# Overview of Predictive Modeling for Actuaries

Michigan Actuarial Society March 18, 2015

Louise Francis, FCAS, MAAA, Francis Analytics and Actuarial Data Mining www.data-mines.com Louise\_francis@msn.com

#### Why Predictive Modeling?

- Better use of data than traditional methods
- Advanced methods for dealing with messy data now available
- New ways to test and validate models



# Real Life Insurance Application – The "Boris Gang"

#### R New York Fraud Ring No Surprise to Russian Drivers

By SABRINA TAVERNISE

New Yorkers may have been shocked by news of an insurance scheme that involved fake car crashes. But in Russia, fraud is a rule of the road.

August 16, 2003 | WORLD | NEWS MORE ON ORGANIZED CRIME AND: FRAUDS AND SWINDLING, FOREIGN BANK ACCOUNTS, AUTOMOBILE INSURANCE AND LIABILITY, STATE FARM INSURANCE COS, NEW YORK CITY, RUSSIA, LONG ISLAND (NY)

#### 🐻 Investigators Say Fraud Ring Staged Thousands of Crashes

By PATRICK HEALY

The ring used Russian immigrants to stage car accidents and then employed its own network of doctors and fake clinics in New York State to bilk an insurance company out of \$48 million.

August 13, 2003 | FRONT PAGE | NEWS

MORE ON ORGANIZED CRIME AND: ACCIDENTS AND SAFETY, FRAUDS AND SWINDLING, FOREIGN BANK ACCOUNTS, CHILDREN AND YOUTH, AGED, WOMEN, AUTOMOBILE INSURANCE AND LIABILITY, SPOTA, THOMAS J, STATE FARM INSURANCE COS, NEW YORK CITY, RUSSIA, WESTCHESTER COUNTY (NY), LONG ISLAND (NY), SWITZERLAND

#### Kinds of Applications

Classification
Target variable is categorical
Prediction
Target variable is numeric

## A Casualty Actuary's Perspective on Data Modeling

- The Stone Age: 1914 ...
  - Simple deterministic methods
    - Use of blunt instruments: the analytical analog of bows and arrows
  - Often ad-hoc
  - Slice and dice data
  - Based on empirical data little use of parametric models
- The Pre Industrial age: 1970 ...
  - Fit probability distribution to model tails
  - Simulation models and numerical methods for variability and uncertainty analysis
  - Focus is on underwriting, not claims
- The Industrial Age 1985 ...
  - Begin to use computer catastrophe models
- The 20<sup>th</sup> Century 1990...
  - European actuaries begin to use GLMs
- The Computer Age 1996...
  - Begin to discuss data mining at conferences
  - At end of 20<sup>st</sup> century, large consulting firms starts to build a data mining practice
- The Current era A mixture of above
  - In personal lines, modeling the rule rather than the exception
    - Often GLM based, though GLMs evolving to GAMs
  - Commercial lines beginning to embrace modeling

#### Data Complexities: Nonlinearities MARS Prediction of Primary Paid Severity



#### Major Kinds of Data Mining

- Supervised learning
  - Most common situation
  - A dependent variable
    - Frequency
    - Loss ratio
    - Fraud/no fraud
  - Some methods
    - Regression
    - CART
    - Some neural networks
    - MARS

- Unsupervised learning
  - No dependent variable
  - Group like records together
    - A group of claims with similar characteristics might be more likely to be fraudulent
    - Ex: Territory assignment, Text Mining
  - Some methods
    - Principal Components
    - K-means clustering
    - Kohonen neural networks

### Methods

- Classical
- Decision Trees
- Neural Networks
- Unsupervised learning
  - Clustering
- Newer Methods
  - Ensemble
  - SVM
  - Deep learning
  - Text Mining



#### **Classical Statistics: Regression**

 Estimation of parameters: Fit line that minimizes deviation between actual and fitted values

$$\min(\sum (Y_i - \widehat{Y})^2)$$



#### Linear Modeling Tools Widely Available: Excel Analysis Toolpak

- Install Data Analysis
   Tool Pak (Add In) that comes with Excel
- Click Tools, Data
   Analysis, Regression

