

# *Reporting Data Quality Issues*

Mike Glasser

Office of Institutional Research  
University of Maryland – Baltimore County

# Agenda

- Introduction
- Elements of Data Quality
- Correcting Issues
- Tables and Procedures
- Reports
- Data Quality Firewall
- Questions

# UMBC

## University of Maryland - Baltimore County

- Located in suburban Baltimore County, between Baltimore, MD and Washington, DC
- One of the three public research campuses in the University of Maryland System
- 10K undergraduate and 2K graduate students
- 2,200 Employees; 715 full-time faculty
- PeopleSoft, SQL Server 2008 R2
- 9 Pan-Am chess championships

# UMBC

## Data Warehouse & Reporting

Headcount	IT	IR	Back Office
Management	1	2	3
Data Warehouse	3	2	0
Reporting	3	5	6
FTE	2.5	3	5

- Student Administration module (incl Fin Aid)
- Finance module
- Learn Analytics module (beta)



# Data Quality

- **Accurate**
  - Data entry
  - Outdated
- **Consistent**
  - Records in one table do not match other
- **Complete**
  - Missing data
- **Within Policy**
  - Violates a policy or practice
- **Reportable**
  - OK in system, but cannot be reported

# Accurate Data

- **Data Entry**
  - Invalid Emplid
  - Invalid Institution
  - Assignment % for Class not 100%
  - Wrong County code for Maryland
- **Outdated**
  - Wrong contact information
  - Wrong birthdate
  - Student no longer in program

# Consistent Data

- Inconsistent within record or tables
  - US Citizenship, F1 visa
  - Withdrawal code, date before classes
  - Class attribute not on catalog
  - StudentPlan without StudentTerm
  - Enroll Total in Class does not match registrations
  - New Plan not in DW setup table

# Complete Data

- Missing data
  - Plan is missing CIP code
  - Class is missing Instructor
- Missing record
  - Class is missing Instructor
  - Course missing graded component
  - New Plan not in DW setup table



# Data Within Policy

- Violation of Policy or Practice
  - Double majors for Grad students
  - Duplicate plans
  - 2<sup>nd</sup> major is Undecided
  - Grad Assistant not enrolled

# Reportable Data

- **Cannot report**
  - Unknown Gender (IPEDS)
  - Invalid major
- **Should not report**
  - Non-degree instead of degree seeking
  - Masters instead of Doctoral
  - Class Section changed after freeze

# Remedies

- Meetings or Email
- Transaction system
  - Edit records manually
  - Change business process
  - Change PS data entry (if lucky)
- Data warehouse
  - “Unknown” values (key = -1)
  - Create fake records
    - Bad majors
    - “Two or More” Ethnicity
  - New fields
    - Gender/IPEDS
  - Tweak data quality check

# Philosophy

- Fix it in the transaction system
- Prevent / Fix it at data entry
- Fix it as soon as possible
- Fix it before IR census
- Fix it in data warehouse

# Carrot / Stick

- Show them the errors
- Explain the impact
- Exposure to more users
- Data Management Cmte
- Provost
- OK, we will fix it in the DW

# Practice

- **Data Quality Team**
  - Part of Data Management Cmte
  - Campus commitment to data quality
  - Identify responsible parties
- **Data Quality Reporting**
  - Identify errors
  - Report errors

# Data Quality Team

- IR Data Administrator (convener)
- Data Managers/Stewards
  - Registrar
  - Scheduling
  - Financial Aid
  - Undergraduate Admissions
  - Graduate School
  - Finance
  - HR

# Data Quality Team

- Kickoff meeting to explain philosophy and processes
- Meetings with back office(s) as needed
- Back office report developers
- Write SQL to identify issues
- Write PS Query in transaction system to identify issues
- Develop DW reports with more detail for specific errors



# Data Quality Reporting

- Identify issue
  - Via meetings, IR, users, back office
- Write SQL
- Nightly procedures
  - Check everything, summarize
- Daily report
  - User subscriptions or on demand
- Fix errors
  - A little prompting

# Tables

- Error Messages
- Daily Errors
- Copy of Yesterday's Errors
- History of Errors
- Exceptions

# Table for Error Messages

- 1 message per error check
- Error Message Number and Text
- Explanation and/or solution
- Module (Admissions, Registrar)
- Table name
- Field name
- Key fields

