

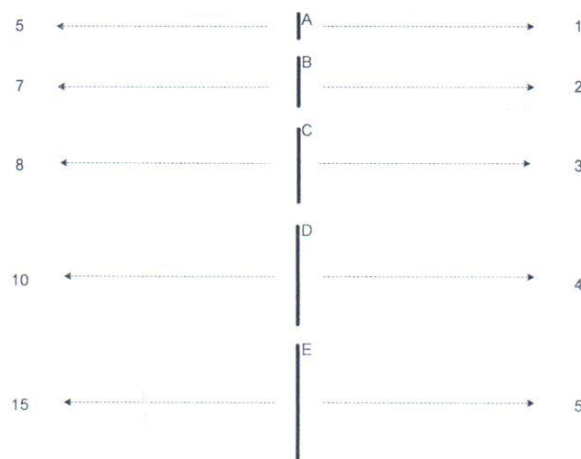
نسبت در درجه و نوع دارد. ما مشخص می‌شود. مثال زده مثلاً طول، ما می‌توانیم طول را اندازه گیری کنیم با مقیاس و واحد اندازه گیری می‌توانیم مشخص کنیم ولی اگر بخواهیم به صورت رسمی این Attribute صفاً صفت طرا مشخص کنیم ما می‌توانیم ۴ مانوع صفت را اندازه گیری کنیم، که می‌توانیم به آن ها بگوییم: ۴ مانوع مقیاس اندازه گیری داریم ۱ - Nominal - ۲ Ordinal - ۳ Interval - ۴ Ratio

Attribute Values

- Attribute values are numbers or symbols assigned to an attribute
- Distinction between attributes and attribute values
 - Same attribute can be mapped to different attribute values
 - ◆ Example: height can be measured in feet or meters
 - Different attributes can be mapped to the same set of values
 - ◆ Example: Attribute values for ID and age are integers
 - ◆ But properties of attribute values can be different
 - ID has no limit but age has a maximum and minimum value

Measurement of Length

- The way you measure an attribute is somewhat may not match the attributes properties.



مثال اسعی = مثلاً گروه خونی هیچ برتری نسبت به هم ندارند، یک اسم است به یک گروه خونی نسبت می دهیم به خاطر همین بهش میگویند که Nominal است.

مثال ترتیبی = مثلاً میگویم رتبه اول، دوم، سوم، رئیس، معاون، کارمند / ... اینها یک ترتیبی دارند بر این برقرار است که با آنها میگویم داده های Ordinal.

Types of Attributes

● There are different types of attributes

اسعی - Nominal

- Examples: ID numbers, eye color, zip codes

ترتیبی - Ordinal

- Examples: rankings (e.g., taste of potato chips on a scale from 1-10), grades, height in {tall, medium, short}

فاصله ای - Interval

- Examples: calendar dates, temperatures in Celsius or Fahrenheit.

نسبی - Ratio

- Examples: temperature in Kelvin, length, time, counts

نسبی Flexible ترین نوع است. آویز می خواهد یک صفت را اندازه گیری کنیم، اگر بتوانیم برای هر مقیاس

نسبی اندازه گیری اش کنیم خیلی بهتر است مثلاً: درم حرات برای س کلوین میگویند سلسیوس کم Ratio نسبت و س کلوین این چنین است باز قان و زمان نسبت است، طول هم همین طور. اکثر اقدار مقیاسی که ما با آنها کار می داریم سلسیوس، درم، سلسیوس هستند. خوبی نسبتی ها این است که خیلی از اعمال ریاضی بر روی آن ها انجام می شود. مقیاس های نسبی صحت در بین هستند، فقط مقیاسی برای بگویم که این دو تا یکی هستند یا خیر.

Properties of Attribute Values

● The type of an attribute depends on which of the following properties it possesses:

- Distinctness: $= \neq$
- Order: $< >$
- Addition: $+ -$
- Multiplication: $* /$

- Nominal attribute: distinctness
- Ordinal attribute: distinctness & order
- Interval attribute: distinctness, order & addition
- Ratio attribute: all 4 properties

