

Automating Business Modeling with the AutoModelR package

Derek McCrae Norton

InterContinental Hotels Group



Outline

1 Introduction

- The Reasoning Behind AutoModelR
- What is AutoModelR?

2 AutoModelR Specifics

- Data Exploration and Reduction
- Modeling / Model Assesment and Selection

Outline

- 1 Introduction
 - The Reasoning Behind AutoModelR
 - What is AutoModelR?
- 2 AutoModelR Specifics
 - Data Exploration and Reduction
 - Modeling / Model Assesment and Selection

Why AutoModelR?

... or should I keep my eyes open for the remainder of this talk?

- The modeling process can be lengthy.
- A standardized process can be defined.
- Many steps can be simple in implementation, if not in scope.
- A standardized report can aid in understanding.

Why AutoModelR?

... or should I keep my eyes open for the remainder of this talk?

- The modeling process can be lengthy.
- A standardized process can be defined.
- Many steps can be simple in implementation, if not in scope.
- A standardized report can aid in understanding.

Why AutoModelR?

... or should I keep my eyes open for the remainder of this talk?

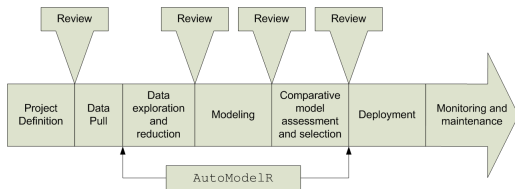
- The modeling process can be lengthy.
- A standardized process can be defined.
- Many steps can be simple in implementation, if not in scope.
- A standardized report can aid in understanding.

Why AutoModelR?

... or should I keep my eyes open for the remainder of this talk?

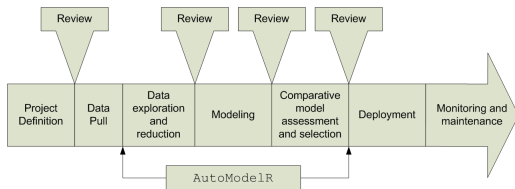
- The modeling process can be lengthy.
- A standardized process can be defined.
- Many steps can be simple in implementation, if not in scope.
- A standardized report can aid in understanding.

Modeling Process



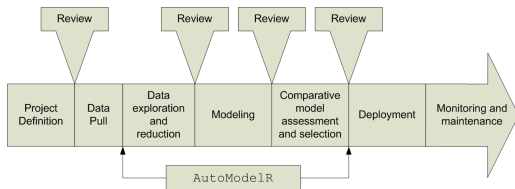
- Process is standardized in a broad sense.
- Process is followed with or without automation.
- When 40+ requests are sitting in queue, any savings of time is *greatly* appreciated.

Modeling Process



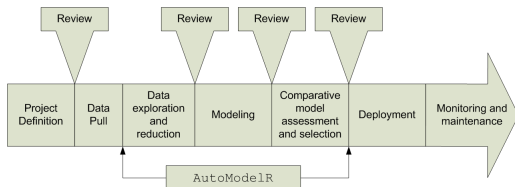
- Process is standardized in a broad sense.
- Process is followed with or without automation.
- When 40+ requests are sitting in queue, any savings of time is *greatly* appreciated.

Modeling Process



- Process is standardized in a broad sense.
- Process is followed with or without automation.
- When 40+ requests are sitting in queue, any savings of time is *greatly* appreciated.

Modeling Process



- Process is standardized in a broad sense.
- Process is followed with or without automation.
- When 40+ requests are sitting in queue, any savings of time is *greatly* appreciated.

Outline

- 1 Introduction
 - The Reasoning Behind AutoModelR
 - What is AutoModelR?
- 2 AutoModelR Specifics
 - Data Exploration and Reduction
 - Modeling / Model Assesment and Selection

What is AutoModelR?

- Automation of three standardized tasks.
 - 1 Data Exploration and Reduction
 - 2 Modeling
 - 3 Model Assesment and Selection
- Most of the focus is on task 1.
- Tasks 2 and 3 are generally first steps to more involved work.

What is AutoModelR?

- Automation of three standardized tasks.
 - 1 Data Exploration and Reduction
 - 2 Modeling
 - 3 Model Assesment and Selection
- Most of the focus is on task 1.
- Tasks 2 and 3 are generally first steps to more involved work.

What is AutoModelR?

- Automation of three standardized tasks.
 - 1 Data Exploration and Reduction
 - 2 Modeling
 - 3 Model Assesment and Selection
- Most of the focus is on task 1.
- Tasks 2 and 3 are generally first steps to more involved work.

What is AutoModelR?

- Automation of three standardized tasks.
 - 1 Data Exploration and Reduction
 - 2 Modeling
 - 3 Model Assesment and Selection
- Most of the focus is on task 1.
- Tasks 2 and 3 are generally first steps to more involved work.

What is AutoModelR?

