# NISS

## Decision-Theoretic Framework for Data Quality

Alan Karr
March 3, 2005

---

## Summary

- Specific context and testbed database
- Specific set of DQ strategies
- Evaluation of strategies
- Predictive statistical models

---

## Notation

- $\mathcal{D}^{\text{true}}$ = true database (flat file of cases × attributes

- $\mathcal{D}^{\text{pre}}$ = database prior to clean-up

- $S$ = clean-up strategy

- $\mathcal{D}^{\text{post}}(S)$ = database resulting from applying $S$ to $\mathcal{D}^{\text{pre}}$

## Measuring Effectiveness

Conceptually,

$$\text{Eff}(S) = d(\mathcal{D}^{\text{post}}(S), \mathcal{D}^{\text{true}}),$$

where $d$ is a data quality metric

## Inference-Based Effectiveness

More meaningfully,

$$\text{Eff}(S, P, \mathcal{D}^{\text{pre}}) = d_P(\mathcal{D}^{\text{true}}, \mathcal{D}^{\text{pre}}) - d_P(\mathcal{D}^{\text{true}}, \mathcal{D}^{\text{post}}(S)),$$

where

- $P$ = inference procedure that can be applied to the data

- $d_P$ = function measuring the difference in the results of $P$ applied to two different databases

## What if Truth is Not Known?

Use

$$\text{Eff}^{\text{naive}}(S, P, \mathcal{D}^{\text{pre}}) = d_P(\mathcal{D}^{\text{pre}}, \mathcal{D}^{\text{post}}(S)).$$

Relevant points:

- $+$ sign for $\text{Eff}^{\text{naive}}(S, P, \mathcal{D}^{\text{pre}})$ may not signal improvement

- Small values of $\text{Eff}^{\text{naive}}(S, P, \mathcal{D}^{\text{pre}})$ mean no improvement

## Prediction

- $\{S(\theta) : \theta \in \Theta\}$ = parameterized family of clean-up strategies

- Goal: solve

$$\theta^* = \arg\max_{\theta} \text{Eff}(S(\theta), P, \mathcal{D}^{\text{pre}})$$

- Problem: only know $\text{Eff}(\theta)$ for a few values of $\theta$

## Predictive Models

Build statistical model

$$\widehat{\text{Eff}}(\theta) = f(\theta) + \text{uncertainty}$$

Challenges:

- "Form" of the model

- Nature of the uncertainties

- What data are necessary to fit the model

- Validation