## 2005 Fall Conference

# Data Mining, Dashboards and Data Quality

**John Rome, Arizona State University**

ARIZONA STATE UNIVERSITY

---

# Review of "BI" Buzzwords

OLAP          De-Normalized            Data Mart

Operation Data Store (ODS)     Bit-Mapped Indexing

Drill-Down     **Data Mining**               ROLAP

Aggregation                    Replication
         MOLAP      XML

                              Facts/Dimensions
                    Metadata

**Data Quality**              Business Intelligence

Star Schema     Multi-dimensional    **Dashboards**

Transformation      SQL    Snowflake Schema
Tools (ETL)

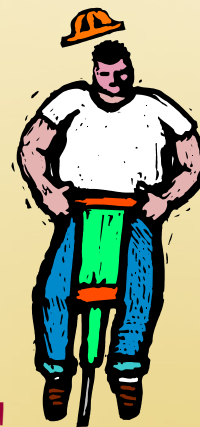ARIZONA STATE UNIVERSITY

# What is Data Mining?

- Analysis of data with the intent to prove a hypothesis or to discover gems of information in the vast quantity of data
- Looking for patterns in a collection of facts or observations
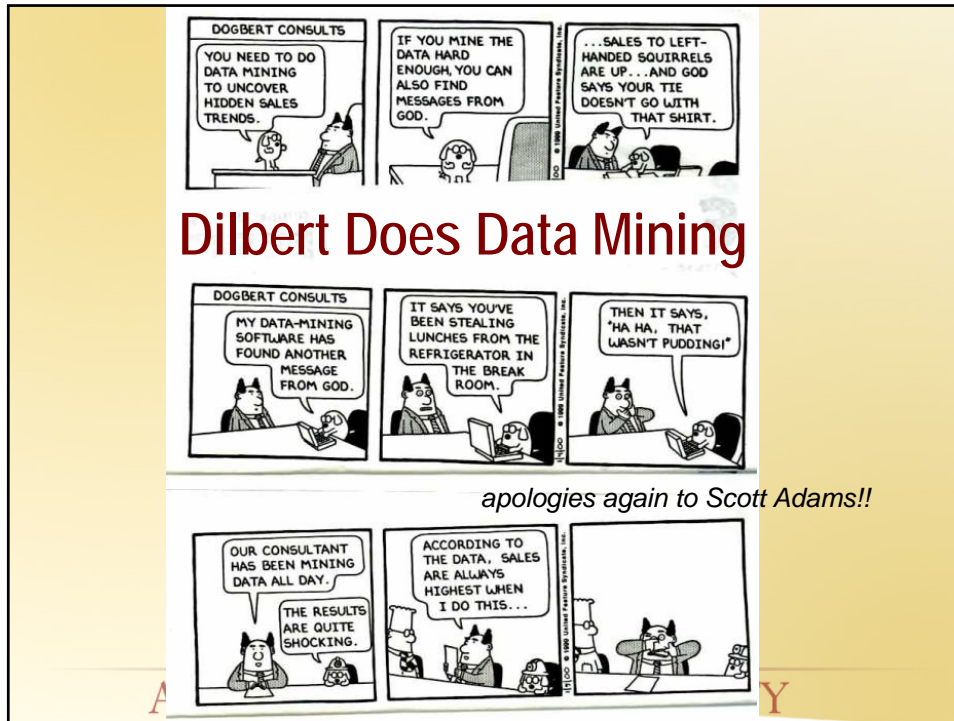- Techniques include Neural Networks, Visualization, Decision Trees, etc.

ARIZONA STATE UNIVERSITY

# Data Mining Observations

- Higher Ed Still Wrestling with Dirty Work of Building Data Warehouse
- How Applicable in Higher Ed?
- Technology Isn't Cheap
- Need to Agree on Common Set of Definitions to Work
- Level of Expectation High
- Don't Feel Guilty if You Don't Do It!!

