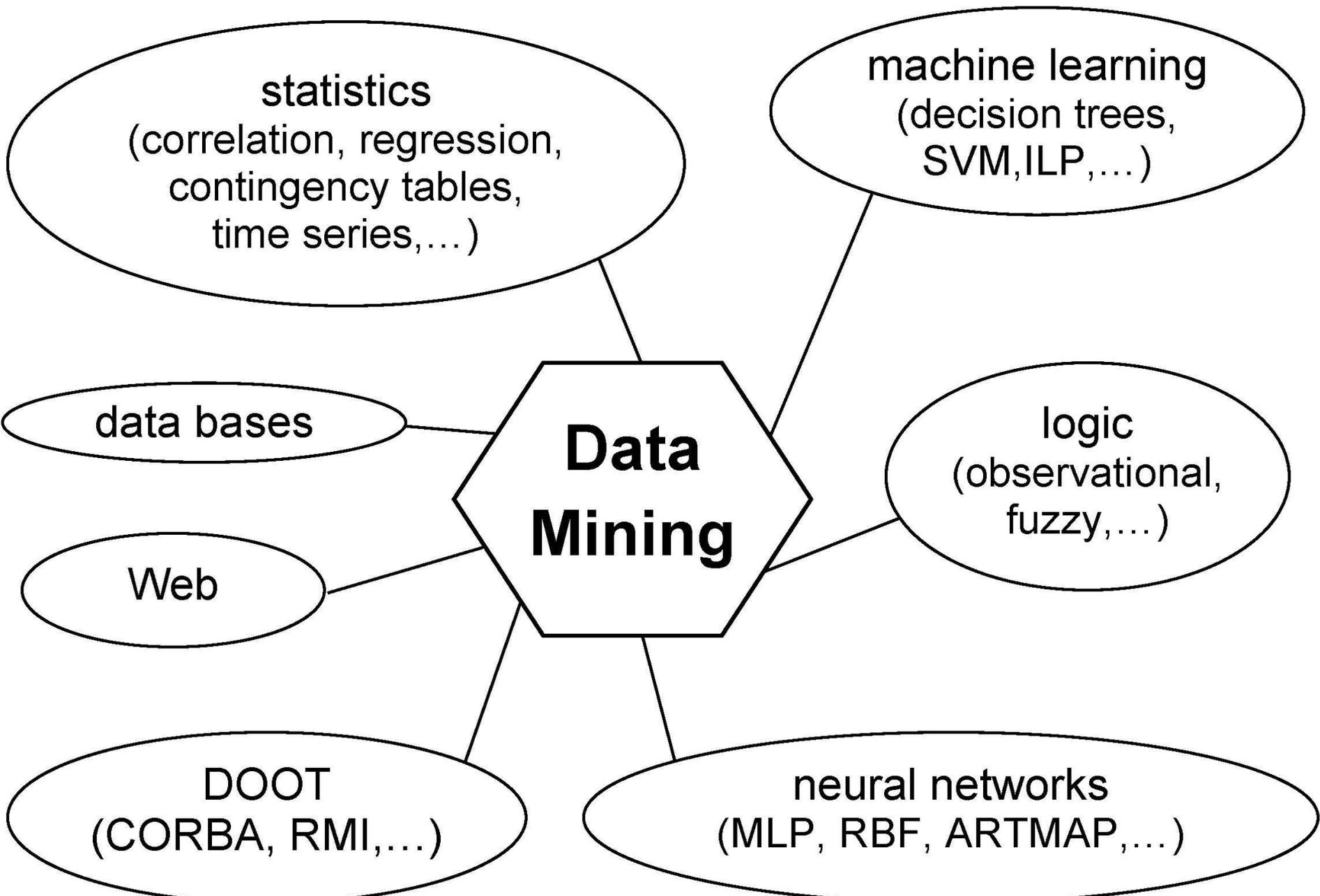


John W. Tukey

EXPLORATORY DATA ANALYSIS





statistics
(correlation, regression,
contingency tables,
time series,...)

machine learning
(decision trees,
SVM,ILP,...)

logic
(observational,
fuzzy,...)

neural networks
(MLP, RBF, ARTMAP,...)

