# Column 1

Regression: modeling ~~the~~ and analysis

Histogram: Divide data into buckets and store average (sum) for each bucket.

PCA: a statistical procedure ~~that~~ uses an orthogonal transformation to convert a set of observations of possibly correlated variables into a set of values of linearly uncorrelated variables called principal components

Binning: 19, 71, 48, 63, 35, 85, 69, 81, 72, 88, 99, 35

equal-f: bin1: 19, 35, 48
bin2: 63, 69, 71

equal-width: width = (99-19)/4 = 20
bin1 = [19,39) 19, 35
bin2 = [39,59) 48

Important Characteristics:
- Dimensionality
- Sparsity — KL: the expected number of extra bits required to code $q(x)$ to $p(x)$ from
- Resolution
- Distribution

Data Warehouse usage:
- information processing - Statistical
- Analytical processing - dimensional
- Data Mining -

shell-fragment:

$$\frac{4}{\_} \mid 6 \qquad 2^6 x$$

4:a = 1
2:b = $\binom{4}{2} \times 2^2 = 24$
3:b = $\binom{4}{3} \times 2 = 8$
4:b = 1

Sup=2
Closed cell $\times 2^6$
$\binom{4}{3} + \binom{4}{3} + \binom{4}{4} = 11$

Mean: $\bar{x} = \frac{1}{N}\sum x_i$, $\mu = \frac{\sum x}{N}$

Median: $L_1 + \left(\frac{N/2 - \sum f}{\sum f}\right)$ width

Mode: mean - mode = 3×(mean - median) unimodal

Var: $s^2 = \frac{1}{n-1}\sum(x-\bar{x})^2 = \frac{1}{n-1}[\sum x^2 - \frac{1}{n}(\sum x)^2]$

$\sigma^2 = \frac{1}{N}\sum(x-\mu)^2 = \frac{1}{N}\sum x^2 - \mu^2$

Data visualization:
- Pixel-oriented
- Geometric: · Direct Visualization
- Scatterplot & scatterplot matrices ($k^2 \div k$)
- Landscape · Projection · Hyperslice
- Parallel coordinates
- Icon-Based: · chernoff Faces
  → Shape coding ( stick Figures
  → color icons → Tile bars
- Hierarchical: · Dimensional stacking
- World-within-World · Tree-Map 2D fine cannot avoid overlap
- Cone Trees · Info Cube
- Tag cloud: font size/color

$Z = \frac{x-\mu}{\sigma}$, $z_{if} = \frac{x_{if} - m_f}{s_f} \frac{1}{n}(|x-m|+\dots)$ robust

Minkowski distance:
$d(i,j) = \sqrt[p]{|x_{i1}-x_{j1}|^p + |x_{i2}-x_{j2}|^p + \dots + |x_{il}-x_{jl}|^p}$
- $d(i,j) \ge 0$ if $i \neq j$ · $d(i,j) = d(j,i)$
- $d(i,j) \le d(i,k) + d(k,j)$
p=1 Manhattan p=2: Euclidean distance
h=3: Supremum → $\max_f^{l=1} |x_{if} - x_{jf}|$
maximum diff between 2 vectors

# Column 2

Proximity Measure for Binary Attributes
categorical object j

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

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

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

Jaccard: $sim_{Ja}(i,j) = \frac{q}{q+r+s}$ → medical test

Proximity Measure for Categorical:
- Simple Matching: $d(i,j) = \frac{p-m}{p}$
- Use a large number of binary

Ordinal Variables:
$r_{if} \in \{1,\dots,M_f\}$  $z_{if} = \frac{r_{if}-1}{M_f-1}$

Mixed type:
$$d(i,j) = \frac{\sum_{f=1}^{p} w_{ij}^{(f)} d_{ij}^{(f)}}{\sum_{f=1}^{p} w_{if}^{(f)}}$$

- numeric: ~~no~~ normalized distance
- binary/nominal: 0, if $x_{if} = x_{jf}$, 1
- ordinal: rank $z_{if} = \frac{r_{if}-1}{M-1}$

