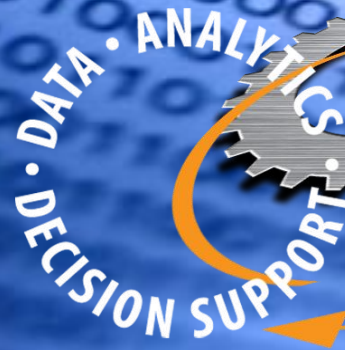


# The Stars are Aligned for P&C Analytics

Karthik Balakrishnan, Ph.D.  
ISO Innovative Analytics



DATA • ANALYTICS •  
DECISION SUPPORT

CADS Spring Meeting – 12 June 2009

Scottsdale, AZ



# Infrastructure Capabilities



- **Improvements in infrastructure capabilities**
  - Increasing computing power
  - Declining cost of storage and memory
  - Advances in parallel and distributed computing
    - E.g., Grid computing
  - Emerging capabilities
    - Floating data centers
    - Cloud computing
      - Hosted data mining
    - Etc.

# GRID Computing



IIA Analysts

## Virtual Machines



Active Directory



VM Server1  
16GB RAM  
2 Quad Core CPU  
146GB Internal



VM Server2  
16GB RAM  
2 Quad Core CPU  
146GB Internal



VM Server3  
16GB RAM  
2 Quad Core CPU  
146GB Internal

## Data Warehousing Environment

Production Database  
SQL 2008 Enterprise

Staging Database  
SQL 2008 Enterprise



Production Server 1  
4 Quad Core CPU  
32GB RAM

Production Server 2  
4 Quad Core CPU  
32GB RAM

Staging Server 1  
4 Quad Core CPU  
32GB RAM

Staging Server 2  
4 Quad Core CPU  
32GB RAM

Development Server  
4 Quad Core CPU  
32GB RAM

Storage Area Network

Backup Network

Production Network

Private Network



Grid Node1  
8GB RAM  
2 Dual Core CPU  
146GB Internal



Grid Node2  
8GB RAM  
2 Dual Core CPU  
146GB Internal



Grid Node3  
8GB RAM  
2 Dual Core CPU  
146GB Internal



Grid Node4  
8GB RAM  
2 Dual Core CPU  
146GB Internal



Grid Node5  
8GB RAM  
2 Dual Core CPU  
146GB Internal



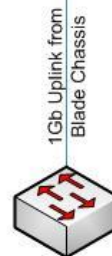
Grid Node6  
8GB RAM  
2 Dual Core CPU  
146GB Internal



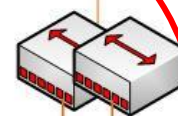
Grid Manager  
GRID Control  
MetaData  
8GB RAM  
2 Dual Core CPU  
146GB Internal



Core Switch 1



Core Switch 2



Fiber Channel  
SAN

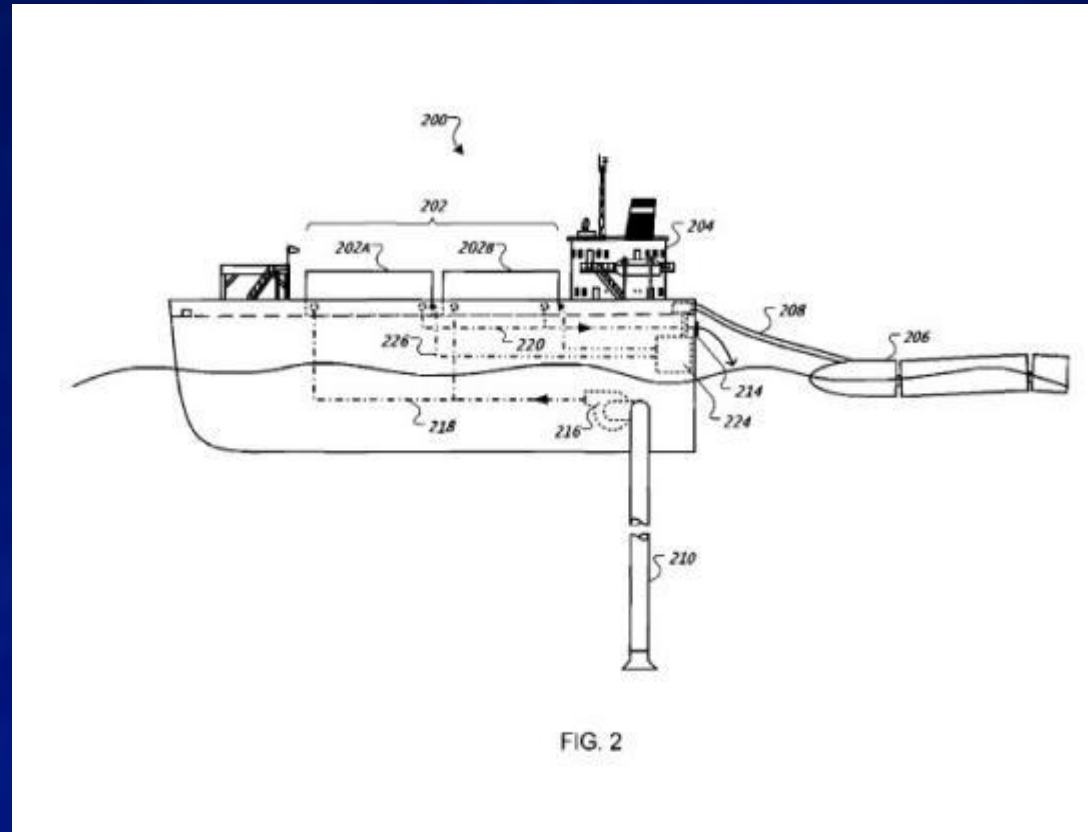
"The Grid"

High-bandwidth  
storage network

# Floating Data Centers



- **Google Patent Filing**
  - Wave-powered
  - Water-cooled
  - Wind turbines
- **International Data Security (IDS)**
  - San Francisco based
  - Refurbished cargo-ships

