

## Outline

# Managing Large Sets Of Models

R. Seger A. Unwin

Department of Computer Oriented Statistics and Data Analysis  
University of Augsburg

UseR 2006

### 1 Introduction

### 2 Managing Models

- Software Requirements
- Model Repository
- Model Comparison

## Model And Model Space

- Models may be fit from a variety of different classes
  - e.g. lm, gam, rpart
- and to variations of the original dataset
  - the (unmodified) dataset
  - subsets of the data
  - versions using transformed variables
  - datasets from sampling (e.g. bootstrapping)

### Conclusion

Model space size is affected by the range of models, by transformations and by sampling

## Model And Model Space

- Models may be fit from a variety of different classes
  - e.g. lm, gam, rpart
- and to variations of the original dataset
  - the (unmodified) dataset
  - subsets of the data
  - versions using transformed variables
  - datasets from sampling (e.g. bootstrapping)

### Conclusion

Model space size is affected by the range of models, by transformations and by sampling

## Model And Model Space

- Models may be fit from a variety of different classes
  - e.g. lm, gam, rpart
- and to variations of the original dataset
  - the (unmodified) dataset
  - subsets of the data
  - versions using transformed variables
  - datasets from sampling (e.g. bootstrapping)

### Conclusion

Model space size is affected by the range of models, by transformations and by sampling

## Model And Model Space

- Models may be fit from a variety of different classes
  - e.g. lm, gam, rpart
- and to variations of the original dataset
  - the (unmodified) dataset
  - subsets of the data
  - versions using transformed variables
  - datasets from sampling (e.g. bootstrapping)

### Conclusion

Model space size is affected by the range of models, by transformations and by sampling

## Model And Model Space

- Models may be fit from a variety of different classes
  - e.g. lm, gam, rpart
- and to variations of the original dataset
  - the (unmodified) dataset
  - subsets of the data
  - versions using transformed variables
  - datasets from sampling (e.g. bootstrapping)

### Conclusion

Model space size is affected by the range of models, by transformations and by sampling

## Model And Model Space

- Models may be fit from a variety of different classes
  - e.g. lm, gam, rpart
- and to variations of the original dataset
  - the (unmodified) dataset
  - subsets of the data
  - versions using transformed variables
  - datasets from sampling (e.g. bootstrapping)

### Conclusion

Model space size is affected by the range of models, by transformations and by sampling

## Model And Model Space

- Models may be fit from a variety of different classes
  - e.g. lm, gam, rpart
- and to variations of the original dataset
  - the (unmodified) dataset
  - subsets of the data
  - versions using transformed variables
  - datasets from sampling (e.g. bootstrapping)

### Conclusion

Model space size is affected by the range of models, by transformations and by sampling

## Model And Model Space

- Models may be fit from a variety of different classes
  - e.g. lm, gam, rpart
- and to variations of the original dataset
  - the (unmodified) dataset
  - subsets of the data
  - versions using transformed variables
  - datasets from sampling (e.g. bootstrapping)

### Conclusion

Model space size is affected by the range of models, by transformations and by sampling

## Model Analysis

- model statistics (for model selection/comparison)
  - global statistics for comparing overall fit
  - measures of variable importance
  - residuals for comparing local fit

### Conclusion

Models yield much useful data for further analysis and you need a tool to organise that flood of data

## Model Analysis

- model statistics (for model selection/comparison)
  - global statistics for comparing overall fit
  - measures of variable importance
  - residuals for comparing local fit

### Conclusion

Models yield much useful data for further analysis and you need a tool to organise that flood of data

- model statistics (for model selection/comparison)
  - global statistics for comparing overall fit
  - measures of variable importance
  - residuals for comparing local fit

## Conclusion

Models yield much useful data for further analysis and you need a tool to organise that flood of data

- model statistics (for model selection/comparison)
  - global statistics for comparing overall fit
  - measures of variable importance
  - residuals for comparing local fit

## Conclusion

Models yield much useful data for further analysis and you need a tool to organise that flood of data