Cosine Similarity: document term vectors in ~~text~~ file
$$\cos(d_1, d_2) = \frac{d_1 \cdot d_2}{\|d_1\| \times \|d_2\|}$$

## Chapter 3: Data Preprocessing

Data Quality: · accuracy
- completeness · Consistency
- Timeliness · Believability
- Interpretability

Incomplete (Missing) Data:
- inference-based: Bayesian formula / decision tree

Noisy Data:
- Binning: sort, partition
  → Smooth: mean, median, boundary
- Regression: regression function
- Clustering: detect & remove outliers
- Semi-supervised: Computer & human

Data Integration: Chi-square Correlation Analysis (categorical)
$$\chi^2 = \sum_{i}^{n} \frac{(O_i - E_i)^2}{E_i}$$ (Null: independent)

row × column = expected ($E_i$)
(#row -1) × (#column -1) = DOF

Covariance for two Variables
$\sigma_{12} = E[(X_1-\mu_1)(X_2-\mu_2)] = E[X_1 X_2] - \mu_1 \mu_2$
$= E[X_1 X_2] - E[X_1]E[X_2]$ [E∞ 00]
if $X_1, X_2$ are independent, $\sigma_{12} = 0$

Correlation between Numerical
$\rho_{12} = \frac{\sigma_{12}}{\sqrt{\sigma_1^2 \sigma_2^2}}$
$\rho_{12} > 0$ positive, $x_1$ as $x_2$
$\rho_{12} = 0$ independent (normal distri) x linear
$\rho_{12} < 0$ negative  $y = x^2$
$[-1, 1]$ · Normalized covariance

# Column 3

Covariance Matrix:
$\sum[(x-\mu)(x-\mu)^T] = \begin{bmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{21} & \sigma_2^2 \end{bmatrix}$ det $\geq 0$

## ✦ Numerosity Reduction:

Parametric: fit some model
Non-Parametric: no model

Regression Analysis: → predictors
dep (response) → indep (explanatory)
- Linear ~ · Nonlinear ~
- Multiple ~ · Log-linear ~

Histogram Analysis: bining→distri
Data → buckets → avg.
- equal-width · equal-frequency

Clustering: similarity, representation effective if clustered.

Sampling: small s represent whole data
Key principle: representative subset
types: · simple random sampling
- Sampling without replacement
- Sampling with replacement
- Stratified sampling

Data Cube Aggregation: The aggregation data for an individual entity of interest.

Data Compression: lossless vs. lossy

Wavelet transform: size ↓
- use · hat-shape filters
- effective removal of outliers
- Multi-resolution
- Efficient · $O(N)$ · only low D

Data Transformation: to a new set of representation data
map to new set. construction
- Smoothing · Attribute/feature ↓
- Aggregation · Normalization
- Discretization

Normalization:
min-max: $v' = \frac{v-\min}{\max-\min}(\text{max}_{new} - \text{min}_{new}) + \text{min}_{new}$
Z-score: $v' = \frac{v-\mu}{\sigma}$
decimal scaling: $v' = \frac{v}{10^j}(\text{Max}|v'| < 1)$

Data Discretization:
- ✦ Data Size ↓ · Supervised vs x
- Split (top-down) vs. merge (bot-up)
- Binning: split · unsup
  - equal-width $W = (B-A)/N$ uniform grid
    → outlier dominant · skewed x
  - equal-depth N interval, save size
    → good scaling, x categorical
- histogram: split · unsup
- clustering: split, merge, unsup
- Decision tree: split · sup (class into)
- correlation: merge, unsup↓

concept hierarchy generation:
- drilling & rolling ↓ Nominal

## ✦ Dimensionality Reduction:
- avoid curse · eliminate noise
- reduce space & time · easier visualization
Feature selection: subset
Feature extraction: transfer highD → lowD

Principal Component Analysis (PCA)
orthogonal transformation, eigenvector
→ Normalize → k ortho vectors → linear combi of k → eliminate weak component
For numeric only ↓ sorted →