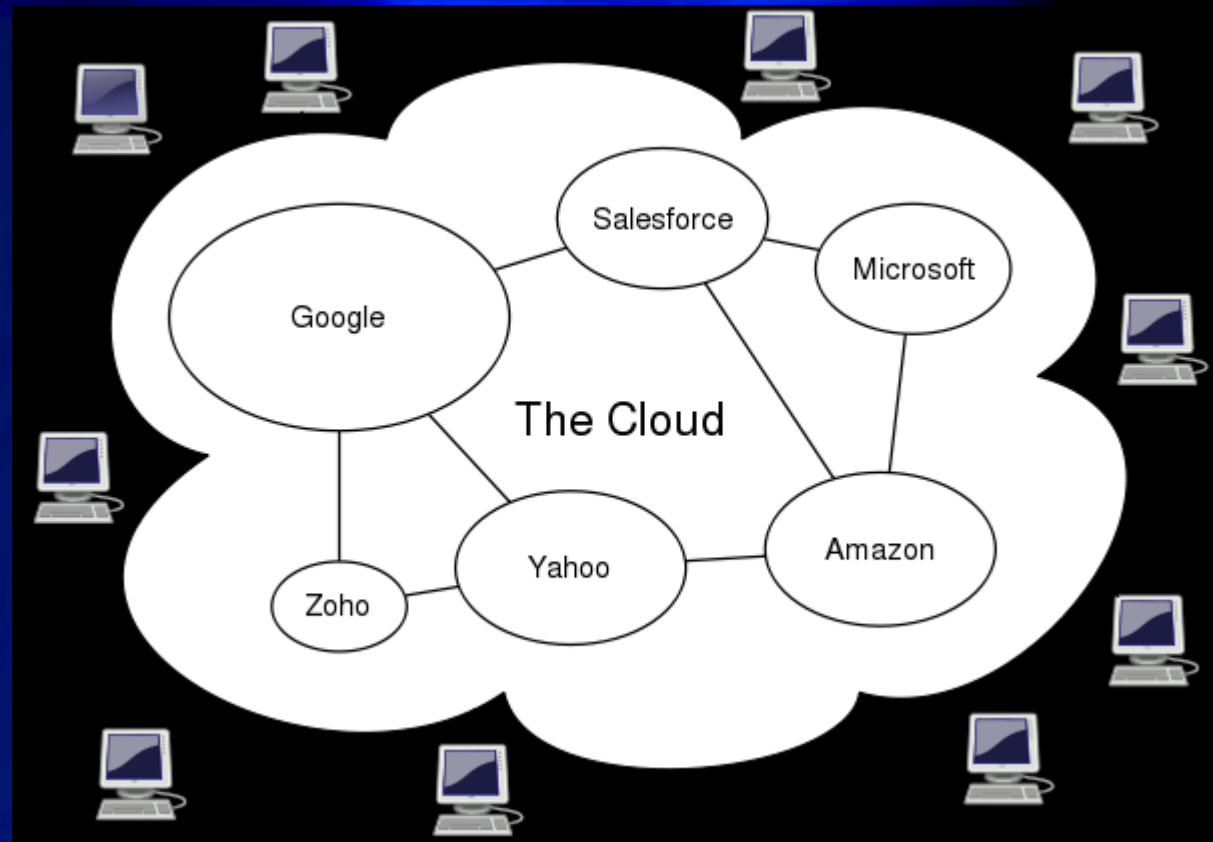


# Cloud Computing



## Infrastructure as a Service (IaaS)

- Scalable
- Virtualized



Cloud computing services usually provide common business applications online that are accessed from a web browser, while the software and data are stored on the servers.

# Availability and Access to Data

- **487 Exabytes ( $10^{18}$ ) data created in 2008**
  - Expected to grow to 2,500 Exabytes by 2012\*
  - In book-form – would stretch to Pluto and back
    - 10 times!
- **“Useful” data is also growing rapidly**
  - Public, e.g., government sources
  - Spurt in fee-based data sources
- **Addressed in various CAS meetings**
  - Albeit, limited to census/geography data

\* Source "Digital Universe" report published by International Data Corp. (IDC)

# Availability and Access to Data



- **Multi-media, rich detail**
  - Text
  - Voice
  - Video
  - Sensor-data (RFID, GPS, etc.)
    - Progressive's MyRate
      - Small device that records speed and time (but not location)
      - Progressive can determine what time of day you tend to drive, how many miles you average and how aggressively you drive