- model statistics (for model selection/comparison)
  - global statistics for comparing overall fit
  - measures of variable importance
  - residuals for comparing local fit

## Conclusion

Models yield much useful data for further analysis and you need a tool to organise that flood of data

- 1 Introduction
- 2 Managing Models
  - Software Requirements
  - Model Repository
  - Model Comparison

## Technical Requirements

- store model forms and results
- commands in R  $\Leftrightarrow$  series of textual input
- model information includes
  - formula
  - global statistics (e.g. deviance)
  - coefficient estimates
  - residuals

## MORET

In order to manage and compare model results a lot of data has to be stored. The software "MORET" uses a relational database for this (see next slide).

## Technical Requirements

- store model forms and results
- commands in R  $\Leftrightarrow$  series of textual input
- model information includes
  - formula
  - global statistics (e.g. deviance)
  - coefficient estimates
  - residuals

## MORET

In order to manage and compare model results a lot of data has to be stored. The software "MORET" uses a relational database for this (see next slide).

## Technical Requirements

- store model forms and results
- commands in R  $\Leftrightarrow$  series of textual input
- model information includes
  - formula
  - global statistics (e.g. deviance)
  - coefficient estimates
  - residuals

## MORET

In order to manage and compare model results a lot of data has to be stored. The software "MORET" uses a relational database for this (see next slide).

## Technical Requirements

- store model forms and results
- commands in R  $\Leftrightarrow$  series of textual input
- model information includes
  - formula
  - global statistics (e.g. deviance)
  - coefficient estimates
  - residuals

## MORET

In order to manage and compare model results a lot of data has to be stored. The software "MORET" uses a relational database for this (see next slide).

# Technical Requirements

- store model forms and results
- commands in R  $\Leftrightarrow$  series of textual input
- model information includes
  - formula
  - global statistics (e.g. deviance)
  - coefficient estimates
  - residuals

## MORET

In order to manage and compare model results a lot of data has to be stored. The software "MORET" uses a relational database for this (see next slide).

# Technical Requirements

- store model forms and results
- commands in R  $\Leftrightarrow$  series of textual input
- model information includes
  - formula
  - global statistics (e.g. deviance)
  - coefficient estimates
  - residuals

## MORET

In order to manage and compare model results a lot of data has to be stored. The software "MORET" uses a relational database for this (see next slide).

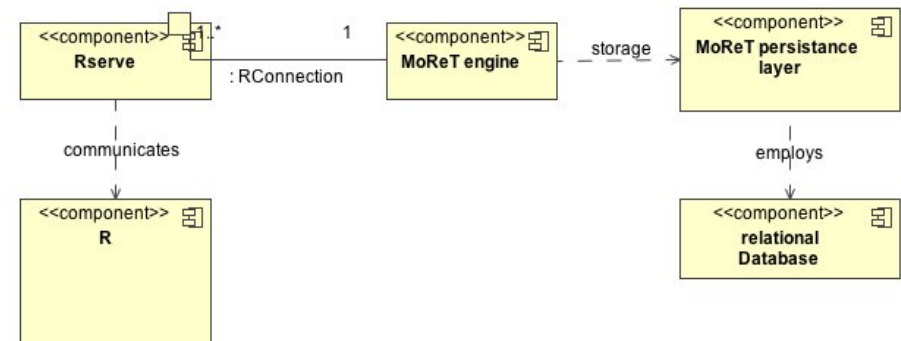
# Technical Requirements

- store model forms and results
- commands in R  $\Leftrightarrow$  series of textual input
- model information includes
  - formula
  - global statistics (e.g. deviance)
  - coefficient estimates
  - residuals

## MORET

In order to manage and compare model results a lot of data has to be stored. The software "MORET" uses a relational database for this (see next slide).

