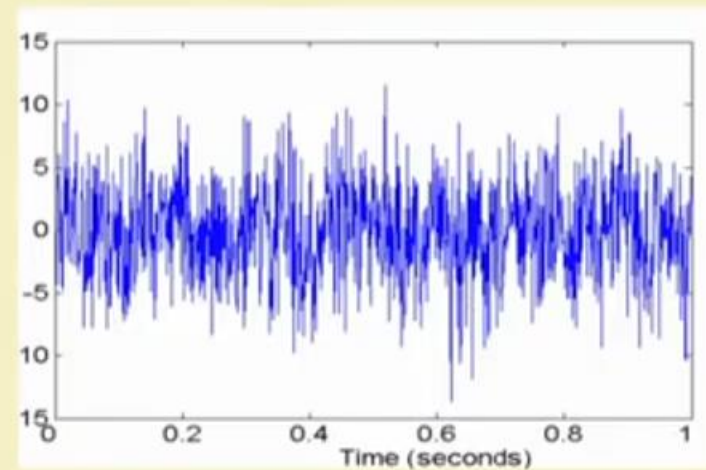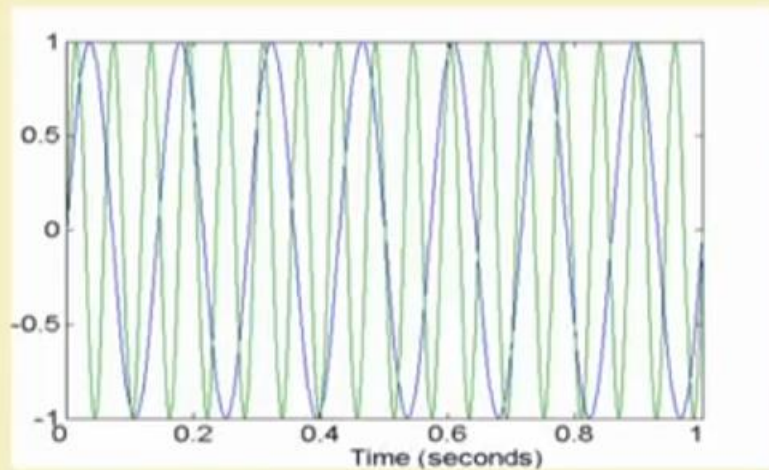# Data Quality

- What kinds of data quality problems?
- How can we detect problems with the data?
- What can we do about these problems?
- Examples of data quality problems:
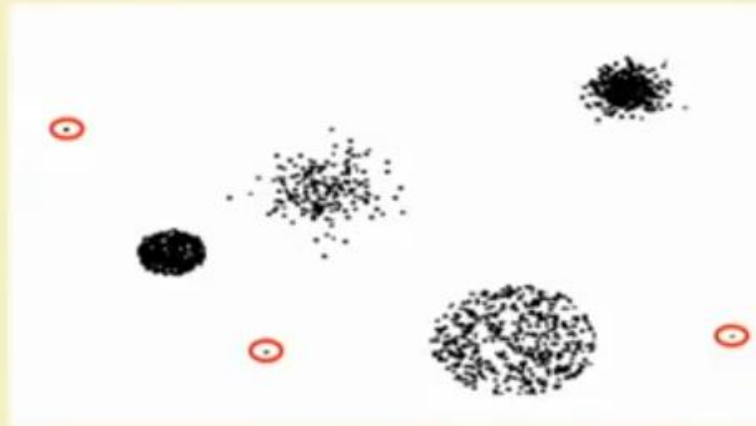  - Noise and outliers
  - missing values
  - duplicate data

# Noise

- Noise refers to modification of original values
  - Examples: distortion of a person's voice when talking on a poor phone and "snow" on television screen

# Outliers

- Outliers are data objects with characteristics that are considerably different than most of the other data objects in the data set