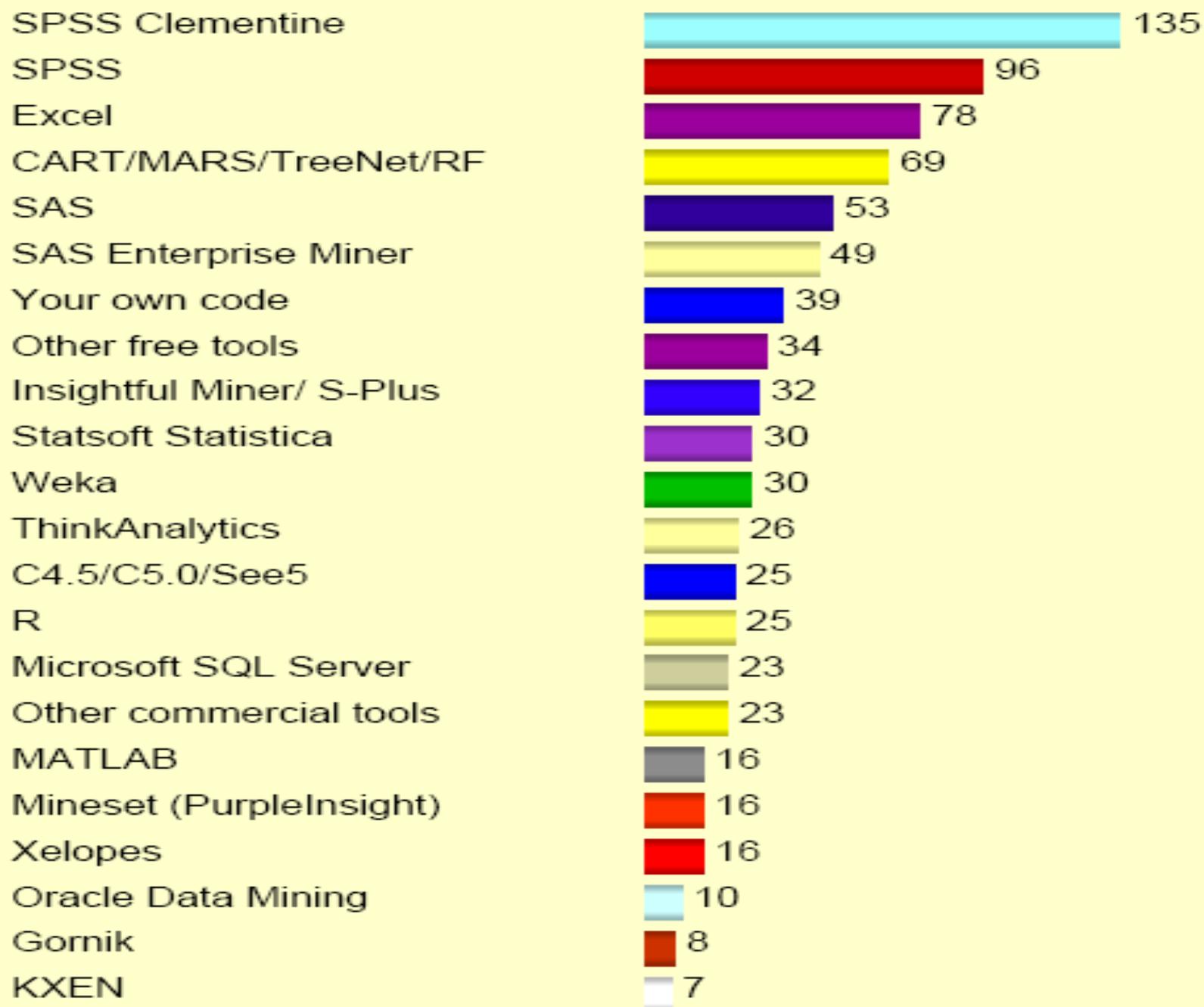
DOOT
(CORBA, RMI,...)

Web

data bases

**Data
Mining**

Data mining/Analytic tools you used in 2005 [376 voters, 860 votes total]



Data Mining Tools Used Poll (May 2009)

What data mining tools have you used for a real project (not just for evaluation) in the past 6 months? (364 voters)

SPSS PASW Modeler (formerly Clementine)
(68 alone, 52 with other tools, 120 total)



RapidMiner (36 alone, 41 w. other tools, 77 total)



SAS (39 alone or with SAS EM; 36 with other tools, 75 total)



Excel (1 alone, 68 total)



SAS Enterprise Miner (39 alone or with SAS; 28 w/ other tools; 67 total)



R (2 alone, 51 total)



Your own code (3 alone, 44 total)



KXEN (25 alone, 31 total)



Weka (now Pentaho) (0 alone, 31 total)



MATLAB (0 alone, 26 total)



Other commercial tools (0 alone, 19 total)



KNIME (1 alone, 18 total)



Other free tools (0 alone, 15 total)



Microsoft SQL Server (1 alone, 15 total)



Zementis (5 alone, 13 total)



Oracle DM (0 alone, 9 total)



Statsoft Statistica (0 alone, 8 total)



Orange (0 alone, 5 total)



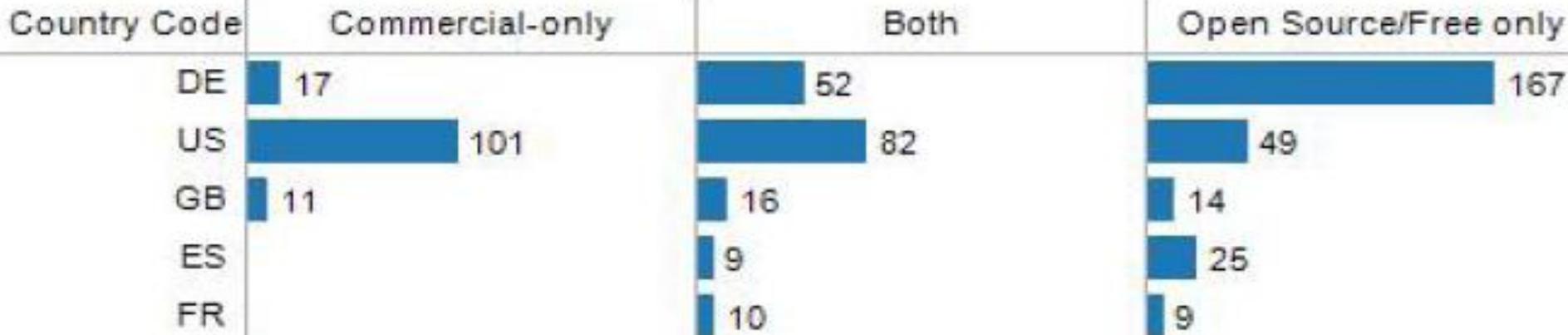
Salford CART, Mars, other (1 alone, 5 total)



