





© 2012 Microsoft Corporation. All rights reserved. Microsoft, Windows, Windows Vista and other product names are or may be registered trademarks and/or trademarks in the U.S. and/or other countries. The information herein is for informational purposes only and represents the current view of Microsoft Corporation as of the date of this presentation. Because Microsoft must respond to changing market conditions, it should not be interpreted to be a commitment on the part of Microsoft, and Microsoft comot guarantee the accuracy of any information provided after the date of this presentation. MICROSOFT MAKES NO WARRANTIES, EXPRESS, IMPLIED OR STATUTORY, AS TO THE INFORMATION IN THIS PRESENTATION. What is Similarity? Slide based on one by Eamonn Keogh



Similarity is hard to define, but... "We know it when we see it"

Similarity and Dissimilarity So, how do we do it using computers? Using the notion of "distance" between two objects If they are "similar", the distance should be around 0 or 1 If they are "dissimilar", the distance should be opposite (1 or 0) Easily computable for a given data/attribute type Appropriate for a given data/attribute type So, we have several "distance" metrics "city block" or "taxi", Euclidan, L-forms, edit distance, similarity matching, Jaccard, Pearson correlation, Spearman, to name a few This also means you have to choose the correct one (burden on the analyst)

Similarity and Dissimilarity

Similarity

- o Numerical measure of how alike two data objects are
- Is higher (close to 1) when objects are more alike
- Often falls in the range [0,1] or [-1 to 1]
- Dissimilarity
 - o Numerical measure of how different are two data objects
 - o Lower when objects are more alike
 - $\circ \hspace{0.1in} \text{Minimum dissimilarity is often } 0$
 - o Upper limit varies
- Proximity refers to similarity or dissimilarity



© 2012 Microsoft Corporation. All rights reserved. Microsoft, Windows, Windows Vista and other product names are or may be registered trademarks and/or trademarks in the U.S. and/or other countries. The information herein is for informational purposes only and represents the current view of Microsoft Corporation as of the date of this presentation. Because Microsoft must respond to changing market conditions, it should not be interpreted to be a commitment on the part of Microsoft, and Microsoft contor guarantee the accuracy of any information provided after the date of this presentation. MICROSOFT MAKES NO WARRANTIES, EXPRESS, IMPLIED OR STATUTORY, AS TO THE INFORMATION IN THIS PRESENTATION.

 \mathbf{k}

	Similarity/Dissimilarity for Simple Attributes p and q are the attribute values for two data objects d(p, q) and s(p, q) are dissimilarity and similarity between p and q.			
	Attribute	Dissimilarity	Similarity	
	Type			
	Nominal	$d = \begin{cases} 0 & \text{if } p = q \\ 1 & \text{if } p \neq q \end{cases}$	$s = \left\{ egin{array}{cc} 1 & ext{if } p = q \ 0 & ext{if } p eq q \end{array} ight.$	
	Ordinal	$d = \frac{ p-q }{n-1}$ (values mapped to integers 0 to $n-1$, where n is the number of values)	$s = 1 - \frac{ p-q }{n-1}$	
	Interval or Ratio	d = p - q	$s = -d, s = \frac{1}{1+d} \text{ or}$ $s = 1 - \frac{d-\min d}{\max d-\min d}$	
Table 5.1. Similarity and dissimilarity for simple attributes			attributes	
9	9 Convide C2007-2017 The University of Texas at Adjuston. All Rights Reserved.			

Euclidean Distance $dist = \sqrt{\sum_{k=1}^{n} (p_k - q_k)^2}$ Where *n* is the number of dimensions (attributes) and *p_k* and *q_k* are, respectively, the kth attributes (components) or data objects *p* and *q*.

Standardization is necessary, if scales differ.

Euclidean Distance





© 2012 Microsoft Corporation. All rights reserved. Microsoft, Windows, Windows Vista and other product names are or may be registered trademarks and/or trademarks in the U.S. and/or other countries. The information herein is for informational purposes only and represents the current view of Microsoft Corporation as of the date of this presentation. Because Microsoft must respond to changing market conditions, it should not be interpreted to be a commitment on the part of Microsoft, and Microsoft comot guarantee the accuracy of any information provided after the date of this presentation. MICROSOFT MAKES NO WARRANTIES, EXPRESS, IMPLIED OR STATUTORY, AS TO THE INFORMATION IN THIS PRESENTATION.