# Algorithms and Tools

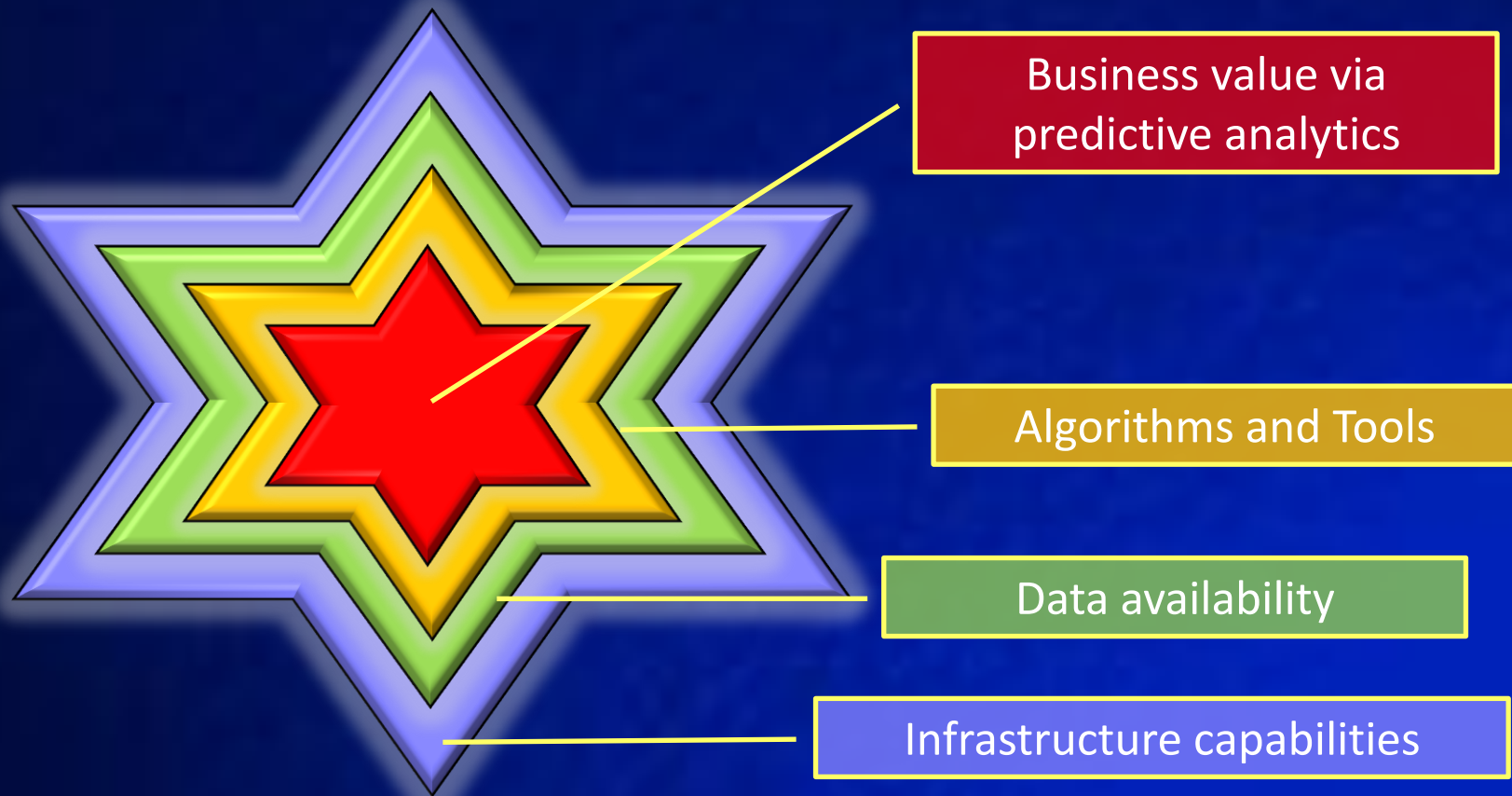
- **Convergence of quantitative disciplines –**
  - Statistics, Machine learning, Econometrics, Actuarial science, etc.
- **Result – a diverse array of algorithms for data manipulation, pattern analysis, and modeling**
  - Efficient algorithms to
    - Discover non-linearities/transformations
    - Identify interactions
    - Bin/group variables
    - Perform variable reduction/selection
    - Visualize data, etc.

# Algorithms and Tools



- **Emerging methodologies**
  - Text mining
  - Ensemble computing
  - Image recognition – OCR, handwriting, pictures, etc.
  - Speech/voice recognition
  - Video analysis, etc.
- **Importantly, tools available in the market**
  - Data Analysis and Modeling
    - R (public domain)
    - Data mining workbenches (SAS, SPSS, Statistica, etc.)
  - Visualization
    - SAS/Graph, R, ArcView, etc.