C4.5/C5.0 (0 alone, 4 total)

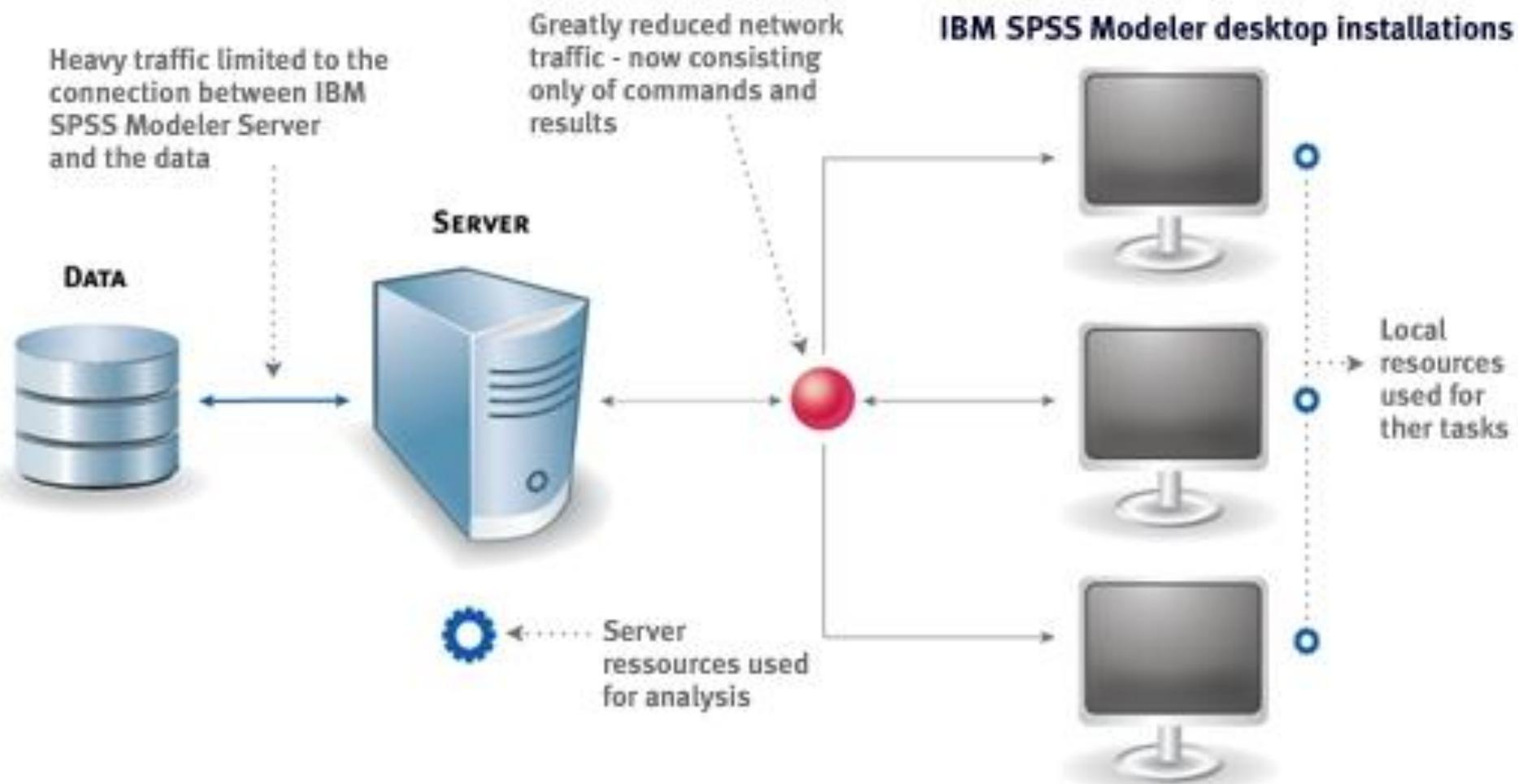


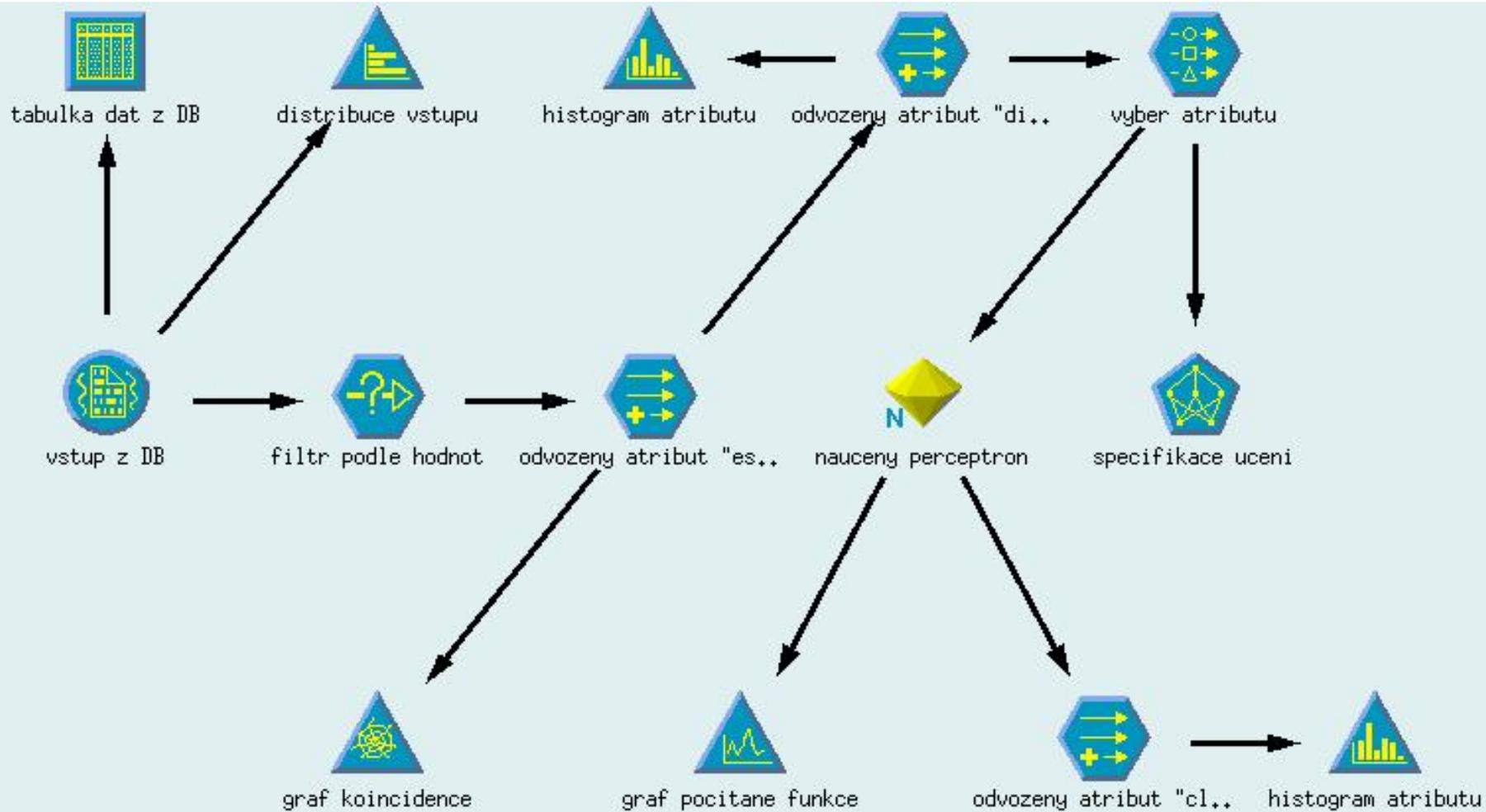
Angoss (0 alone, 4 total)

Our customers include:

- All 50 U.S. state governments
- 22 of the top 24 global commercial banks
- 100% of the top U.S. universities
- The 18 top property and casualty insurance companies in the U.S.
- The leading 12 global pharmaceutical companies





Pentaho Data Mining Enterprise Edition

Data Mining - Classify Panel

Weka Explorer

Preprocess **Classify** Cluster Associate Select attributes Visualize

Classifier: Choose **Logistic -R 1.0E-8 -M -1**

Test options:
 Use training set
 Supplied test set (Set...)
 Cross-validation (Folds: 10)
 Percentage split (%: 66)
More options...

(Nom) class: (Nom) class

Start Stop

Result list (right-click for options):
15:13:01 - functions.Logistic

Weka Classifier Visualize: ThresholdCurve. (Class value good)

X: False Positive Rate (Num) Y: True Positive Rate (Num)
Colour: Threshold (Num) Select Instance

Reset Clear Open Save Jitter

Plot (Area under ROC = 0.7849)

Class colour: 0.031 0.52 1

Classifier output:

```
num_dependents
om_telephone
foreign_worker

Time taken to build classifier: 0.001 seconds

=== Stratified cross-validation ===
=== Summary ===

Correctly Classified Instances 891
Incorrectly Classified Instances 99
Kappa statistic 0.891751
Mean absolute error 0.10825
Root mean squared error 0.328876
Relative absolute error 12.1213%
Root relative squared error 1.10825
Total Number of Instances 1000

=== Detailed Accuracy By Class ===
```

	TP Rate	FP Rate	Precision	Recall	F-Measure	ROC Area	Class
good	0.864	0.51	0.798	0.864	0.83	0.785	good
bad	0.49	0.136	0.607	0.49	0.542	0.785	bad
Weighted Avg.	0.752	0.398	0.741	0.752	0.744	0.785	

```
=== Confusion Matrix ===
 a  b  <-- classified as
605 95 | a = good
153 147 | b = bad
```

Status: OK

Log x0

Pentaho

In addition to ad hoc analysis, Pentaho Data Mining and a suite of algorithms that may otherwise be available. These can be used and also export predictive analytics. There's a high level of integration based on an