Common Properties of a Distance

- Distances, such as the Euclidean distance, have some well known properties.
 - 1. $d(p, q) \ge 0$ for all p and q and d(p, q) = 0 only if p = q. (Positive definiteness)
 - 2. d(p, q) = d(q, p) for all p and q. (Symmetry)
 - 3. $d(p, r) \le d(p, q) + d(q, r)$ for all points p, q, and r. We know it is true for scalars.

(Triangle Inequality) but it is also true for vectors! where d(p, q) is the distance (similarity) between points (data objects), p and

- q.
- ✤ A distance that satisfies these properties is a metric

Slide based on one by Eamonn Keogh

Intuitions behind desirable distance measure properties

D(A,B) = D(B,A) Otherwise you could claim "Alex looks like Bob, but Bob looks nothing like Alex."

D(A,A) = 0Otherwise you could claim "Alex does not look like Alex."

© 2012 Microsoft Corporation. All rights reserved. Microsoft, Windows, Windows Vista and other product names are or may be registered trademarks and/or trademarks in the U.S. and/or other countries. The information herein is for informational purposes only and represents the current view of Microsoft Corporation as of the date of this presentation. Because Microsoft must respond to changing market conditions, it should not be interpreted to be a commitment on the part of Microsoft, and Microsoft comparison that et et accuracy of any information provided after the date of this presentation. MICROSOFT MAKES NO WARRANTIES, EXPRESS, IMPLIED OR STATUTORY, AS TO THE INFORMATION IN THIS PRESENTATION.

 \bigstar

Slide based on one by Eamonn Keogh

Intuitions behind desirable distance measure properties (continued)

D(A,B) = 0 Iff A=B Otherwise there are objects in your world that are different, but you cannot tell anart

 $D(A,B) \le D(A,C) + D(B,C)$ Otherwise you could claim "Alex is very like Carl and Bob is very like Carl, but Alex is very unlike Bob."

Common Properties of a Similarity

- Similarities, also have some well known properties.
 - 1. s(p, q) = 1 (or maximum similarity) only if p = q.
 - 2. s(p, q) = s(q, p) for all p and q. (Symmetry)

where s(p, q) is the similarity between points (data objects), p and q.

Similarity Between Binary Vectors

- Common situation is that objects, *p* and *q*, have only binary attributes ۰.
- Compute similarities using the following quantities ٠.
- $\begin{array}{l} \label{eq:simple Matching and Jaccard Coefficients (SMC and JC) \\ \mbox{SMC} = number of matches / number of attributes \\ = (M_{11} + M_{00}) / (M_{01} + M_{10} + M_{11} + M_{00}) \end{array}$ \diamond
- $J = number of 11 matches / number of not-both-zero attributes values = (M_{11}) / (M_{01} + M_{10} + M_{11})$
 - can also be calculated for sets using intersection and union



© 2012 Microsoft Corporation. All rights reserved. Microsoft, Windows, Windows Vista and other product names are or may be registered trademarks and/or trademarks in the U.S. and/or other countries. The information herein is for informational purposes only and represents the current view of Microsoft Corporation as of the date of this presentation. Because Microsoft must respond to changing market conditions, it should not be interpreted to be a commitment on the part of Microsoft, and Microsoft canot guarantee the accuracy of any information provided after the date of this presentation. MICROSOFT MAKES NO WARRANTIES, EXPRESS, IMPLIED OR STATUTORY, AS TO THE INFORMATION IN THIS PRESENTATION.