# The Stars are Aligned!



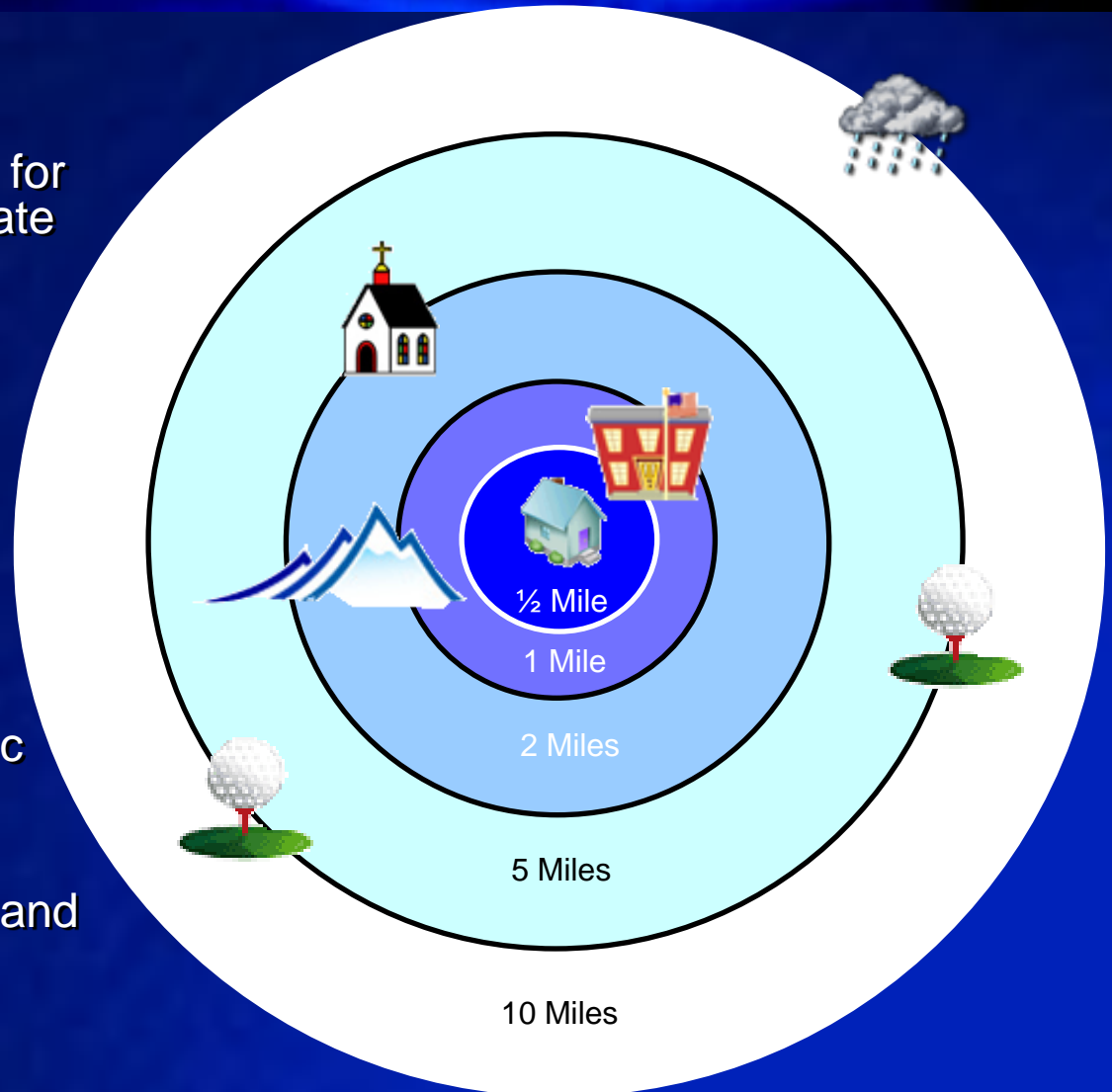
# RiskAnalyzer® Homeowners

- **Goal**

- Produce highly-refined prediction of Loss Costs for HO risks using multivariate modeling techniques

- **Model Structure**

- Loss Cost =  
Frequency \* Severity
- Frequency
  - probability of loss modeled with logistic regression
- Severity
  - GLM with a log link and Gamma error



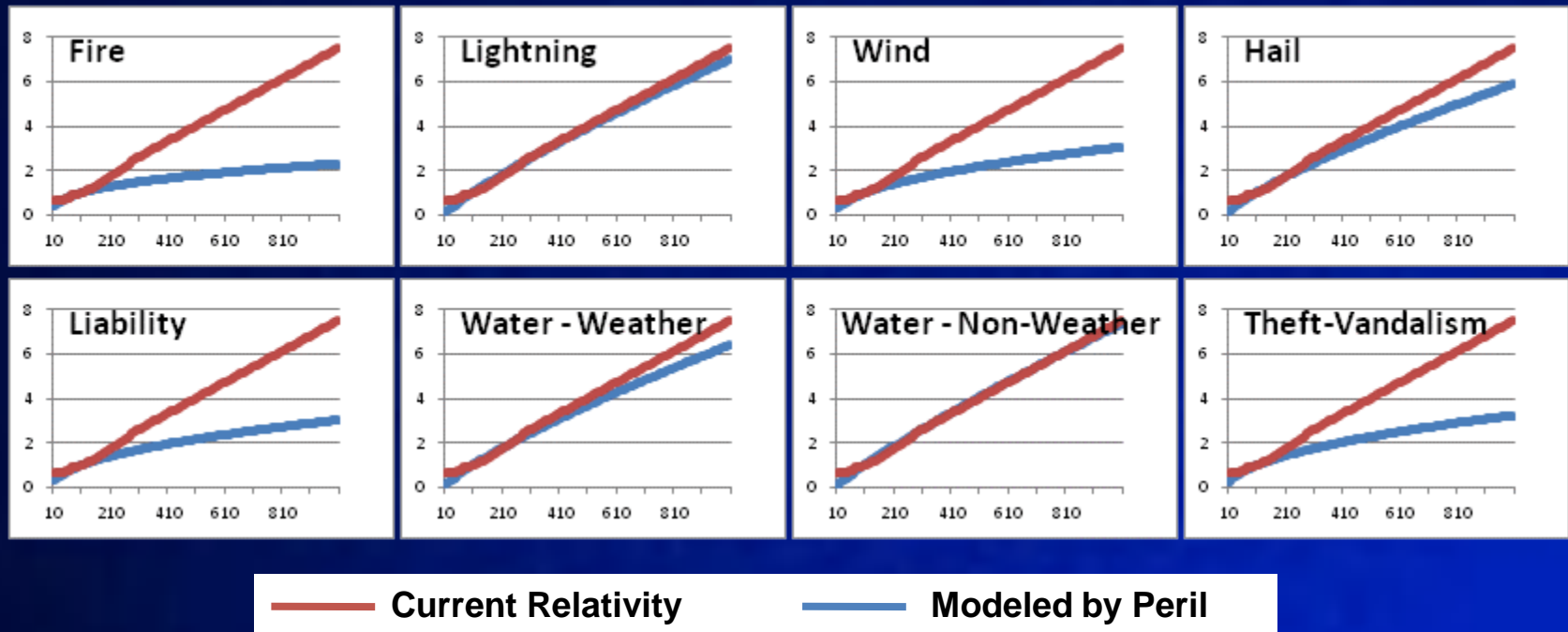
# Modeling at a Granular Level

1



**Decompose HO losses and model by peril to produce “tighter” models**

# AOI Relativities by Peril



- Significant variation by peril
- Source of lift

# Explore Detailed Data



- **North American Regional Reanalysis (NARR)**
  - “Best/most accurate North American weather and climate dataset”
- **Data Range – 1979 – 2007**
- **Granularity – 32 x 32 km grid**
- **8 daily readings (every 3 hrs)**
  - Accumulated precipitation
  - Air temperature at 2 meters
  - Rain
  - Wind
  - Relative humidity
  - Snow depth
  - etc.
- **Data Size ~ 150 GB**

<http://www.emc.ncep.noaa.gov/mmb/rrean/>

# Derive Novel Data Features

3

- **Temperature**
  - Mean
  - Max deviation from mean
  - # consecutive days below freezing, etc.
- **Wind**
  - # days with High wind, etc.
- **Precipitation**
  - # days with severe precipitation
  - # days without precipitation, etc.
- **Interactions**
  - Days without precipitation, high temperature, and high wind, etc.
- **2 person-years of effort**
- **80+ derived predictors**