Regression			? 🛛
Input Input <u>Y</u> Range: Input <u>X</u> Range:	\$H\$11:\$H\$23 \$J\$11:\$J\$23		OK Cancel
□ Labels □ □ Confidence Level: 9	Constant is <u>Z</u> ero		Help
Output options	\$5\$4	3	
Residuals           Residuals           Standardized Residuals           Normal Probability           Vormal Probability	✓ Residual Plots ✓ Line Fit Plots		

## Goodness of Fit

SUMMARY OUTPUT					
Regression Stati	stics				
Multiple R	0.303				
Adjusted R Square	0.092				
Standard Error	1.206				
Observations	1818				
ANOVA					
	df	SS	MS	F	Significance F
Regression	1	266.29	266.29	183.07	0.00
Residual	1816	2641.50	1.45		
Total	1817	2907.79			

#### **Classical Model: Discriminant Analysis**



# Generalized Linear Models (GLMs)

- Relax normality assumption
  - Exponential family of distributions
- Models some kinds of nonlinearity

# Similarities with GLMs

#### Linear Models

- Transformation of Variables
- Use dummy coding for categorical variables
- Residual
- Test significance of coefficients

#### <u>GLMs</u>

Link functions

- Use dummy coding for categorical variables
- Deviance
- Test significance of coefficients

## Linear Model vs GLM

• Regression:

• GLM:

 $Y_{i} = \sim_{i} + \vee$   $\sim_{i} = X ' B$   $\vee \sim N(0, †^{2})$   $Y = h(\sim_{i}) + \vee$   $h(\sim) = X ' B$   $\vee \sim \text{exponential family}$ h is a link function

# **Estimating** Parameters

- As with nonlinear regression, there usually is not a closed form solution for GLMs
- A numerical method used to solve for parameters
- For some models this could be programmed in Excel but statistical software is the usual choice
- If you can't spend money on the software, download R for free

## GLM fit for Poisson Regression

>devage<-as.factor((AGE)

>claims.glm<-glm(Claims~devage, family=poisson)</p>

- >summary(claims.glm)
- Call:
- glm(formula = Claims ~ devage, family = poisson)
- Deviance Residuals:
- Min 1Q Median 3Q Max
- ► -10.250 -1.732 -0.500 0.507 10.626
- Coefficients:
- Estimate Std. Error z value Pr(> | z | )
- (Intercept) 4.73540 0.02825 167.622 < 2e-16 \*\*\*</p>
- devage2 -0.89595 0.05430 -16.500 < 2e-16 \*\*\*</p>
- devage3 -4.32994 0.29004 -14.929 < 2e-16 \*\*\*</p>
- devage4 -6.81484 1.00020 -6.813 9.53e-12 \*\*\*
- -
- Signif. codes: 0 `\*\*\*' 0.001 `\*\*' 0.01 `\*' 0.05 `.' 0.1 ` ' 1
- (Dispersion parameter for poisson family taken to be 1)
- Null deviance: 2838.65 on 36 degrees of freedom
- Residual deviance: 708.72 on 33 degrees of freedom
- AIC: 851.38



# Data Complexities: Missing Data

- It is not uncommon for one third of the possible predictors to contain records with missing values
- Possible solutions:
  - A data mining method such as CART that uses a statistical algorithm to find an alternative parameterization in the presence of missing data
  - A statistical method such as expectation maximization or data imputation to fill in a value

# Data

- Data Management
- Data quality
  - Francis, "Dancing With Dirty Data", CAS forum, <u>www.casact.org</u>
  - CAS Working Party, "Actuarial IQ", www.casact.org
- Big Data



# **Examples of Applications**

- Claim Frequency, Claim Severity
  - Use features of data to predict
  - Chapter in Predictive Modeling book
  - www.casact.org, "Intro to GLMs"
- Liklihood a claim will occur (life insurance)
- Insurance Fraud
  - Derrig and Francis "Distinguishing the Forest from the Trees", Variance, 2008
- Financial Crisis
  - Could the defaulting mortgages have been predicted?
  - Francis and Prevosto, "Data and Disaster: The Role of Data in the Financial Crisis"

# The Questionable Claims Study Data

- 1993 AIB closed PIP claims
- Simulated data based on research performed on original data
- Dependent Variables
  - Suspicion Score
  - Expert assessment of likelihood of fraud or abuse
- Predictor Variables
  - Red flag indicators
  - Claim file variables

Francis Analytics and Actuarial Data Mining, Inc.

