

ANTITRUST Notice



The Casualty Actuarial Society is committed to adhering strictly to the letter and spirit of the antitrust laws. Seminars conducted under the auspices of the CAS are designed solely to provide a forum for the expression of various points of view on topics described in the programs or agendas for such meetings.

Under no circumstances shall CAS seminars be used as a means for competing companies or firms to reach any understanding – expressed or implied – that restricts competition or in any way impairs the ability of members to exercise independent business judgment regarding matters affecting competition.

It is the responsibility of all seminar participants to be aware of antitrust regulations, to prevent any written or verbal discussions that appear to violate these laws, and to adhere in every respect to the CAS antitrust compliance policy.

Machine Learning

Insurance Applications and Recent Advances

Prof. Paul Beinart
Research and Development, EagleEye Analytics and
Centre For Quantum Computation and Intelligent Systems
University of Technology Sydney

Value Proposition – What's Different


Statistical Science
2001, Vol. 16, No. 3, 199–231
Statistical Modeling: The Two Cultures
Leo Breiman

"Abstract. There are two cultures in the use of statistical modeling to reach conclusions from data. One assumes that the data are generated by a given stochastic data model. The other uses algorithmic models and treats the data mechanism as unknown.

The statistical community has been committed to the almost exclusive use of data models. This commitment has led to irrelevant theory, questionable conclusions, and has kept statisticians from working on a large range of interesting current problems.


Algorithmic modeling, both in theory and practice, has developed rapidly in fields outside statistics. It can be used both on large complex data sets and as a more accurate and informative alternative to data modeling on smaller data sets. If our goal as a field is to use data to solve problems, then we need to move away from exclusive dependence on data models and adopt a more diverse set of tools."

Leo Breiman (January 27, 1928 – July 5, 2005) was a distinguished [statistician](#) at the [University of California, Berkeley](#). He was the recipient of numerous honors and awards, and was a member of the [United States National Academy of Science](#). Breiman's work bridged the gap between statisticians and computer scientists, particularly in the field of [machine learning](#).




Machine Learning

- Why do it?
- Types of problem classes
- Data issues
- Methods
- Insurance Applications
 - Induction
 - Continuous Induction
 - Multivariate Fusion




Why do Machine Learning?

- Database resources are largely going to waste
 - MIS, DSS are not machine learning
- Increasing rate of data collection
- Knowledge Based Systems
 - Domain Experts have limited knowledge
 - Knowledge Acquisition is a slow process




Machine Learning Myths

- Machine learning tools need no guidance.
- Machine learning requires no data analysis skill.
- Machine learning eliminates the need to understand your business and your data.
- Machine learning tools are “different” from statistics.




Types of Problems

- Classification
 - Fraud detection, loan approval, etc.
- Regression
 - Insurance rating, stock price prediction, etc.
- Clustering (Pattern Detection)
 - Customer clustering, time series, etc
- Mixtures



Data

- Data Types
- Reliability
 - Missing data
 - Skewed data
 - Noisy Data
- Variability
 - Data changes over time
- Sufficiency
 - Sample size
- Testing implications
 - Validation data



Classification Learning

- Supervised Learning, Inductive Learning
- Dataset contains examples of input attributes and their corresponding classification

How Does It Work?


- Examine every input attribute
 - Find the best split that can be made
- Select the best attribute, with the best split
- Build a branch for it and assign the appropriate subset of the data
- Determine if any branch is a leaf node
- Send that subset back to the start

Limitations

- No guaranteed way to find optimal solution
- Induction must compromise
 - Specificity
 - Simplicity
- OR
 - Accuracy
 - Description length


Good and Bad of Induction

- | | |
|--|--|
| <ul style="list-style-type: none">• Many input attributes• Easy to understand results• Relatively fast to run• Relatively easy to use• Several algorithms to choose from<ul style="list-style-type: none">◦ IDS, C4.5, CART, O-Btree | <ul style="list-style-type: none">• No continuous non-linear effects• No optimal way of partitioning numeric attributes• Poor performance on noisy data• No algorithm selection rules found• Regression uncommon |
|--|--|




Knn - K Nearest Neighbor

- Take k nearest instances to make estimate
- Issues
 - How many is k
 - How far away can neighbors be
 - CBR



How Many in k

- Static - determined at development time
- Dynamic
 - Statistically based
 - Heuristically based
 - Composites
 - With/without Gaussian biases
 - With/without directional imperatives



How Far Away for Neighbours

- Euclidian distance
 - Difficult with categorical data
 - Development of biases for axes
- Heuristic distance
 - Encodes domain knowledge
- Non-linear
 - Interaction between variables