# Visual Data Analysis (VDA)



Explore - PVA.PVA\_RAW\_DATA

File Edit View Actions Window

INCOME\_GROUP

TARGET\_B

Graph8

PVA.PVA\_RAW\_DATA

DONOR_A...	FILE_AVG...	INCOME_G...	LIFETIME...	MEDIAN_H...	Target Vari...	URBA
	6.29		44	535		1T
58	8.88	2	142	266	0?	0S
48	10.86	2	152	1221		0S
85	15.33	2	368	340	1T	
	8.86	2	328	789	1R	
59	21	7	210	2573	0C	
47	20.14	7	141	3770	0T	
	7.38		96	1821	1S	
76	4.36		144	303	0?	
80	26.25	5	315	711	1T	

DONOR\_AGE

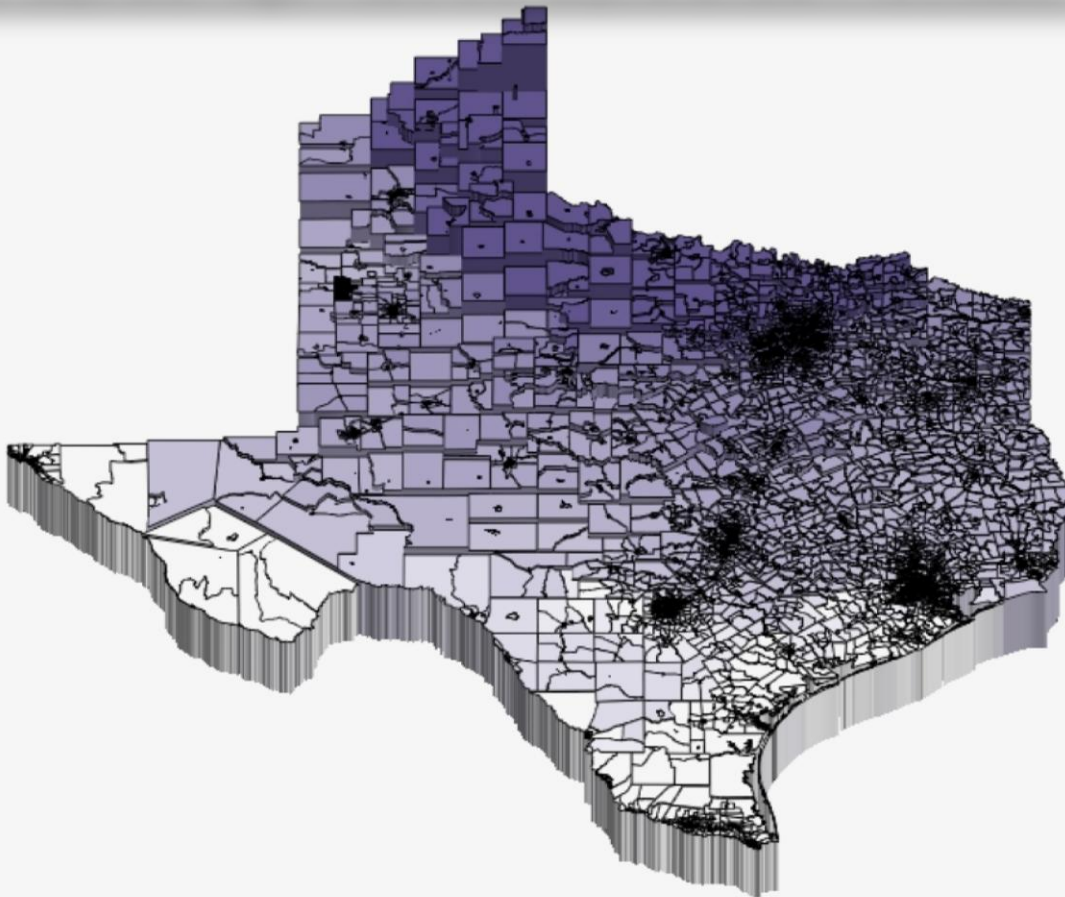
Sample Properties

Property	Value
Data	
Sample Method	Top
Sample Details	
Random Seed	12345
Top Rows	1000
Random Rows	1000
Percent Size	10.0

Apply Plot...

# Visualizing Aids Understanding

% of days with High < 32 x % of days with Low > 72 (Texas)



value

□ -1.05 - -1.01	□ -1.01 - -0.99	□ -0.99 - -0.99	□ -0.99 - -0.98	□ -0.98 - -0.97	□ -0.97 - -0.91	□ -0.91 - -0.79
■ -0.79 - -0.65	■ -0.65 - -0.52	■ -0.52 - -0.34	■ -0.34 - -0.26	■ -0.26 - -0.20	■ -0.20 - -0.13	■ -0.13 - 1.20

Positive coefficient in  
Wind Frequency  
model

Spatial visualization  
shows it is "Tornado  
Alley"