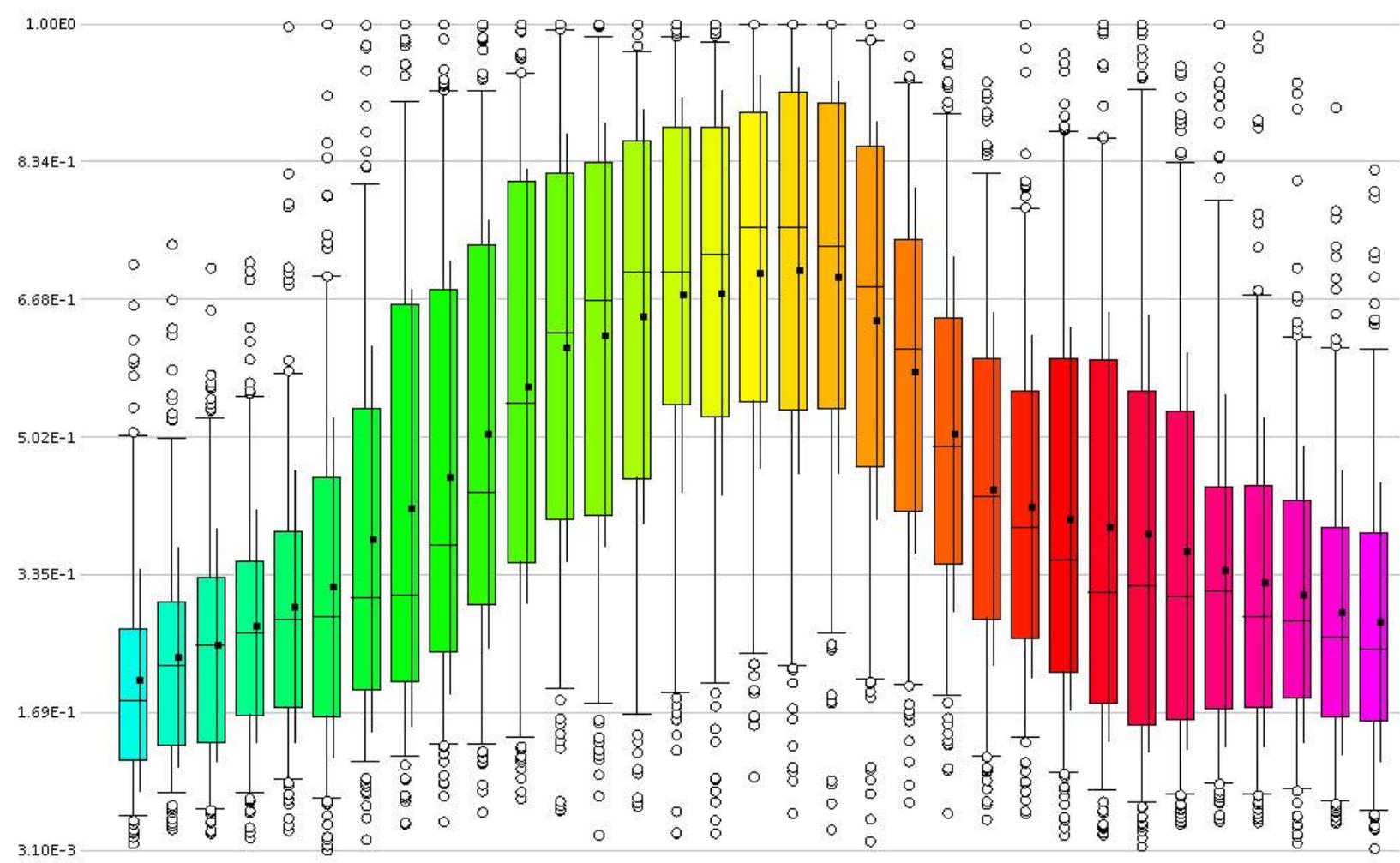
Pentaho Data Mining compliant native tight integration including reporting Suite Enterprise

Pentaho

- An out-of-the-box analyst results Dashboard
- A set of Java OData portals.
- Together



- Plotter
- Quartile
- Dimensions
- attribute_7
 - attribute_8
 - attribute_9
 - attribute_10
 - attribute_11
 - attribute_12
 - attribute_13
 - attribute_14
 - attribute_15
 - attribute_16
 - attribute_17
 - attribute_18
 - attribute_19
 - attribute_20
 - attribute_21
 - attribute_22
 - attribute_23
 - attribute_24
 - attribute_25
 - attribute_26
 - attribute_27
 - attribute_28
 - attribute_29
 - attribute_30
 - attribute_31
 - attribute_32
 - attribute_33
 - attribute_34
 - attribute_35
 - attribute_36
 - attribute_37
 - attribute_38
 - attribute_39
 - attribute_40
 - attribute_41
- Save Image...