# Used Components



1 Introduction

2 Managing Models

- Software Requirements
- Model Repository
- Model Comparison

- load datasets with the MORET-GUI
- fit models to the data
- use MORET to manage the database of the models.
  - export model data
  - use the Model Explorer



- load datasets with the MORET-GUI
- fit models to the data
- use MORET to manage the database of the models.
  - export model data
  - use the Model Explorer

- load datasets with the MORET-GUI
- fit models to the data
- use MORET to manage the database of the models.
  - export model data
  - use the Model Explorer



## Working With Moret

- load datasets with the MORET-GUI
- fit models to the data
- use MORET to manage the database of the models.
  - export model data
  - use the Model Explorer

## Working With Moret

- load datasets with the MORET-GUI
- fit models to the data
- use MORET to manage the database of the models.
  - export model data
  - use the Model Explorer

## Working With Moret/2

```
Model REpository CLI
File Data Options
Welcome to MORET
read.table('/Users/ralfseger/Documents/workspace/StatisticalModels/data/election.txt',header=TRUE);
All Objects removed from workspace.
election.txt<-read.table('/Users/ralfseger/Documents/workspace/StatisticalModels/data/election.txt',header=TRUE);
save(file="/Users/ralfseger/Documents/workspace/StatisticalModels/binary/data/election.txt",election.txt)
file "/Users/ralfseger/Documents/workspace/StatisticalModels/data/election.txt" added to database. Caching binary data file.
attach(election.txt)

dataset in workspace is 'election.txt'
```

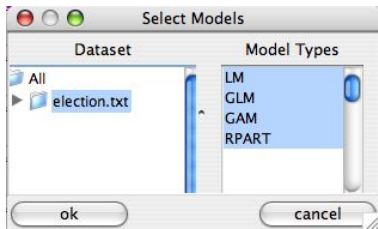
## Outline

- 1 Introduction
- 2 Managing Models
  - Software Requirements
  - Model Repository
  - Model Comparison



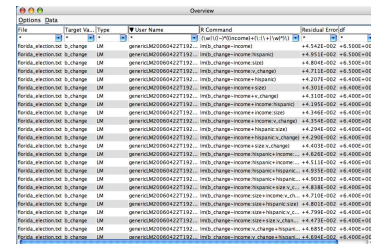
# Model Selection

# Model Selection



## Preselect Models

- Models space is huge.
- Work with selected subsets of all models.



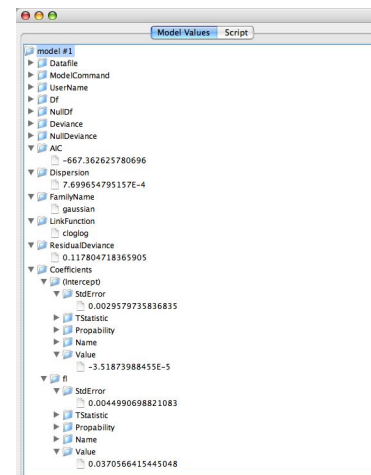
## Preselect Models

- Models space is huge.
- Work with selected subsets of all models.