# Missing Values

- Reasons for missing values
  - Information is not collected
    (e.g., people decline to give their age and weight)
  - Attributes may not be applicable to all cases
    (e.g., annual income is not applicable to children)

- Handling missing values
  - Eliminate Data Objects
  - Estimate Missing Values
  - Ignore the Missing Value During Analysis
  - Replace with all possible values (weighted by their probabilities)

# Duplicate Data

- Data set may include data objects that are duplicates, or almost duplicates of one another
  - Major issue when merging data from heterogenous sources

- Examples:
  - Same person with multiple email addresses

- Data cleaning
  - Process of dealing with duplicate data issues

# Data Preprocessing

- Aggregation
- Sampling
- Dimensionality Reduction
- Feature subset selection
- Feature creation
- Discretization and Binarization
- Attribute Transformation

# Aggregation

- Combining two or more attributes (or objects) into a single attribute (or object)

- Purpose
  - Data reduction
    - Reduce the number of attributes or objects
  - Change of scale
    - Cities aggregated into regions, states, countries, etc
  - More "stable" data
    - Aggregated data tends to have less variability

# Sampling

- Sampling is the main technique employed for data selection.
  - It is often used for both the preliminary investigation of the data and the final data analysis.

- Statisticians sample because **obtaining** the entire set of data of interest is too expensive or time consuming.

- Sampling is used in data mining because **processing** the entire set of data of interest is too expensive or time consuming.

# Sample Size



8000 points      2000 Points      500 Points

# Sampling ...

- The key principle for effective sampling is the following:
    - using a sample will work almost as well as using the entire data sets, if the sample is representative
    - A sample is representative if it has approximately the same property (of interest) as the original set of data

# Types of Sampling

- **Simple Random Sampling**
  - There is an equal probability of selecting any particular item

- **Sampling without replacement**
  - As each item is selected, it is removed from the population

- **Sampling with replacement**
  - Objects are not removed from the population as they are selected for the sample.
    - In sampling with replacement, the same object can be picked up more than once

- **Stratified sampling**
  - Split the data into several partitions; then draw random samples from each partition
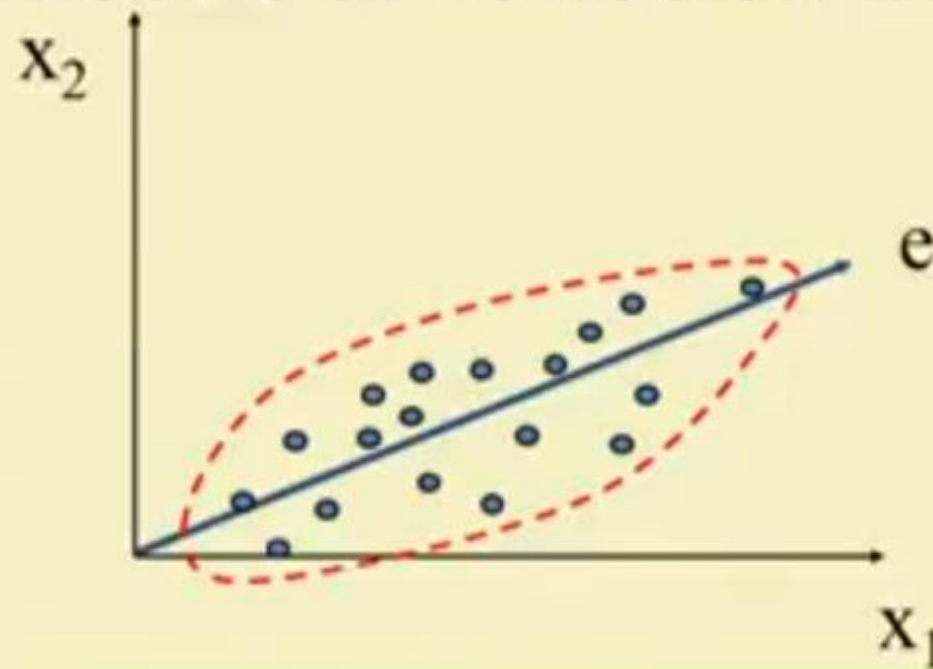
# Curse of Dimensionality