 \bigstar





© 2012 Microsoft Corporation. All rights reserved. Microsoft, Windows, Windows Vista and other product names are or may be registered trademarks and/or trademarks in the U.S. and/or other countries. The information herein is for informational purposes only and represents the current view of Microsoft Corporation as of the date of this presentation. Because Microsoft must respond to changing market conditions, it should not be interpreted to be a commitment on the part of Microsoft, and Microsoft canot guarantee the accuracy of any information provided after the date of this presentation. MICROSOFT MAKES NO WARRANTIES, EXPRESS, IMPLIED OR STATUTORY, AS TO THE INFORMATION IN THIS PRESENTATION. Cosine Similarity We also know that cos(0) is 1 and cos(20) is 0, and cos(180) is -1 The general formula for non-right-angled triangles is $|| a b ||^2 = || a ||^2 + || b ||^2 - 2 || a || || b || Cos(0)$ When this is simplified for Cos(0), you get $(a \cdot b) / ||a|| ||b||$ (try simplifying!) If the vectors are orthogonal, the value is 0 showing they are perpendicular (more separation) If the vectors overlap (have 0 degree between them) the value is 1, showing complete similarity! If the value is -1, they are opposite vectors > Easy to generalize to n-dimensions

Very widely used similarity metric!

>



* Cosine Similarity: why we ignore magnitude? For example, vectors b and c below have 0 angle between them, but different magnitudes So, they are similar! However, vectors a and b look similar in magnitude, but *a* [5,3] have a large angle between them a-b They are considered dissimilar! Why? C [3, 6] Think of a, b, and c as document vectors b [2, 4] Doca 5 3 Doc b 2 4 HW: what if the vectors are of different length? How do you 6 compute cos(theta)? Doc c 3



© 2012 Microsoft Corporation. All rights reserved. Microsoft, Windows, Windows Vista and other product names are or may be registered trademarks and/or trademarks in the U.S. and/or other countries. The information herein is for informational purposes only and represents the current view of Microsoft Corporation as of the date of this presentation. Because Microsoft must respond to changing market conditions, it should not be interpreted to be a commitment on the part of Microsoft, and Microsoft contor guarantee the accuracy of any information provided after the date of this presentation. MICROSOFT MAKES NO WARRANTIES, EXPRESS, IMPLIED OR STATUTORY, AS TO THE INFORMATION IN THIS PRESENTATION.



Variance and standard deviation

- > If you want the data spread, variance and standard deviation can be used. Variance If you with the average degree to which each point differs from the mean (variability) We define the variance to be $r^2 = \frac{1}{n-1} \frac{1}{2n} (\alpha - 3^2)$ (why n-1 and not n?)
- The units of variance is not the same as that of data
- > A variance of 0 means all data points are identical!
- A high variance indicates that data points are very spread out from the mean and from one another.
- A distinction is made between sample mean and population mean represented as mu (μ) (n-1) is used fro sample mean and n is used for population mean;
- > There is also the notion of average deviation or Mean absolute deviation
- Uses absolute values instead of squaring to circumvent negative differences
 Used less frequently because the use of absolute values makes further calculations more complicated



© 2012 Microsoft Corporation. All rights reserved. Microsoft, Windows, Windows Vista and other product names are or may be registered trademarks and/or trademarks in the U.S. and/or other countries. The information herein is for informational purposes only and represents the current view of Microsoft Corporation as of the date of this presentation. Because Microsoft must respond to changing market conditions, it should not be interpreted to be a commitment on the part of Microsoft, and Microsoft canot guarantee the accuracy of any information provided after the date of this presentation. MICROSOFT MAKES NO WARRANTIES, EXPRESS, IMPLIED OR STATUTORY, AS TO THE INFORMATION IN THIS PRESENTATION.

Use of standard deviation

- Use: One of the most common methods of determining the risk an investment poses is standard deviation.
- Standard deviation helps determine market volatility or the spread of asset prices from their average price.
- When prices move wildly, standard deviation is high, meaning an investment will be risky. Low standard deviation means prices are calm, so investments come with low risk.
- Average spending in a restaurant can be better estimated by using standard deviation rather than the mean

Use of standard deviation, percentile

- The most common definition of a percentile is a number where a certain percentage of scores fall below that number.
- > Percentile: Let p be any integer between 0 and 100. The p^{tb} percentile of data set is the data value at which p percent of the value in the data set are less than or equal to this value Remember GRE, SAT, LSAT
- > You might know that you scored 67 out of 90 on a test. But that figure may not be helpful unless you know what percentile you fall into. If you know that your score is in the 90th percentile, that means you scored better than 90% of people who took the test.