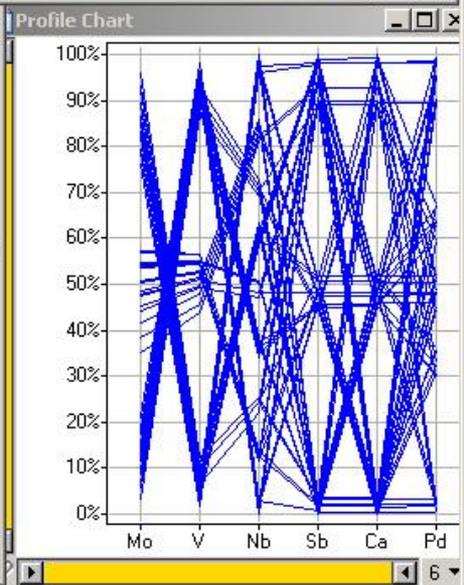
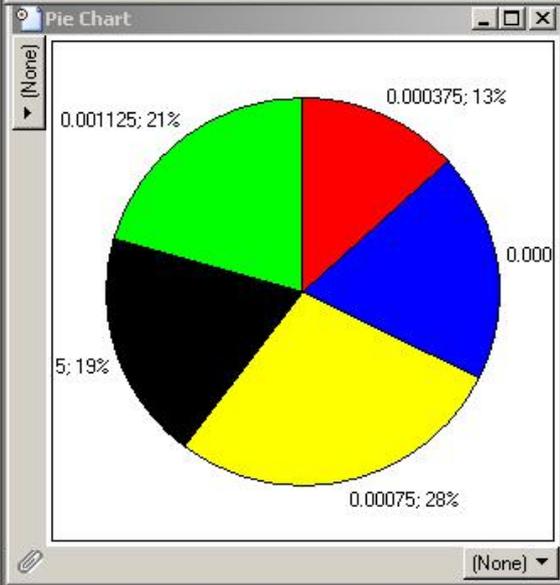
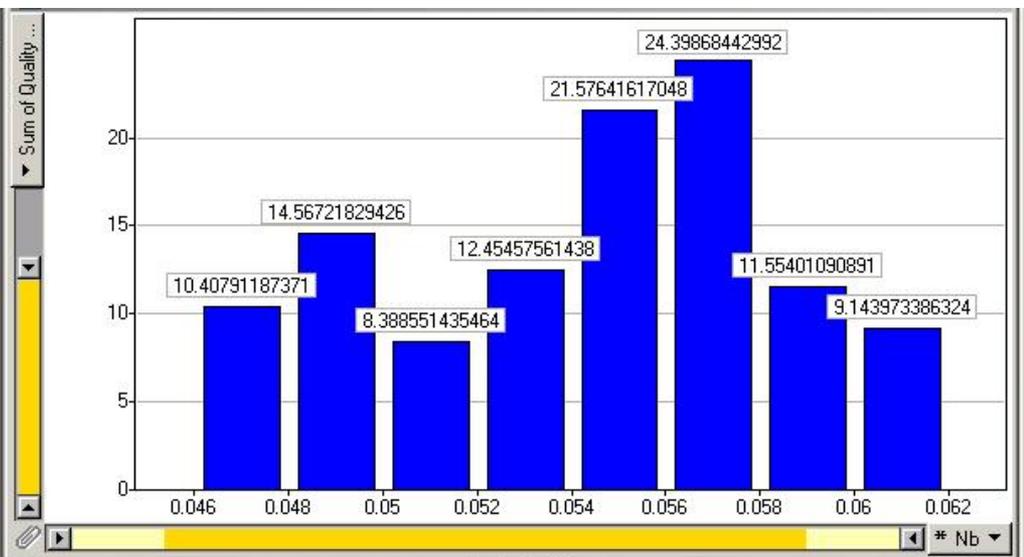
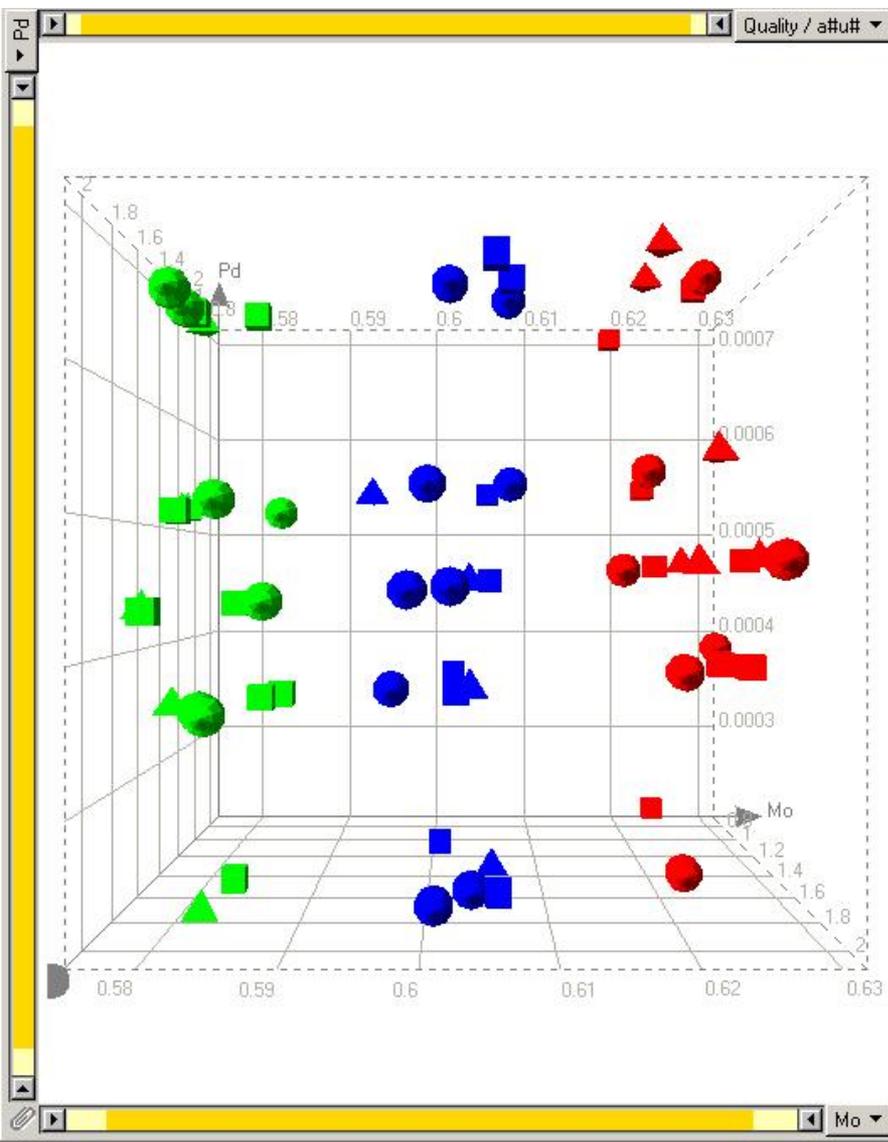