- When dimensionality increases, data becomes increasingly sparse in the space that it occupies

- Definitions of density and distance between points, which is critical for clustering and outlier detection, become less meaningful

# Dimensionality Reduction

- Purpose:
  - Avoid curse of dimensionality
  - Reduce amount of time and memory required by data mining algorithms
  - Allow data to be more easily visualized
  - May help to eliminate irrelevant features or reduce noise

- Techniques
  - Principle Component Analysis
  - Singular Value Decomposition
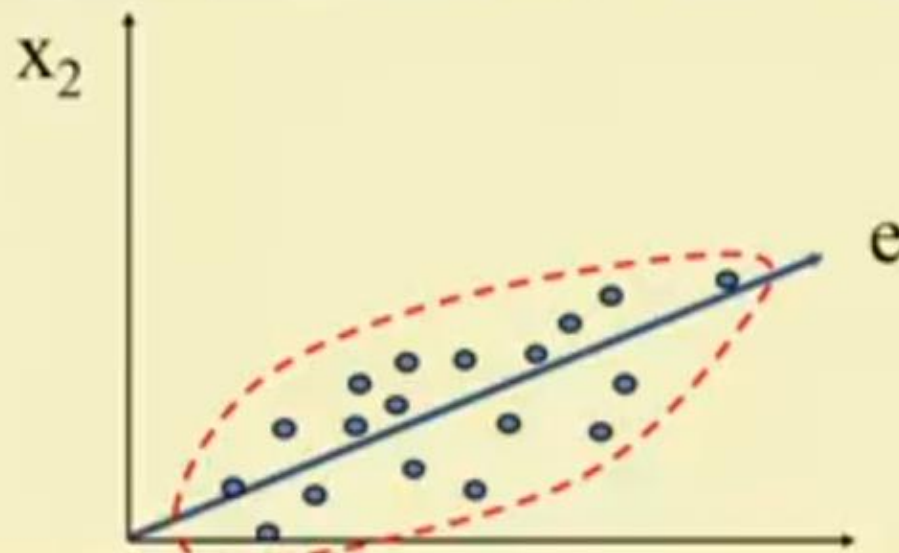  - Others: supervised and non-linear techniques

# Dimensionality Reduction: PCA

- Goal is to find a projection that captures the largest  amount of variation in data

# Dimensionality Reduction: PCA

- Find the eigenvectors of the covariance matrix
- The eigenvectors define the new space

# Feature Subset Selection

- Another way to reduce dimensionality of data

- Redundant features
  - duplicate much or all of the information contained in one or more other attributes
  - Example: purchase price of a product and the amount of sales tax paid

- Irrelevant features
  - contain no information that is useful for the data mining task at hand
  - Example: students' ID is often irrelevant to the task of predicting students' GPA
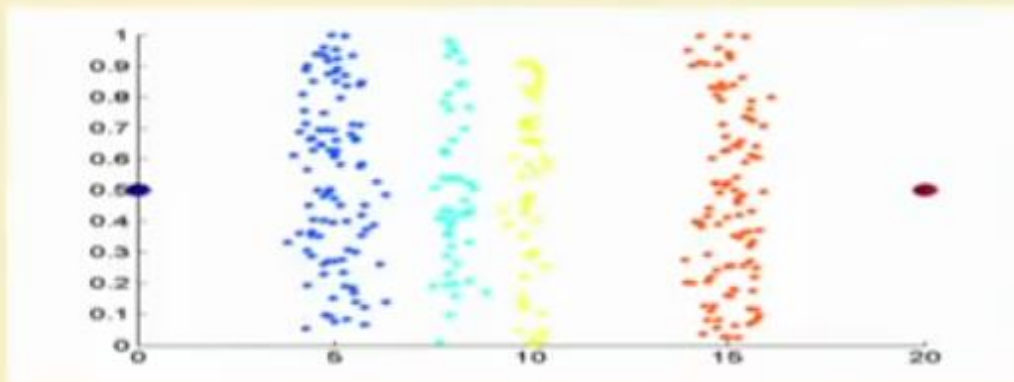
# Feature Subset Selection