Using SAS/Graph

# Enable Serendipitous Discoveries

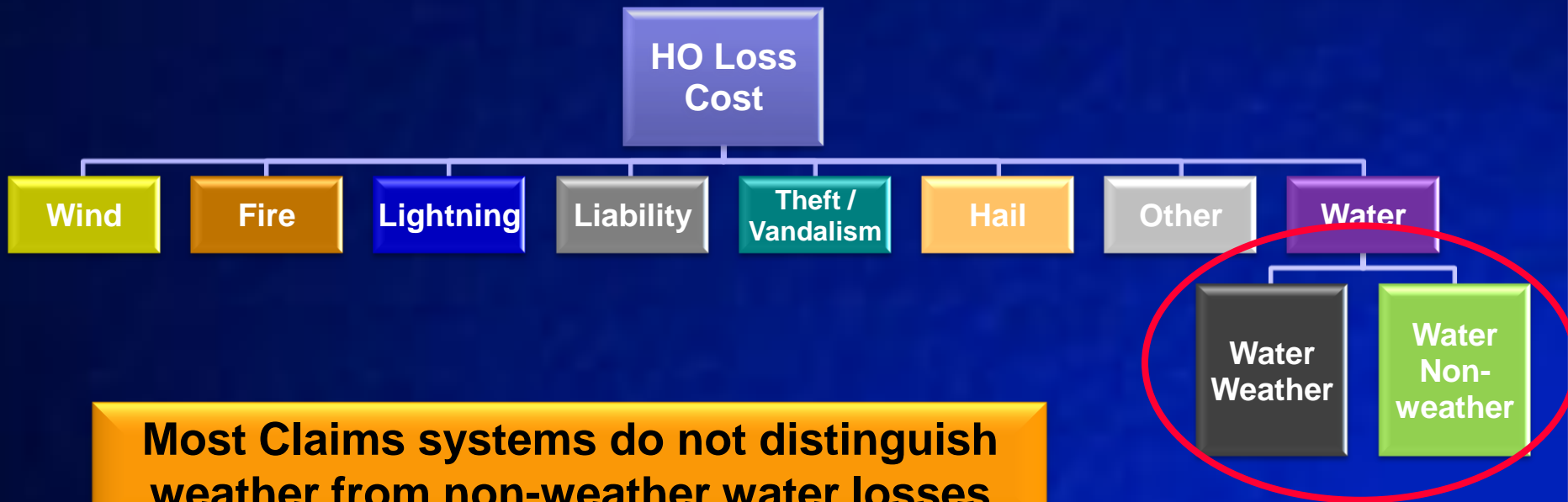
5

Weather & Elevation	FIRE	LIGHT	WIND	HAIL	WW	LIAB	THEFT
Elevation							
Temperature							
Precipitation							
Relative Humidity							
Snow							
Wind							
Ice Pellets							

Ellen Cohn. "Weather and Crime". *The British Journal of Criminology* 30:51-64 (1990)

# Exploit Novel Technologies

6



**Text Mining to the rescue!**

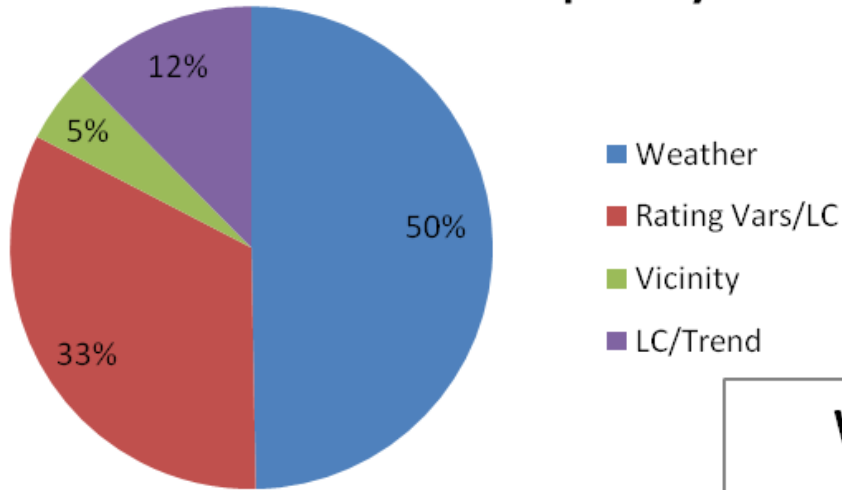
# Text Mining for “Cause of Loss”

- Rich information buried in unstructured data, such as loss descriptions or adjuster notes
  - Challenge – typos, abbreviations, poor structure, etc.
- Text mining loss descriptions

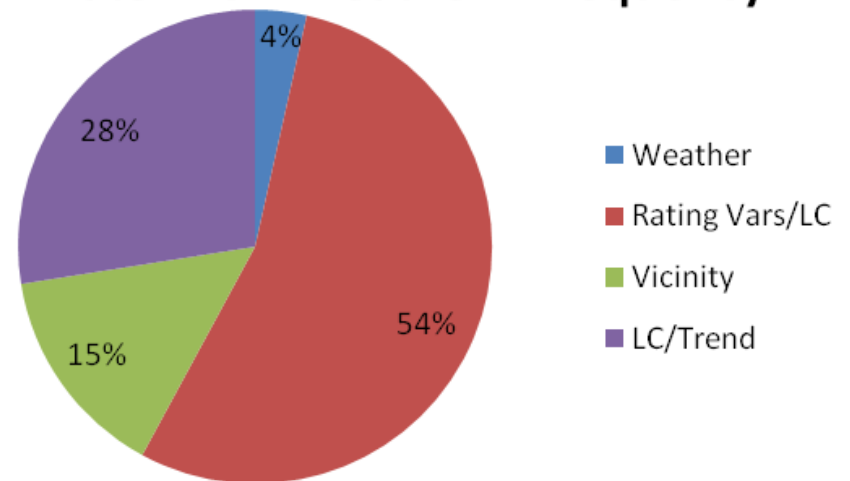


# Tighter and Relevant Predictors

## Water Weather - Frequency



## Water Non-Weather - Frequency



# Premium Audit Overview

- **Issued premium based on estimated payroll or sales**
- **Review policyholder's records and operations**
  - Determine the correct premium based on actual experience (accurate risk exposures)
  - Policyholder contractually obligated to comply
- **AP – Additional premium – insured → carrier**
- **RP – Return premium – carrier → insured**
- **Premium Audit Process**
  - Physical Audit — On-site/location audit by a person
  - Phone Audit — Audit via phone
  - Mail Audit — Form sent to Insured to complete and return

# Building a WC Premium Audit Model

- **Type of Problem**
  - Classification vs Value-Prediction
  - Few “reliable” audits to model actual dollar results
- **Modeled as a three class problem**
  - Returned Premium: audit result  $\leq$  \$0
  - Low Additional Premium:  $\$0 < \text{audit result} < \$625^*$
  - High Additional Premium : audit result  $\geq$  \$625
- **Aligns with potential business strategies**