در مقیاس ترتیبی میتوانیم کوچکتر یا بزرگتر بودن را تعیین کنیم، فاصله ای میتوانیم با هم جمع یا تفریق کنیم و نسبتی تمام کارهایی که میتوانیم برای انجام داده ها داریم برقرار است. سوال: اگر داده ها Nominal باشند میتوان بر روی آن ها کاری انجام داد؟ آری در اینجا نه، ماضی از کارهایی را که با فاصله نسبتی اندازه گیری میکنند میتوانیم آن کارها را روی داده های Nominal انجام دهیم مثلاً ما وقتی یک داده بیان می دهیم می توانیم حدس بزنیم که این داده ای است که رتبه بین کارهایی را دارد 1 تمام داده ها با انواع مختلف

رابطه توان مدار برای محاسبه کرد. entropy را می توان برای این محاسبه کرد، می توان آنهارا در یک جدول توافقی، جدول
 contingency (contingency) خلاصه کنیم و correlation آن را یا رابطه آنها با هم را با یک بسنجیم و این خیلی خوب
 است پس داده های Nominal در یک جدول داده های نیستند، نمودار و نشان کاری انجام داده خیلی از کارهایی که ما روی داده های
 نسبت انجام می دهیم می توانیم روی داده های Nominal نیز انجام دهیم.

Attribute Type	Description	Examples	Operations
Nominal	The values of a nominal attribute are just different names, i.e., nominal attributes provide only enough information to distinguish one object from another. ($=$, \neq)	zip codes, employee ID numbers, eye color, sex: {male, female}	mode, entropy, contingency correlation, χ^2 test
Ordinal	The values of an ordinal attribute provide enough information to order objects. ($<$, $>$)	hardness of minerals, {good, better, best}, grades, street numbers	median, percentiles, rank correlation, run tests, sign tests
Interval	For interval attributes, the differences between values are meaningful, i.e., a unit of measurement exists. (+, -)	calendar dates, temperature in Celsius or Fahrenheit	mean, standard deviation, Pearson's correlation, t and F tests
Ratio	For ratio variables, both differences and ratios are meaningful. (*, /)	temperature in Kelvin, monetary quantities, counts, age, mass, length, electrical current	geometric mean, harmonic mean, percent variation

entropy یک معیاری است که می آید تسام و عدم تسام یا وابستگی را به نحوی که ما فواید خیلی مفهومی در نظر به اطلاعات است.
 می توان از entropy در clustering به عنوان یک فاصله در نظر گرفت، با هاش می آید داده ها را خوشه بندی می کنیم برای داده های Nominal.
 در Ordinal علاوه بر دو می توانیم median حساب کنیم (ده دانه، ده دانه دوم و ... چهارم و ... هفتم و ...) میان بدنی داده ها را

Attribute Level	Transformation	Comments
Nominal	Any permutation of values	If all employee ID numbers were reassigned, would it make any difference?
Ordinal	An order preserving change of values, i.e., $new_value = f(old_value)$ where f is a monotonic function.	An attribute encompassing the notion of good, better best can be represented equally well by the values {1, 2, 3} or by {0.5, 1, 10}.
Interval	$new_value = a * old_value + b$ where a and b are constants	Thus, the Fahrenheit and Celsius temperature scales differ in terms of where their zero value is and the size of a unit (degree).
Ratio	$new_value = a * old_value$	Length can be measured in meters or feet.

به دو دسته تقسیم کنیم، در مورد
 Nominal ساینز نه هفت median
 حساب کرد. انا برای Ordinal
 هفت میس. همچنین می توانیم
 رتبه بندی داده ها را با رتبه Rank
 اصلی محاسبه کرد. همچنین می توانیم
 در مورد قرار گرفتن داده ها را
 اطلاعاتی کسب کنیم به روش Run tests
 این را مقایسه می کنیم که در آن کار به عمل می آید
 برای داده های Ordinal
 صفت در مورد نشان صفت کرد. (۱۶۷)

حالا ما اگر ترتیب داشته باشیم، ما می توانیم اینها را به عنوان یک دنباله sequence در نظر بگیریم و برای آن این اولین آزمون برآورد است، داده های که
 سبب می شود هم هستند یا نیستند اینها با اعداد صفر و یک می نویسند به هم، به صورت یک sequence binary درش می آوریم، این sequence binary
 را می توان تست کرد که تعداد دفعی است یا نه.
 در Interval می توانیم mean بگیریم، standard deviation بگیریم، اینها دیگر می توانیم Pearson's correlation محاسبه کنیم به جایی ...

