

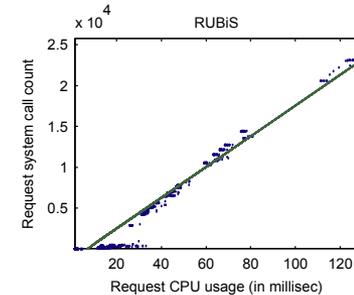
Regression

Kai Shen

Dept. of Computer Science, University of Rochester

Regression

- What is regression?
 - predict a target random variable as a function of one or more input variables
- Linear regression
 - linear functional relationship
 - $y = b_0 + b_1x$



2/28/2008

URCS 573 - Spring 2008

2

Regression Error

- Least-square fitting
 - y_1, \dots, y_n are observed samples for target variable while y_1^*, \dots, y_n^* are corresponding predictions
 - minimize mean square error
 - $((y_1^* - y_1)^2 + \dots + (y_n^* - y_n)^2) / N$
- Problems with mean square errors
 - doesn't give a true sense of modeling difficulty
 - coefficient of determination: $1 - E_{\text{model}} / E_{\text{mean}}$
 - F-test - must also consider degrees of freedom
 - trivial to find a regression function with two parameters that fit perfectly with two sample observations
 - too much weight for outliers - mean absolute error

2/28/2008

URCS 573 - Spring 2008

3

Non-linear Regression

- Curvilinear
 - nonlinear relationships that can be converted to linear forms with variable transformations
- Exponential
 - $y = b a^x$
 - $\ln y = \ln b + x \ln a$
- Powerlaw
 - $y = b x^a$
 - $\ln y = \ln b + a \ln x$
 - Powerlaw relationships in practice

2/28/2008

URCS 573 - Spring 2008

4



Powerlaw Fitting

- Loglog-transformations and then linear fitting
- Problems
 - Heavy weight on tails after transformations
 - Enlarged errors on non-tails

2/28/2008

URCS 573 - Spring 2008

5



Other Prediction Models

- Nearest matching
- Decision trees
- Bayesian networks

2/28/2008

URCS 573 - Spring 2008

6



Other Issues

- Outliers
- Confidence beyond measured range
- Intuitive explanation of derived models
- Correlation vs. Dependency

2/28/2008

URCS 573 - Spring 2008

7



Disclaimer

- Most materials in these slides were developed from the book "The Art of Computer Systems Performance Analysis", R. Jain, 1991, Wiley.

2/28/2008

URCS 573 - Spring 2008

8