ARIZONA STATE UNIVERSITY

# Dilbert Does Data Mining

*apologies again to Scott Adams!!*

# Santa Claus Uses Data Mining…

Here's some news, and brace yourselves
There's trouble hiring Santa's elves
The little guys put down their hammers
And took new job as C programmers

But don't you worry, Santa's finding
Salvation in **data mining**
To prep his list, his query would be…

SELECT * FROM KIDS, STATUS = 'GOOD'

*-Anonymous*

# What are Dashboards?

A dashboard is a graphical display that compares performance against predefined goals.



# Types of Dashboards?

- Operational
- Strategic
- Tactical

**-Wayne Eckerson**



Types of Dashboards in use

Operational dashboards 23%
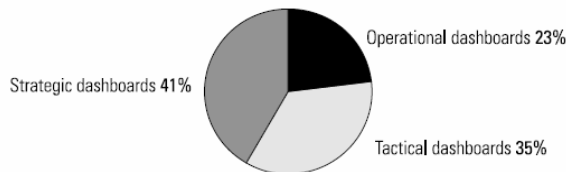
Strategic dashboards 41%

Tactical dashboards 35%

*Illustration 32. A larger percentage of organizations are implementing strategic dashboards than any other type of dashboard. Based on 240 respondents who said their group has deployed a dashboard.*
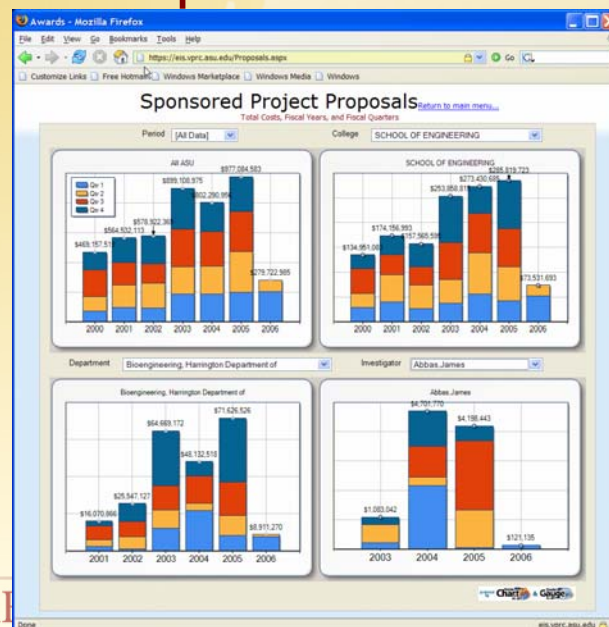
# Dashboard Tools…

How Organizations Are Building:

- 41% - BI Tools
- 22% - Custom Code
- 13% - Microsoft
- 17% - Purchased from Vendor

ARIZONA STATE UNIVERSITY

# Example of Dashboard

# Example of MDD (Dashboard?)



# What is Data Quality (DQ)?

Data Having:
- Accuracy
- Integrity
- Consistency
- Completeness
- Validity
- Timeliness
- Accessibility

## Data Quality Impact

*"The Business costs of nonquality data, including irrecoverable costs, rework of products or services, workarounds, and lost and missed revenue may be as high as **10 to 25 percent of revenue or total budget** of an organization."*

-Larry English
"Father of Data Quality"

## Data

It's good,
it's bad,
and it's ugly!