Covariance

- Covariance is a statistical tool that is used to determine the relationship between the
- movement of two data sets or random variables (e.g., stock prices) Covariance measures the directional relationship between the two <u>data sets</u>. A
- positive covariance means that asset returns move together while a negative covariance means they move inversely.
- Used for diversifying portfolios in an investment!
- For choosing a portfolio where all stocks that do not behave the same way in prices/returns (not putting all eggs in one basket!) $Cov(x,y) = SUM [(x_i - x_m) * (y_i - y_m)] / (n - 1)$
- While the covariance does measure the directional relationship between two assets, it does not show the strength of the relationship between the two assets; the coefficient of correlation is a more appropriate indicator of this strength. Has units. Does not measure the strength or the dependency between variables!
- Strength is measured using a coefficient; dependency by regression!

Covariance and correlation

- > On the other hand, correlation measures the strength of the relationship between variables. Correlation is the scaled measure of covariance. It is dimensionless. In other words, the correlation coefficient is always a pure value and not measured in any units.
- > The relationship between two concepts can be expressed using the formula $\rho(X,Y) = \frac{Cov(X,Y)}{\sigma_X \sigma_Y} \quad (\text{Pearson Correlation coefficient})$

> Where:

- r(X,Y) the correlation between the variables X and Y
- Cov(X,Y) the covariance between the variables X and Y
- > σX − the standard deviation of the X-variable
- > σ Y the standard deviation of the Y-variable

 \bigstar





Correlation coefficient

- It is a value between -1 and 1
- Both -1 or 1 indicate very strong (or perfect) positive or negative correlation
- > Values between 0.3 and 0.7 (-0.3 and -0.7) indicate a moderate positive
- (negative) linear relationship
- The relationship between two variable is considered strong when their r value is large than $0.7 \ (\text{-}0.7)$
- A r value of 0 (or near 0) indicates no linear relationship!
- When the r value is 0, non-linear relationships could still exist. ≻
- The book has an example. Take a look at it. \triangleright
- This is used when both variables being studied are normally distributed. This
- correlation is affected by extreme values
- Otherwise, use Spearman's rank correlation coefficient which is better if one or both are not normally distributed!
- Look up Spearman's coefficient on your own



© 2012 Microsoft Corporation. All rights reserved. Microsoft, Windows, Windows Vista and other product names are or may be registered trademarks and/or trademarks in the U.S. and/or other countries. The information herein is for informational purposes only and represents the current view of Microsoft Corporation as of the date of this presentation. Because Microsoft must respond to changing market conditions, it should not be interpreted to be a commitment on the part of Microsoft, and Microsoft contor guarantee the accuracy of any information provided after the date of this presentation. MICROSOFT MAKES NO WARRANTIES, EXPRESS, IMPLIED OR STATUTORY, AS TO THE INFORMATION IN THIS PRESENTATION.





Selecting the right proximity/similarity measure

- For many types of dense and continuous data, metric distance measures (e.g., Euclidian distance) are used
- Scale issues need to be dealt with using normalization and weighting attributes
 For sparse data, which often consist of asymmetric attributes similarity measures that ignore 0-0 matches are used. Cosine, Jaccard are appropriate for
- measures that ignore 0-0 matches are used. Cosine, Jaccard are appropriate for such data
- Invariance to scaling (multiplication), translation (addition) makes cosine, Euclidian, and correlation useful for sparse data
- Correlation works better for time series where both scaling and translation are
- important
 In some cases, normalization and transformation of data may be needed (for
- periodicity data)
- In summary, proper choice of a proximity measure can be a time consuming task requiring both domain knowledge and the purpose for which the measure is being used!



© 2012 Microsoft Corporation. All rights reserved. Microsoft, Windows, Windows Vista and other product names are or may be registered trademarks and/or trademarks in the U.S. and/or other countries. The information herein is for informational purposes only and represents the current view of Microsoft Corporation as of the date of this presentation. Because Microsoft must respond to changing market conditions, it should not be interpreted to be a commitment on the part of Microsoft, and Microsoft comparison that et et accuracy of any information provided after the date of this presentation. MICROSOFT MAKES NO WARRANTIES, EXPRESS, IMPLIED OR STATUTORY, AS TO THE INFORMATION IN THIS PRESENTATION.

*