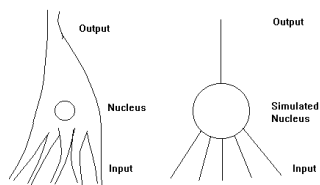
Classification


- Easy when all neighbors agree
- Problems
 - Sparsity of data locally
 - Neighbors not in agreement
 - Probabilistic democratic classification
 - Weighted probabilities base on distance
 - Edges - all neighbors on one side

Regression

- Central estimate based on neighbours
- Estimate biased by distance
- Certainty
 - Sample sufficiency
 - Distribution
- Edges
 - Incorporate multi-dimensional linear trends
- Care with skewed distributions


Neurons - Real v Artificial






Artificial Neural Networks

- Fundamentals
- Training
- Problems with NN's
- NN Variants




Fundamentals

- Each Neuron
 - A Number of Inputs
 - Each Input Has a Weight Associated With It
 - The Weight Moderates Sensitivity To That Input
 - One Output
 - The Output Can be Connected to Many Neurons
 - As Input
 - A Mathematical Function Inside
 - Controls Firing




Training

- Initialization
 - Weights Must be Randomized
- Learning Function
 - Back propagation
 - Weight Space and Error Function
- Data
 - Training Set
 - Test Set
 - OverFitting




Using Neural Nets

- Set learning parameters
 - Learning rate, momentum
 - Tolerance
- Train
- Test
- Iterate and tune




Problems?

- Local Minima!
 - Shaking
 - Genetic Algorithms
- OverFitting
 - Needs validation data to control
- Topology
 - Brittleness to Change (Sine Function NN)




More Problems

- Black Box
 - No Insight
 - No Expert Validation
- No Domain Knowledge
 - No Ability to Use the Obvious




Adaptive Networks

- Cascade Correlation
 - Start With Inputs Connected to Outputs - Train
 - Add One Hidden Node - Set Orthogonal to Existing Nodes - Retrain Other Nodes
 - Go Back to Step 2
 - Stop When Adding a Node Does Not Improve Fit
 - Can Produce Mapping with Abrupt Transitions



Optimal Brain Damage

- Start With Overly Populated NN – Train
- Select Hidden Layer Node with Least
- Sensitivity to Input Nodes
- Delete that Node - Re-Train
- Repeat
- Stop When Next Iteration Does Not Improve Fit



Radial Basis Functions

- Hidden Units Use Gaussian Activation
- Activation Governed by Euclidian Distance Between Weight and Input Vectors
- Hidden Units Represent Clusters in Data
- Gradient Descent Learning
- Input Attribute Numbers Cause Geometric Explosion of Hidden Units

Kohonen Self Organising Maps

- Unsupervised, Winner Take All Learning (No Output Values)
- Clustering Problem
- Units Arranged in 2 Dimensional Lattice
- Weight Vector Same Dimensionality as Input Vectors
- Randomize Weight Vectors Initially

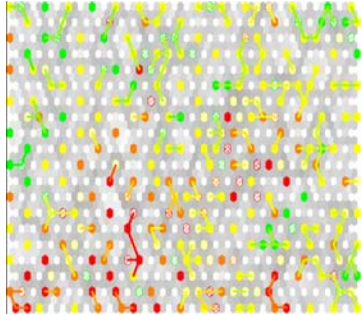
SOMs

- For Each Input Vector
 - Determine Which Unit Wins
 - Change Its Weight Vector Towards Input Vector (by Learning Rate)
 - Change All Its Neighbors (Limited by Neighborhood Size) Towards Input Vector
- As Training Proceeds
 - Reduce Neighborhood Size
 - Reduce Learning Rate

SOM Network View



SOM Data View



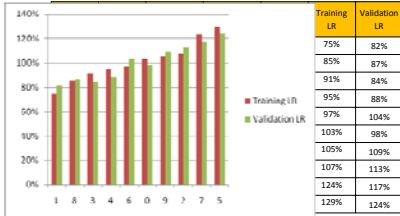
GLM Related Results

- Methodology
 - Use machine learning inductive techniques to search the residuals of the GLMs
 - Medium size auto portfolio
 - Pure premiums set by consulting actuaries using GLM tools on 4 years of data
 - Pure premiums derived from ALL of the data
 - We divide data into training and validation
 - Search for difference between pure premiums and claims signal
 - Test on validation data

Results

- Methodology
 - Use only policy attributes
 - One year for training, subsequent year for validation
 - Derive point estimates for segments of the portfolio using a minimum data requirement (number of claims)
 - Derive continuous estimates
 - Test accuracy of estimates on validation data

Consulting Actuaries A Results



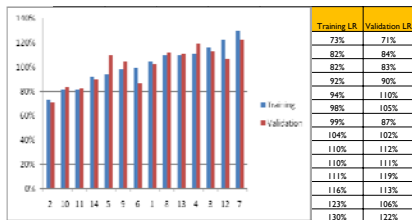
How well does it do?