3/20/2015

22

#### The Fraud Red Flags

- Binary variables that capture characteristics of claims associated with fraud and abuse
- Accident variables (acc01 acc19)
- Injury variables (inj01 inj12)
- Claimant variables (ch01 ch11)
- Insured variables (ins01 ins06)
- Treatment variables (trt01 trt09)
- Lost wages variables (Iw01 Iw07)

#### The Fraud Problem

from: www.agentinsure.com

X< RETURN TO ARTICLES LIST

#### In Florida, Cops are Cracking Down on Car Insurance Fraud Rings

Recent raids by police in the Tampa Bay, Florida area are shedding light on a serious problem that's plaguing the car insurance industry and having a negative impact on the ability of law abiding citizens to buy auto insurance at low prices.

According to statistics, incidents of staged car accidents in Tampa Bay are much higher than they are anywhere else. One of the reasons for this is said to be Florida's existing PIP law, which stands for Personal Injury Protection, and guarantees as much as \$10,000 in medical payments for every person injured in a car accident, no matter who's to blame for it.

Francis Analytics and Actuarial Data Mining, Inc.

# Fraud and Abuse

- Planned fraud
  - Staged accidents
- Abuse

25

- Opportunistic
- Exaggerate claim
- Both are referred to as "questionable claims"

Francis Analytics and Actuarial Data Mining, Inc.

#### Neural Networks

- Theoretically based on how neurons function
- Can be viewed as a complex non-linear regression

See Francis, "Neural Networks Demystified",
 CAS Forum, 2001, www.casact.org



## Hidden Layer of Neural Network (Input Transfer Function)



# Assessing Results

Confusion Matrix

ROC Curve

		Predicted			
Sample	Observed	1	2	Percent Correct	
Training	1	608	95	86.5%	
	2	52	278	84.2%	
	Overall Percent	63.9%	36.1%	85.8%	
Testing	1	271	36	88.3%	
	2	24	136	85.0%	
	Overall Percent	63.2%	36.8%	87.2%	

Classification

Dependent Variable: Suspicion



#### **Regression** Trees

- Tree-based modeling for continuous target variable
  - most intuitively appropriate method for loss ratio analysis
- Find split that produces greatest separation in

#### ∑[**y – E(y)**]²

- i.e.: find nodes with minimal within variance
  - and therefore greatest between variance
  - like credibility theory i.e.: find nodes with minimal within variance
- Every record in a node is assigned the same expectation → model is a step function





# Different Kinds of Decision Trees

- Single Trees (CART, CHAID)
- Ensemble Trees, a more recent development (TREENET, RANDOM FOREST)
  - A composite or weighted average of many trees (perhaps 100 or more)
  - There are many methods to fit the trees and prevent overfitting
    - Boosting: Iminer Ensemble and Treenet
    - Bagging: Random Forest

#### The Methods and Software Evaluated

- 1) TREENET 5) Iminer Ensemble
- 2) Iminer Tree 6) Random Forest
- 3) SPLUS Tree 7) Naïve Bayes (Baseline)
- 4) CART 8) Logistic (Baseline)

#### **Ensemble Prediction of Total Paid**





## The Fraud Surrogates used as Dependent Variables

- Independent Medical Exam (IME) requested
- Special Investigation Unit (SIU) referral
- IME successful
- SIU successful
- DATA: Detailed Auto Injury Claim Database for Massachusetts
- Accident Years (1995-1997)

# **Results for IME Requested**

Area Under the ROC Curve – IME Decision					
	CART	S-PLUS			
	Tree	Tree	Iminer Tree	TREENET	
AUROC	0.669	0.688	0.629	0.701	
Lower Bound	0.661	0.680	0.620	0.693	
Upper Bound	0.678	0.696	0.637	0.708	
	Iminer	Random	Iminer		
	Ensemble	Forest	Naïve Bayes	Logistic	
AUROC	0.649	703	0.676	0.677	
Lower Bound	0.641	695	0.669	0.669	
Upper Bound	0.657	711	0.684	0.685	