# Column 4

- Attribute Subset Selection:
  - Redundant attributes
  - Irrelevant attributes
  - Heuristic Search: 2d
    → best single attri → best step-wise select
    → best step-wise eliminate → best combi
    → optimal branch & bound

- Attribute Creation (Feature Generation) new more effective set
  - Attributes extraction
  - Mapping data to new space
  - Attribute construction

Data Integration: · Entity identification · Remove redundant · Detect inconsistency

Data reduction: · Dimensionality · Numerosity · Data Compression

## Chapter 4: OLAP & data warehouse

- decision Support processing, separately from the organization's operational data
- Support information processing by providing a solid platform of consolidated, historical data for analysis

Data warehouse: · subject-oriented · integrated · time-variant · nonvolatile

subject-oriented: organized around major subjects · modeling & analysis → decision · simple & concise: excluding not useful.

Integrated: · integrating multiple, heterogenous data source · Data cleaning & data integration

Time Variant: · longer time horizon · operational: ~current time - no key time · warehouse: historical, key → time element

Nonvolatile: · independent → physical Separation store. · Static: no update of data ⟹ 1. initial loading 2. access of data

OLTP: Online transactional processing · Day to Day · application-oriented · repetitive · small size · more users

OLAP: online analytical processing · historical · complex · large

Why separate: · High performance · Different functions & data: missing data → data consolidation · data quality

Architecture: · Top tier: Front End Tools · Middle: OLAP · Bottom: Data Warehouse · Data

- Enterprise warehouse: subject, entire org. · Data Mart: specific group { indep, dep (direct fun) · Virtual warehouse: views, summary

Extraction, transformation, loading (ETL)
- Data extraction · Data cleaning
- Data transformation · load · refresh
Metadata: · description of structure · operational meta-data → data lineage → currency of data → monitoring info · algorithm · summary · mapping · business · data related to performance
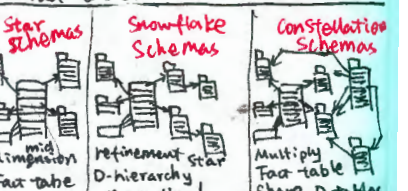
Data Cube: a lattice of cuboids
Base cuboid: n-D base cube
apex cuboid: 0-D highest summary

Data Warehouse: multidimensional · Dimension tables: item, time · Fact table: measures, key

| Star Schemas | Snowflake Schemas | Constellation Schemas |
|---|---|---|
|  |  |  |
| mid dimension Fact table | refinement star D-hierarchy normalized | Multiply Fact table share D tables |

## Data Cube Measure: 3 Categories:

- **Distributive:** it the result derived by applying the function n aggregate values is the same as that derived by applying the function on all the data without partition
- **Algebraic:** it can be computed by any algebraic function with M arguments, each of which is obtained by applying a distributive aggregate
  - $avg(x) = sum(x)/count(x)$ · max min
- **Holistic:** it there is no constant bound on the storage size needed to describe a subaggregate
  - median( ), mode( ), rank( )

## Data Cube



## Typical OLAP Operations:

- **Roll up (drill-up):** summarize data
  - climbing hierarchy · D reduction
- **Drill down (roll down):** reverse
- **Slice & dice:** project & selection
- **Pivot (rotate):** reorient than, fact table
- **Drill across:** involving more, fact table
- **Drill through:** bottom level to b-e SQL

## Design: Business Analysis Framework

- **Top-down:** select of relevant data
- **Data source view:** expose info. uper.sys.
- **Data warehouse:** fact table, D table
- **Business query:** perspective, end-user

## Data Cube Computation: $T = \prod_{i=1}^{n}(L_i + 1)$

**Bitmap index:** · value → vector
- len(vector) = # of records

**Join indices:** · $J(R,S) \to R(R-) \bowtie S(s-)$
- relates the value of dimensions → row in the fact table

- **ROLAP:** relational OLAP: extended-r. DBMS optimization of DBMS. → Greater Scalability
- **MOLAP:** multidimensional : · sparse-array based · fast indexing → pre-computed
- **Hybrid (HOLAP):** flexibility