- Techniques:
  - Brute-force approch:
    - Try all possible feature subsets as input to data mining algorithm
  - Embedded approaches:
    - Feature selection occurs naturally as part of the data mining algorithm
  - Filter approaches:
    - Features are selected before data mining algorithm is run
  - Wrapper approaches:
    - Use the data mining algorithm as a black box to find best subset of attributes
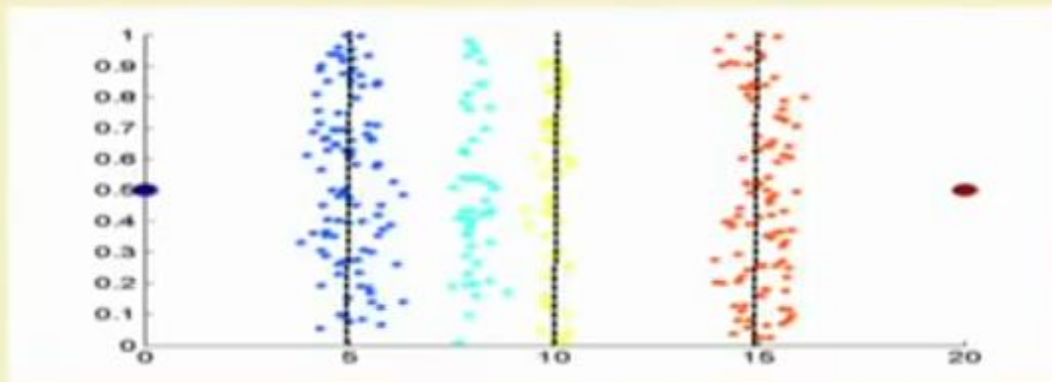
# Feature Creation

- Create new attributes that can capture the important information in a data set much more efficiently than the original attributes

- Three general methodologies:
  - Feature Extraction
    - domain-specific
  - Mapping Data to New Space
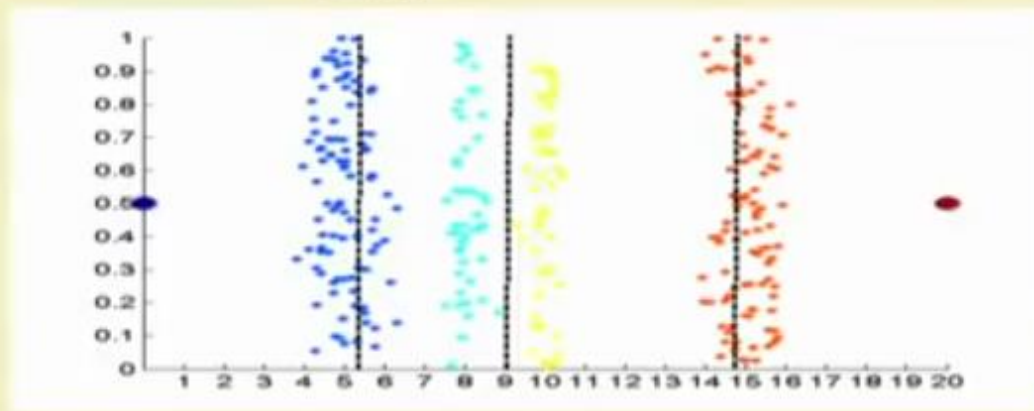  - Feature Construction
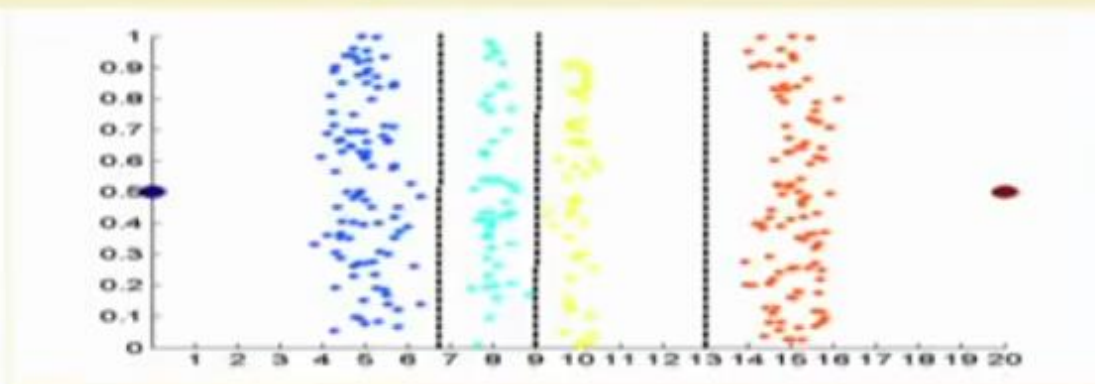    - combining features