- Automation of three standardized tasks.
 - 1 Data Exploration and Reduction
 - 2 Modeling
 - 3 Model Assessment and Selection
- Most of the focus is on task 1.
- Tasks 2 and 3 are generally first steps to more involved work.

What is AutoModelR?

- Automation of three standardized tasks.
 - 1 Data Exploration and Reduction
 - 2 Modeling
 - 3 Model Assessment and Selection
- Most of the focus is on task 1.
- Tasks 2 and 3 are generally first steps to more involved work.

Outline

- 1 Introduction
 - The Reasoning Behind AutoModelR
 - What is AutoModelR?
- 2 AutoModelR Specifics
 - Data Exploration and Reduction
 - Modeling / Model Assessment and Selection

General Steps

- Drop variables with a high percentage of missing values.
- Create indicator variables for presence in variables with a medium percentage of missing values.
- Impute variables with a low percentage of missing values.
- Drop variables with zero variation.
- Add log transformed variables.
- Conduct basic filter variable selection based on relevance and redundancy.

General Steps

- Drop variables with a high percentage of missing values.
- Create indicator variables for presence in variables with a medium percentage of missing values.
- Impute variables with a low percentage of missing values.
- Drop variables with zero variation.
- Add log transformed variables.
- Conduct basic filter variable selection based on relevance and redundancy.

General Steps

- Drop variables with a high percentage of missing values.
- Create indicator variables for presence in variables with a medium percentage of missing values.
- Impute variables with a low percentage of missing values.
- Drop variables with zero variation.
- Add log transformed variables.
- Conduct basic filter variable selection based on relevance and redundancy.

General Steps

- Drop variables with a high percentage of missing values.
- Create indicator variables for presence in variables with a medium percentage of missing values.
- Impute variables with a low percentage of missing values.
- Drop variables with zero variation.
- Add log transformed variables.
- Conduct basic filter variable selection based on relevance and redundancy.

General Steps

- Drop variables with a high percentage of missing values.
- Create indicator variables for presence in variables with a medium percentage of missing values.
- Impute variables with a low percentage of missing values.
- Drop variables with zero variation.
- Add log transformed variables.
- Conduct basic filter variable selection based on relevance and redundancy.

General Steps

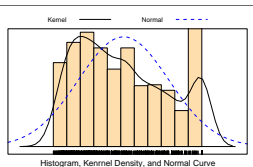
- Drop variables with a high percentage of missing values.
- Create indicator variables for presence in variables with a medium percentage of missing values.
- Impute variables with a low percentage of missing values.
- Drop variables with zero variation.
- Add log transformed variables.
- Conduct basic filter variable selection based on relevance and redundancy.

General Steps II

- Create summary tables.
- Create plots. (Thanks to Jim Porzak)

IHG_TOTAL_SHARE

MIN	MEDIAN	MAX	MEAN	STD DEV
0.0125	0.4286	1	0.4798	0.3008
N	UNIQUE	MISSING	SKEWNESS	KURTOSIS
1000	326	0	0.3837	-1.066



MI_EMAIL_PREF_VAL_CD

	Label	Frequency	Barchart
1	25	31.6	
2	9	18.6	
3	8	15.7	
4	0	15.2	
5	17	4.7	
6	1	4.2	
7	24	4	
8	57	1.6	
9	13	1	
10	** Others (23 Levels) **	3.4	

This Space For Rent
Apply Inside

Outline

- 1 Introduction
 - The Reasoning Behind AutoModelR
 - What is AutoModelR?
- 2 AutoModelR Specifics
 - Data Exploration and Reduction
 - Modeling / Model Assessment and Selection

Modeling Goals

- As yet incomplete.
- Create a portfolio of models for particular types of dependent variables
 - Continuous
 - Binary Categorical
 - Multicategory Categorical

Modeling Goals

- As yet incomplete.
- Create a portfolio of models for particular types of dependent variables
 - Continuous
 - Binary Categorical
 - Multicategory Categorical

Modeling Goals

- As yet incomplete.
- Create a portfolio of models for particular types of dependent variables
 - Continuous
 - Binary Categorical
 - Multicategory Categorical

Summary

- The initial goal of AutoModelR was to automate as much as possible of the Exploratory Data Analysis in the modeling process.
 - This is achieved with a Sweave document that is intelligent enough to make choices similar to mine.
- The goal is now to expand this to include automated initial modeling and model comparison.

Outlook

- Continue to add initial models to be fit.
- Consider a rewrite specifically for time series data.

Summary

- The initial goal of AutoModelR was to automate as much as possible of the Exploratory Data Analysis in the modeling process.
 - This is achieved with a Sweave document that is intelligent enough to make choices similar to mine.
- The goal is now to expand this to include automated initial modeling and model comparison.

Outlook

- Continue to add initial models to be fit.
- Consider a rewrite specifically for time series data.