# **Results for SIU Referral**

Area Under the ROC Curve – SIU Decision					
	CART	S-PLUS			
	Tree	Tree	<b>Iminer Tree</b>	TREENET	
AUROC	0.607	0.616	0.565	0.643	
Lower Bound	0.598	0.607	0.555	0.634	
Upper Bound	0.617	0.626	0.575	0.652	
	Iminer	Random	Iminer		
	Ensemble	Forest	Naïve Bayes	Logistic	
AUROC	0.539	0.677	0.615	0.612	
Lower Bound	0.530	0.668	0.605	0.603	
Upper Bound	0.548	0.686	0.625	0.621	

#### Volumes 1 and 2, Book Project

- Predictive Modeling Applications in Actuarial Science Volume 1
  - The first volume contains an introduction to predictive modeling methods used by actuaries
  - It was published in 2014
- Predictive Modeling Applications in Actuarial Science Volume 2
  - The second volume would be a collection of applications to P&C problems, written by authors who are well aware of the advantages and disadvantages of the first volume techniques but who can explore relevant applications in detail with positive results.



W. Frees, Richard A. Derrig

# Focus on Using R for Applications

2.• 🖯	k • 🔒 🔒 🗁 🕐	Go to file/function			
0 Untitl	ed5* × R data sets ×	Test ChainLadder.R ×	R data sets ×	PRIDIT - Cluster3.R ×	PRIDIT - »
15		0 /-		Run	Source +
1	t Code to calcula	To DIDITS and DDIDI	Ta on Quastic	unable Claime Data	
2	# Code to Carcula	nable claims data_	uestic	Maple Claims Data	-
4	+ this varsion of	data has dependent	war and is 1	000 lines_	
4	# mydatal<-read c	sv/"C:/ClusterData	Sim PIPPridDa	ta cav" header=TR	UF) -
5	# this version of	data has dependent	var and is 1	500 lines-	5 - / · · ·
6	mydata1<-read.cay	("C:/ClusterData/Si	mPTP.csv", hea	der=TRUE)-	
7	names (mydata1)	,,			
8	nrow(mydata1)-				
9	ncol(mydata1)-				
10	table (Suspicion, 1	egalrep)-			
11					
11 12	mydata=mydata1[,3	:27]¬			
11 12 13	mydata=mydata1[,3 mydata[1:5,]¬	:27]¬			
11 12 13 14	mydata=mydata1[,3 mydata[1:5,]¬	:27]¬			
11 12 13 14 15	mydata=mydata1[,3 mydata[1:5,]¬	:27]¬			•
11 12 13 14 15 1:1	mydata=mydata1[,3 mydata[1:5,]¬ 	:27]¬ III			R Scrip
11 12 13 14 15 1:1	mydata=mydata1[,3 mydata[1:5,]¬ ( (Top Level) \$ e ~/ &	:27]¬ '''	-		R Scrip
11 12 13 14 15 1:1 Console	<pre>mydata=mydata1[,3 mydata[1:5,]¬  ( ( (Top Level) \$ a ~/ \$ 1 -2.9678803</pre>	:27]¬ ''' -1.1013521	-		R Scrip
11 12 13 14 15 1:1 Console 4 5	¬ mydata=mydata1[,3 mydata[1:5,]¬ ( (Top Level) ‡ e ~/  1 -2.9678803 2 -4.9972007	:27]¬ 	-		R Scrip
11 12 13 14 15 1:1 Console 4 5 6	<pre>mydata=mydata1[,3 mydata[1:5,]¬  ( ( (Top Level) \$</pre>	:27]¬ 			R Scrip
11 12 13 14 15 1:1 Console 4 5 6 7	<pre>mydata=mydata1[,3 mydata[1:5,]¬  ( ( (fop Level) \$</pre>	-1.1013521 -1.0641477 -2.0426494 0.8856234			R Scrip
11 12 13 14 15 1:1 <b>Console</b> 4 5 6 7 8	<pre>mydata=mydata1[,3 mydata[1:5,]¬  ( ( (fop Level) \$</pre>	:27]¬ 			R Scrip
11 12 13 14 15 1:1 <b>Consol</b> 4 5 6 7 8 9	<pre>mydata=mydata1[,3 mydata[1:5,]¬  ( ( (fop Level) \$</pre>	:27]¬ 			R Scrip

# R Libraries

41

- Code is provided with book
- The "cluster" library from R used
  - Many of the functions in the library are described in the Kaufman and Rousseeuw's (1990) classic book on clustering, Finding Groups in Data.
- randomForest R library used to get dissimilarity matrix
- prcomp, princomp and factanal used for PRIDITs
- Some custom coding needed

# Dependent Variable Problem: Unsupervised Learning

- Insurance companies frequently do not collect information as to whether a claim is suspected of fraud or abuse
- Even when claims are referred for special investigation
- Solution: unsupervised learning

Francis Analytics and Actuarial Data Mining, Inc.