- Correlation 0.93
- Lift
 - Using exposure weighted standard deviation of loss ratios
 - Training 16.4%
 - Validation 14.5%

	Deviance	Squared error	Chi squared error
GLM Premiums	34.15388	1463.838	5.47708
Estimated Premiums	15.02064	243.8955	0.946702

- Result
 - Consistent and large signal present in GLM residuals
 - Better fitting premiums

Consulting Actuaries B Training Results



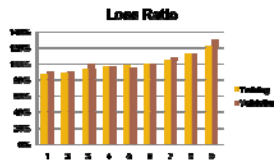
How well does it do?

- Correlation 0.87
- Lift
 - Using exposure weighted standard deviation of loss ratios
 - Training 16.3%
 - Validation 15.2%

	Deviance	Squared error	Chi squared error
GLM Premiums	41.60947	2241.3318	7.9576
Estimated Premiums	18.30203	737.3713	1.814276

- Result
 - Consistent and large signal present in GLM residuals
 - Better fitting premiums

Signal in GLM Residuals

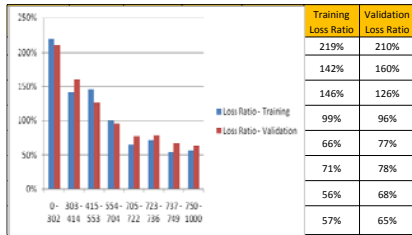


Loss Ratio		Frequency		Severity	
Training	Validation	Training	Validation	Training	Validation
87.64%	90.52%	10.75%	11.02%	\$1,775	\$1,799
89.24%	90.41%	14.18%	14.01%	\$1,864	\$1,895
93.09%	100.12%	14.82%	15.19%	\$2,045	\$2,123
97.06%	95.87%	11.73%	11.59%	\$1,880	\$1,886
97.62%	95.24%	13.95%	14.03%	\$2,080	\$2,038
99.03%	99.45%	17.99%	17.43%	\$2,509	\$2,613
104.44%	107.57%	16.17%	16.70%	\$2,239	\$2,220
113.41%	112.19%	19.09%	18.04%	\$2,611	\$2,655
122.37%	129.69%	15.36%	15.44%	\$2,449	\$2,582
Correlation 0.962		Correlation 0.986		Correlation 0.989	

Continuous Estimates

- Modify induction to produce continuous estimates
- Estimates made of exposures based on a 0 to 1000 range
 - 0 is best loss ratio
 - 1000 is worst loss ratio
 - An insurance score – parallel to credit score

Training Results



How well does it do?

- Correlation 0.98
- Lift
 - Using exposure weighted standard deviation of loss ratios
 - Training 30.4%
 - Validation 28.9%

	Deviance	Squared Error	Chi Squared Error
GLMs	44.75	5722	21.74
Output Ranges	18.23	528	1.99

- Result
 - Consistent and large signal present in GLM residuals
 - Better fitting premiums

Current Research

- Domain Driven Algorithms
 - Fraud detection for electronic payment transactions
 - Market manipulation detection
- Fusion Algorithms
 - Combining two or more existing algorithms

Fusion Algorithms

$$F(\mathbf{x}) = \hat{a}_0 + \sum_{k=1}^K \hat{a}_k r_k(\mathbf{x}) + \sum_{j=1}^n \hat{b}_j l_j(x_j)$$

r = rule identity function

- Combination of:
 - Rules (hundreds to thousands)
 - Linear model
 - Additive

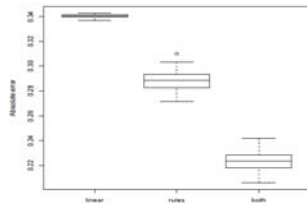
“Predictive Learning via Rule Ensembles”
Jerome Friedman and Bogdan Popescu

Fusion Algorithms

$$\{(\hat{a}_k)_k^N, (\hat{b}_j)_j^n\} = \arg \min_{\{a_k\}_k^N, \{b_j\}_j^n} \sum_{i=1}^N L \left(y_i, a_0 + \sum_{k=1}^K a_k r_k(\mathbf{x}_i) + \sum_{j=1}^n b_j l_j(x_{ij}) \right) + \lambda \cdot \left(\sum_{k=1}^K |a_k| + \sum_{j=1}^n |b_j| \right)$$

- While minimizing the loss function L
 - L can be any loss function
- Note
 - Lasso term
 - Controls complexity

Performance



$$\left\{ y_i = 10 \cdot \prod_{j=1}^6 e^{-2x_{ij}} + \sum_{j=6}^{20} \pi_{ij} + \epsilon_i \right\}_{i=1}^N$$

$x_{ij} \sim U(0,1)$ and $\epsilon \sim N(0, \sigma^2)$ Gaussian with 2.0 signal to noise

What about Real Data?