# Model Overview

# Model Overview Options

File	Target Va...	Type	User Name	R Command	Residual Error	df	#
*	*	*	*	(\w \( ~)*((income)+(?: + \w)*)	*	*	*
florida_election.txt	b_change	LM	genericLM20060422T192...	lm(b_change~income)	+4.542E-002	+6.500E+001	2
florida_election.txt	b_change	LM	genericLM20060422T192...	lm(b_change~income+hispanic)	+4.951E-002	+6.500E+001	6
florida_election.txt	b_change	LM	genericLM20060422T192...	lm(b_change~income+size)	+4.804E-002	+6.500E+001	7
florida_election.txt	b_change	LM	genericLM20060422T192...	lm(b_change~income+v_change)	+4.711E-002	+6.500E+001	8
florida_election.txt	b_change	LM	genericLM20060422T192...	lm(b_change~income+hispanic)	+4.207E-002	+6.400E+001	22
florida_election.txt	b_change	LM	genericLM20060422T192...	lm(b_change~income+size)	+4.301E-002	+6.400E+001	23
florida_election.txt	b_change	LM	genericLM20060422T192...	lm(b_change~income+v_change)	+4.310E-002	+6.400E+001	24
florida_election.txt	b_change	LM	genericLM20060422T192...	lm(b_change~income+hispanic)	+4.195E-002	+6.400E+001	25
florida_election.txt	b_change	LM	genericLM20060422T192...	lm(b_change~income+income+size)	+4.346E-002	+6.400E+001	26
florida_election.txt	b_change	LM	genericLM20060422T192...	lm(b_change~income+income+v_change)	+4.354E-002	+6.400E+001	27
florida_election.txt	b_change	LM	genericLM20060422T192...	lm(b_change~income+size)	+4.294E-002	+6.400E+001	28
florida_election.txt	b_change	LM	genericLM20060422T192...	lm(b_change~income+hispanic+v_change)	+4.290E-002	+6.400E+001	29
florida_election.txt	b_change	LM	genericLM20060422T192...	lm(b_change~income+size+v_change)	+4.403E-002	+6.400E+001	30
florida_election.txt	b_change	LM	genericLM20060422T192...	lm(b_change~income+hispanic+income+size)	+4.626E-002	+6.400E+001	52
florida_election.txt	b_change	LM	genericLM20060422T192...	lm(b_change~income+hispanic+income+size)	+4.511E-002	+6.400E+001	53
florida_election.txt	b_change	LM	genericLM20060422T192...	lm(b_change~income+hispanic+hispanic+size)	+4.935E-002	+6.400E+001	54
florida_election.txt	b_change	LM	genericLM20060422T192...	lm(b_change~income+hispanic+hispanic+size)	+4.903E-002	+6.400E+001	55
florida_election.txt	b_change	LM	genericLM20060422T192...	lm(b_change~income+hispanic+size+v_change)	+4.838E-002	+6.400E+001	56
florida_election.txt	b_change	LM	genericLM20060422T192...	lm(b_change~income+size+income+v_change)	+4.710E-002	+6.400E+001	57
florida_election.txt	b_change	LM	genericLM20060422T192...	lm(b_change~income+size+hispanic+size)	+4.801E-002	+6.400E+001	58
florida_election.txt	b_change	LM	genericLM20060422T192...	lm(b_change~income+size+hispanic+v_change)	+4.799E-002	+6.400E+001	59
florida_election.txt	b_change	LM	genericLM20060422T192...	lm(b_change~income+size+size+v_change)	+4.473E-002	+6.400E+001	60
florida_election.txt	b_change	LM	genericLM20060422T192...	lm(b_change~income+v_change+hispanic+size)	+4.685E-002	+6.400E+001	61
florida_election.txt	b_change	LM	genericLM20060422T192...	lm(b_change~income+v_change+hispanic+size)	+4.694E-002	+6.400E+001	62