# Mainly Ugly

## Problems I've Seen…

- Building named "China"
- High number of +95 year old students
- Is Pat "M" or "F"?
- Negative SQFT
- Addresses
- Gender Identity
- Macintosh worth $6

ARIZONA STATE UNIVERSITY

# What's Wrong With This Data?

| Name | Address | City | State | Zip | Country |
|---|---|---|---|---|---|
| John Rome | 1235 North Sunnyvale #47 | Mesa | AZ | 85205-4347 | USA |
| Tiger Woods | PO Box 871203 | Boulder | Colo | 12345 | US |
| John Porter | 123 Oak Street | Temple | AZ | 852011223 | US |
| Foyt, A.J. | 555 Temper Blvd. | Indianapolis | INDIANA | 33045 | |
| Evander Holyfield | Missing Ear Drive | Atlanta | GA | 11111-1234 | US |
| John Rome | 145 N. Greenfield Rd. | Mesa | AZ | 85203 | US |
| Sir Charles Barkley | 4545 S. Scotsdale Blvd | Paradise Valley | AZ | 85288 | US |

ARIZONA STATE UNIVERSITY

# Data Warehouse Helps Identify…

1. Objectionable Outliers
2. Nuisance Nulls
3. Disorderly Domains
4. Downright Bad Data

ARIZONA STATE UNIVERSITY

# Objectionable Outlier Example

| FACILITY_SERVICE_SQFT | FACILITY_MECHANICAL_SQFT | FACILITY_WALL_SQFT |  |
|---|---|---|---|
| 822 | 1303 | | -40307 |
| 168 | 249 | | -6825 |
| 412 | 302 | | -6344 |
| 310 | 260 | | -6117 |
| 281 | 261 | | -3396 |
| 595 | 95 | | -594 |
| [NULL] | [NULL] | [NULL] | |
| [NULL] | [NULL] | [NULL] | |
| [NULL] | [NULL] | [NULL] | |
| [NULL] | [NULL] | [NULL] | |
| [NULL] | [NULL] | [NULL] | |
| [NULL] | [NULL] | [NULL] | |
| [NULL] | [NULL] | [NULL] | |
| [NULL] | [NULL] | [NULL] | |
| 33 | 5 | 117 | |
| 228 | [NULL] | [NULL] | |
| 23 | 120 | 233 | |
| 862 | 7 | 239 | |
| 296 | 7 | [NULL] | |
| 296 | 7 | [NULL] | |
| 71 | 29 | 240 | |
| 122 | [NULL] | [NULL] | |
| 41 | [NULL] | [NULL] | |
| [NULL] | [NULL] | [NULL] | |
| [NULL] | [NULL] | [NULL] | |
| 77 | 157 | 246 | |
| 69 | [NULL] | [NULL] | |
| 510 | 214 | 272 | |

ARIZONA STATE UNIVERSITY

# Objectionable Outlier Example

**AGE**

**BIRTHDATE**

| | | Total |
|---|---|---|
| 1 | 05/10/2002 | 1 |
| | 06/04/2002 | 1 |
| | 07/10/2002 | 1 |
| 2 | 05/09/2001 | 1 |
| 86 | 06/27/1917 | 1 |
| 99 | 01/01/1905 | 3 |
| 100 | 07/03/1903 | 1 |
| 101 | 07/23/1902 | 1 |
| 102 | 04/23/1901 | 1 |
| Total | | 11 |

Sort Age by Lab

# What Wrong with These Birthdates?

Distribution Birth Days Clients

# Objectionable Outlier Techniques

- Min/Max Functions (SQL)
- Standard Deviation
- Data Visualization (after sorting results)

```
Query   Results
  1:    select MIN (FACILITY_GROSS_SQFT)    as MIN_GROSS_SQFT,
  2:           MAX (FACILITY_GROSS_SQFT)    as MAX_GROSS_SQFT
  3:    from FACILITY
```

| | MIN_GROSS_SQFT | MAX_GROSS_SQFT |
|---|---|---|
| 1 | 0 | 318030 |

# Nuisance Nulls Example

```
Query   Results
  1:  SELECT    MAJOR_CODE       as MAJOR_CODE,
  2:            COUNT(ASU_ID)    as COUNT
  3:  FROM      STUDENT
  4:  WHERE     MAJOR_CODE in (null,' ')
  5:  GROUP BY MAJOR_CODE
```

```
Query   Results
  1:  SELECT    MAJOR_CODE       as MAJOR_CODE,
  2:            COUNT(MAJOR_CODE) as COUNT
  3:  FROM      STUDENT
  4:  WHERE     MAJOR_CODE in (null,' ')
  5:  GROUP BY MAJOR_CODE
```