## Chapter 5: Data Cube Technology

**Iceberg Cube:** having count(*)>=min
- no need to save value below threshold
- avoid computing the un-needed
- avoid explosive growth

**close cube:** no cell ∈ descendant of c
→ d has same measure value as c
**cube shell:** cuboid, small # D

## General Heuristic: Sorting, hashing, grouping

- **Smallest-child:** smallest pre cuboid
- **Cache-results:** caching result & reduce disk I/o
- **Amortize-Scans:** as many as possible
- **Share-sorts:** share sort cross multi cuboid
- **Share-partition:** hash-based.

## Multiway array aggregation (bottom up)

- multi-D chunks · simultaneous aggre
- intermedia aggre value reused
- Cannot Apriori pruning× iceberg×
- Efficient for small

- Partition chunks: small subcube, memory
- compressed sparse array addressing
- aggregation in multiway, visiting the cube cells → min the # of visit, reduce memory access & storage cost



X pruning A→0
Bottom-up
X apriori

---

## BUC: top-down:
- partition & iceberg pruning
- it partition < min, pruning child
- No simultaneous aggregation
- Entire data cannot fit in memory
- Sort distinct values
- Continue processing · optimizations
  - Partitioning · ordering D · ×duplicate
  → cannot holistic aggregates ×



## High-D OLAP with Shell-fragment

- Trade off between the amount of pre-computation and the speed of online computation
- Reduce high-D ⇒ low-D cubes shell fragment
- Online re-construction of h-D space
- Lossless reduction ← pre-compute

| tid | A | B | C | D |
|-----|-----|-----|-----|-----|
| 1 | $a_1$ | $b_1$ | $c_1$ | $d_1$ |
| 2 | $a_1$ | $b_2$ | $c_1$ | $d_2$ |
| 3 | $a_1$ | $b_2$ | $c_1$ | $d_1$ |
| 4 | $a_2$ | $b_1$ | $c_1$ | $d_1$ |
| 5 | $a_2$ | $b_1$ | $c_1$ | $d_1$ |

| Value | TID | size |
|-----|-----|-----|
| $a_1$ | 123 | 3 |
| $a_2$ | 45 | 2 |
| $b_1$ | 145 | 3 |
| $b_2$ | 23 | 2 |
| $c_1$ | 12345 | 5 |
| $d_1$ | 1345 | 4 |
| $d_2$ | 2 | 1 |

| cell | ∩ | TID | size |
|-----|-----|-----|-----|
| $a_1 b_1$ | 123∩145 | 1 | 1 |
| $a_1 b_2$ | 123∩145 | 23 | 2 |
| $a_2 b_1$ | 45∩145 | 45 | 2 |
| $a_2 b_2$ | 45∩23 | ∅ | 0 |

Shell-fragment

## Size & design:
- T: tuples  D: dimension
- F: shell-fragment size.
→ Cube space: $O\left(T\lceil\frac{D}{F}\rceil(2^F-1)\right)$
- × disjoint

## General Form: $\langle a_1, a_2, \dots, a_n : M \rangle$
3 values: · instantiated value
- Aggregate * function · inquire ? function

## Complex Aggregation at Multiple Granularities
- **Multi-Feature cubes:** compute complex queries involving multiple dependant aggregates at multiple granularities ← cube
- **Discovery-driven Exploration of Data:**
  - huge cube space · pre-compute
  - Indicating Exception

## Chapter 6: Frequent Pattern

**Pattern:** itemset $X = \{x_1, \dots, x_k\}$
(absolute) support: frequency of number
· relative support: fraction

**Associate Rule:** $X \to Y (s,c)$
- Support $\frac{X \cup Y}{total}$ · Confidence $= \frac{Sup(X \cup Y)}{Sup(X)}$

**closed pattern:** no $Y \supset X$ w/ same sup.
$P_1\{a_1-a_{100}\}: 2 \quad P_2\{a_1..a_{100}\}: 1$
- lossless compression