Jan 20, 2010 10:18:24 AM INFO: No filename given for result file, using stdout for logging results!

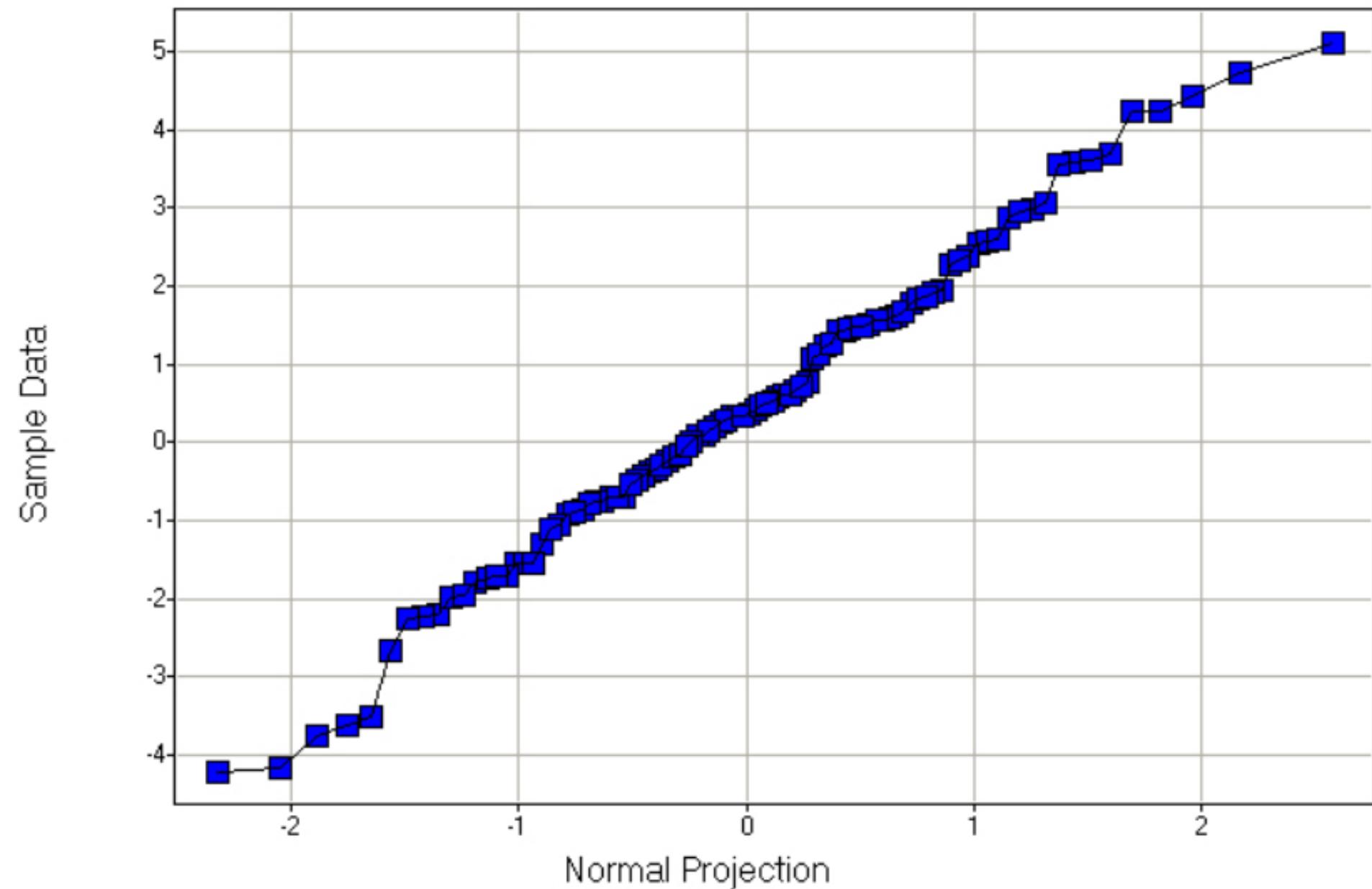
Jan 20, 2010 10:18:24 AM INFO: Process starts