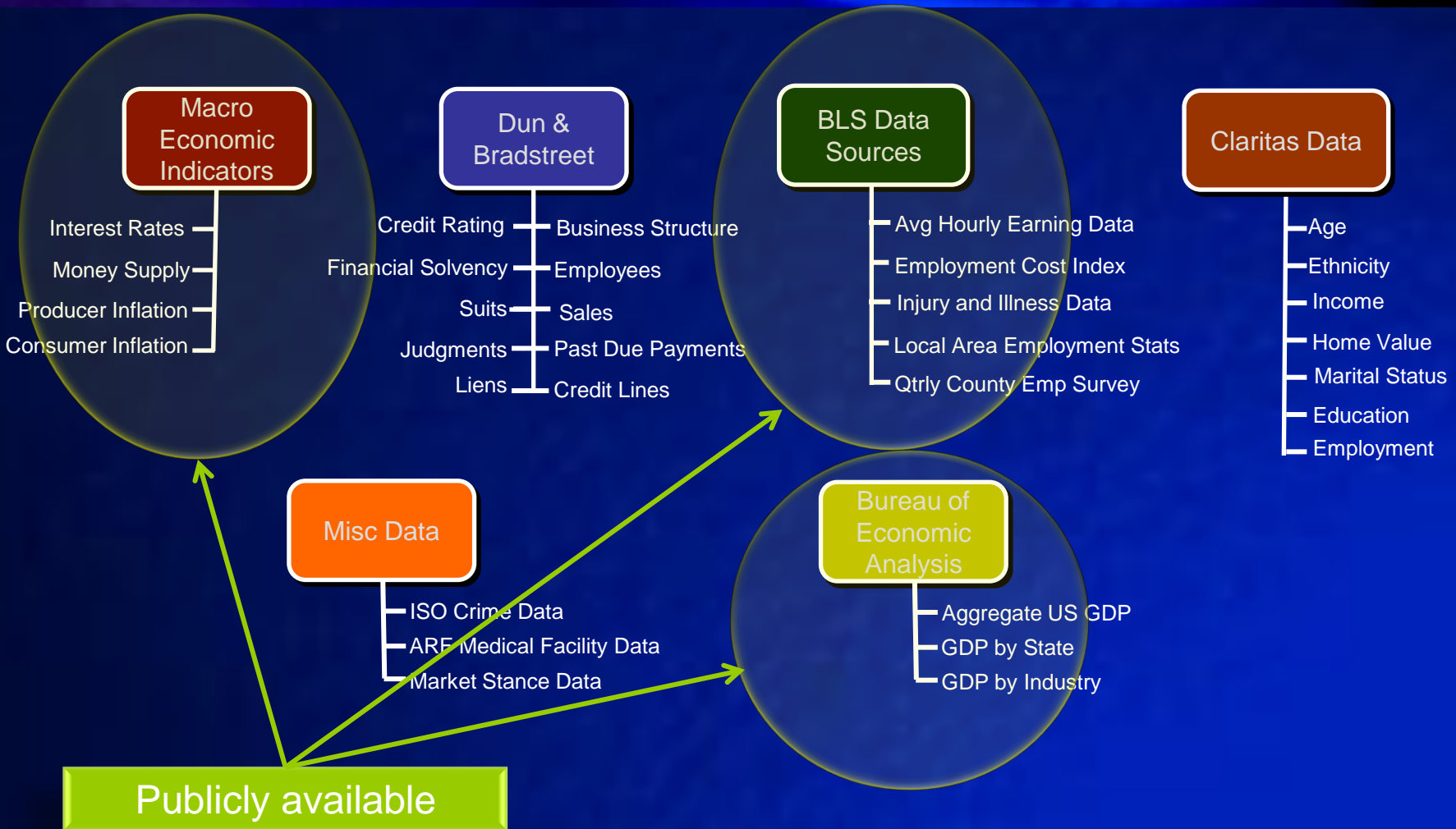
Model Prediction	Business Strategy
Returned Premium	Mail Audit
Low AP	Telephone Audit
High AP	Physical Audit

\* Illustrative break-even cost of a physical audit



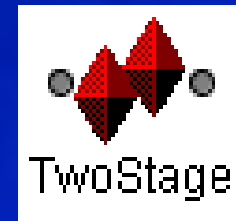
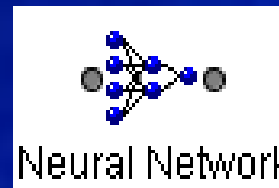
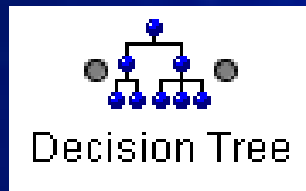
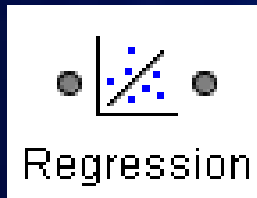
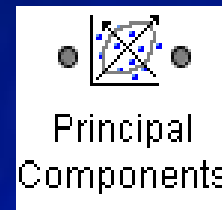
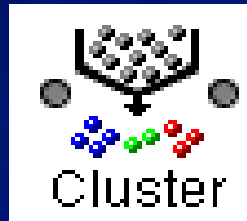
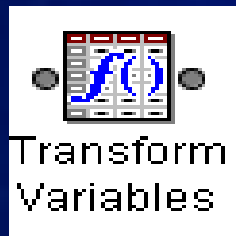
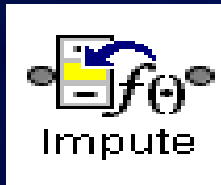
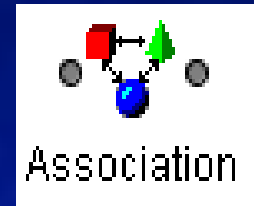
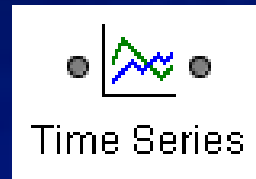
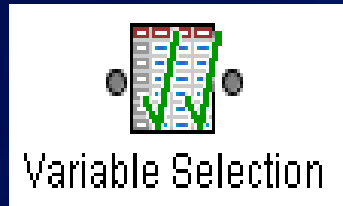
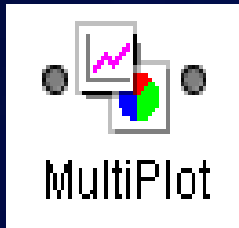
# Consider Diverse Data