و می‌توانیم بصورت آماری خیلی از ست‌های مثل \sum و σ را برای به کار ببریم.
 برای داده‌های فشرده می‌توانیم با بیان را از این مباحثی که گفتیم فزاینده‌تر کنیم به جای این می‌توانیم حساب کنیم، geometric mean
 می‌گیریم (مثلاً وقتی می‌خواهیم از سرعت دو ماشین میانگین بگیریم، فاصله میانگین معمولی گرفته، geometric mean بدرد می‌خورد).
 harmonic mean می‌گیریم همین‌طور... پس دست‌ان خیلی بازتر است. دقت کنید که برای ratio تمام کارهای بالا از سر خود انجام ندهید.

Discrete and Continuous Attributes

- Discrete Attribute
 - Has only a finite or countably infinite set of values
 - Examples: zip codes, counts, or the set of words in a collection of documents
 - Often represented as integer variables.
 - Note: binary attributes are a special case of discrete attributes
- Continuous Attribute
 - Has real numbers as attribute values
 - Examples: temperature, height, or weight.
 - Practically, real values can only be measured and represented using a finite number of digits.
 - Continuous attributes are typically represented as floating-point variables.

© Tan, Steinbach, Kumar Introduction to Data Mining 4/18/2004 9

* ما می‌توانیم روی داده‌هایی که با
 مقیاس‌های مختلف اندازه‌گیری شده‌اند یک تبدیل بگیریم مثلاً برای داده‌هایی
 اسکی یک جایگاه‌نشانی در نظر بگیریم به عنوان یک
 تبدیل و برای داده‌هایی
 ترتیبی، می‌توانیم یک تایی بودن یا کار می‌ریم که آن تابع
 لینک‌ها باشد مثلاً با اصول آنزولی.
 برای Interval می‌توانیم یک تبدیل پیدا کنیم به هم
 برای ratio می‌توانیم بخش
 کنیم و ضرب کنیم.
 * این صفاتی مگر ما گفتیم دو بخش تقسیم می‌شوند Discrete و Continuous.

* پس صفات ۳ نوع شدند اما می‌توان آنها را به دو بخش عمده تقسیم کرد: Discrete و Continuous
 * معمولاً قضاوتی نمی‌توانیم بکنیم Discrete و Continuous به ندرت می‌توانیم پیدا کرد و در مورد نسبتی بین continuous و Discrete.

Types of data sets

- Record
 - Data Matrix
 - Document Data
 - Transaction Data
- Graph
 - World Wide Web
 - Molecular Structures
- Ordered
 - Spatial Data
 - Temporal Data
 - Sequential Data
 - Genetic Sequence Data

© Tan, Steinbach, Kumar Introduction to Data Mining 4/18/2004 10

* Discrete یعنی چه؟ یعنی مقدار شمارا (اعداد طبیعی D است) یا این
 را آنقدر براب می‌کنیم تا به عدد ۱ برسیم، ممکن است همان دفعه اول یا بارهایش باشد.
 اگر یک‌بار هم فضای نفوذ همچنین آن را به ۱ می‌رسانیم تا به نهایت ادامه دارد اما تعداد اعضایش
 شمارا است، پس یک مجموعه شمارا است. پس Discrete گسترده (مقدور است)
 این است که تعداد اعضای مجموعه را می‌توان با انگشتان دست و یا بشماریم یا تعداد اعضایش
 اعضایش تا به نهایت ادامه داشته‌اند اما تا به نهایت نمی‌توان یک نقطه‌ای به دست آورد
 مجموعه اعداد طبیعی پرکاربرد. صفاتی به صورت خیلی زیاد هستند ۸ وزن، مقدار درجه حرارت
 اینها با یک مقیاس فیزیکی اندازه‌گیری می‌کنیم و مقیاس‌های پیوسته این‌ها مقدار پیوسته
 است. در واقع Continuous Attribute یک Superset از اعداد حقیقی است.

* نوع‌های داده به چه شکلی هستند؟ می‌توانند به شکل Record باشند یعنی ثبت شده، می‌توانند بصورت گراف باشند و Network
 که در وب است، انواع و اقسام گراف‌هایی را می‌توانیم به صورت Data در نظر بگیریم همچنین داده‌ها می‌توانند بصورت ordered
 باشند یعنی ترتیبی.
 Sequential data مثل داده‌های DNA که داده‌های ژنتیکی به صورت یک دنباله می‌باشند.