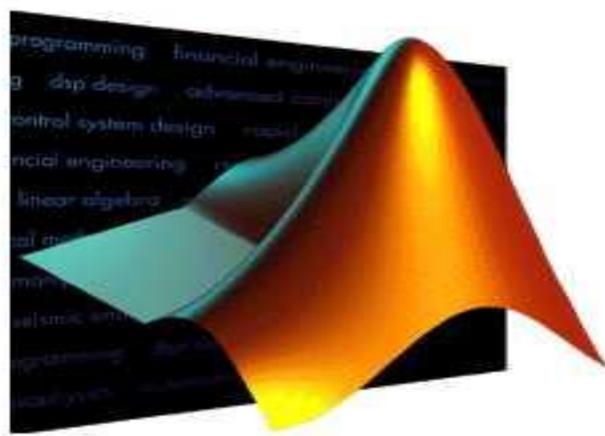
Jan 20, 2010 10:18:24 AM INFO: Loading initial data

Max:	1002 MB
Total:	1002 MB



Normal Probability Plot





REVENUE

- \$500 million in 2009, with more than 50% from outside the United States
- Profitable every year since its founding

FAST FACTS

- Founded in 1984
- Privately held
- There are more than:
 - » 1 million users of MATLAB worldwide
 - » 350 third-party solutions that build on MATLAB and Simulink
 - » 1200 MATLAB based books in 26 languages

Workspace

Name	Value
MessageHandle	<1x1 struct>
gui_Singleton	1
gui_State	<1x1 struct>
varargin	<1x1 cell>

Current Directory Workspace

Command History

```

for kk = 2: size(theta,1)
yy = [yy Ymn];
end
    
```

Array Editor - gui_State

Field	Value
gui_Name	'gamessage'
gui_Singleton	1
gui_OpeningFcn	@gamessage_OpeningFcn
gui_OutputFcn	@gamessage_OutputFcn
gui_LayoutFcn	[]

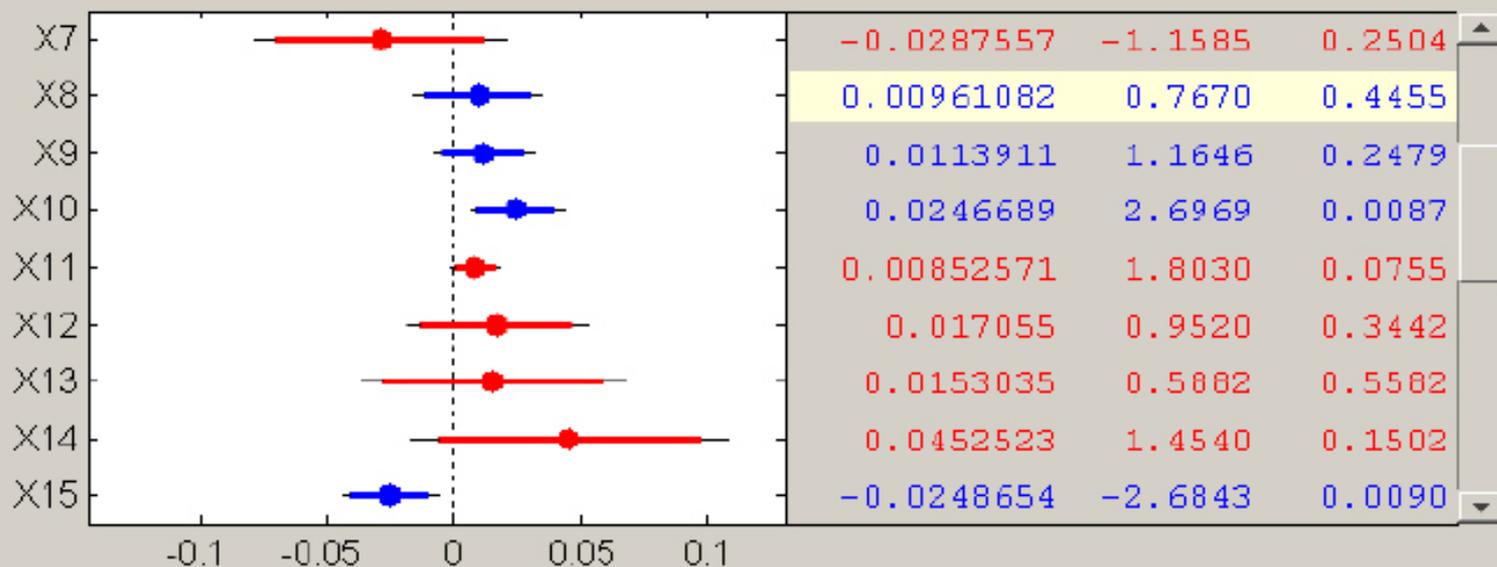
TwoLayers x gui_State

Command Window

```

% Apply spherical coordinate e
r = rho.*sin(theta);
x = r.*cos(phi); % spherica
y = r.*sin(phi);
    
```

Coefficients with Error Bars



Next step:

Move X8 out

Next Step

All Steps

Export ...

Intercept = 0.291283

R-square = 0.316512

F = 8.56706

RMSE = 0.405153

Adj R-sq = 0.270331

p = 9.75464e-006

Model History