42

# Grouping Records

43



Francis Analytics and Actuarial Data Mining, Inc.



# Clustering

- Hierarchical clustering
- K-Means clustering
- Most frequent is k-means

Francis Analytics and Actuarial Data Mining, Inc.

#### **Cluster** Plot

45

plot(clara(x = ClusterDat1, k = 2, metric = "manhattan", stand clusplot( keep.data = TRUE))



These two components explain 100 % of the point variability.

Francis Analytics and Actuarial Data Mining, Inc.

# The Mortgage Crisis

Could simple descriptive statistics have predicted the meltdown?

#### Time Series of Loan-to-Value



Data from Demyanyk and Van Hemert, "Understanding the Subprime Mortgage Crisis", 2008

## Subprime Loan Volume and Size



Data from Demyanyk and Van Hemert, 2008

# Balloon Payments and Completed Documentation



Data from Demyanyk and Hemert, 2008

# **Observations from HMDA**

- HMDA indicates lower income applicants tend to have a higher loan to income ratio
- HMDA cross-state comparison indicates states with a foreclosure problem have consistently higher loan to income ratios compared to states not experiencing a foreclosure problem

#### The Data

HMDA Data

- LISC ZIP Foreclosure Needs Score
  - Subprime component
  - Foreclosure component
  - Disclosure component

http://www.housingpolicy.org/foreclosure-response.html

Zip Code Demographic Data

# CART Subprime Tree

n

%



# **CART** Foreclosure Variable Ranking

Independent Variable	Importance	Normalized Importance
Denial Percent	.027	100.0%
Mean Denial Score	.027	99.9%
PctApprove	.024	88.5%
ZipCodePopulation	.020	72.6%
PctPropNot1-4Fam	.019	69.5%
Median Rate Spread	.017	61.6%
PInCom	.016	60.5%
HouseholdsPerZipcode	.015	56.1%
Mean LTV Ratio	.014	52.7%

#### Results of Applying Clustering to HMDA Data

 K-means clustering applied to loan characteristics but not result data (i.e., approval)

	Cluster		
	1	2	3
Avg Loan Amount	297.23	566.96	163.80
Average Income	165.71	356.66	87.26
Mean LTV <sup>[2]</sup> Ratio	2.53	2.38	2.48
Rate Spread - mean	4.84	4.54	5.05
Median LTV Ratio	2.29	2.09	2.31
Median Rate Spread	4.40	3.95	4.67
Percent Applicants High LTV	4.4	3.8	4.5
Pct Applicants High Rate Spread	4.7	4.5	5.6
Percent Manufactured, Multi Family Houses	1.9	.4	6.1
Pct Home Improvement	57.8	56.5	65.6
Percent Refinance	52.4	52.5	57.3
Pct Owner Occupied	18.1	28.4	13.5

# Library for Getting Started

- Dahr, V, Seven Methods for Transforming Corporate into Business Intelligence, Prentice Hall, 1997
- Berry, Michael J. A., and Linoff, Gordon, Data Mining Techniques, John Wiley and Sons, 1997, 2003
- Derrig and Francis, "Distinguishing the Forest from the Trees", Variance, 2008
- If you use R, get a book on doing analysis in R. See www.r-project.org
- Francis, L.A., Neural Networks Demystified, Casualty Actuarial Society Forum, Winter, pp. 254-319, 2001. Found at www.casact.org
- Francis, L.A., "Taming Text: An Introduction to Text Mining", CAS Winter Forum, March 2006, www.casact.org
- Francis, L.A., Martian Chronicles: Is MARS better than Neural Networks? Casualty Actuarial Society Forum, Winter, pp. 253-320, 2003.
- Frees, Derrig and Francis, Predictive Modeling Applications in Actuarial Science, vol 1, Cambridge, 2014
- James, Witten, Hastie and Tibshirani, An Introduction to Statistical Learning with applications in R, Springer