**Max Pattern:** no $Y \supset X$.
$P_1\{a_1..a_{100}\}: 1$ do not care sup subpat
- lossy compression

**Downward closure (Apriori):** Any Subset of a frequent itemset must be frequent

**Apriori pruning principle:** it there is any itemset which is infrequency, no more super-set will be generated
→ self-join → prune

**Partition: reduce passes:**
- local frequent → global frequent
- **DHP: Direct Hashing & Pruning**
  - reduce number of candidates
  - whose corresponding hash below threshold

**ECLAT: Exploring Vertical Data Format:**
- t(x) = t(xy), together
- t(x) ⊂ t(y): x, always Y.
→

---

## FP Growth: Pattern Growth.

100 ftacdjimp $\Rightarrow$ ft,c,a,m,p,t
200 abtfjmo $\Rightarrow$ ft,c,a,b,m
300 bfhjowt $\Rightarrow$ ft,b
conditional



a fc:2
b f:1 conditional

## Parallel projection: each frequency item
## Partition project: in order

**closed:** it Y appears in every occurence of X, then Y is merged with X

## Interestingness: subjective vs objective
- **Lift:** $Lift(B,C) = \frac{c(B \to C)}{s(C)} = \frac{s(B \cup C)}{s(B) \times s(C)}$
  - $= 1$: independent  $< 1$: 
- **$\chi^2$:** $\chi^2 = \sum \frac{(O-E)^2}{E}$  $\begin{cases} = 0 & indep \\ > 0 & dep \end{cases}$
- **Null invariant:** value × change → # Null
- **IR (Imbalance Ratio):** $IR(A,B) = \frac{|S(A) - S(B)|}{S(A) + S(B) - S(AB)}$
- **Kulczynski:** $Kul(A,B) = \frac{1}{2}\left(\frac{S(A \cup B)}{S(B)} + \frac{S(A \cup B)}{S(B)}\right)$

## Chapter 7: Advanced FP

**Multiple-Level:** uniform / reduced min
- Redundant Rule · group based min
**Multiple-D:** · 1-D. · Multi-D
- inter-D (no repeat) · Hybrid-D (repeat)
**Quantitative:** numerical · static: concept
- Dynamic: distribution · Clustering: distance
- Deviation: $G = Fe \Rightarrow$ usage: mean=? (covered)
- **Extraordinary:** Rule: confirm · sub: highlight
- **Rare Pattern:** group-based min-sup
- **Negative:** · Sup-based $\Rightarrow sup(A \cup B) << sup(A) \cdot sup(B)$ → × null
  - kulcsy: $(P(A|B) + P(B|A))/2 < \varepsilon$
- **Compressed:** $D(P_1, P_2) = 1 - \frac{|T(P_1) \cap T(P_2)|}{|T(P_1) \cup T(P_2)|}$
**Constrains:** · Knowledge type · Data
- D / Level · Rule (Pattern) · interesting
  all autonomously × data mining query

## Pattern Space Pruning
- **Anti-monotic:** Violate → terminate i.e. Apriori is anti-monotic
- **Mono tone:** Satisfies
  e.g. · $sum < v \to$ anti  $> v$ mono
  - $range < v \to$ anti  $> v$ mono
  - $min < v$ mono · $Sup > \sigma$ anti

## Data Space Pruning:
- **Anti-monotone:** it every cannot satisfy pattern p under c.
  e.g. $sum > v$, $min < v$, $range > v$.
  - check recursively
- **Succinctness:** it the constrain c can be enforced by directly manipulating the data
  - no i → remove i · with i → i-projected
  - $min < v$ ✓ · $sum > v$ × → increase

## Convertible Constraints:
- tough → (anti) monotone by ordering

**Core Pattern:** subpatterns of α that cluster around α by sharing a similar sup
**Robustness:** $\left|\frac{D\alpha}{D\beta}\right| > \tau \quad 0 < \tau \leq 1$
(d,τ) robust it d is max number of items that can removed from α
- $\sqrt{2}(2^d)$ core Pattern · colossal → core
- **Dense ball:** $D(\alpha, \beta) = 1 - \frac{|D\alpha \cap D\beta|}{D\alpha \cup D\beta|}$