## Working with the Model Table

- inspect a single model
- view the creation script
- export models

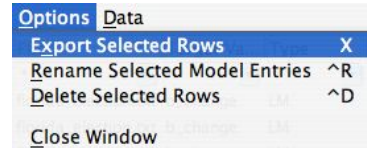
# Model Overview Options

# Model Overview Options



### Working with the Model Table

- inspect a single model
- view the creation script
- export models

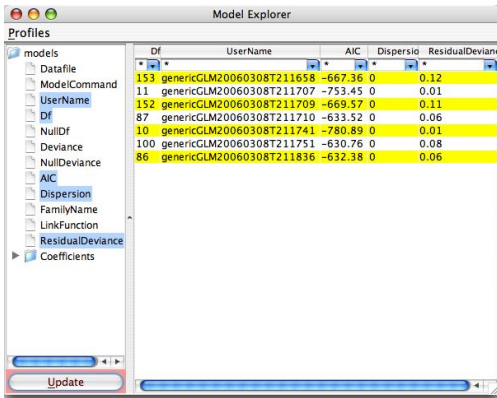


### Working with the Model Table

- inspect a single model
- view the creation script
- export models

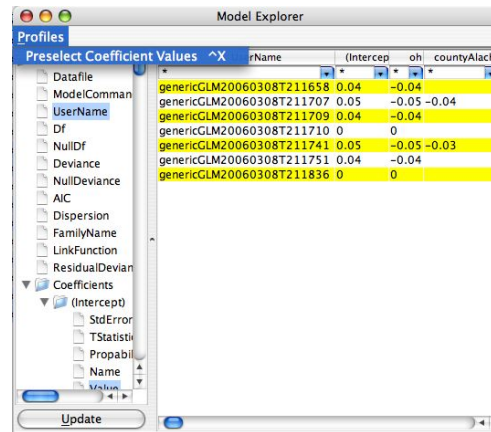
# Model Explorer

# Model Explorer



### Actions

- dynamically select attributes of the stored models
- use predefined profiles (e.g. for all coefficient values)



### Actions

- dynamically select attributes of the stored models
- use predefined profiles (e.g. for all coefficient values)

# Using External Applications

# Using External Applications

## External Interfaces

- export models to XML or CSV
- copy and paste to a spreadsheet

## External Interfaces

- export models to XML or CSV
- copy and paste to a spreadsheet



# Using External Applications

# Summary

## External Interfaces

- export models to XML or CSV
- copy and paste to a spreadsheet

A	B	C	D	E	F	G	H	I	J	K	L	M	
1	UserName	(Intercept)	fi	sh	countyAlachua	countyAllen	countyAshland	countyAshTabul	countyAthens	countyAurlice	countyBaker	countyBay	countyBelmont
2	generisGLM200	-3.518739884	0.0370566415445048										
3	generisGLM200	0.0370214541456593	-0.0370566415445048										
4	generisGLM200	-4.792809486391E-4		0.0147326598	-0.005589544	-0.013851523	-0.008180469	-0.0637169770	0.0298029184	0.0796021819	0.0454324483	0.0293666720	
5	generisGLM200	0.0095431208610535											
6	generisGLM200	-3.518739884	0.0370566415445048										
7	generisGLM200	-4.792809486	0.0513064290086428		-0.036573770	-0.005589544	-0.013851523	-0.008180469	-0.0637169770	0.0298029184	0.0282957529	-0.005873980	0.0293666720
8	generisGLM200	-3.518739884	0.0095782082598991										
9	generisGLM200	0.050827148060004		-0.051306429	-0.036573770	-0.005589544	-0.013851523	-0.008180469	-0.0637169770	0.0298029184	0.0282957529	-0.005873980	0.0293666720
10	generisGLM200	0.0095431208610535		-0.0095783082598991									
11	generisGLM200	-4.792809486399E-4		0.0147326598	-0.005589544	-0.013851523	-0.008180469	-0.0637169770	0.0298029184	0.0796021819	0.0454324483	0.0293666720	
12	generisGLM200	-4.792809486	0.0513064290086428		-0.036573770	-0.005589544	-0.013851523	-0.008180469	-0.0637169770	0.0298029184	0.0282957529	-0.005873980	0.0293666720
13	generisGLM200	-3.518739884	0.0095782082598991										
14	generisGLM200	-4.792809486	0.0634346604347228		-0.048702001	-0.005589544	-0.013851523	-0.008180469	-0.0637169770	0.0298029184	0.0161675214	-0.0180022120	0.0293666720
15	generisGLM200	0.0629553794860839		-0.063434660	-0.048702001	-0.005589544	-0.013851523	-0.008180469	-0.0637169770	0.0298029184	0.0161675214	-0.0180022120	0.0293666720
16	generisGLM200	-4.792809486	0.0634346604347228		-0.048702001	-0.005589544	-0.013851523	-0.008180469	-0.0637169770	0.0298029184	0.0161675214	-0.0180022120	0.0293666720
17													
18													
19													
20													
21													
22													

- **Managing large sets of models** is possible using MORET
- **MORET** can be used as link to other applications for further analysis

- Further developments
  - Integration of more model alternatives
  - Better control of transformed variables

- **Managing large sets of models** is possible using MORET
- **MORET** can be used as link to other applications for further analysis
  
- Further developments
  - Integration of more model alternatives
  - Better control of transformed variables



R. Seger

*project homepage.*

<http://stats.math.uni-augsburg.de/software/>.