| | MAJOR_CODE | COUNT |
|---|---|---|
| 1 | [NULL] | 63254 |
| 2 | | 209590 |

| | MAJOR_CODE | COUNT_ |
|---|---|---|
| 1 | [NULL] | 0 |
| 2 | | 209590 |

ARIZONA STATE UNIVERSITY

# Nuisance Nulls Techniques

- WHERE Column_Name in (null, ' ')
- Data Visualization (after sorting results)

Nulls

# Disorderly Domains Example
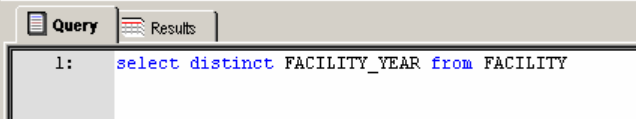
| | Facility Location Code | Facility Location Description |
|---|---|---|
| 1 | 12 | UNKNOWN |
| 2 | DT | DOWNTOWN CENTER |
| 3 | EC | EAST CAMPUS |
| 4 | MC | MAIN CAMPUS |
| 5 | OC | OFF CAMPUS |
| 6 | RP | RESEARCH PARK |
| 7 | TZ | CAMP TONTOZONA |
| 8 | WC | WEST CAMPUS |

# Disorderly Domains Techniques

- Select Distinct (SQL)
- Outer joins to find missing values



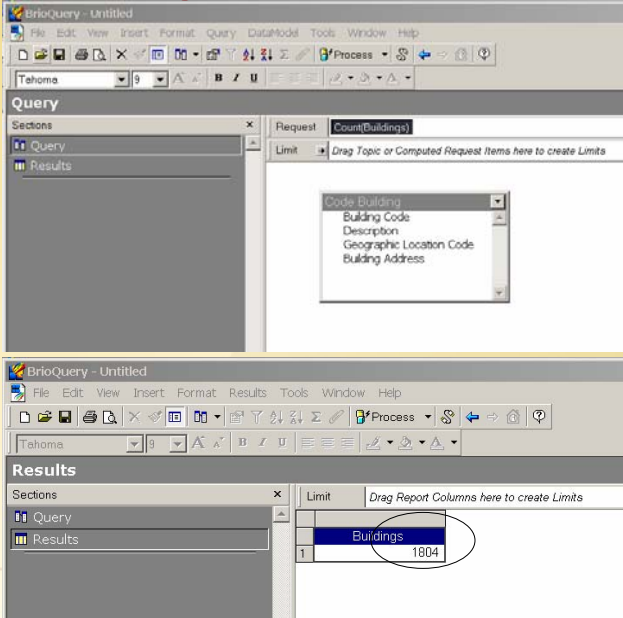ARIZONA STATE UNIVERSITY

# Downright Bad Data Example

## Some of the Buildings We Found…

| | Building Code | Description |
|---|---|---|
| 308 | CHHS | CHOLLA HIGH SCHOOL |
| 309 | CHILD | CAMPUS CHILDRENS CENTER |
| 310 | CHINA | PEOPLE'S REPUBLIC OF CHINA |
| 311 | CHKEE | CHEROKEE |
| 312 | CHLGR | CHALLENGER ELEMENTRAY SCHOOL |
| 313 | CHNDH | CHANDLER HS |
| 314 | CHNDL | CHANDLER DISTRICT 80 |
| 315 | CHNDR | CHANDLER, AZ |
| 316 | CHNDS | CHANDLER JHS, CHANDLER |
| 317 | CHNLB | CHINLE BOARDING SCHOOL |
| 318 | CHNLE | CHINLE,AZ |
| 319 | CHNLH | CHINLE HS |
| 320 | CHNLJ | CHINLE JHS,CHINLE |
| 321 | CHOLA | CHOLLA APARTMENTS |

