# SOA Predictive Analytics Initiatives

**KEN GUTHRIE**
**Managing Director, Education**
**Society of Actuaries**

SOCIETY OF
ACTUARIES

---

# Predictive Analtyics

- Also known as
  - Data Science
  - Big Data
  - Statistics, but with better marketing

# What has Changed? New Data Sources

- Information about insured individuals
  - Driving patterns via telematics
  - Exercise habits via fitness monitors

- Ancillary data
  - Social media
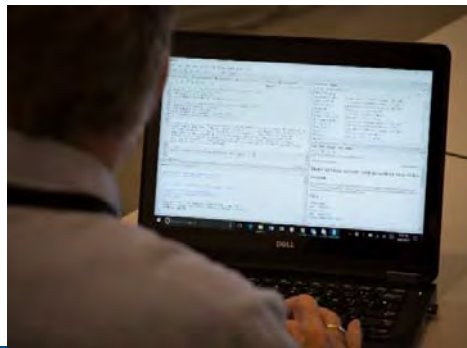  - Credit reports

# What has Changed? New Questions

- Previous
  - How many people do we expect to die/lapse/be hospitalized…?
- Now
  - What factors relate to persistence?
  - What factors relate to purchasing additional products?
  - What factors drive claims experience?

# What Has Changed? Software

- Previous
  - Expensive
  - Proprietary
- Now
  - Free
  - Open Source
  - Specialized (e.g., ChainLadder package)
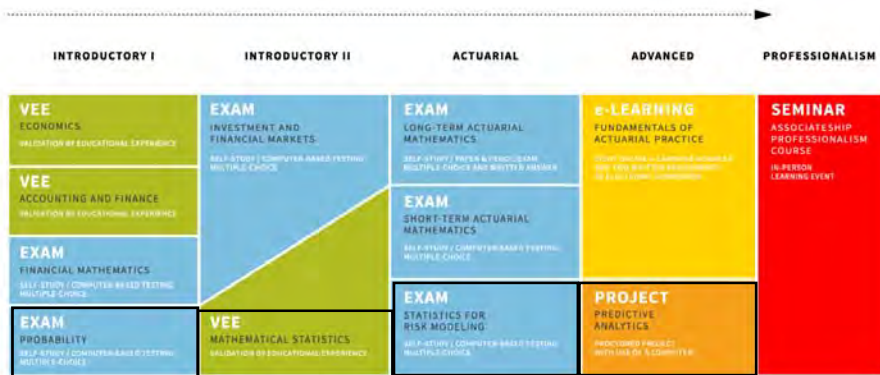
# Challenges for Actuarial Societies

- Methods/Models to Cover
  - GLM, PCA, SVM, GAM, MARS, …
- Related skills
  - Data preparation
  - Exploratory data analysis
  - Feature selection
  - Communication

# Challenges for Actuarial Societies

- Educating candidates and members
  - Read a book?  OR
  - Have directed practice
- Assessing candidates and members
  - Need to assess more than recollection of facts
  - Must verify that an individual can conduct an analysis from start (data and understanding the problem) to finish (a formal report of findings)

# Society of Actuaries Solution - Candidates

# Society of Actuaries Solution - Candidates

- Probability Exam
- Validation by Educational Experience: Mathematical Statistics
  - Ensures fundamentals of estimation and hypothesis testing covered
- Statistics for Risk Modeling Exam
  - Multiple choice
  - Models: Generalized Linear Model, Time Series, Principal Components, Decision Trees, Clustering
  - Methods: Cross-validation and regularization

# Society of Actuaries Solution - Candidates

- Predictive Analytics Exam
  - Instruction via e-Learning modules.
    - Additional information on models and methods
    - Data preparation and understanding, and
    - Communication.
  - Proctored project
    - Five hours computer-based exam at testing center
    - Computer equipped with R/Rstudio, Word and Excel
    - Candidates presented with business problem and dataset
    - Must write report that describes their analysis and states their findings

# Society of Actuaries Solution - Candidates

- Predictive Analytics Exam (cont')
  - Fully proctored
    - No internet access
    - Same security measures as for all SOA exams
  - Fully graded
    - Same protocol as fellowship written-answer exam

# Society of Actuaries Solution - Members

Certificate Program
- Open to any credentialed actuary
- Six e-Learning modules
  - 30-40 hours of study time per module
- Two-day seminar with project based assessment
- Self-study portion
  - Discussion forum interactions
  - Exercises
  - Offline Readings
  - Practice
  - End of Module Tests

# Society of Actuaries Solution - Members

**Certificate Program (cont')**

- e-Learning content
  - Predictive Analytics Tools
  - Effective Problem Definition and Project Management
  - Data Design, Transformation and Visualization
  - Data Exploration
  - Feature Generation and Selection
  - Model Development and Validation

# Certificate Program Sample Content

# Certificate Program Sample Content



# Certificate Program Sample Content

# Society of Actuaries Solution - Members

Professional Development Offerings

- **Predictive Analytics Seminar**, Kuala Lumpur, Malaysia,
27 August 2018

- **Predictive Analytics Seminar**, Hong Kong,
29 August 2018

- **Predictive Analytics Seminar**, Taipei, Taiwan,
31 August 2018

- **ERM/Big Data and Predictive Analytics**, Hangzhou,
China, September 2018 (Date TBD)

- **Big Data Seminar**, Jakarta, Indonesia,
13 November 2018