# Table for Error Messages

- Create table DW.Data\_Quality\_Error\_Messages (
  - Error\_Msg\_Nbr int IDENTITY(1, 1),
  - Error\_Msg\_Text varchar(100),
  - Error\_Msg AS (((Error\_Msg\_Text + ' [') + CONVERT(varchar,Error\_Msg\_Nbr)) +']'),
  - Error\_Msg\_Severity varchar(10),
  - Error\_Explanation varchar(500),
  - DQ\_Module varchar(3),
  - Addl\_Recipients varchar(200),
  - Table\_Name varchar(128),
  - Field\_Name varchar(128),
  - Process\_Name varchar(128),
  - DW\_Load\_Dttm datetime,
  - Key\_Fields varchar(500),
  - DQ\_Sub\_Module varchar(5))

# Table for Daily Errors

- Error Message Number
- Value of the field with error
- Value of fields to identify record
- Date/time
- New error indicator
- Current error indicator

# Table for Daily Errors

- create table DW.Data\_Quality\_Daily (
  - Error\_Msg\_Nbr int,
  - Fieldvalue varchar(100),
  - Key1 varchar(100),
  - Key2 varchar(100),
  - Key3 varchar(100),
  - Key4 varchar(100),
  - Key5 varchar(60),
  - Key6 varchar(60),
  - Key7 varchar(60),
  - Key8 varchar(60),
  - Key9 varchar(60),
  - DW\_Load\_Dttm datetime,
  - Key\_Values varchar(1000),
  - New\_Error\_Yn varchar(1),
  - Exception\_YN varchar(1),
  - Current\_Yn varchar(1))

# SQL

- SELECT statement to identify the errors
- Has to return the key fields and the field with error (if applicable)
- Convert to INSERT statement for table **Data\_Quality\_Daily**
- Written by me or back office report developer

# SQL

## Academic Plan has invalid CIP code

```
INSERT ( DW_Load_Dttm, Fieldvalue, Error_Msg_Nbr
        , Current_YN, Key1, Key2)
SELECT getdate(), a.CIP_Code,          180, 'Y'
       a.ACAD_PLAN, a.EFFDT
FROM iPSSA.Source.PS_ACAD_PLAN_TBL A
LEFT JOIN iPSSA.Source.PS_CIP_CODE_TBL B
WHERE b.CIP_CODE IS NULL
```

\* Removed effective dating logic for simplicity



# Table for Error History

- Summarize daily error messages
- One record per message per day
  - Errors
  - Unique values
  - Exceptions
  - New errors
  - Current errors

# Table for Error History

- create table DW.Data\_Quality\_History (
  - Error\_Msg\_Nbr int,
  - ETL\_Load\_Dttm datetime,
  - ETL\_Date varchar(10),
  - Error\_Count int,
  - Unique\_Values\_Count int,
  - Exception\_Count int,
  - DW\_Load\_Dttm datetime,
  - New\_Error\_Count int,
  - Current\_Count int)

# Table for Exceptions

- Some errors can be warnings
- Exceptions are on individual case
- Still identified during the checks
- Exceptions are deleted from daily errors after summary, so exceptions can be counted

# Table for Exceptions

- create table DW.Data\_Quality\_Exceptions (
  - **Error\_Msg\_Nbr** int,
  - **Fieldvalue** varchar(100),
  - **Key1** varchar(100),
  - **Key2** varchar(100),
  - **Key3** varchar(100),
  - .....
  - **Key8** varchar(100),
  - **Key9** varchar(100)

Error_Msg_Nbr	Fieldvalue	Key1	Key2	Key3
179	%	SPSE PBC	%	%

Error Message 179 : Academic Plan is missing Degree code

# Procedures

- **Nightly Data Quality**
  - Run Data Quality Checks
  - Cleanup error table
  - Identify errors as new
  - Summarize errors
  - Delete exceptions
- **Data Quality Checks SA**
  - SQL for individual DQ checks

# Reports

- Microsoft Reporting Services
- Summary report for any date
  - Count of errors, current and new
  - Broken down by module
  - Links to details of error message
- Detail report for last night
  - Which records had the errors
- List of possible error messages
  - With explanation and relevant data

# Summary Report

Error Messages from Oct 9 2012 5:46AM

Data Quality Module	Error Message	Errors	Unique Values	Except	New Errors	Current
<b>AD</b>	<b>1 messages</b>	<b>3</b>	<b>1</b>	<b>0</b>	<b>0</b>	<b>3</b>
AD	<a href="#">Student has invalid Plan in Admissions stack</a>	3	1	0	0	3
<b>CC</b>	<b>1 messages</b>	<b>4</b>	<b>4</b>	<b>0</b>	<b>0</b>	<b>4</b>
CC	<a href="#">Invalid EMPLID</a>	4	4	0	0	4
<b>IR</b>	<b>2 messages</b>	<b>28</b>	<b>24</b>	<b>0</b>	<b>0</b>	<b>2</b>
IR	<a href="#">Academic Plan not found in OIR.DW.UW_Acad_Plan_Tbl</a>	2	2	0	0	2
IR	<a href="#">Class Section changed after Day 10</a>	26	22	0	0	0
<b>SR</b>	<b>18 messages</b>	<b>1290</b>	<b>154</b>	<b>13</b>	<b>0</b>	<b>506</b>
CORE	<a href="#">Academic Plan has an invalid CIP code</a>	1	1	0	0	1
DEG	<a href="#">Degrees in different Plan, same HEGIS</a>	24	5	0	0	0
DEG	<a href="#">More than one degree with same Academic Plan</a>	17	10	0	0	0