- Does real data exhibit linear and nonlinear components?
- Credit card applications from major bank
 - Bank model continuously developed
 - Compound variables identified
 - Derived variables developed
- Recast algorithm within logistic regression framework

Credit Card Applications Bank Results

Score Band	Accepts	Goods	Bads	%Accepts	%Goods	%Bads	Bad Rate	OML	Current Accepts	Current Goods	Current Bads	Diagonal	Gain Calc.
1	11,310	9,760	1,551	85%	69%	49%	13.72%	6.29	4,971	4,171	41.17%	43%	0.010
2	11,528	10,430	610	90%	89%	10%	5.69%	11.22	9,771	9,091	57.92%	57%	0.022
3	11,316	10,879	438	95%	95%	5%	3.87%	24.64	11,491	11,371	69.14%	69%	0.020
4	11,848	11,366	482	95%	95%	5%	2.98%	28.14	10,971	10,761	77.87%	78%	0.020
5	11,176	10,929	248	97%	97%	3%	1.72%	37.28	10,821	10,711	82.71%	83%	0.021
6	10,822	10,650	172	98%	97%	3%	1.61%	41.21	10,721	10,611	86.41%	87%	0.022
7	11,258	11,428	170	97%	97%	3%	0.92%	44.64	11,271	11,411	87.91%	88%	0.021
8	11,262	11,464	106	98%	98%	2%	0.67%	116.01	10,771	10,751	92.21%	92%	0.019
9	10,476	10,421	55	99%	99%	1%	0.53%	180.47	10,271	10,251	93.81%	94%	0.017
10	12,465	12,549	84	99%	99%	1%	0.37%	276.77	12,071	12,051	95.21%	95%	0.011
11	11,442	11,396	46	99%	99%	1%	0.39%	247.77	11,671	11,651	96.41%	96%	0.011
12	10,966	10,968	41	99%	99%	1%	0.37%	260.63	10,971	10,971	97.51%	98%	0.011
13	11,822	11,793	29	99%	99%	1%	0.24%	348.42	11,871	11,871	98.21%	98%	0.010
14	11,443	11,434	9	99%	99%	1%	0.08%	1250.24	11,471	11,471	99.71%	99%	0.001
15	10,211	10,202	9	100%	100%	0%	0.13%	977.41	10,271	10,271	99.91%	99%	0.001
16	12,431	12,399	32	100%	100%	0%	0.26%	1110.93	12,471	12,471	99.71%	99%	0.002
17	12,832	12,800	32	100%	100%	0%	0.25%	2470.03	12,871	12,871	99.91%	99%	0.002
18	11,111	11,079	32	100%	100%	0%	0.29%	1776.76	11,151	11,151	100.00%	100%	0.001
19	12,142	12,101	41	100%	100%	0%	0.34%	1210.13	12,181	12,181	100.00%	100%	0.001
20	241,258	227,483	13,775	100%	100%	100%	1.83%	48.24	100.00%	100.00%	100%	0.000	
Total													GINI 74.83%