ARIZONA STATE UNIVERSITY

## Downright Bad Data Techniques

- Min/Max/Avg/Count/Sum,etc. (Aggregate)
- Data Visualization
- Numeric, String, Date functions
- Subselects, other SQL syntax, etc

```
Query    Results
 1:  SELECT AL1.AFFILIATE_ID,
 2:         AL1.FACILITY_CODE,
 3:         AL1.ROOM,
 4:         AL1.PRIMARY_ROOM_FLAG
 5:  FROM   FACILITY_ROOM_EMPLOYEE AL1
 6:  WHERE  AL1.PRIMARY_ROOM_FLAG = 'Y'
 7:  GROUP BY AL1.AFFILIATE_ID
 8:
 9:  HAVING COUNT (AL1.AFFILIATE_ID) > 1
10:  ORDER BY AFFILIATE_ID
```

*Or is the problem definitional?*

ARIZONA STATE UNIVERSITY

## Data Cleansing Techniques

### AKA - Data Hygiene

- Focus on high-payoff data elements
- Interrogate data elements individually and collectively
- Standardization on national codes
- Conduct data audit for conformity of domain
- Document transformation rules and test
- Go back to the source if necessary

ARIZONA STATE UNIVERSITY

## Case For Data Warehousing

- Perform queries/reporting on servers not used by transaction processing system
- Use data models or server technologies that speed up query and reporting not appropriate for transaction processing
- Having an environment where you don't need a programmer to get your information
- Having an environment that contain history or where you can generate data "as of"
- To provide data that is "cleaned up"

ARIZONA STATE UNIVERSITY

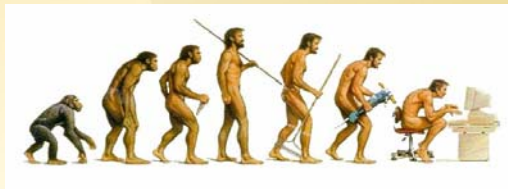# Case <u>Against</u> Data Warehousing

- Historical data often has limited value
- Data warehouses complicate the business processes
- All the data you may need is in your operational system
- Data warehouses required a great deal of maintenance which many can't support
- Might not "take" in the user community
- Cost might be too great

ARIZONA STATE UNIVERSITY

# Data Warehouse Evolution

- Access to Data
- Reporting (Query Tools)
- Analysis (OLAP)
- "Operationalize" (Applications)
- Prediction (Data Mining)

ARIZONA STATE UNIVERSITY

## Warehousing Biggest Obstacles

- Political Issues
- Poorly defined goals
- Lack of resources
- Technical limits
- Poor understanding of legacy data
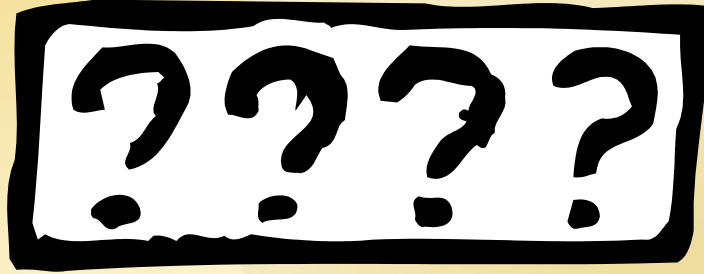- Poor data quality
- Lack of end-user support

## IR Challenges and Next Steps?

- Learn more about your Institution's DW or start discussion
- Learn to Leverage your DW and get data you need (have a a section for you!!)
- Become More Technical (modeling/DBA,etc.)
- Educate Yourself on DW (books/Web/blogs)
- Visit a couple of organizations that have had warehousing systems in production
- http://dheise.andrews.edu/dw/DWData.htm

Questions

ARIZONA STATE UNIVERSITY