* این خصوصیت هایی که در داده ها وجود دارد خیلی مهم است، یعنی داده ها را که ما جمع آوریم می توانیم به آن اشاره کنیم. یعنی رکورد داده ها است، چون اگر بعد داده ها خیلی زیاد باشد ما با کامپیوتر هیچ کاری نمی توانیم بکنیم و آن ها انجام دهیم. مثلاً فرض کنید PCA می تواند انجام دهد روی بردار داده هایی که دهه اند می بینیم، اصلاً با کامپیوتر نمی توانیم در همین یک ماتریس 1000×1000 را به صورت Real Time حساب کرده و عملاً با اینجایا مشکل مواجه می شود. با بعد داده ها با وجود اینکه

Important Characteristics of Structured Data

1 - Dimensionality

- ◆ Curse of Dimensionality

داریم، یکی این است از نظر ریاضی ما مشکلی نداریم ولی از نظر عملیاتی مشکل داریم. مثلاً PCA در پس این است که تعداد ابعاد را کم می کند.

2 - Sparsity

- ◆ Only presence counts

یعنی اینست داده های ما داده های اسپارسی باشند، داده های پراکنده ای باشند. اینقدر دراز هم باشند که ما نتوانیم در مورد آنها هیچ تصمیمی بگیریم، هیچ ایده ای در هیچ ردیفی برای این برایش داریم. برای علم بر اسپارسی دوباره روش های زیادی است.

3 - Resolution

- ◆ Patterns depend on the scale

یعنی اینست داده های ما Resolution آن خوب باشد، Resolution تغییر پیدا می کند، یعنی مواقع کارهای ما را خوب می کند، اطلاعات بر روی ماس (مقیاس) دارد و اگر جمع آوری کردن با وقت خوبی جمع آورده شده ما نتوانیم عملاً از آنها استفاده کنیم.

* وقتی که ما میگوییم Record of Data منظوریان این است که داده های ما درون فایل مثل Excel یا برای SQL ضبط شده که معمولاً هر ردیف نام تغییراتی ما است و بعد اطلاعاتی در مورد حالات و غیره این.

Record Data

- Data that consists of a collection of records, each of which consists of a fixed set of attributes

Tid	Refund	Marital Status	Taxable Income	Cheat
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

* ممکن است داده ما همان به صورت ماتریس باشد یعنی از یک بعدی باشد، ممکن است از دو بعد هم باشد.

Data Matrix

- If data objects have the same fixed set of numeric attributes, then the data objects can be thought of as points in a multi-dimensional space, where each dimension represents a distinct attribute
- Such data set can be represented by an m by n matrix, where there are m rows, one for each object, and n columns, one for each attribute

Projection of x Load	Projection of y load	Distance	Load	Thickness
10.23	5.27	15.22	2.7	1.2
12.65	6.25	16.22	2.2	1.1

* این مثال همان Document است که برای آن گفته

Document Data

- Each document becomes a 'term' vector,
 - each term is a component (attribute) of the vector,
 - the value of each component is the number of times the corresponding term occurs in the document.

	team	coach	play	ball	score	game	win	lost	timeout	season
Document 1	3	0	5	0	2	6	0	2	0	2
Document 2	0	7	0	2	1	0	0	3	0	0
Document 3	0	1	0	0	1	2	2	0	3	0

* اسلاید ۱۵: داده های Transaction است که این داده ها را می توان کتابخانه ها استفاده می کنند.

Transaction Data

- A special type of record data, where
 - each record (transaction) involves a set of items.
 - For example, consider a grocery store. The set of products purchased by a customer during one shopping trip constitute a transaction, while the individual products that were purchased are the items.

Transaction

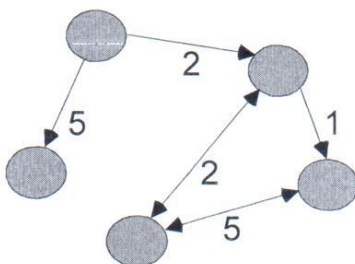
TID	Items
1	Bread, Coke, Milk
2	Beer, Bread
3	Beer, Coke, Diaper, Milk
4	Beer, Bread, Diaper, Milk
5	Coke, Diaper, Milk