# Discretization



Data

Equal interval width

Equal frequency

K-means

# Attribute Transformation

- A function that maps the entire set of values of a given attribute to a new set of replacement values such that each old value can be identified with one of the new values
  - Simple functions: $x^k$, $\log(x)$, $e^x$, $|x|$
  - Standardization and Normalization

# Similarity and Dissimilarity

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

# Similarity/Dissimilarity for Simple Attributes

$p$ and $q$ are the attribute values for two data objects.

| Attribute Type | Dissimilarity | Similarity |
|---|---|---|
| Nominal | $d = \begin{cases} 0 & \text{if } p = q \\ 1 & \text{if } p \neq q \end{cases}$ | $s = \begin{cases} 1 & \text{if } p = q \\ 0 & \text{if } p \neq q \end{cases}$ |
| Ordinal | $d = \frac{\lvert p - q \rvert}{n-1}$ (values mapped to integers $0$ to $n-1$, where $n$ is the number of values) | $s = 1 - \frac{\lvert p - q \rvert}{n-1}$ |
| Interval or Ratio | $d = \lvert p - q \rvert$ | $s = -d$, $s = \frac{1}{1+d}$ or $s = 1 - \frac{d - min\_d}{max\_d - min\_d}$ |

**Table** Similarity and dissimilarity for simple attributes
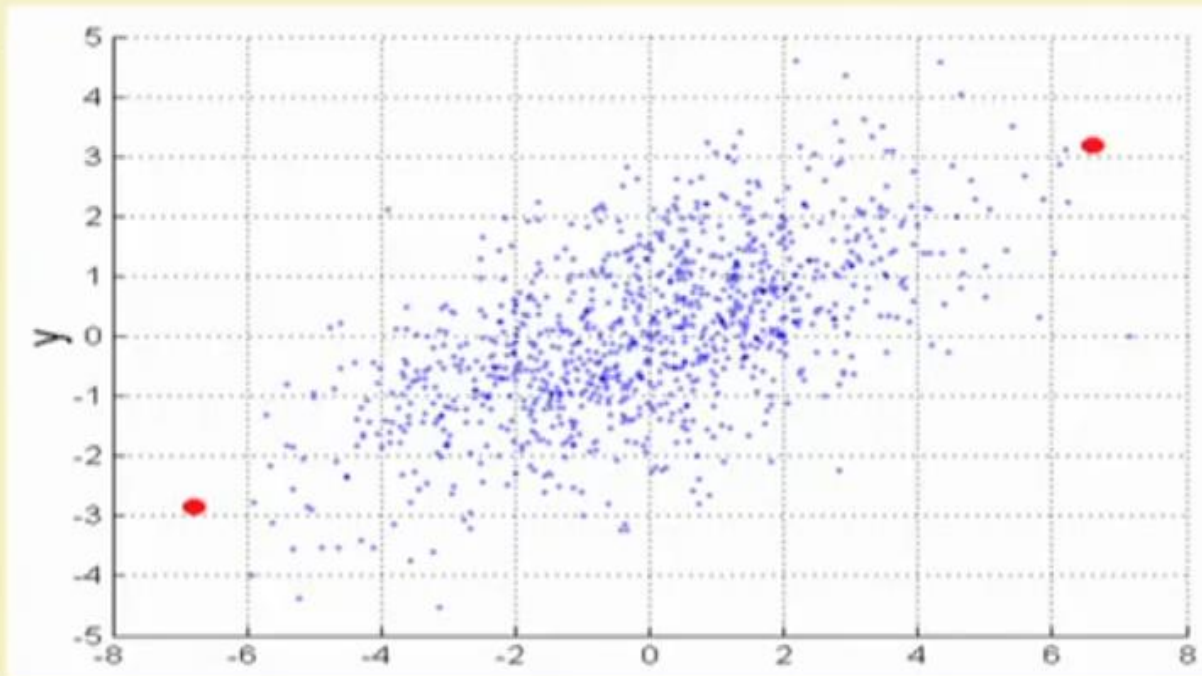
# Euclidean Distance

- 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 $k^{th}$ attributes (components) or data objects $p$ and $q$.

- Standardization is necessary, if scales differ.

# Mahalanobis Distance

$$mahalanobis(p, q) = (p - q) \Sigma^{-1} (p - q)^{T}$$