Summary

- The initial goal of AutoModelR was to automate as much as possible of the Exploratory Data Analysis in the modeling process.
 - This is achieved with a Sweave document that is intelligent enough to make choices similar to mine.
- The goal is now to expand this to include automated initial modeling and model comparison.

Outlook

- Continue to add initial models to be fit.
- Consider a rewrite specifically for time series data.

Screenshots

EDA Report for Data Frame x.fact

August 1, 2008

Emerging Technologies

© 2008 IBM Corp. All rights reserved. IBM, the IBM logo, and the e-business logo are trademarks of International Business Machines Corporation. Other trademarks and registered trademarks are the property of their respective owners.

Dataset Information EDA Report for Data Frame x.fact

1 Dataset Information

1.1 Independent Variable Assignments

The variables from the data frame x.fact have been mapped to the following planes:

Numeric Variables: 4 variables
Categorical Variables: 8 variables
Binary Variables: 1 variable
Variables With More Than 100 Missing Values: 8 variables
Variables With Between 101 and 1000 Missing Values: 8 variables

1.2 Data Demographics

Initial Data Summary:

- Total Variables: 20
- Numeric Variables: 12
- Categorical Variables: 8
- Total Observations: 8000

Data Summary After Data Cleaning and Variable Selection:

- Numeric Variables: 12
- Categorical Variables: 8
- Binary Variables: 0

© Version 3.1.0 (2008-04-01)
File Name: 080400001.rpt Page: 1 of 1

1.3 Categorical Variable Summaries

EDA Report for Data Frame x.fact

1.1 Independent Variable Information

1.1.1 Numeric Variable Summaries

Variable Name	Min	Q1	Q2	Q3	Max	Missing	Percent
CR_CREDIT_RAT	762	869	917	1000	1000	0	0
CR_CREDIT_RAT2	236	869	917	1000	1000	0	0
CR_CREDIT_RAT3	0	869	917	1000	1000	0	0
CR_CREDIT_RAT4	1074	869	917	1000	1000	0	0
CR_CREDIT_RAT5	869	917	1000	1000	1000	0	0
CR_CREDIT_RAT6	7	0	1000	1000	1000	0	0
CR_CREDIT_RAT7	0	0	1000	1000	1000	0	0
CR_CREDIT_RAT8	0	0	1000	1000	1000	0	0
CR_CREDIT_RAT9	0	0	1000	1000	1000	0	0
CR_CREDIT_RAT10	0	0	1000	1000	1000	0	0
CR_CREDIT_RAT11	0	0	1000	1000	1000	0	0
CR_CREDIT_RAT12	0	0	1000	1000	1000	0	0
CR_CREDIT_RAT13	0	0	1000	1000	1000	0	0
CR_CREDIT_RAT14	0	0	1000	1000	1000	0	0
CR_CREDIT_RAT15	0	0	1000	1000	1000	0	0
CR_CREDIT_RAT16	0	0	1000	1000	1000	0	0
CR_CREDIT_RAT17	0	0	1000	1000	1000	0	0
CR_CREDIT_RAT18	0	0	1000	1000	1000	0	0
CR_CREDIT_RAT19	0	0	1000	1000	1000	0	0
CR_CREDIT_RAT20	0	0	1000	1000	1000	0	0

1.1.2 Categorical Variable Summaries

Variable Name	Min	Q1	Q2	Q3	Max	Missing	Percent
CR_CREDIT_RAT	0	0	0	0	0	0	0
CR_CREDIT_RAT2	0	0	0	0	0	0	0
CR_CREDIT_RAT3	0	0	0	0	0	0	0
CR_CREDIT_RAT4	0	0	0	0	0	0	0
CR_CREDIT_RAT5	0	0	0	0	0	0	0
CR_CREDIT_RAT6	0	0	0	0	0	0	0
CR_CREDIT_RAT7	0	0	0	0	0	0	0
CR_CREDIT_RAT8	0	0	0	0	0	0	0
CR_CREDIT_RAT9	0	0	0	0	0	0	0
CR_CREDIT_RAT10	0	0	0	0	0	0	0
CR_CREDIT_RAT11	0	0	0	0	0	0	0
CR_CREDIT_RAT12	0	0	0	0	0	0	0
CR_CREDIT_RAT13	0	0	0	0	0	0	0
CR_CREDIT_RAT14	0	0	0	0	0	0	0
CR_CREDIT_RAT15	0	0	0	0	0	0	0
CR_CREDIT_RAT16	0	0	0	0	0	0	0
CR_CREDIT_RAT17	0	0	0	0	0	0	0
CR_CREDIT_RAT18	0	0	0	0	0	0	0
CR_CREDIT_RAT19	0	0	0	0	0	0	0
CR_CREDIT_RAT20	0	0	0	0	0	0	0

© Version 3.1.0 (2008-04-01)
File Name: 080400001.rpt Page: 1 of 1