# Report of Daily Errors

## Data Quality Details for Last Night

**Academic Plan has an invalid CIP code [180]**

DQ Module : **Student Records ( CORE )**

Table Name : **PS\_ACAD\_PLAN\_TBL**

Field with Error : **CIP\_Code**

Keys in Table : **Acad\_Plan ~ Effdt**

Total Errors : 1

<b>Error Value</b>	<b>Keys</b>	<b>New</b>	<b>Current</b>
	ACCT UDC~Jan 2 1901 12:00AM	N	Y

- Error Message with Number
- Field with error
- Keys in table
- Value of field with error
- Value of keys to identify record
- Error explanation on hover



# Subscriptions

- Microsoft Reporting Services
- Setup email recipients
  - Can be anyone
  - Don't need access to report
- Any schedule
  - Recommend daily new, weekly all
- Set report parameters
- Choice of formats
- Link to report online

# Email Summary

Error Messages from Oct 9 2012 5:46AM

Data Quality Module	Error Message	Errors	Unique Values	Except	New Errors	Current
<b>AD</b>	<b>1 messages</b>	<b>3</b>	<b>1</b>	<b>0</b>	<b>0</b>	<b>3</b>
AD	<a href="#">Student has invalid Plan in Admissions stack</a>	3	1	0	0	3
<b>CC</b>	<b>1 messages</b>	<b>4</b>	<b>4</b>	<b>0</b>	<b>0</b>	<b>4</b>
CC	<a href="#">Invalid EMPLID</a>	4	4	0	0	4
<b>IR</b>	<b>2 messages</b>	<b>28</b>	<b>24</b>	<b>0</b>	<b>0</b>	<b>2</b>
IR	<a href="#">Academic Plan not found in OIR.DW.UW_Acad_Plan_Tbl</a>	2	2	0	0	2
IR	<a href="#">Class Section changed after Day 10</a>	26	22	0	0	0
<b>SR</b>	<b>18 messages</b>	<b>1290</b>	<b>154</b>	<b>13</b>	<b>0</b>	<b>506</b>
CORE	<a href="#">Academic Plan has an invalid CIP code</a>	1	1	0	0	1
DEG	<a href="#">Degrees in different Plan, same HEGIS</a>	24	5	0	0	0
DEG	<a href="#">More than one degree with same Academic Plan</a>	17	10	0	0	0

# HR Data Quality

- Process built prior to BbA
- Similar to SA process, but not same tables or procedures
- Email summary sent daily with NEW errors, weekly with ALL
- Email created with SQL, not RS
- Details reported with Crystal
- 180 error checks

# HR Email

Comparing yesterday's errors with today, the following NEW errors were found ...

Count Error Message

-----

- 1 Eligible retirement code not found [151]
- 1 EMPL\_CLASS inconsistent with EEO6CODE [113]
- 1 Unable to find Benefit\_Plan for employee [81]
- 6 EMPLID not found in UM\_Person\_Info table [103]

=====

9 NEW errors  
1085 Total errors

Email produced by Email\_New\_ETL\_Errors\_Sp on Oct 9 2012 12:55AM

# Data Quality Firewall

- **Procedures**
  - missing data
  - foreign key discrepancies
  - duplicate source keys
- **Warnings**
  - Usually missing data
  - Loaded as “Unknown” (Key = -1)
- **Critical Errors**
  - Usually duplicate source keys
  - Only the first is loaded

# Data Quality Firewall

## Why don't we use it?

- I already had a system in place for HR data
- I did not know much, if anything, about it
- I did not know XML
- Maybe we could, if time were invested

# Recap

- Identify people responsible for data quality in each area
- Back office commitment
- Identify issues and resolutions
- Use DW to capture and report issues
- Report issues to appropriate people
- Impact of data quality
- Accept that some things are wrong only for reporting

# Wrap Up

Any Questions?

**Mike Glasser**

University of Maryland - Baltimore County

[mglasser@umbc.edu](mailto:mglasser@umbc.edu)

(410) 455-3577

Source code is available upon request