7



# Use a “Toolkit” of Algorithms

8



# Use Multiple Algorithms

**Enterprise Miner - Donations**

File Edit View Actions Options Window Help

Sample Explore Modify Model Assess Utility

**Donations**

- Data Sources
  - Market Baskets
  - Donor Score
  - Purchase Data
  - Historical Loans
  - Donor Data
- Diagrams
  - Churn Analysis
  - Donor Analysis
  - Fraud Analysis
  - Loan Analysis
- Model Packages
  - Regression Results
  - Tree Results
- Users
  - EMUser

Property	Value
Name	Donations
Start-Up Code	...
Exit Code	...
Server	EMDDEV SAS 9.1
Path	d:\users\projects5
Max. Concurrent T:	Default

**Donor Analysis**

Donor Data

StatExplore → MultiPlot

Impute

Data Partition

Filter

Transform Variables

Variable Selection

Decision Tree

Regression

Rule Induction

AutoNeural

Model Comparison

Score

Donor Score

Market Baskets → Association → Cluster

Purchase Data → Merge → Data Partition (2)

Data Partition (2) → MBR

Data Partition (2) → Neural Network

MBR → Model Comparison (2)

Neural Network → Model Comparison (2)

Model Comparison (2) → Score

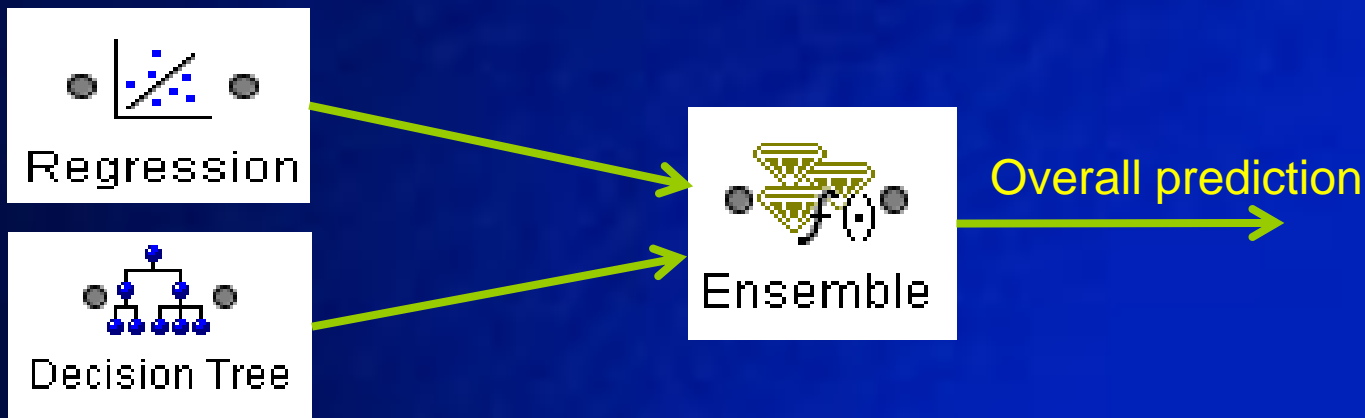
Model Package generated

Connected to EMDDEV SAS 9.1 - Logical Workspace Server

# Explore “Ensemble” Models

9

- **Several modeling techniques explored**
  - Multinomial Logistic Regression
  - CART Decision Tree, Boosted Trees, Random Forests
  - Support Vector Machines (SVM)
  - Neural Networks
- **Best model – Ensemble of Logistic and Tree**
  - Where predictions match, high confidence it is correct
  - Where predictions do not match select model/prediction with the higher confidence



# Where can Analytics be Applied?

## Operations Analytics

- Claims
  - Subrogation
  - Fraud
  - Litigation
  - IME
  - etc.
- Premium Inadequacy
  - Premium Audit WC/GL
  - Cov A ITV (PL)
- Loss Control
- Attrition Scoring
- etc...

## Marketing Analytics

- Strategic Market Dev.
  - Target Mkt
  - Niche identification
- Channel Optimization
  - Segmentation & LTV
- Product Innovation
  - Ideation support
- Customer Optimization
  - Segmentation & LTV
- Targeted Marketing Campaigns
  - Acquisition
  - X-sell/Up-sell
- etc.

## Operations Analytics

## Marketing Analytics

Insurance Lifecycle

## U/W & Actuarial Analytics

## Actuarial Analytics

- New Binning for factors
- Novel Rating Factors
- Novel Pricing Models
- Enhancing Reserving Models
- New Product/Coverage Pricing
- etc.

## U/W Analytics

- Risk Understanding
  - Causes of Loss
  - U/W sweet-spots
- Risk Qualification rules
- Risk Scoring Models
- Risk Tiering/Subsidy Models
- Renewal Scoring
- etc.

# In Sum...

- **“Perfect storm” created by advances in**
  - Infrastructure capabilities
  - Data availability and access
  - Methodologies and Tools
- **...has opened up tremendous opportunities for Analytical solutions within P&C**
- **If not doing so already, exploit the timing, leverage the opportunities, and create successes!**

**Thank you!**

**Karthik Balakrishnan**  
**[kbalakrishnan@iso.com](mailto:kbalakrishnan@iso.com)**