Fusion Algorithm Result

Score Band	Accepts	Goods	Bads	%Accepts	%Goods	%Bads	Bad Rate	OML	Current Accepts	Current Goods	Current Bads	Diagonal	Gain Calc.
1	11658	9814	1794	84%	67%	47%	18.43%	5.67	4,371	41,461	11,441	0.02441	
2	11609	10,948	856	94%	87%	7.99%	13	10.04%	9,041	75,171	10,171	0.03041	
3	11504	11,085	418	96%	97%	3.85%	28	15.01%	13,301	81,801	81,801	0.03881	
4	11588	11,985	310	98%	98%	2.07%	36	20.03%	18,801	93,021	93,021	0.04781	
5	11642	11,886	146	99%	99%	1.45%	78	25.08%	23,021	93,961	93,961	0.04641	
6	11594	11,828	76	99%	99%	1%	0.61%	191	30,271	28,161	95.71%	0.05081	
7	11164	11,131	33	99%	99%	1%	0.30%	217	36,081	34,061	95.58%	0.05271	
8	12292	12,247	45	99%	99%	1%	0.29%	437	40,391	39,441	97.24%	0.05451	
9	12105	12,076	29	99%	99%	1%	0.24%	594	49,021	48,701	98.15%	0.05614	
10	10719	10,158	561	94%	87%	11%	0.29%	508	50,021	49,221	98.68%	0.04971	
11	9961	9,964	10	100%	100%	0%	0.17%	874	66,371	66,371	99.73%	0.04761	
12	10008	10,071	10	100%	100%	0%	0.10%	1007	60,431	60,391	99.39%	0.03021	
13	6792	6,747	45	100%	100%	0%	0.64%	2248	63,351	62,701	99.47%	0.02974	
14	9748	9,746	2	100%	100%	0%	0.02%	470	67,191	66,991	99.82%	0.03061	
15	10961	10,965	5	100%	100%	0%	0.03%	3653	71,871	71,411	99.62%	0.03101	
16	8478	8,472	6	100%	100%	0%	0.07%	2038	78,071	77,641	99.70%	0.03134	
17	12719	12,712	7	100%	100%	0%	0.06%	1819	83,001	82,231	99.37%	0.02943	
18	12719	12,718	1	100%	100%	0%	0.01%	12718	89,021	88,821	100.00%	0.00046	
19	12719	12,719	0	100%	100%	0%	0.00%	#DIV/0!	94,021	94,411	100.00%	0.00002	
20	12719	12,719	0	100%	100%	0%	0.00%	#DIV/0!	100,00%	100,00%	100.00%	0.00000	
Total													GINI 82.99%

Insurance Fusion Algorithm


- Use personal auto data set with GLM derived premiums
- Adapt fusion algorithm
 - Multiplicative rather than additive
 - Recast as maximizing likelihood
- Can it improve GLM premiums?
 - Use test and validation data (random 70%/30%)
 - Measure Tweedie deviance against data set ($p=1.72$)

Insurance Results - Summary

	Training	Validation
Null	19,165,615	8,213,623
Existing GLM	18,659,759	7,965,385
Improvement	2.64%	3.02%
Fusion Model	18,551,066	7,949,261
Improvement	3.21%	3.22%


Insurance Results - Relativities

Variable	Rel	Variable	Rel	Variable	Rel	Variable	Rel	Variable	Rel
a1	3.73	b1	0.95	b21	1.06	e1	1.23	f1	0.97
a2	1.19	b2	1.00	b22	1.21	e2	1.00	n2	1.00
a3	1.19	b3	1.23	b23	1.12			n3	1.01
a4	1.06	b4	1.17	b24	1.12	f1	1.18	n4	1.08
a5	1.07	b5	1.09	b25	1.15	f2	1.00	n5	1.32
a6	1.09	b6	0.99	b26	1.02			n6	1.41
a7	1.00	b7	1.05	b27	1.14	g1	1.16	n7	1.47
a8	1.00	b8	1.09	b28	1.13	g2	1.06	n8	1.54
a9	1.12	b9	1.09	b29	1.22	g3	0.98	n9	1.20
a10	1.00	b10	1.12	b30	1.24	g4	1.00		
a11	1.12	b11	1.21	b31	1.20	g5	1.01	o1	0.83
		b12	1.27	b32	1.11	g6	1.06	o2	1.00
c1	0.86	b13	1.06	b33	1.10	g7	1.13	o3	0.76
c2	1.00	b14	1.12	b34	1.16			o4	0.97
		b15	1.00	b35	1.24	h1	1.00	o5	0.79
d1	0.93	b16	1.04	b36	1.10	h2	0.86		
d2	0.94	b17	1.21	b37	1.39	h3	1.02		
d3	0.90	b18	1.33						
d4	0.96	b19	1.18						
d5	0.86	b20	1.63						
d6	0.92								
d7	0.89								
d8	1.00								
d9	0.98								
d10	1.00								



Summary

- Point estimates (a piecewise constant function) derived from training data fit the validation data much better than the GLM alone
 - Improvement in fit is very significant
- Regardless of who fits the GLM
- Continuous estimates (scores) also fit validation data better than the GLM alone
 - Have greater lift
 - More lift and better fit than credit scores
- Fusion Algorithms
 - Contribute significantly more variables, and rules
 - Produce better fitting premiums



Conclusion

- GLMs
 - Overfit their beta values (validation)
 - Underfit the signal in the data
- Machine learning methods can be a valuable supplement to GLMs
 - Can find extra lift
 - Can improve fit of pure premiums
- More accuracy
 - Better risk selection
 - Reliable price optimization
