = 0.94

# K Nearest Neighbor

- Introduction
- kNN
- Similarity and Dissimilarity
- Summary

# Summary

- Instance-Based Learning

  - Lazy learning vs. eager learning; K-nearest neighbor algorithm; Similarity / dissimilarity measures