نوعی از داده (فهرست خرید) می باشد
خرید ۱

Graph Data

داده های گراف

- Examples: Generic graph and HTML Links



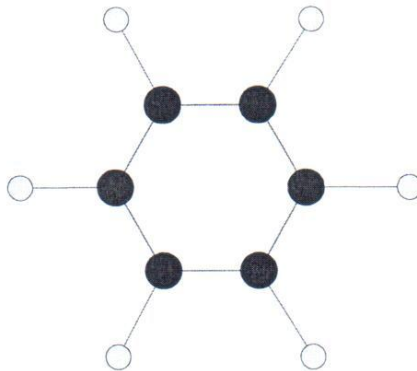
```

<a href="papers/papers.html#bbbb">
Data Mining </a>
<li>
<a href="papers/papers.html#aaaa">
Graph Partitioning </a>
<li>
<a href="papers/papers.html#aaaa">
Parallel Solution of Sparse Linear System of Equations </a>
<li>
<a href="papers/papers.html#ffff">
N-Body Computation and Dense Linear System Solvers
  
```


Chemical Data

جاده ها معین است chemical با سبز داده های (شیمی)

- Benzene Molecule: C_6H_6

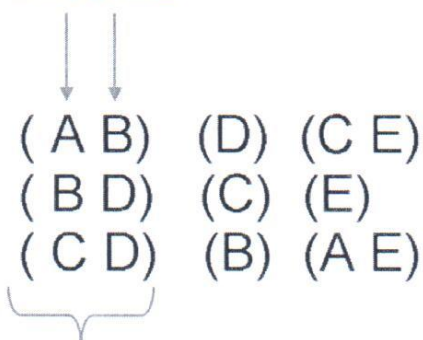


Ordered Data

جاده ها مرتب با سبز

- Sequences of transactions

Items/Events



An element of
the sequence

Ordered Data

● Genomic sequence data

GGTTCCGCCTTCAGCCCCGCGCC
CGCAGGGCCCCGCCCCGCGCCGTC
GAGAAGGGCCCGCCTGGCGGGCG
GGGGGAGGCGGGGCCGCCCGAGC
CCAACCGAGTCCGACCAGGTGCC
CCCTCTGCTCGGCCTAGACCTGA
GCTCATTAGGCGGCAGCGGACAG
GCCAAGTAGAACACGCGAAGCGC
TGGGCTGCCTGCTGCGACCAGGG

DNA
این Sequence است.

این داده‌های ما است می‌خواهیم

آینده را Predict کنیم.

یعنی آمار این ژن‌ها بیا بیندازیم

چه اتفاقی خواهد افتاد و چه خواهیم

داشت؟!

* داده‌های مکانی =

* زمانی = می‌خواهیم در زمان‌های مختلف گزینش کردن ژن‌ها را بررسی کنیم.

Ordered Data

● Spatio-Temporal Data

Average Monthly
Temperature of
land and ocean



اسلام 25، ما وقتی داده‌ها را جمع کرده می‌بینیم، داده‌های تکراری نبولیم این Duplicate & Redundancy است، اینها کیفیت داده‌ها را پایین می‌آورند. مثلاً یک نفر ممکن است در یک جای چند بار نامش می‌آید در این سایت‌های اینترنتی. این باعث می‌شود داده‌ها مار orip کند و ایراد دارد.

Duplicate Data

- Data set may include data objects that are duplicates, or almost duplicates of one another
 - Major issue when merging data from heterogeneous sources
- Examples:
 - Same person with multiple email addresses
- Data cleaning
 - Process of dealing with duplicate data issues

اسلام 26 = Data Preprocessing چند بخش دارد که در زیر آمده است. 92, 8, 7, 6, 5, 4, 3, 2, 1

Data Preprocessing

- 1- ● Aggregation
- 2- ● Sampling
- 3- ● Dimensionality Reduction
- 4- ● Feature subset selection
- 5- ● Feature creation
- 6- ● Discretization and Binarization
- 7- ● Attribute Transformation

Aggregation همان طور که از اسمش مشخص است، هدف از آن Data Reduction است یا تغییر اسکیل (Scale) است. یا پایداری یا ثبات بیشتر در داده‌ها است. در واقع Aggregation کاری که انجام می‌دهیم آنست که داده‌ها را از این Attribute ها را با هم جمع می‌کنیم مثلاً: ما می‌خواهیم میزان برق مصرفی را در نردگانه‌ها پیش‌بینی کنیم، لیکن راه این است که اطلاعاتی که به ما می‌دهند بگونه است؟! مصرف برق خانگی، مصرف برق معابر، مصرف برق صنعتی، چندتا از این معیار برق را ما می‌دهند. آنرا اینجا را با هم جمع می‌کنیم.