$\Sigma$ is the covariance matrix of the input data $X$

$$\Sigma_{j,k} = \frac{1}{n-1} \sum_{i=1}^{n} (X_{ij} - \overline{X}_j)(X_{ik} - \overline{X}_k)$$

# Cosine Similarity

- If $d_1$ and $d_2$ are two document vectors, then
$$\cos(d_1, d_2) = (d_1 \bullet d_2) / ||d_1||\ ||d_2||\ ,$$
  where $\bullet$ indicates vector dot product and $|| d ||$ is the length of vector $d$.
- Example:

$$d_1 = 3\ 2\ 0\ 5\ 0\ 0\ 0\ 2\ 0\ 0$$
$$d_2 = 1\ 0\ 0\ 0\ 0\ 0\ 0\ 1\ 0\ 2$$

$d_1 \bullet d_2 = 3^*1 + 2^*0 + 0^*0 + 5^*0 + 0^*0 + 0^*0 + 0^*0 + 2^*1 + 0^*0 + 0^*2 = 5$

$||d_1|| = (3^*3+2^*2+0^*0+5^*5+0^*0+0^*0+0^*0+2^*2+0^*0+0^*0)^{0.5} = (42)^{0.5} = 6.481$

$||d_2|| = (1^*1+0^*0+0^*0+0^*0+0^*0+0^*0+0^*0+1^*1+0^*0+2^*2)^{0.5} = (6)^{0.5} = 2.245$

$$\cos(d_1, d_2) = .3150$$

# Similarity Between Binary Vectors

- Common situation is that objects, *p* and *q*, have only binary attributes

- Compute similarities using the following quantities
  $M_{01}$ = the number of attributes where p was 0 and q was 1
  $M_{10}$ = the number of attributes where p was 1 and q was 0
  $M_{00}$ = the number of attributes where p was 0 and q was 0
  $M_{11}$ = the number of attributes where p was 1 and q was 1

- Simple Matching and Jaccard Coefficients
  SMC = number of matches / number of attributes
  $= (M_{11} + M_{00}) / (M_{01} + M_{10} + M_{11} + M_{00})$

  J = number of 11 matches / number of not-both-zero attributes values
  $= (M_{11}) / (M_{01} + M_{10} + M_{11})$