**Sequential Pattern:** set of sequences ⇒ complete set of frequency subsequences $\langle(ab)c\rangle$ · an element may contain a set of items

---

## GSP: Apriori-Based SPM:
- initial → scan → generate

## SPADE: Vertical Data Format
- map to $\langle SID, EID \rangle$
- Grow the subsequences one item at a time by apriori candidate gene

## Pattern Growth: Prefix Span
- Given $\langle a(abc)(ac)d(cf)\rangle$
- Prefix $\langle a \rangle, \langle aa \rangle, \langle a(ab) \rangle, \langle a(abc) \rangle$
- Suffix: Prefixes based projection

## CloSpan: There exists no super-P
Sub $= S' \supset S$, s' & s has same sup.
→ reduce # of Pattern → Same expressive power
- $S \supset S_1$, $S_1$ is closed itt same size



## Constraint-Based SPM:
- Anti-monotonic · Monotonic
- Data Anti-monotonic · Succinct
- Convertible

## Timing-Based Constraints:
- order constraints · min-gap const
- Max span : time diff 1st & last
- Window size: events in one ele do not occur at same time

## Episodes, Episode SPM:
- **Regular Expression:** serial $A \to B$
- Parallel: $A|B$ ← partial order relationship
- Regular Exp: $(A|B)C^*(D \to E)$
- **Method:** Variation of GSP

## Frequent Graph Patterns:
$D = \{G_1, G_2, \dots G_n\}$
$D_g = \{G_i | g \subseteq G_i, G_i \in D\}$
$Support(g) = |D_g|/|D| \geq min\_sup$
- **Apriori-Based:** candidate generation → pruning → sup counting → cand eliminate
- **Vertex Growing vs. Edge Growing:** new graph with new vertex/edge
- **Pattern-Growth:** (duplicate)

$(a_1, a_2, a_3, a_4, \dots a_{10})$
$(b_1, a_2, a_3, a_4, \dots a_{10})$  3,96
$(c_1, c_2, a_3, a_4, \dots a_{10})$
$(d_1, d_2, d_3, a_4, \dots a_{10})$  11
$(D_1, \underline{\quad} a_4 - a_{10})$ } $4 \times 2^9$
$(*, D_2, \underline{\quad} a_4 \leftarrow a_{10})$  $2 \times 2^8$
$(*, *, D_3 \underline{\quad} a_4 - a_{10})$  $2 \times 2^7$
$(*, *, *, \underline{\quad} a_4 \leftarrow a_{10})$  $2^7$

| | | | | |
|-----|-----|-----|-----|-----|
| $\chi^2(A,B)$ | $\sum \frac{(o-E)^2}{e(a;b)}$ | | $[0,\infty]$ | No |
| Lift(A,B) | $\frac{s(A \cup B)}{s(A) \times s(B)}$ | | $[0,\infty]$ | No |
| AllConf(A,B) | $\frac{s(A \cup B)}{max\{s(A), s(B)\}}$ | $[0,1]$ | Yes |
| Jaccard(A,B) | $\frac{s(A \cup B)}{s(A)+s(B)-s(A,B)}$ | $[0,1]$ | Yes |
| Cosine(A,B) | $\frac{s(A \cup B)}{\sqrt{s(A) \cdot s(B)}}$ | $[0,1]$ | Yes |
| Kulczynski(A,B) | $\frac{1}{2}(\frac{s(A \cup B)}{s(A)} + \frac{s(A \cup B)}{s(B)})$ | $[0,1]$ | Yes |
| maxconf(A,B) | $max\{\frac{s(A \cup B)}{s(A)}, \frac{s(A \cup B)}{s(B)}\}$ | $[0,1]$ | Yes |

$Sup(A) = A/total$

## Bitmap: 3D



3×4 vectors   4 distinct values