# Correlation
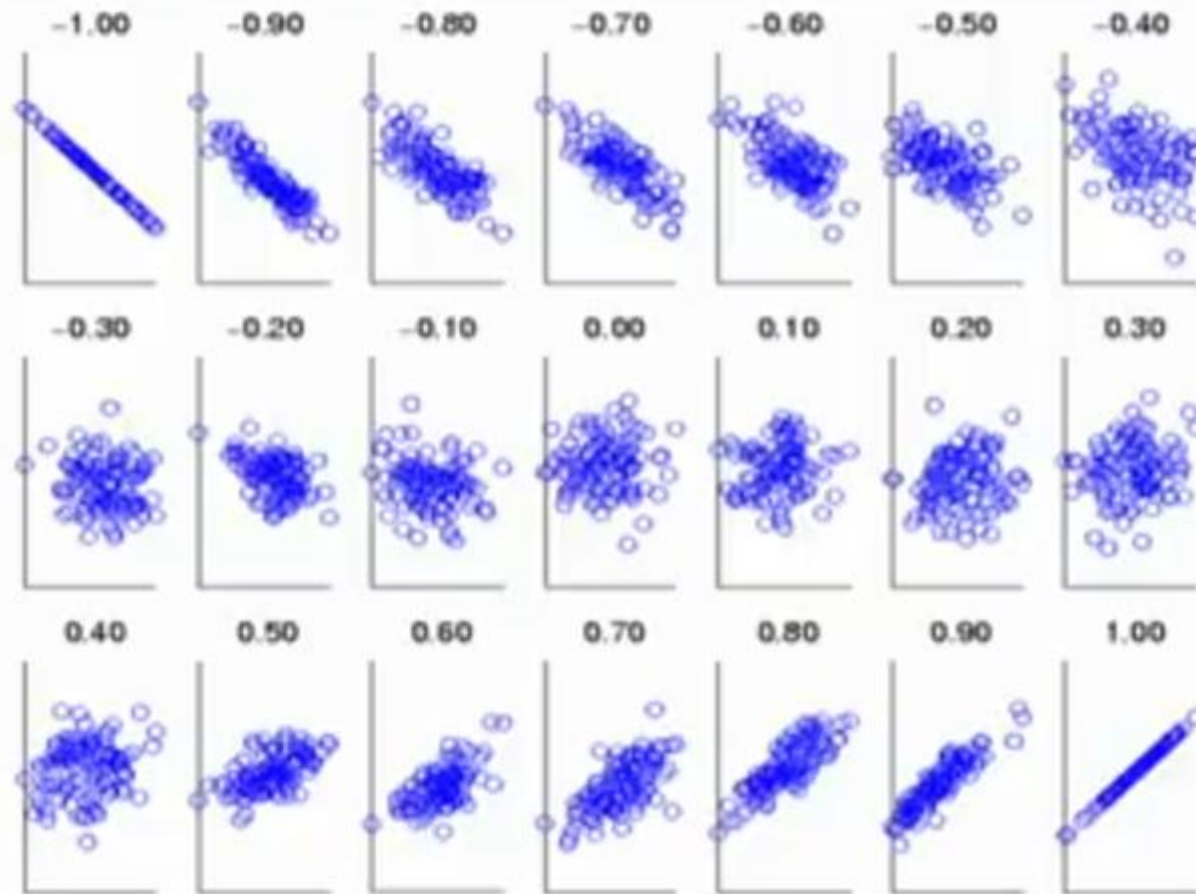
- Correlation measures the linear relationship between objects
- To compute correlation, we standardize data objects, p and q, and then take their dot product

$$p'_k = (p_k - mean\ (p)) / std\ (p)$$

$$q'_k = (q_k - mean\ (q)) / std\ (q)$$

$$correlation(p, q) = p' \bullet q'$$

# Visually Evaluating Correlation



Scatter plots showing the similarity from −1 to 1.